

Efficient Temporal Synopsis of Social Media Streams

by

Younes Abouelnagah

A thesis
presented to the University of Waterloo
in fulfillment of the
thesis requirement for the degree of
Master of Math
in
Computer Science

Waterloo, Ontario, Canada, 2013

© Younes Abouelnagah 2013

I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.

Abstract

Search and summarization of streaming social media, such as Twitter, requires the ongoing analysis of large volumes of data with dynamically changing characteristics. Tweets are short and repetitious – lacking context and structure – making it difficult to generate a coherent synopsis of events within a given time period. Although some established algorithms for frequent itemset analysis might provide an efficient foundation for synopsis generation, the unmodified application of standard methods produces a complex mass of rules, dominated by common language constructs and many trivial variations on topically related results. Moreover, these results are not necessarily specific to events within the time period of interest. To address these problems, we build upon the Linear time Closed itemset Mining (LCM) algorithm, which is particularly suited to the large and sparse vocabulary of tweets. LCM generates only closed itemsets, providing an immediate reduction in the number of trivial results. To reduce the impact of function words and common language constructs, we apply a filtering step that preserves these terms only when they may form part of a relevant collocation. To further reduce trivial results, we propose a novel strengthening of the closure condition of LCM to retain only those results that exceed a threshold of distinctiveness. Finally, we perform temporal ranking, based on information gain, to identify results that are particularly relevant to the time period of interest. We evaluate our work over a collection of tweets gathered in late 2012, exploring the efficiency and filtering characteristic of each processing step, both individually and collectively. Based on our experience, the resulting synopses from various time periods provide understandable and meaningful pictures of events within those periods, with potential application to tasks such as temporal summarization and query expansion for search.

Acknowledgements

I would like to thank all the people who made this possible. First and foremost, my supervisor Charles Clarke for giving me the opportunity to join such a great university, and for providing a very encouraging learning environment. Second, my family who gave me a lot of support – and the free voice chat services that kept me close to home. Last but not least, my colleagues in this journey of development with whom I worked and grew.

Dedication

This is dedicated to the people struggling for freedom and democracy in the Arab world, while I comfortably study in beautiful Canada.

Table of Contents

List of Tables	ix
List of Figures	x
1 Introduction	1
1.1 Social Media Text	1
1.2 Frequent Itemset Mining	2
1.3 Motivation and Contributions	3
1.4 Thesis outline	4
1.5 Terminology and Notation	5
2 Background and Related Work	6
2.1 Frequent Itemset Mining	6
2.1.1 The Apriori Algorithm	6
2.1.2 The FP-Growth Algorithm	7
2.1.3 Solution Space Algorithms	8
2.2 Election of Representative or Interesting Itemsets	11
2.3 Frequent Itemsets in Search and Summarization	15

2.4	Social Media as a Source of Information	17
2.5	Topic Detection and Tracking for Social Media	19
3	Mining Social Media Text	23
3.1	The Linear-Time Closed Itemset Mining Algorithm	23
3.1.1	Implementation Details	27
3.2	Using LCM for Mining Social Media Text	29
3.2.1	Dataset	29
3.2.2	Timeseries Model of the Dataset	29
3.2.3	Dynamic Support Threshold	34
3.2.4	Sliding Window Model	34
3.3	Filtering Language Constructs	37
4	Selecting and Ranking Itemsets	42
4.1	Social Media Dynamics	42
4.2	Distinct Itemsets	44
4.3	Strongly Closed Itemsets	48
4.3.1	Bounding Space and Time Requirements	55
4.4	Temporal Ranking	55
5	Evaluation	58
5.1	Efficiency	58
5.2	Effectiveness	60
5.2.1	US presidential elections day	61
5.2.2	Google Trends	65

6 Conclusion and Future Work	72
6.1 Contributions	72
6.2 Results	74
6.2.1 The Complete Results of Mining the Twitter Dataset	75
6.3 Future Work	75
APPENDICES	76
A Unigram Tokenization	77
References	78

List of Tables

3.1	Generation of closed itemsets by Prefix Preserving Closure Extension . . .	26
5.1	Top 3 (or 5) <i>strongly closed</i> itemsets for hours in US presidential elections day	63
5.2	The most interesting itemsets out of the top 10 picked by the MTV algorithm for hours in US presidential elections day. The actual rank given to each itemset by MTV is in column 3 and the rank of its equivalent <i>strongly closed</i> itemset is in column 4.	64
5.3	Top 10 itemsets in an hour in the morning of the day of MTVEMA	68
5.4	Top 30 itemsets from 18:30 EST on the day of American Music Awards (AMAs)	70

List of Figures

2.1	A few steps in the FP-Growth algorithm as depicted by Han et al. [31]	8
2.2	Example of a Galois lattice adapted from Pasquier et al. [66]	10
3.1	Seasonal and noise components of the volume velocity	32
3.2	Stochastic level of the volume velocity	33
3.3	Mean runtime of LCM at different epoch spans	36
3.4	Mean support corresponding to minimum support 10, at different epoch spans	36
3.5	Average of token frequency percentiles at different times of day and maximum N-gram length	40
3.6	Effect of changing the maximum N-gram length on mining hour long epoches	41
4.1	Closed, distinct and strongly closed sets	43
4.2	Scatter plot of <i>distinct</i> itemsets from the hour starting at 10 PM EST on 6 Nov. 2012	47
4.3	Itemsets containing ‘barack’ or ‘romeny’ that are closed but not <i>distinct</i>	48
4.4	The hierarchy of itemsets in the cluster pertaining to the “sham and travesty” tweet	53
4.5	Several clusters for itemsets pertaining to the MTV Europe Music Awards votes	54

5.1	Effect of changing κ on mining results	59
5.2	Runtime of two itemset mining algorithms alone and of the LCM algorithm with filtering and clustering, at different epoch spans	60
5.3	Runtime of the MTV algorithm [57] at different resultset sizes (top k) . . .	65
5.4	Google Trends' volume with time plot for the queries "Barack Obama" and "Mitt Romney", the top 2 in the Top Chart for the <i>people</i> category in November 2012	66
5.5	Google Trends' volume with time plot for the queries about pop artists in the Top Chart for the <i>people</i> category in November 2012	67

Chapter 1

Introduction

1.1 Social Media Text

Text posted on social media outlets, such as Twitter, is a good and timely source of information about real world events [72]. It also includes posts about other topics attracting the collective attention of user communities, such as Internet memes or large scale discussions [60]. Nevertheless, the collective stream from all users is overwhelmed by personal updates and non-informative chatter [39, 35]. Furthermore, the attention span given to the majority of topics is quite short [45]. Finding posts about topics of interest in a timely fashion requires an efficient mining algorithm.

To realize the benefits of social media in giving a voice to ordinary people, the algorithm should mine the content of posts without explicitly assigning weights based on popularity of users. However, the nature of text in social media poses a challenge when applying traditional text mining algorithms. Text in social media is usually short, undermining the effectiveness of within document frequency counting. It lacks context, since posts are standalone and most platforms don't allow explicit linkage. It also lacks structure and other useful formatting cues. One type of mining algorithms that can be applied to such data is frequent itemset mining.

1.2 Frequent Itemset Mining

Frequent itemset mining algorithms count the number of times each combination of “items” appear together. This is called the support of the itemset. An itemset is considered frequent if its support is above a threshold. When mining textual data, an item is usually taken to be a unigram. A set of items are considered to be appearing together if they appear within the same document or paragraph. The frequent itemset mining family of algorithms is fast and efficient, however it is not readily suited for application on text. Following are a number of challenges faced when applying frequent itemset mining to textual data. In this thesis we address those problems, adapting frequent itemset mining to social media text without degrading its efficiency:

1. The number of itemsets mined grows with the number of distinct items – which is particularly high in the text domain. In social media, the problem is complicated by the continuous creation of new tokens [53], in the form of hashtags or usernames, and due to shortening of words because of the length limit of some platforms.
2. The number of frequent itemsets is generally high. Frequent itemset mining was originally proposed as a preliminary stage for association rules mining, which sifts through the numerous itemsets and produces a smaller number of rules associating combinations of items with each other. To reduce the number of itemsets, they may be limited by setting a high support threshold. This is not possible in text mining because the frequency of most items is low, and only a few items have high frequencies; that is, frequencies of items follow a long-tailed Zipfean distribution.
3. Function words are among the few items having high frequencies. This leads to mining a large number of itemsets which are uninformative language constructs. Even if a maximum frequency threshold is set, incurring the risk that we will filter out important itemsets, many non-English constructs will be mined because the proportion of posts in English is much higher than other languages. Notice that we do not remove function words to allow mining itemsets made up solely of function words, and to make the mining results easier to understand at the user level.

4. There is considerable redundancy in frequent itemsets caused by trivial differences in the language used.

1.3 Motivation and Contributions

Frequent itemset mining is suitable to the dynamic nature of social media streams. It does not require prior knowledge of the distribution of items, nor does it require selecting a few items to monitor or a number of topics to mine. It is fast, efficient and scalable. It is also robust against data sparsity, and can actually exploit this sparsity for faster computation. Unlike trending topics¹ [58], the results of frequent itemset mining include itemsets that have high frequency because of sustained interest, as well as a spike of interest.

The mined itemsets provide the vocabulary associated with events and can be used as a preliminary step for search and summarization. For example, the collection of mining results from different epochs of time can be used for temporal query expansion and document expansion [20, 24]. Expansion methods can be specially developed based on the mining results, but it is also possible to use existing methods. The results from each epoch can be treated as a document, and relevant “documents” can be used to create a *temporal profile* [40] of the query or document being expanded. For summarization, frequent itemsets can provide a good foundation for summary creation. While the frequent itemsets themselves are not summaries, since they lack qualitative properties such as coherence and cohesion, the results are understandable at the user level. As we shall see in later examples, the top ranked itemsets cover a variety of open topics, and within one topic different opinions are reported as separate contrastive itemsets.

Previous work has proposed methods tailored for the use of frequent itemsets (patterns) to improve search performance. Itemsets provide better semantics than keywords, and have better statistical properties than phrases [87]. By using carefully selected itemsets and filtering out “meaningless” ones, *pattern-based* approaches achieve better performance than keyword-based approaches [88]. However, the choice of “interesting” or “informative” itemsets and filtering out “meaningless” ones remains an area of active research. In Li et

¹<http://blog.twitter.com/2010/12/to-trend-or-not-to-trend.html>

al. [51] itemsets are mined from paragraphs of newswire text, and are used to determine term weights for query expansion. Improvements in performance have been achieved by using itemsets taken from a training set of related documents, as well as ones from unrelated documents. In a more recent work, an unsupervised version that relies on co-occurrence of itemsets within the same paragraph has been proposed [4].

In this thesis, we propose efficient methods for filtering out non-informative frequent itemsets, reducing redundancy in the mining results, and ranking the selection according to temporal novelty. Our methods exploit the dynamics of social media and make use of the collaborative filtering that users naturally undergo on social media, by sharing interesting posts and participating in conversations. Our main contributions are as follows:

- We propose the use of variable length term N-grams as items to avoid mining uninformative language constructs as itemsets.
- We propose a condition for selecting informative itemsets, and a clustering scheme for reducing redundancy in the mining results.
- We propose a formula for ranking itemsets according to temporal novelty.

The effectiveness of our methods is shown quantitatively, in terms of the drop in the number of itemsets at each processing step. The quality of the final outcome is verified by showing mining results from various time periods.

1.4 Thesis outline

The next two chapters provide necessary background. We start by discussing related work in chapter 2, then we explain the algorithm on which we build our work in chapter 3. In chapter 3, we also introduce the dataset we use, and the data preprocessing necessary for running a frequent itemset mining algorithm on social media text effectively. The most important contribution in this chapter is a simple method for flattening the Zipfean distribution of items, and thus reduce the number of itemsets made up of function words. Chapter 4 presents our main contributions for creating temporal synopses from frequent

itemsets. We present in it our proposed condition for selecting itemsets, the clustering scheme for grouping together itemsets according to their topic, and a ranking formula that is suitable for ranking clusters according to temporal novelty (based on the collective Information Gain of itemsets in them). In chapter 5, we empirically evaluate the efficiency and the effectiveness of our proposed methods. Finally, we conclude our presentation and suggest future directions in chapter 6.

1.5 Terminology and Notation

Classically, frequent itemset mining is applied to a *database of transactions* made at a retail store. This terminology is suitable for market basket data and we retain it out of convention, even though we are mining text where the terms “corpus” and “document” are normally used. Because of the dynamic nature of social media, rather than giving the whole database as input to mining algorithms, the input is an *epoch* of data; data with timestamps within a certain period of time. The epoch’s *span* is the length of this period in hours, and the *volume velocity* at this epoch is the number of transactions in the epoch divided by its *span*.

We now introduce the notation used in this thesis:

- $W = \{w_1, w_2, \dots, w_n\}$: The set of all items occurring within the epoch of data being mined. Items can be terms or term N-grams in this thesis.
- $t_a = \{w_{a_1}, \dots, w_{a_m}\}$: A transaction made up of a set of items. Each transaction has a sequential id derived from its timestamp (denoted by the subscript letter).
- $E^{span} = \langle t_a, t_b, \dots, t_v \rangle$: An epoch of data of a certain span, such as an hour, made up of the sequence of the transactions created within this hour.
- $s \subset W$: An itemset; any possible combination of items.
- $T_s = \{t : t \in E \text{ and } s \subseteq t\}$: All transactions containing itemset s . We refer to it as the itemset’s postings list, as is common in the Information Retrieval (IR) literature.

Chapter 2

Background and Related Work

2.1 Frequent Itemset Mining

Frequent itemset mining comprises a large body of work that goes back to the early 1990s [2]. We categorize frequent itemset mining algorithms into 1) algorithms that operate in the item space, 2) algorithms that operate in the transaction space, and 3) algorithms that operate in the solution space. While a comprehensive survey is beyond the scope of this thesis, we discuss at least one representative from each category of algorithms.

2.1.1 The Apriori Algorithm

Agrawal et al. [3] proposed the Apriori algorithm as part of the first efficient solution to the problem of association rules mining [2]. In the domain of market basket data, an association rule is an implication that an item is likely to be purchased in the same transaction with a certain set of items. Association rules are mined by finding itemsets that co-occur together frequently, then producing a rule with each individual item as a consequent to the purchase of the rest of the itemset. The Apriori algorithm finds frequent itemsets of increasing lengths iteratively. It starts by counting frequencies of individual items (itemsets of length 1). Then it iteratively increases the length of itemsets. At each

iteration, it generates candidate itemsets of length k (k -itemsets) by merging frequent itemsets of length $(k-1)$ ($(k-1)$ -itemsets) that differ in only 1 item, usually the last item given a certain total ordering of items. By using only the frequent $(k-1)$ -itemsets for generating candidate k -itemsets, many possible k -itemsets are implicitly pruned, based on the anti-monotonicity between itemsets' support and length: all subsets of a frequent itemset have to be frequent. However, each candidate has to be explicitly checked to verify that it does not have any infrequent subsets.

Apriori and algorithms based on it suffer performance degradation and large increases in memory requirement when the number of distinct items is high. These limitations are caused by the candidate generation bottleneck. A large number of candidates can be generated, especially in early iterations of the algorithm. Consider, for example, the generation of candidate 2-itemsets from a database. This generation requires producing all unordered pairs of 1-itemsets (terms), after pruning out rare ones with frequency less than the support threshold. In many domains, including text mining, the number of frequent 1-itemsets is large enough to prohibit generating a number of candidates in the order of this number squared. In text mining, a rather low support threshold has to be used, because the frequency of terms follow a long-tailed Zipfean distribution.

The Apriori based family of algorithms operates in the item space in the sense that an itemset can be expanded by any item regardless of whether there is any support for the combination. Operating in the transaction space avoids this blind enumeration of unsupported candidates, and the subsequent DB scan required to count the support of those candidates. This involves the creation of a succinct representation of the DB that is well suited for itemset mining.

2.1.2 The FP-Growth Algorithm

The first algorithm proposed based on the idea of avoiding candidate generation was FP-Growth [30]. The database is represented as a frequency ordered prefix tree called the FP-tree. Items with low support are removed, so that any node in the FP-tree is necessarily an item with enough support. An FP-tree imposes the invariant that within each branch the frequency is non-increasing from the root down to the leaves. This increases the chance



Figure 2.1: A few steps in the FP-Growth algorithm as depicted by Han et al. [31]

of finding common prefixes. An FP-tree also contains links between nodes representing the same item. Those links are used to create the *conditional* FP-tree for an item, which is the FP-tree made only from branches containing this item, excluding any items with lower frequency. The algorithm proceeds by iterating over items in increasing order of frequency. Each item is first output as an itemset, and then its *conditional* FP-Tree is created and recursively mined. At each recursive call, the current item is appended to a growing prefix, which is prepended to all itemsets mined from the conditional FP-tree. Figure 2.1, borrowed from Han et. al [31], shows a few steps in the FP-Growth algorithm.

The navigation of a representation of the database avoids the candidate generation bottleneck, by directly deriving itemsets from the transaction space. However, the creation of such representation is a bottleneck itself. Building conditional pattern bases and sub-FP-trees becomes very time and space consuming as the recursion goes deep and the number of patterns goes large [85].

2.1.3 Solution Space Algorithms

The last category of frequent itemset mining algorithms comprises algorithms that operate in the solution space. It can be argued that algorithms that traverse a representation of the DB are operating in the solution space rather than the transaction space – given that they avoid generating candidates with low support. However, in our categorization

we focus on the logic behind how itemsets are generated or pruned rather than the data structure used to generate them. After all, any algorithm must calculate the support of itemsets, and this requires the algorithm to either maintain a representation of the DB or scan the transactions. Algorithms that operate in the solution space prune itemsets that do not possess certain properties. The limited number of itemsets that possess the desired property are sufficient to deduce the rest of the solution (other frequent itemsets and their support). The reduction in the number of itemsets mined improves the performance of such algorithms, even if they rely on candidate generation to enumerate possible itemsets before checking for the desired property. It is not always necessary to produce the full solution since the properties usually filter out only redundant itemsets, and in many cases the elect itemsets can be used directly.

The *closure property* [65, 66, 96] prunes an itemset if it has the same support as its supersets. It is easy to see that all itemsets and their support can be derived from the set of *closed* itemsets. A formal proof is given by Pasquier et al. [65]. The smallest itemset of an equivalence class of itemsets having the same support is called a *generator* [44] or a *free set* [12]. It is also possible to keep only the generator of each equivalence class, but in this case some infrequent itemsets has to be kept in order to be able to calculate the support of all frequent itemsets [44]. A different selection of itemsets can be made according to the Non-Derivable property [17, 18], which is linked to the closed property.

Finding closed itemsets can be done by arranging possible itemsets into a Galois lattice, then traversing the lattice using the Galois connection between itemsets and transactions in which they appear. Figure 2.2 shows an example of a Galois lattice. The Galois connection is a pair of operators, one to map an itemset to transactions in which it appears, and another to map a set of transactions to the largest itemset that appears in all of them. Applying the two operators in cascade grows an itemset directly to its closed superset. More details can be found in Pasquier et al. [66] and Zaki et al. [96].

Some algorithms take advantage of the closure property to reduce the search space more dramatically, such as CLOSET [67] and LCM [81]. Our work is based on LCM (which stands for Linear-time Closed itemset Mining) [81]. As a starting point for our work, we use the implementation of LCM submitted to the workshop for Frequent Itemset Mining Implementations (FIMI) in 2004 [41], which was the workshop’s award winner. The

TID	Items
1	A B C
2	B C E
3	A B C E
4	B E
5	A B C E

The transaction database



Figure 2.2: Example of a Galois lattice adapted from Pasquier et al. [66]

LCM algorithm is robust against data sparsity, and has low memory requirements since it does not store intermediate results. We will explain the closure property and describe LCM in detail in section 3.1. We proceed by discussing other areas of related work.

2.2 Election of Representative or Interesting Itemsets

The problem of having many itemsets, with redundancy and noise, can be addressed by picking out representative itemsets that satisfy a certain condition. The representative itemsets are not necessarily sufficient to deduce the rest of the solution (itemsets and their support information). In this thesis we propose a condition for selecting interesting itemsets, which we call *distinct* itemsets. Related approaches to picking out representative itemsets are 1) choosing a set that will minimize the reconstruction error of the full solution, and 2) limiting the number of itemsets to a user specified number.

The δ -free [13], δ -covered [90], and *master* itemsets [92] are examples of conditions motivated by providing a compressed representation of itemsets through sacrificing support information. The δ -free condition selects an itemsets if its support is different from the support of all its subsets at least by δ . This selection of itemsets can be used to approximate the support of other frequent itemsets with a guarantee on the error, but only if some infrequent itemsets are included in the selection. The δ -free condition implies setting an upper bound on the strength of association rules that can be derived from selected itemsets. It is suitable for compressing mining results from dense datasets, but it is exactly the opposite of *distinct* itemsets we propose in this thesis.

The δ -covered condition is also the opposite of the *distinct* condition, as it “relaxes the closure condition to further reduce pattern set size” [55]. On the other hand, the *distinct* condition strengthens the closure condition. The δ -covered condition sets an upper bound on the confidence of the rule that an itemset implies a superset; if an itemset implies its superset with confidence greater than δ , then it is considered to be covered by the superset and can be pruned. On the other hand, the *distinct* condition sets a lower bound on the confidence of a derived association rule. Nevertheless, the clustering scheme Xin et al. [90] propose (δ -clusters) is similar to the *strongly closed itemset* clusters proposed in this

thesis. The efficient version of the clustering algorithm, RPLocal, limits the search space for clustering candidates using a similar reasoning to the reasoning used in our algorithm – exploiting the sequence in which frequent itemsets are found. The main difference between δ -clusters and *strongly closed itemset* clusters is that δ -clusters are soft clusters, while *strongly closed itemset* clusters are hard clusters. Also, the distance measure used in δ -clusters is the Jaccard distance, while the *strongly closed itemset* clusters are based on association rule confidence.

Yan et al. [92] proposed a method to choose k itemsets as representatives of the itemsets in the full solution, which they call *master* itemsets. Rather than leaving k as a user specified parameter, a method is proposed for choosing k by minimizing the restoration error of support information. A *master* pattern is the union of all itemsets in a cluster, similar to our proposed *strongly closed itemsets*. While the use of clustering is similarly motivated by the trivial difference between itemsets, the k *master* itemsets have to cover all itemsets. On the other hand, *strongly closed itemsets* actually avoid clustering together itemsets representing different topics or different opinions within a topic. Furthermore, Yan et al. improve performance by approximating the calculation of the distance between itemsets, to avoid accessing the data. We can achieve high performance without approximation because our mining results include the postings list of each itemset.

Another approach to choosing itemsets is according to their interestingness. The interestingness of itemsets (patterns) has been an area of active research for a long time. According to Silberschatz and Tuzhilin [75] interestingness of an itemset can be *objective* or *subjective*. Commonly used *objective* measures include (a more extensive coverage is provided by Geng et al. [26] and Tan et al. [77]):

1. *Support*: The number (or fraction) of transactions that contain the itemset.
2. *Confidence*: The confidence of a rule derived from the itemset. A rule is derived from the itemset by treating one item, w , as the consequent of the rule while the rest of items, $s \setminus w$, are the antecedent. The confidence of the rule is then defined as:

$$conf(s \setminus w \rightarrow s) = \frac{|T_s|}{|T_{s \setminus w}|} \quad (2.1)$$

Since different rules with varying confidence can be derived from an itemset, either the minimum (*all-confidence*) or the maximum (*any-confidence*) is used [63].

3. *Lift* [16]: The probability (support) of the itemset over the product of the probabilities of all items in the itemset. This is a measure of dependence and can be linked to Pointwise Mutual Information [14]. Such a measure is suitable for Knowledge Discovery or Exploratory Search tasks, where the goal is to extract correlated item sets or significant phrases.
4. *Bond (or Coherence)* [63]: The support of an itemset over the number of transactions containing *any* of its items.

An important property of a measure, upon which all Apriori based mining algorithms are built, is the downward closure property; the value of a downward closed measure cannot increase as the size of the itemset increases. All the above measures are downward closed [9] with the exception of *lift* which is upward closed (cannot decrease with the increase of the size of the itemset), and *any-confidence* which is neither upward nor downward closed. Algorithms that mine itemsets based on the above measures were proposed [47, 42], but if an application specific measure is to be used then it has to be evaluated in a post mining stage or within the Redundancy-Aware Top-K Framework [89].

The Redundancy-Aware Top-K Framework [89] selects the top k itemsets according to an interestingness criterion. The major difference between their approach and all of the previous work is that it emphasizes both interestingness and redundancy when selecting the top k itemsets. The itemset interestingness is defined in the context of the application; summarizing the whole collection of itemsets is not the goal. Our work fits in this framework: we reduce redundancy through filtering and clustering, and we propose a ranking that orders itemsets according to some definition of interestingness. If only the top k itemsets are desired, itemsets ranked at positions higher than k can be omitted. An explicit number of itemsets needs not be set, however. This makes our computational model more efficient, since we do not need to solve a constrained optimization problem. In their experiments, Xin et al. [89] used two application specific interestingness measures: they used one that is based on TF-IDF for extracting interesting itemsets from a text corpus,

and another that is tailored for predicting block fetches in storage systems [52]. We rank itemsets according to their novelty using a measure based on Information Gain. Preliminary experimentation with various other measures of interestingness ruled them out as ineffective for our purposes.

Subjective measures of interestingness select itemsets that are actionable or unexpected [76]. Selecting unexpected itemsets entails building a probabilistic model to use it for estimating itemsets probabilities. Jaroszewicz and Simovici [38] use a Bayesian Network as a model, which is suitable for cases where items denote discrete values of a limited number of attributes. A commonly used model is the maximum entropy model [83, 78, 57]. The Maximally informaTiVe itemsets (MTV) algorithm [57] uses a greedy algorithm to select k itemsets whose probabilities diverge the most from the estimate given by a maximum entropy model. The maximum entropy model is initialized as a uniform distribution. After selecting each of the k itemsets it is updated to fit the selected itemsets actual frequencies using the *iterative scaling* procedure [23]. If k is set to infinity, the algorithm will generate the best model according to the Bayesian Information Criterion (BIC). Since BIC incorporates a penalty term for the number of parameters (itemsets included in the model), this selection will minimize redundancy while maximizing surprise (divergence from model). The model can be initialized using itemset frequencies known from prior knowledge. It is also possible to use the Minimum Description Length (MDL) as a criterion for selecting itemsets. We compare our selection of itemsets to ones obtained from MTV using the implementation provided by the authors¹. However, the value of k has to be kept low (at most 100) for the algorithm to finish without an error, even though it is running without a limit on its resource usage on a machine with 256 GB of RAM. In our work we focus on efficiency, allowing the summary to include larger numbers of itemsets than the values of k for which MTV finishes in a reasonable time (at most 10). This is crucial for scalability to the volumes of data typical in social media streams.

Another approach to picking out itemsets is choosing ones that can be used to compress the data (not the itemsets). The KRIMP algorithm [82] is a good example of methods that follow this approach. Our goal is different because we aim to filter out itemsets pertaining to personal updates, which make up a large portion of social media data.

¹<http://www.adrem.ua.ac.be/implementations>

2.3 Frequent Itemsets in Search and Summarization

Different ways to use frequent itemset mining in search have been proposed. Possas et al. [70] proposed using closed itemsets in place of terms, using them in exactly the same way as terms are used. In the tokenization process, they use CHARM [96] to mine closed itemsets of terms from each document. They then add the document to the postings lists of each closed itemset it contains. Queries are tokenized similarly and results are ranked using TF-IDF. This simple approach achieves increased precision on various test collections, possibly because frequent itemsets take into account collocation information. The approach is simplistic, however. The overlap between itemsets was not taken into account while tokenizing documents and queries into itemsets. Also, the mining algorithm was not modified to accommodate the high number of terms and the greater length of documents compared to the length of transactions typical to market basket data. Finally, the number of itemsets is higher than the number of terms so the index size will increase.

Wu et al. [88, 87] proposed the use of sequential closed patterns in place of terms, in what they call a *pattern based approach* to information retrieval. They consider patterns as surrogates to phrases that have better statistical properties than phrases, thus achieving a break through in the performance of *phrase based approaches*. “Although phrases contain less ambiguous and narrower meanings than individual words, the likely reasons for the discouraging performance [50] from the use of phrases are: (1) phrases have inferior statistical properties to words, (2) they have low frequency of occurrence, and (3) there are a large number of redundant and noisy phrases among them.” [87] They addressed the problem of overlap between itemsets, and they developed an algorithm specifically tailored for mining closed sequential patterns from text. Documents are broken into paragraphs such that a paragraph is the unit of mining and retrieval. This alleviates the problem that documents are usually much longer than traditional market basket transactions. Albathan et al. [4] extended the work of Wu et al. by proposing a method for further reducing the number of closed patterns. Lau et al. [46] used related methods for term weighting in pseudo-relevance feedback for Twitter search, and achieved substantial improvements over a baseline.

In an earlier work [1], we have also explored the viability of using frequent itemsets for query expansion in microblog search. We could achieve improvements over a strong baseline on the TREC 2011 microblog track queries. As a baseline, we used Okapi BM25 [73] with its parameters tuned to achieve the best performance on short documents. The baseline’s performance ranked high compared to the TREC 2011 submissions. After query expansion, the performance was better than all of the TREC 2011 submissions. Queries were expanded using terms from itemsets that are most relevant to the query, according to the same relevance ranking formula used in the baseline. We used frequent itemsets mined by an unmodified mining algorithm that returns the top k itemsets according to their support². However, query expansion achieved no improvement over the baseline on the TREC 2012 microblog track queries. This indicates that there was not enough relevant itemsets, because of the use of support to limit the number of itemsets. The methods proposed in this thesis can be used for mining a short ranked list of relevant itemsets, which can then be used by methods such as our earlier work, or the work of Wu et al.

This work also complements the work done by Yang et al. [93] for using frequent itemsets as a temporal summary. They have proposed a framework for storing mining results of temporally consecutive intervals (batches of data) using a pyramidal time window. Their framework allows for fast execution of temporal queries, and for tracking the evolution of topics. However, the itemsets stored for each interval are not selected for providing the best summary. Instead, they are selected according to their utility in compressing the data. The temporal summary is created from transactions relevant to a query in a separate step. First, the pyramidal time window is queried using keywords and a specific time intervals (for example, “world cup” on the day before a certain match). Then, non-negative matrix factorization is used to extract topics from the returned transactions. The pyramidal time window structure can be used to store itemsets selected by our methods, thus directly providing a temporal summary without the need for the matrix factorization step. Our methods exploit the nature of social media to choose topically relevant itemsets without the need for a query. Social media has been shown to be a good source of real-time information, as we will discuss in the next section.

²<https://cwiki.apache.org/MAHOUT/parallel-frequent-pattern-mining.html>

2.4 Social Media as a Source of Information

The problem of finding information in social media posts is motivated by the uniqueness and variety of information that can be found in these media. Some stories are reported mainly on social media because of censorship on journalism, such as the 2009 Iranian elections protests (nicknamed “Iran’s Twitter Revolution”³). Other stories are reported mainly on social media because the content creator happens to be at the right place at the right time. This could be by chance, like in the case of the emergency landing of a US Airways plane on the Hudson river in New York⁴. This could also be an intentional act of citizen journalism. Citizen journalists are “the people formerly known as the audience,” who “were on the receiving end of a media system that ran one way, in a broadcasting pattern, with high entry fees and a few firms competing to speak very loudly while the rest of the population listened in isolation from one another – and who today are not in a situation like that at all.”⁵

In a study of the use of Twitter as a news medium during the 2011 Tunisia and Egypt revolutions, Lotan et al. [56] categorized 963 accounts that made influential posts according to the *type of actor*. A post is considered influential if it is among the top 10% of posts in terms of the number of near duplicates made out of it; posts that had at least 16 near duplicates in the Tunisia data, or at least 19 near duplicates in the Egypt data. The types of actors included *mainstream media organizations* and *mainstream new media organizations*, which are news and media organizations with and without offline presence, respectively. It also included *journalists*, who are individuals employed by media organizations or who regularly work as freelancers for them.

In both datasets, media organizations make up about 11% only of the 963 accounts (approximately 7% for ones with offline presence and 4% for ones that exist solely online). Journalists make up another 15%, so the total of mainstream media affiliated accounts is

³<http://www.theatlantic.com/technology/archive/2010/06/evaluating-irans-twitter-revolution/58337/>

⁴<http://www.telegraph.co.uk/technology/twitter/4269765/New-York-plane-crash-Twitter-breaks-the-news.html>

⁵http://archive.pressthink.org/2006/06/27/ppl_frmr.html#more

approximately 26%. On the other hand, more than half of the accounts from which influential tweets originated belong to individuals not affiliated to the government or mainstream media (55% for Tunisia and 56% for Egypt). Out of those individual, almost half are ordinary people who don't self identify and cannot be identified by assessors as bloggers or activists.

The rate at which an account makes posts is approximately 16 per hour for mainstream media affiliated accounts, and approximately 10 for individuals accounts (thus a mainstream media account makes 60% more posts). Also the probability that a post will be echoed is 0.88 for mainstream media accounts, 0.55 for journalists, 0.4 for bloggers and activists, and 0.31 for individuals (this is the probability that any follower will echo the post, the average number of followers of an account belonging to an ordinary person is much lower than that of an account belonging to an actor of any of the other types). However, collectively the fraction of highly echoed threads originating from ordinary people is on a par with the fraction of threads originating from a journalist. Threads originating from ordinary people, bloggers and activists together make up more than half of highly echoed threads, followed by threads originating from journalists (around 20%), then threads originating from mainstream media (around 10%, divided equally between media with and without offline presence).

All this emphasizes the influence citizen journalists and ordinary people now have on the information creation and dissemination process. User created content is rich with unique content tackling a variety of topics from different points of view. However, it is dominated by non-informative chatter and personal updates. In a study comparing the use of Twitter's search feature to Web search, Teevan et al. [79] found that Twitter is mostly used for getting timely information: updates about events, news about topics of interest (such as technology news), and realtime local updates (such as traffic jams). It was noted that Twitter search could be the main way of accessing Twitter, since posts from the followed network could be outside of the user's topical interests. This is supported by findings from an earlier study by Naaman et al. [61], in which 3,379 tweets from 350 accounts belonging to individuals were studied. The tweets were hand coded into different categories, then the activity of users in terms of what type of tweets they make was analyzed. It was found that 40% of the studied tweets are personal updates (about

“*me now*”), and that users can be clustered into a big cluster of “*meformers*” (80% of the users, half of their posts are personal updates) and a small cluster of *informers* (20% of the users, half of their posts are for sharing information). The study also found that most people engage in some form of *meformer* activity; on average, users had 41% of their messages in the “*me now*” category. Therefore, filtering out posts by *meformers*, or giving higher weights to posts by *informers* would not be sufficient for filtering posts about “*me now*” and finding informative posts made by ordinary people.

To tap into social media as a source of information, social media users collaboratively select good content. On Twitter, users can *retweet* an interesting post, which means forwarding it to their followers. The dynamics of retweeting was studied by Boyd et al [15]. One of the findings was that posts being retweeted are frequently modified to accommodate additional text, or at least the notation meant to indicate that the post is a retweet. The limit on tweets length might be the reason behind this modification, but it can also be a way to emphasize what the retweeter finds most interesting. The selection of what to keep from the post is actually a form of collaborative filtering itself, because the most interesting parts has to be retained. Parts that are selected by many people are likely to be good summaries of what is interesting. Using frequent itemset mining, the words in those selections would be mined as itemsets. In chapter 4 we discuss how to exploit such social media dynamics to create temporal synopses. In the next section, we discuss other methods for finding emerging topics in social media.

2.5 Topic Detection and Tracking for Social Media

The problem of Topic Detection and Tracking (TDT) was introduced by Allan et al. in 1998 [6]. Topic detection is the discovery of previously unknown topics, where a topic was first defined as an “event” then its definition was broadened to include the triggering event as well as other events and activities directly related to it. Our work can be categorized as a method for topic detection from social media, or more specifically “event” detection since we make no effort to group events into topics. Ever since it was introduced, TDT has been an area of active research. Besides the fact that there is always room for improvement,

TDT was originally targeting the domain of newswire data, and the solutions developed for this domain are not directly applicable to data from other domains.

In the domain of social media, Petrović et al. [68] used an adaptation of Locality Sensitive Hashing (LSH) [36] to perform first story detection at scale. The computation framework is the same as the one used in traditional systems [94, 7]:

- The distance between a tweet and its nearest neighbour from the preceding tweets is calculated.
- If the distance is less than a threshold, then the tweet is added to the topic represented by the cluster containing the nearest neighbour.
- If the distance is greater than the threshold, a new cluster is created. The new cluster represents a new topic, with the current tweet as its first story.

The adaptation of LSH that Petrović et al. propose reduces the number of times a tweet is considered far from all of its neighbours while it actually is not. In case of failure to find a near neighbour, a tweet is compared with the last 1000 tweets, updating the minimum distance if necessary. The performance of the modified LSH on a traditional TDT task (TDT-5) was comparable to the UMass FSD system [7]. On the other hand, the run times of the two systems are 2 hours and 28 hours respectively. Another experiment was done on a dataset of a 160 million tweets collected over a period of 6 months from Twitter. The UMass system failed to terminate, while the LSH based system could process each tweet in bounded space and time. One problem remains, however. Their system does not achieve high precision in determining if a tweet cluster pertains to an event or not. Actually, their system can detect spam clusters with high precision (*average precision* = 0.963), but not event clusters (*average precision* = 0.34). The problem of identifying event clusters was further studied by Beker et al. [10], where they used a Support Vector Machine (SVM) classifier and achieved a higher precision. However, the test set used was limited to 5 hours of Twitter data and the precision decreases as the number of clusters increases.

Another approach to event detection follows the framework of bursty keywords detection proposed by Kleinberg [43]. In this framework, a set of keywords are tracked – possibly

all tokens in the dataset. A state machine is used to indicate the state of each keyword, where consecutive states correspond to higher levels of *burst* intensity. Each state transition is associated with a cost, and a keyword is moved between states such that the overall cost of its transitions is minimized. A simplified version of this framework is used by Parikh and Karlapalem [64]. In their work, there are only two states for each token: the normal state, and the bursty state. A token is moved to the state of being bursty if the increase in its frequency exceeds a threshold, and it stays in this state for a fixed time interval. The timeline is divided into fixed length intervals, and a token is totally ignored in intervals when it is in the normal state. Tokens are then clustered according to the similarity of the sets of bigrams with which they co-occur, and the pattern of their *appearance* (moving to the bursty state). The clusters are considered to be events, and they are ranked according to their sizes. The clustering they propose has to process data from a long interval, so that the similarity of appearance can be calculated. This approach is not suitable for stream processing. Besides, the number of clusters produced by their methods is small compared to the number of events one would expect to have happened in the interval of time processed (23 in 20 days and 15 in a week).

Another method that tracks the frequency of all tokens is EDCoW [86], where the wavelet transformation is used to convert counts of occurrences from the time domain to the frequency domain then auto-correlation is used to detect burstiness. However, the use of wavelet transformation is not justified specially for the enormous number of tokens in case of social media. Other methods based on burst detection track a small set of keywords. Lehmann et al. [48] track hashtags and study their usage pattern. Chakrabarti and Punera [19] track tokens appearing in tweets tagged by the name of an American football team, and use them to create a summary of a game.

Our work can be seen in the light of meme-tracking [49], where a particularly interesting quotation (meme) is tracked across news and blog posts. The quotation being tracked is not quoted verbatim in all posts, but rather parts of it are selected and it could be slightly paraphrased. This is similar to modifications made during retweeting an interesting tweet. Leskovec et al. [49] tracked quotations made by American politicians during the 2008 U.S. presidential elections. Phrases made up of 4 words or more are tracked if they got repeated 10 times or more in blogs and news articles, where at most 25% of the repetitions could be

on the same domain name. A directed acyclic graph is constructed where vertices represent phrases and edges are weighted according to the edit distance between the phrases, as well as the number of occurrences of the longer phrase. The graph is then divided into clusters where each cluster contains one *root* that has no outgoing edges. Each vertex is assigned to the cluster with which it has the highest number of edges. Clusters are then ranked according to the number of news articles and blog posts containing any phrase from the cluster. Leskovec et al. then use the clusters and their associated news articles and blog posts to perform an analysis of the temporal variation of the volume of posts about each meme. They propose a model that fits the data well, and study the difference between blog posts and news articles in terms of how fast a meme is picked and how long it remains in focus before getting dropped. In our work, we are also using the volume of posts containing phrases and their variations (itemsets). We also do clustering but in a different way, and we work in a totally different scale of time.

O’Connor et al. [62] also start from phrases to summarize tweets about a user specified query. They select *significant* phrases which are statistically unlikely, then cluster them. Frequent itemsets are a better starting point than phrases, since there are much more phrases than frequent itemsets.

Event detection in Twitter can also be done using natural language processing to identify *event* parts of speech [72], or to identify action phrases [69]. Of course, it is also possible to apply topic modelling by Latent Dirichlet Allocation (LDA) [33] or an online variant [32, 95]. Our method tackles the problem in a totally different way than these two approaches.

Chapter 3

Mining Social Media Text

3.1 The Linear-Time Closed Itemset Mining Algorithm

In section 2.1.3 we have discussed frequent itemset mining algorithms that overcome the bottleneck of *candidate generation* by limiting the solution to itemsets with a certain property, reducing the size of the solution and possibly providing a way to quickly navigate the solution space. In this thesis, we expand on the Linear-time Closed itemset Mining (LCM) algorithm [81], which mines only *closed itemsets* [65]. The LCM algorithm is distinctively fast because it also takes hints from the transaction space during candidate generation. This also makes it resilient against data sparsity, since it will not consider a candidate that never appears in the data. The ability to efficiently mine sparse data makes LCM particularly suitable for mining social media text. In this section we explain in detail how LCM works.

Since LCM makes use of the properties of closed itemsets, we begin our presentation by discussing these properties. Informally, a closed itemset contains any item that is present in all the transactions containing this itemset. A formal definition of closed itemsets is given in equation 3.1:

$$\mathcal{C} = \{s_c : s_c \subset W \text{ and } \nexists s_d \text{ where } s_c \subset s_d \text{ and } |T_{s_c}| = |T_{s_d}|\} \quad (3.1)$$

The properties of closed itemsets are as follows:

1. Adding an item to a closed itemset reduces its support.
2. A subset of a closed itemset is not necessarily closed, but one or more closed subset must exist for any itemset (formally this could be the empty set, given that any item that appears in all transactions is removed in a preprocessing step).
3. If a closed k -itemset can be extended any further then one of its supersets will be closed, however not necessarily a $(k+1)$ superset. Itemsets that cannot be extended any further are called *maximal itemsets*, which form a subclass of closed itemsets.

Besides being much smaller than the solution space of frequent itemsets, the solution space of closed itemsets can be navigated efficiently. By using an arbitrary total ordering of items, any closed itemset can be considered an extension of exactly one of its subsets. Thus, only this subset is extended during candidate generation. All other subsets do not need to be extended by items that would lead to the longer closed itemset. This property is called *prefix preserving closure extension (PPC-Extension)* and it was proposed and formally proved by Uno et al. [81].

PPC-Extension is achieved by following three rules, which we state after a few definitions to facilitate their statement. First, an item is *larger/smaller* than another item if it comes later/earlier in the total ordering. This terminology comes from the fact that LCM is most efficient if the items are ordered in ascending order of their frequency. Second, the *suffix* of an itemset is one or more items which have to be removed to get an itemset with greater support. Notice that they are necessarily the last items added to the itemset, regardless of the total ordering. Finally, we call the first item added to the suffix of the itemset its *suffix head*. With this terminology, the rules for *PPC-Extension* are as follows:

1. An itemset must be extended by every item that occurs in $T_{itemset}$, except items which are *smaller* than its *suffix head*; extending by *smaller* items will lead to closed

itemsets already generated in an earlier step. Extension items are added to the itemset in turn in the order given by the total ordering, recursively applying PPC-Extension to the extended itemset before adding the next extension item.

2. After adding an extension item, w_e , to an itemset, s , we add all other items that appear in all transaction containing $s \cup \{w_e\}$ – that is, items whose frequency within $T_{s \cup \{w_e\}}$ is equal to $|T_{s \cup \{w_e\}}|$. The newly added items become the *suffix*.
3. If all items in the *suffix* are *larger* than the *suffix head* then add the itemset to the solution. Otherwise, prune this solution branch; all closed itemsets within this branch have already been generated, when processing the smallest suffix member.

Table 3.1 is an example of how *PPC-Extension* is used to generate closed itemsets starting from the 1-itemset $\{\text{'barack'}\}$. The upper table enumerates $T_{\{\text{'barack'}\}}$. The lower table shows steps of itemsets generation. The current itemset along with its frequency is in column 2. Itemsets marked by an (*) are the closed itemsets that are part of the solution. The suffix of the itemset is shown in *italic*. All possible extension items and their frequencies are in column 3. Extension items that are *smaller* than the *suffix head* are shown with a line striked through them. For each itemset, the extension items are kept so that it is known which extension item is next in turn to be added. An item is bolded when its turn to be added has come. After adding each item, a pass is done on $T_{itemset}$ to enumerate and count possible extension items. To enforce a support threshold infrequent extension items would be removed after counting, but in this example there is no such threshold. Finally, column 4 is a comment explaining each step.

Table 3.1 shows that the number of steps is linear in the number of closed itemsets, and the only additional storage required, besides storage for the documents, is that required for possible extension items. Of course, this is a simplified example, but it shows in essence how LCM achieves its low run time and memory requirements. We refer the interested reader to Uno et al. [81] for a theoretical proof that the algorithm runs in linear time in the number of closed itemsets, and that this number is quadratic in the number of transactions. Performance on a real data set is shown in section 3.2. We proceed by describing how to implement this algorithm using an inverted index.

Doc. Id	Document	Doc. Id	Document
a	barack & mitt	b	brack obama & mitt romney
c	brack obama & romney	d	brack obama

Documents (two per row)

Step	Current Itemset	Possible Extension Items	Comments
1	{barack} (4)	mitt (2), obama (3), romney (2)	Items in $T_{\{\text{'barack'}\}}$ are counted. There is not an item appearing in all transactions.
2	{ <i>barack</i> } (4)*	mitt (2), obama (3), romney (2)	Rule 3: suffix is ordered, add itemset to solution. This is not shown as a separate step in the rest of the example.
3	{ <i>barack</i> } (4)	mitt (2), obama (3), romney (2)	Items are ordered lexicographically. Adding ‘mitt’ to itemset.
4	{barack, <i>mitt</i> } (2)*	obama (1), romney (1)	Extension items reenumerated & counted.
5	{barack, <i>mitt</i> , obama} (1)	romney (1)	Rule 2: ‘romney’ appears in all $T_{itemset}$.
6	{barack, mitt, <i>obama</i> , <i>romney</i> } (1)*		Rule 3: ‘obama’ is the <i>suffix head</i> .
7	{ <i>barack</i> } (4)	mitt (2), obama (3), romney (2)	Nothing more to add, back to {‘barack’}.
8	{barack, <i>obama</i> } (3)*	mitt (1), romney (2)	Rule 1: skipping ‘mitt’, adding ‘romney’
9	{barack, obama, <i>romney</i> } (2)*	mitt (1)	Rule 1: Nothing more to add.
10	{ <i>barack</i> } (4)	mitt (2), obama (3), romney (2)	Back to {‘barack’}, adding ‘romney’.
11	{ <i>barack</i> , romney} (2)	mitt (1), obama (2)	Rule 2: add ‘obama’ after ‘romney’.
12	{barack, <i>romney</i> , <i>obama</i> } (2)	mitt (1)	Rule 3: suffix is not ordered, prune.
13	{ <i>barack</i> } (4)	mitt (2), obama (3), romney (2)	Back to {‘barack’}, all possible extension items were added. Done.

Closed itemsets containing ‘barack’

Table 3.1: Generation of closed itemsets by Prefix Preserving Closure Extension

3.1.1 Implementation Details

We show in algorithm 1 how to implement LCM and PPC-Extension using an inverted index. The algorithm assumes the presence of a search infrastructure, with the ability of performing boolean queries efficiently. Besides the high likelihood that such an infrastructure already exists, it is also a good representation of textual data. Other representations that are designed for databases in general might not be well suited for textual data. For example, the memory requirements of the FP-tree suffers from the sparsity of data, since the data structure is succinct only if it can find common prefixes within the constraints of its invariant. The possibility of using an inverted index in the implementation of LCM is yet another reason why it is well suited for textual data.

The algorithm takes as input an epoch of data and a support threshold. It outputs the closed itemsets with support above the threshold. Along with each itemset in the solution, it also outputs the transactions in which it occurs – which is represented as $\langle items, T_{itemset} \rangle$. The symbol \succ denotes that the lefthand side succeeds the righthand side in the total ordering. In the implementation shown, it is not necessary that the inverted index’s tokens list ($X.tokens$) follow the total ordering; all itemsets of length 1 will be considered anyway. Each 1-itemset is then expanded according to the PPC-Extension rules. Lines 11-18 are the enumeration and occurrence counting of possible extension items, except line 15 which forms the closed itemset according to PPC-Extension rule number 2. Lines 19-20 prune a solution branch according to rule number 3, and line 21 adds the itemset to the solution according to the rule’s complement. Lines 22-27 extend the itemset according to rule number 1, where line 23 (along with line 5) enforces the support threshold.

The algorithm lends itself to distributed implementations. For example, a map/reduce implementation is straightforward since the only operations are counting (line 15) and projection (line 24). However, the fast execution time and the low memory requirements of the algorithm makes it possible that a distributed implementation will cause unnecessary overhead for all but the largest data sets. We experimented with Hadoop¹ and the overhead was in the order of minutes. Besides, many problems start to arise when the programming model is complicated, and Hadoop was particularly problematic.

¹<http://hadoop.apache.org>

Input: *support*: Support threshold

Data: *E*: Epoch of data

Result: *C*: Closed itemsets occurring in at least *support* records

```

1  $C \leftarrow \{\langle \emptyset, E \rangle\}$ ;    //  $\emptyset$  is a closed itemset. This is skipped in practice
2  $X \leftarrow$  Inverted index of E;
3 foreach  $w \in X.tokens$  do
4    $T_{\{w\}} \leftarrow X.postingsList[w]$ ;
5   if  $|T_{\{w\}}| \geq support$  then // Support threshold enforcement on 1-itemsets
6      $LCM(\{w\}, w, T_{\{w\}})$ ;
7 end
8 return C;
9 Function  $LCM(s: \text{Current itemset}, w_{sh}: \text{Suffix head},$ 
10    $T_s: \text{Transactions (tweets) containing } s)$  is
11    $frequency[1 \dots w_n] \leftarrow 0$ ;
12    $suffix \leftarrow \{w_{sh}\}$ ;
13   foreach  $t \in T_s$  do
14     foreach  $w \in t$  do
15        $frequency[w]++$ ;
16       if  $frequency[w] = |T_s|$  then  $suffix \leftarrow suffix \cup \{w\}$ ;
17     end
18   end
19   if  $\exists v \in suffix : w_{sh} \succ v$  then
20     return; // Prune according to PPC-Extension Rule 3
21    $C \leftarrow C \cup \{\langle s \cup suffix, T_s \rangle\}$ ;
22   foreach  $v \succ w_{sh}$  and  $v \notin suffix$  do
23     if  $frequency[v] \geq support$  then
24        $T_{s \cup \{v\}} \leftarrow T_s \cap T_{\{v\}}$ ; // Results of query s AND v
25        $LCM(s \cup suffix \cup \{v\}, v, T_{s \cup \{v\}})$ 
26     end
27   end
28 end

```

Algorithm 1: The LCM frequent itemset mining algorithm

3.2 Using LCM for Mining Social Media Text

3.2.1 Dataset

Throughout this paper we use data collected from the Twitter public stream² since October 1st, 2012. We use only tweets written in the Latin script to facilitate tokenization using white space and other word boundaries. We collect only the tweet text to avoid reliance on any features specific to a certain social medium, and to make the algorithms applicable to other media where text is short such as Facebook or Google+ status updates. YouTube comment data is a particularly promising candidate since comments are limited to 500 characters; the itemsets could provide textual highlights of non-textual content.

We remove duplicate original tweets (not retweets) using a Bloom filter [59]. This filtering removes spam tweets sent by botnets, averaging at 2.86% of the stream. There were two disruptions during data gathering. One in late October because of Hurricane Sandy which made the Twitter service inaccessible. The other was in late December because of technical problems on the server gathering the data.

The average number of tweets per hour (*volume velocity*) in the data is 119,035.49 tweets ($n = 2,546$ hours, standard error = 491), and the average number of tweets per day is 2,705,931.83 tweets ($n = 112$ days, standard error = 55,155). These are simple moving averages calculated by summing the number of tweets in non-overlapping epochs and dividing by the number of epochs. A more accurate estimate of the average number of tweets at different hours of the day can be acquired using a timeseries model of the activity on the social network.

3.2.2 Timeseries Model of the Dataset

Arrivals entering a system are commonly modelled as a Poisson stochastic process. A Poisson process can be characterized by a random variable $N(t)$ representing the aggregate number of arrivals that has happened up to time t . This random variable has a Poisson

²<https://dev.twitter.com/docs/streaming-apis/streams/public>

probability mass function with a rate λ (as well as two other necessary properties: independent and stationary increment). When the rate is high enough (more than 1000), the Normal distribution can be used as an approximation to the Poisson distribution. Since social media has arrival rates larger than 1000 posts per minute (the time unit is arbitrary but must be long enough to have more than 1000 arrivals), it is therefore possible to assume a normal distribution and use a timeseries model to describe the activity on the social network. One class of timeseries models which is particularly useful for analysis is the class of state space models.

Values of a variable in timeseries data is known to be correlated, and specialized models are developed to deal with such explicit correlation. The explicit autocorrelation (correlation between values of the same variable) is usually uninteresting, since it is known and expected. For example, the data can be known to have a trend such as a chronic increase in the volume velocity. Also, a cyclic increase and decrease in the volume velocity according to the time of day can be expected – this is called the seasonal effect in timeseries analysis. State space methods provide an explicit structural framework for the decomposition of timeseries [22]. Unlike the more popular Box-Jenkins ARIMA models, trend and seasonal effects are not treated as nuisance parameters, and it is not necessary to remove the trend and seasonal effects from the series before the analysis can begin. Instead, a state space model has an equation to capture each type of effect that is suspected to contribute to the values of the observed variable. Different models can be used to describe the same observed variable, and the best model is selected using the Akaike information criterion which penalizes extra parameters.

After experimentation with different models, we found that the model that best describes the volume velocity of Twitter is the following:

1. A deterministic seasonal cycle that is 1 day long, which captures the effect of the hour of day on the number of posts.
2. A stochastic level component, which represents the volume velocity specific to each hour. The stochasticity is represented as a variance attached to the level component.

There was no need to add a trend component, maybe because the sampling used by the Twitter API public stream endpoint emits a stable volume of Tweets, and maybe because there was actually no significant growth in the use of Twitter during the months of data collection. In our model we use the timestamps of tweet creation, and we did not keep the timestamps of when the API returned each tweet as part of the sample.

Figures 3.1 and 3.2 show different components of the model when it is fitted to the hourly number of tweets in the month of October 2012. Figure 3.1 shows the daily cycle component (green) and the noise component (black). The noise component must be present in all state space models, to capture small variabilities in the observed value so that the model does not over-fit the data. The cycle is perfectly aligned with the days in the Hawaii Standard Time (HST) time zone, and thus we always use this timezone (GMT-10) to get the date from a timestamp. Figure 3.2 shows the level component (red) along with a scatter plot of the hourly volume velocity. The blue continuous line is the simple moving average and its confidence bands are the blue dashed lines. The simple moving average is clearly not a good estimate of the hourly volume velocity. A better estimate of the number of tweets in an upcoming hour can be predicted from the model using a Kalman filter [80].

The peaks in the level component coincide with real world events. The peaks on the 3rd, 16th and 22nd coincide with the presidential debates³, and the peak on the 28th coincide with the emergency declaration for the northeastern states of the USA in anticipation of Hurricane Sandy⁴. The use of such a model to indicate the presence of interesting information can be investigated further, however it will probably require modelling the number of occurrences of many keywords. This is not scalable because handling correlation between different random variables increases the number of parameters of the model quadratically in the number of variables. The correlation between random variables is captured by a variance-covariance matrix, and if two variables are not known to be independent a parameter must be added in the cell of the intersection of their row and column. We will not use timeseries models any further than this exposition.

³<http://www.uspresidentialelectionnews.com/2012-debate-schedule/>

⁴<http://www.fema.gov/hurricane-sandy-timeline#oct28>

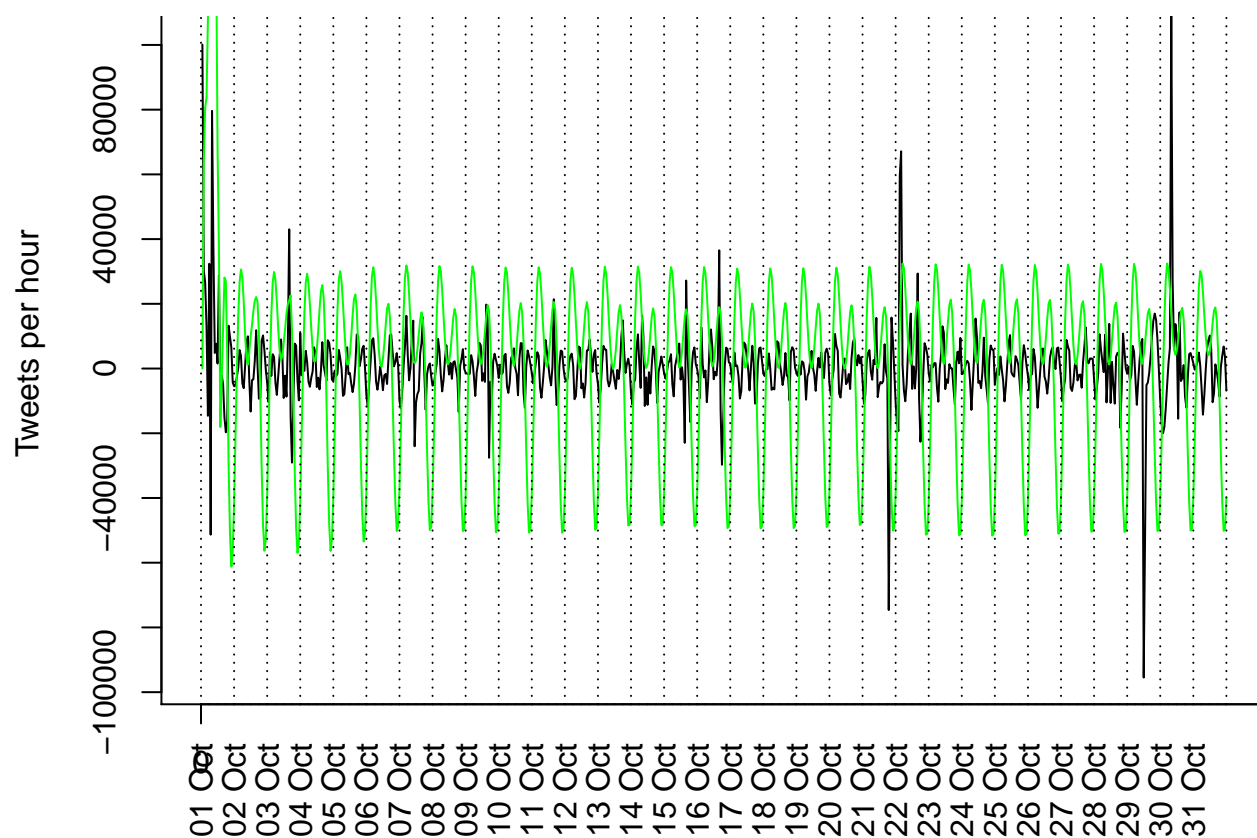


Figure 3.1: Seasonal and noise components of the volume velocity

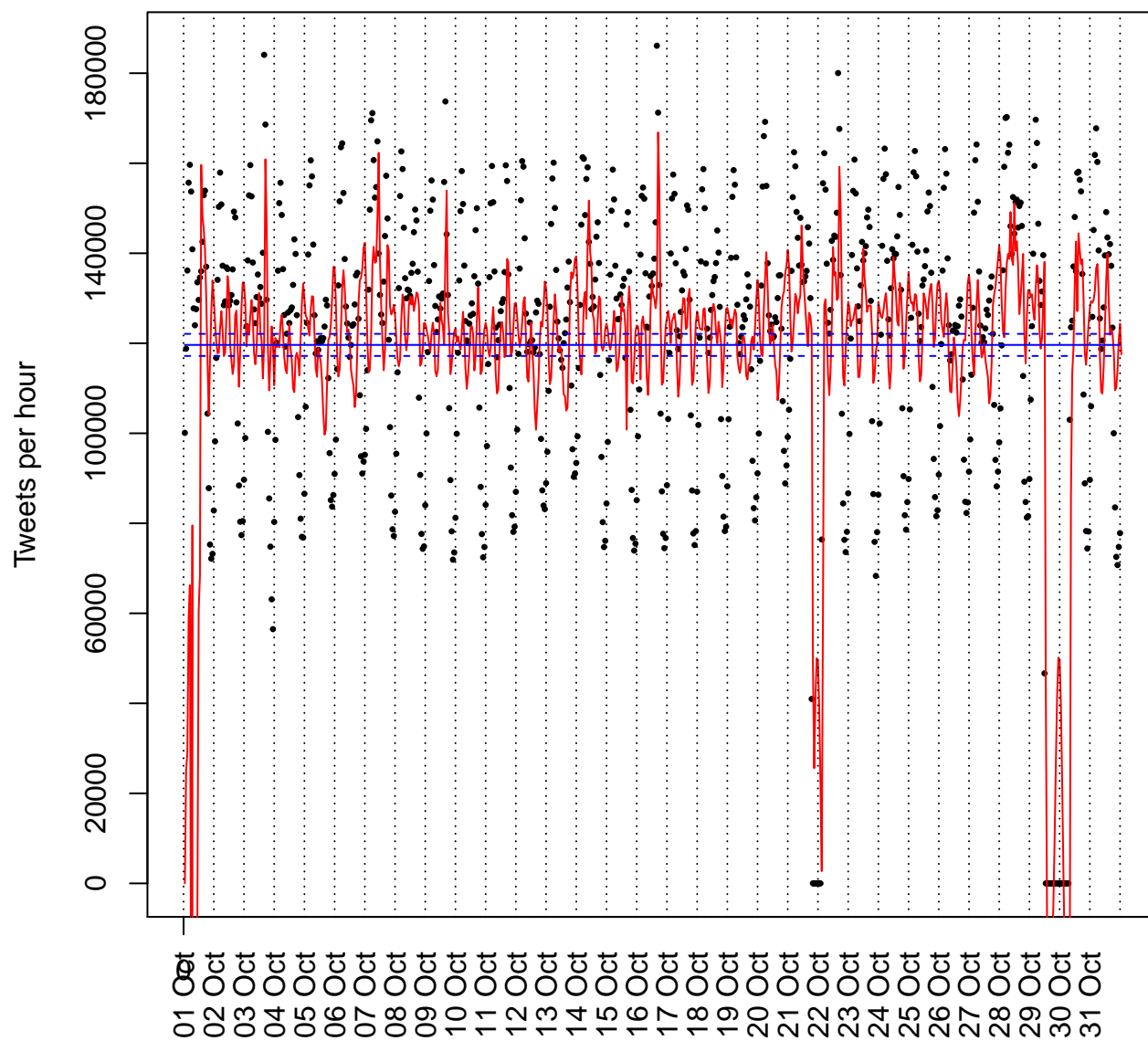


Figure 3.2: Stochastic level of the volume velocity

3.2.3 Dynamic Support Threshold

It is necessary to adapt the support threshold to the dynamicity of the volume velocity at different times of the day. If a ratio support threshold is used then this would not be needed, however determining a support ratio is difficult and it is more desirable to be able to specify the threshold as an absolute number. We therefore define the *minimum support threshold* as the threshold at the hour of the least volume velocity during the day. The *minimum support threshold* is supplied as an absolute number and then converted to a ratio, α , by dividing it by the average of the number of tweets in the hour of minimum volume velocity in different days. The actual support threshold used for mining any given epoch, E , is thus α multiplied by the number of tweets in the epoch:

$$\text{dynamic_support_threshold} = \alpha \times |E| \quad \text{where } \alpha : \text{support ratio} \quad (3.2)$$

$$= \frac{\text{minimum_support_threshold}}{\text{avg}(\min_{\text{day}}(\text{volume velocity}))} \times |E| \quad (3.3)$$

Calculating the support threshold dynamically as such makes it intuitive to select a minimum support threshold, but it would lead to a high support threshold in epochs including a burst of tweets. The high support might lead to mining only itemsets about the topic that is causing the burst. Mining a sliding window can alleviate this problem.

3.2.4 Sliding Window Model

One of the prominent stream processing models is the sliding window model, where a window of data *slides* forwards through the stream. The window keeps a fixed amount of history data (fixed in terms of age not size); data with a timestamp older than a certain time is removed from the window as new data is added. The data in the window is processed every time the window slides. Thus, the sliding window model sets a cut-off age after which a datapoint has no effect at all on the results of the stream processing algorithm. This is the simplest way to overcome the problem that an anomaly in the data, such as a burst of tweets about a certain topic, could affect the mining results more than recent data. More elaborate methods include dynamically changing the sliding window size [11], smoothing the data [28], and weighting data according to its age [27].

We apply our algorithms to epochs of data, similar to the *block evolution* stream mining approach [25, 98]. However, our algorithms are not strictly stream processing algorithms. Stream processing algorithms are supposed to process transactions one at a time as they arrive, looking at each transaction only once – in other words, they should make one pass on the data then discard it. This complicates the model of computation and there is no justification to process social media data as a stream. The stream processing model of computation is justified when the data is not supposed to be stored (such as data from sensor networks), or when the processing time has to be as short as possible such that looking at historic data is an overhead (such as in the case of processing stock ticker data, where making a decision a few milliseconds before/after competitors can translate into profits/losses). In the case of social media text, the data has to be stored anyway and the perception of the human users makes processing times in the order of seconds acceptable for a *real time* system.

In essence, processing epochs of data is similar to mining a sliding window that is moved forward by time steps of short span. The time step must be longer than the time needed to mine an epoch of data, and the performance of our algorithms makes it possible to use a time step as short as a few seconds for epochs up to a day long. The window must be moved such that the combinations of transactions not processed together in a batch is minimal. This is achieved by using a time step as short as possible, within the limits of available resources (such as storage space for mining results, if they are retained). Moreover, the time step must be less than the epoch length to avoid missing mining results from a spike happening at the boundary between two epochs.

Figure 3.3 shows the runtime of LCM on epochs of increasing length, and we will show in section 5.1 that our extensions do not degrade its performance. The times reported in figure 3.3 are averages across all epochs of the specified length in the last 3 months of 2012, using a time step that is half the epoch length. The variance is very low and the confidence bands are not shown because they appear as dots. The *minimum support threshold* used throughout this thesis is 10 occurrences in the hour of day in which the volume velocity is the minimum. Figure 3.4 shows the averages of the actual support threshold to which this minimum support threshold translates at different epoch spans.

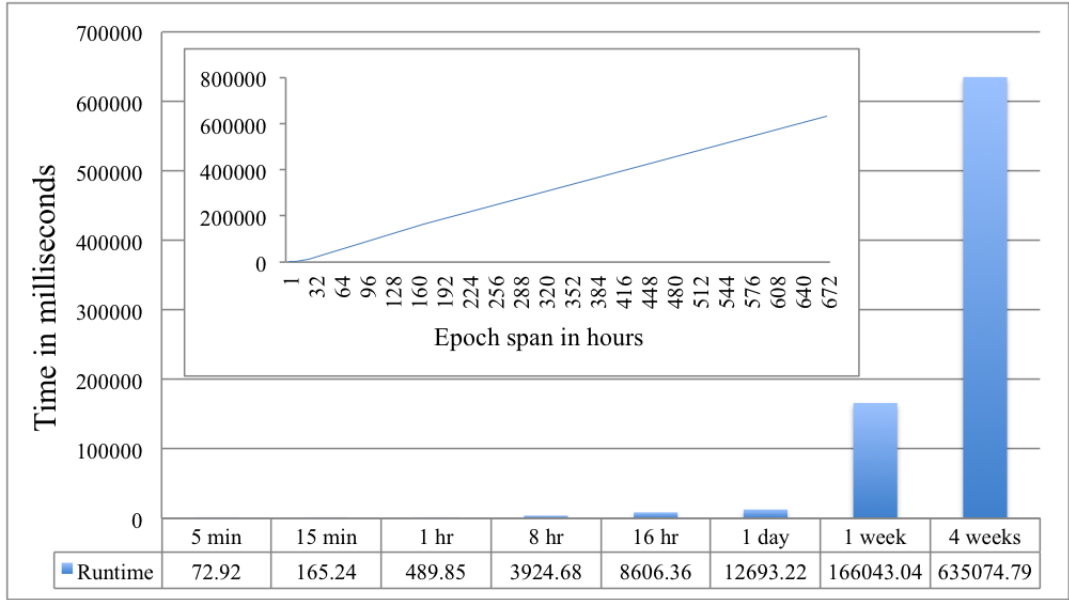


Figure 3.3: Mean runtime of LCM at different epoch spans

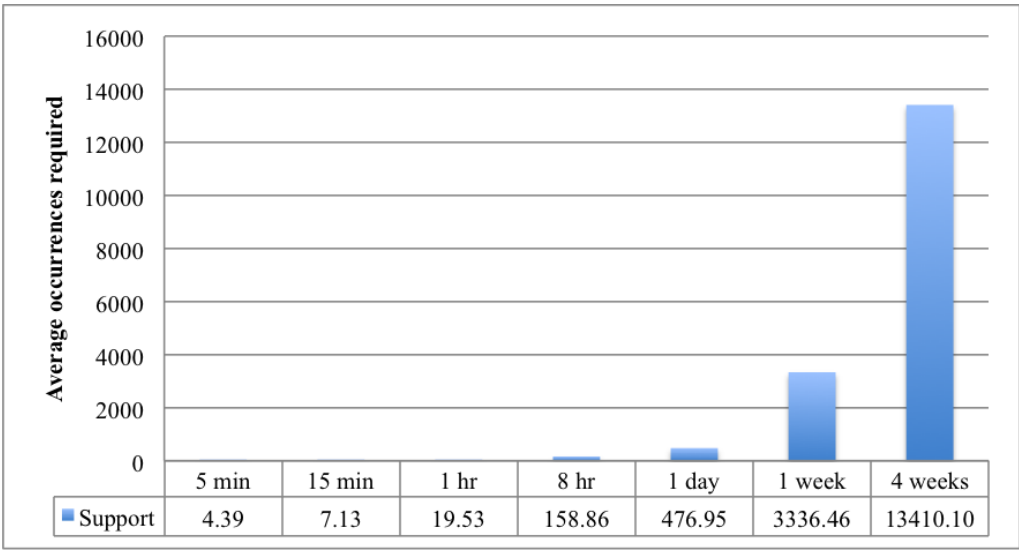


Figure 3.4: Mean support corresponding to minimum support 10, at different epoch spans

In the rest of this paper we mine epochs of 1 hour span. The reason behind this choice is an observation that the number of closed itemsets mined from epochs of span 1 hour or more, at the same minimum support threshold, remains the same. This indicates that itemsets mined from shorter epochs are not included in the results of mining longer epochs. Therefore, the epoch span should be minimized. However, when the epoch span is shorter than an hour the frequency required to surpass the support threshold becomes very low, and the number of mined itemsets increases, with many noise itemsets appearing in the results. For example, the number of itemsets mined from a 15 minutes epoch is double that mined from an hour long epoch.

Regardless of the length of the epoch, many mined itemsets are combinations of function words. In the next section, we outline how we reduce the number of itemsets and eliminate the effect of function words by N-gram filtering.

3.3 Filtering Language Constructs

A large number of itemsets are language constructs that bear no information, such as “such as”. By treating sequential language constructs, and any other multiword expression, as one item we eliminate a large number of such itemsets. We can also eliminate itemsets that are made up of all the different fragments of the language construct along with other items; for example, the itemset {we, did, it, #teamobama} can produce 10 other combinations of length 2 or more. There are many measures of association that can be used to detect multiword expressions, but each measure is good only under certain conditions [71]. After preliminary experiments with various measures, we determined that the best performance could be obtained by tokenizing the documents into term N-grams of varying length.

We use term N-gram tokens such that N-grams of high probability are replaced by (N+1)-grams, resulting in a distribution with no high peaks. An N-gram is considered to have a high probability if its maximum likelihood estimate, from a background model, is higher than a threshold η_N . A background language model built from a long epoch of data from the same stream is used for probability estimation. The language model is simply a hash map of the counts of N-grams that appeared within the last month, for N up to 7.

The tokenization of a tweet starts by tokenizing it into unigrams, as described in appendix A. Then, each unigram of high probability is replaced by two term bigrams, by attaching to it the unigrams before and after it. We keep replacing N-grams of high probability by two (N+1)-grams until there are no more such N-grams.

The threshold of high probability is different for each value of N. The threshold for unigrams is determined as follows: We pick a topical term that is known to steadily appear with a rather high frequency, and is talked about in all languages; for example, ‘obama’. The maximum likelihood estimate of the probability of the term ‘obama’ within the whole collection of tweets is 0.0001. We use $\eta_1 = P(\text{“obama”}) = 0.0001$. At each N, the probability threshold is adjusted to account for the increase in the number of tokens and the overall increase in the grand sum of counts (caused by overlap). The adjusted η_N is:

$$\eta_N = \eta_1 \times \frac{\sum_{\{w: w \in W \text{ and } w.length \leq N\}} freq(w)}{\sum_{\{v: v \in W \text{ and } v.length = 1\}} freq(v)} \quad (3.4)$$

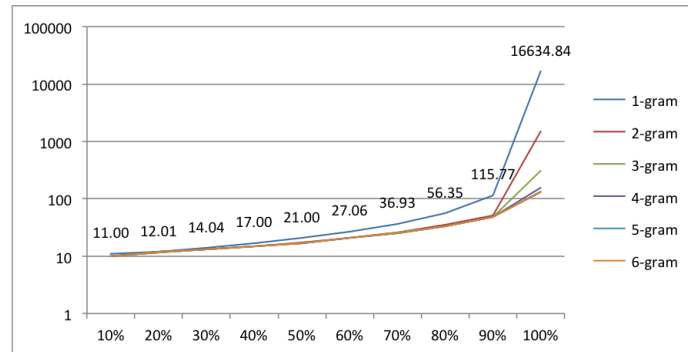
The effect of increasing the maximum length of N-grams is shown in figures 3.5 and 3.6. Figure 3.5 shows the flattening of the distribution by plotting the average decile values at three different times of day. N-grams appearing less than 10 times in an hour (the minimum support threshold we use) are excluded, because they don’t contribute to the distribution of items that the algorithm has to mine. The peakedness of the head of the Zipfean distribution is reflected in the peakedness of the 100th percentile. Increasing the maximum N-gram length reduces the peakedness as well as the variance in the peakedness between different hours of day. Figure 3.6 shows how the maximum N-gram length affects the mining algorithm; it shows the number of tokens in the input, the number of closed itemsets, and the runtime of mining one hour of data. The values shown are averages across all one-hour epochs in the month of November 2012. The value of η_1 used is 0.0001. Figure 3.6(a) shows that the number of distinct items increases substantially as the maximum N-gram length increases from 1 to 2, then continues increasing slightly until it starts decreasing at $N \leq 5$. The decrease happens because all 4-grams with probability above the threshold are parts of tweets from services that use the same text and append a URL, such as tweets reporting scores from Game Insight⁵. Such tweets are tokenized into more 4-grams than 5-grams,

⁵<http://www.game-insight.com/>

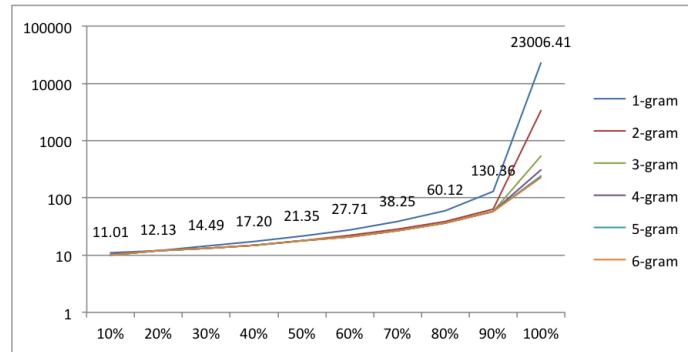
and the 4-grams appearing in them do not appear elsewhere. Thus, some adjacent 4-grams are replaced by fewer 5-grams. Figure 3.6(b) shows that the number of itemsets continues to decrease as expected, with the biggest reduction when the maximum N-gram length increases from 1 to 2. Figure 3.6(c) shows that runtime also decreases as the maximum N-gram length increases from 1 to 5, since LCM’s runtime is proportional to the number of closed itemsets, and it can take advantage of the sparsity of data. The runtimes in this figure are slightly less than those in figure 3.3 because they do not include the time taken for writing the posting list of each itemset.

The numbers of itemsets reported in the plots are counts of distinct unigram sets. After mining itemsets of term N-grams we flatten the itemsets to sets of unigrams again. This is necessary since an itemset will have different N-gram set representations, and its postings list is the union of those of the different representations. This also removes overlap between N-grams of the same itemset, making it easier to reason about how itemsets relate to each other. We will use the relation between an itemset and its subsets to select interesting itemsets in the next chapter.

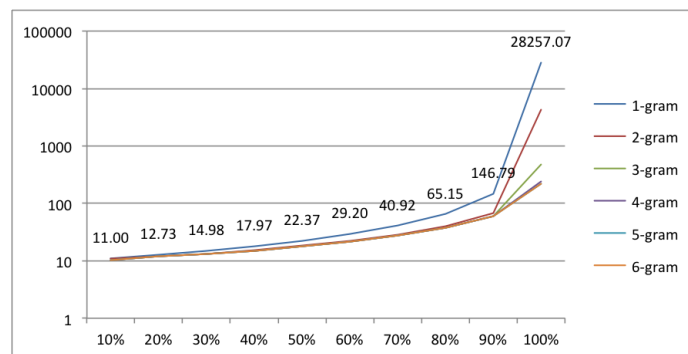
The condition for selecting itemsets that we propose in the next chapter cannot select itemsets of length 1 (frequent unigrams). Thus, to assess its filtering characteristics we have to compare the number of itemsets it selects to the number of itemsets of length 2 or more before its application. At $N \leq 5$, the number of *closed itemsets* of length 2 or more that are mined from an hour long epoch averages at 2,439.17. We also note that the number of *maximal itemsets* of length 2 or more averages at 1,831.92. This high number of *maximal itemsets* show that they cannot substitute the *strongly closed itemsets*, which we propose in the next chapter.



(a) Hour of Day: 00-01 HST (5-6 AM EST)

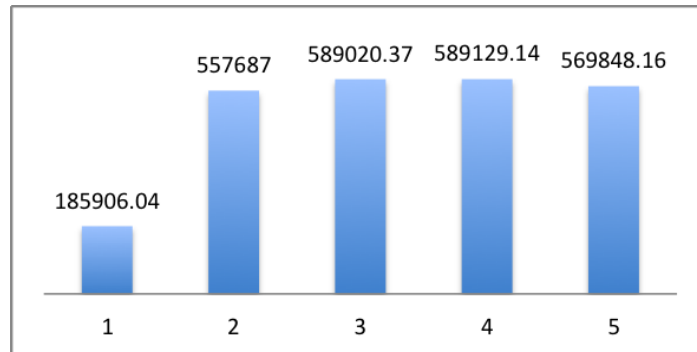


(b) Hour of Day: 08-09 HST (1-2 PM EST)

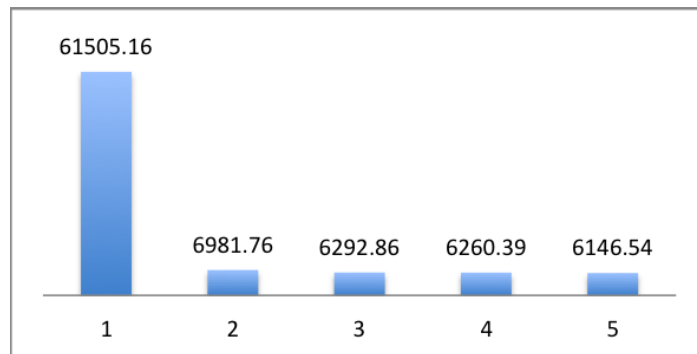


(c) Hour of Day: 16-17 HST (9-10 PM EST)

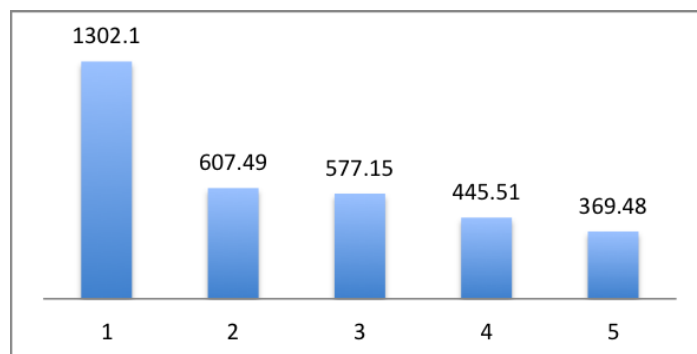
Figure 3.5: Average of token frequency percentiles at different times of day and maximum N-gram length



(a) Mean number of distinct items at different values of maximum N



(b) Mean number of itemsets at different values of maximum N



(c) Mean runtime in milliseconds at different values of maximum N

Figure 3.6: Effect of changing the maximum N-gram length on mining hour long epoches

Chapter 4

Selecting and Ranking Itemsets

In section 3.3, we discussed our handling of function words, using a technique that exploits LCM’s tolerance to sparsity. After applying this technique, the average number of closed itemsets mined from an hour of twitter data drops from 61,505 to 6,146. In this chapter we propose two conditions that reduces the number of itemsets even further. The first condition filters out some itemsets, and the second condition merges together similar itemsets to reduce redundancy. We also propose a ranking function that scores itemsets according to their novelty. Both conditions and the ranking function exploit the dynamics of social media, therefore we start by discussing it.

4.1 Social Media Dynamics

As discussed earlier, one use case of social media is to share content that a user finds important with her network. The user may share an original status message about the topic of interest, or *re-share* content posted by other users. In twitter, re-sharing is done by retweeting, which can be done in several ways. One way is by using the “retweet” feature of the Twitter client, which simply forwards the tweet to the user’s network. The “retweet” feature keeps a reference to the original tweet in the forwarded tweet, but it does not allow editing it – it automatically prepends “RT: @*original_poster_username*”. Since

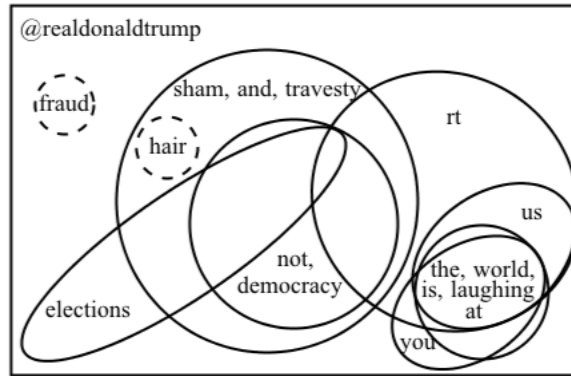


Figure 4.1: Closed, distinct and strongly closed sets

editing is not supported by the “retweet” feature, many users chose to retweet manually to be able to add their own thoughts. Recently this can be done using the “reply” feature, which allows a user to write a tweet that references the original tweet thus maintaining the hierarchy of the conversation. However, manually retweeting is still a prevalent way to re-share a tweet along with the user’s comment.

There are several *styles* of manual retweeting [15]. They all basically quote the original tweet along with its poster’s username and a retweet indicator, and prepend the user’s comment. Due to the 140 characters length limit of tweets the quotation is usually edited to be as short as possible by selecting only the most discriminative words. This selection is an act of collaborative filtering, but it results in many trivially different subsets from the original tweet, and all subsets with enough support will be mined as itemsets. The additions of retweeters also form many different supersets of the of the original tweet, and some additions may represent opinions that are supported enough to be mined as itemsets.

Figure 4.1 illustrates closed itemsets related to Donald Trump’s famous tweets in reaction to Obama’s victory in the 2012 US elections¹. Each area in the figure represents the transactions containing the itemset formed by concatenating the items in all intersecting ellipses. In the figure, the most discriminative words are “sham, and, travesty” which are quoted along with Donald Trump’s user name in most of the retweets. Other people choose

¹http://www.huffingtonpost.com/2012/11/07/donald-trump-election-revolution_n_2085864.html

to also include “not, democracy” and/or “elections”, and in most of the cases the retweet indicator “rt” is added. A few people also added something about Donald Trump’s “hair” while retweeting, and some people talked about “fraud”.

The figure illustrates how the closed property of an itemset is easily violated by modifying one transaction that contains the itemset and removing one of its items. While an update operation is not supported in the model of frequent itemset mining, a similar effect happens when people retweet, or even when they are writing original posts about a certain fine grained topic. In case of retweeting, people actually start from the original tweet and then remove parts of it to make space for their content. We can imagine that for each fine grained topic there is also a base post that represent the information or ideas within the topic. People make posts about the fine grained topic by selecting parts of the base post and adding their own thoughts. This results in many closed itemsets about the topic that are trivially different from each other.

4.2 Distinct Itemsets

We have seen that the dynamics of posting on social media result in redundant closed itemsets. To reduce redundancy, we propose the *distinct* condition which is not as easily violated as the closed condition and use it for selecting itemsets. The *distinct* condition builds on the concept of *association rule confidence*. Confidence is the basic property used for association rules mining, and it is used in the definition of δ -free sets [13]. Mining itemsets based on the confidence of rules they induce has long been recognized as a method for finding “interesting patterns” [21], but since this property is not anti-monotone it cannot be directly used to mine itemsets. It has to be used in a post-processing phase. The confidence of an association rule that the presence of an itemset, s_j , implies the presence of another itemset, s_i , is defined as:

$$conf(s_j \rightarrow s_i) = \frac{|T_{s_i} \cap T_{s_j}|}{|T_{s_j}|} \quad (4.1)$$

The *distinct* condition is a novel strengthening of the closed condition so that it is not easily satisfied by any superset of a closed itemset. A closed itemset of length k , call it s_p ,

can have as many frequent supersets of length $k + 1$ as its support, where s_p will be the intersection of the supersets. From the properties of closed itemsets, each superset that has enough support will result in one or more closed itemsets; if a superset is a maximal itemset then it will result in only one closed itemset, and if it can be extended further it will result in more closed itemsets. Now consider that s_p has a particularly high support, as is the case for itemsets made up of the most discriminative words of a certain topic. In this case, we can filter out closed supersets that have support high enough to satisfy the support threshold but much lower than s_p 's support. These itemsets have enough support to be mined as frequent itemsets and they are closed because they have a closed subset, but they do not represent a *distinctive* piece of information or opinion within the topic that s_p represents. This is similar in essence to mining itemsets at multiple minimum supports [54, 74, 84, 91], where a different support threshold is used for each itemset according to the frequency of its components, or its length, or some other function that calculates the *suitable* support threshold. However, we filter or keep itemsets based on a very simple condition that is easy to implement efficiently and intuitive to reason about.

We define a *distinct* itemset as a closed itemset whose frequency comprises more than a certain proportion of the frequency of its least frequent subset. This condition chooses closed itemsets which substantially violate the closed condition of their subsets. The proportion is a parameter, κ , that controls the selectivity of the condition, or the *distinctiveness* of itemsets selected by the condition. This can be interpreted as selecting itemsets which are implied by a subset with confidence greater than κ . Formally, the set of *distinct* itemsets, \mathcal{D} , is defined as follows:

$$\mathcal{D} = \{s : s \in \mathcal{C} \text{ and } \exists s_p \subset s \text{ where } \frac{|T_s|}{|T_{s_p}|} \geq \kappa\} \quad (4.2)$$

In figure 4.1, *distinct* itemsets are illustrated in solid lines, and closed itemsets that do not satisfy the distinctiveness condition are illustrated in dashed lines. This is not an accurate illustration that is only meant to give an intuition of how *distinct* itemset are different from closed itemsets that are not *distinct*.

Figure 4.2 is a more realistic example. It shows some itemsets (and their support, next to them) from the hour when the US presidential election results were becoming clear.

The itemsets shown are *distinct* at $\kappa = 0.25$. The figure is produced by the graph drawing algorithm proposed by Yifan Hu [34], using the implementation provided in the graph visualization platform Gephi². We run the algorithm on a graph with a node for each itemset and an edge between similar itemsets. Similarity between itemsets is determined based on association rule confidence, as explained in the next section. The similarity is used as the edge weight, thus it determines the distance between itemsets. The figure shows a lot of interesting itemsets. Itemsets pertaining to Donald Trump’s tweets appear on the right side of the page, where the itemsets from the “sham and travesty” tweet are in the middle. Close to the bottom left of the page there are itemsets showing the number of electoral votes each candidate got so far ($\{191, 238\}$ and $\{191, 249\}$). It is clear from the figure that considerable redundancy remains in *distinct* itemsets. In the next section we propose overcoming this redundancy by clustering similar itemsets together.

To assess the filtering power of the *distinct* condition we compare the sizes of the set *distinct* itemsets ($\kappa = 0.25$) with the set of closed itemsets of length 2 or more. In the mining results of hour long epochs from the first two weeks of November 2012, the number of *distinct* itemsets averages at 911.71 itemsets as compared to 2,428.3 closed itemsets (37.55%). To make sure that no important information is filtered out, we look at itemsets that are filtered out. To facilitate the investigation, we look at itemsets containing known topical terms. In the mining results of hour long epochs from the first two weeks of November 2012, the number of closed itemsets containing at least one of the unigrams ‘obama’ or ‘romney’ but are not *distinct* averages at 7.2 itemsets. (an epoch is included in the average only if the keywords occur in it). The hour of the 2012 US elections victory announcement happens to be the hour with the most number of such itemsets, because of the large number of posts about ‘obama’ and ‘romney’. Figure 4.3 shows the filtered itemsets from this hour. Most of the itemsets in the figure are not interesting. Furthermore, interesting ones have surrogates in the set of *distinct* itemsets, either a subset or a superset or an alternative with similar (sometimes more accurate) information. The *distinct* surrogate for each filtered itemset that we find interesting is shown below it.

²<https://gephi.org/>

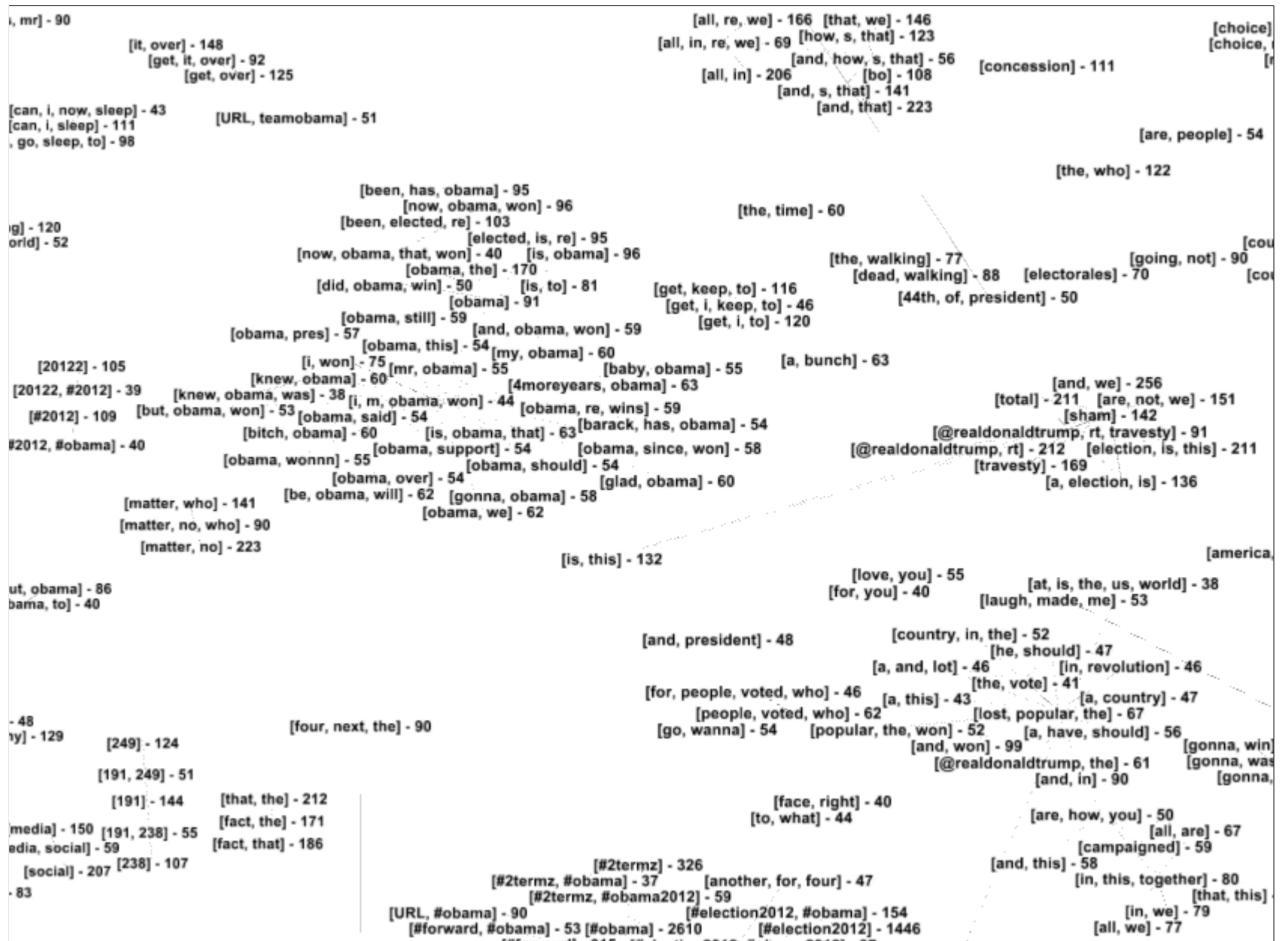


Figure 4.2: Scatter plot of *distinct* itemsets from the hour starting at 10 PM EST on 6 Nov. 2012

1 [romney, that]	24 [romney, 244]	47 [did, obama, it]
2 [still, obama, is]	\hookrightarrow [obama, 244]	48 [for, obama, it]
\hookrightarrow [still, obama]	25 [romney, has]	49 [que, obama]
3 [romney, what]	26 [romney, got]	50 [URL, obama]
4 [romney, with]	27 [romney, can]	51 [can, obama]
5 [will, romney]	28 [romney, but]	52 [obama, all]
6 [romney, if, wins]	29 [romney, was]	53 [so, obama, won]
7 [romney, for, mitt]	30 [romney, the]	\hookrightarrow [wonn, obama]
8 [romney, lost]	31 [voted, for, obama, who]	54 [yes, obama]
9 [romney, like]	32 [2012, obama]	55 [obama, is, the]
10 [romney, campaign]	33 [romney, aint]	56 [supporters, obama]
11 [romney, is, mitt]	34 [obama, is, barack]	57 [obama, when]
12 [romney, still]	35 [gana, obama]	58 [with, obama]
13 [romney, supporters]	36 [obama, gano]	59 [obama, team]
14 [romney, a]	37 [michelle, obama]	60 [obama, has, won]
15 [romney, i]	39 [years, more, obama, 4]	\hookrightarrow [obama, we, won]
16 [s, romney]	\hookrightarrow [#4moreyears, obama]	61 [this, got, obama]
17 [romney, of]	40 [romney, for, vote]	62 [yeah, obama]
18 [romney, on]	41 [de, obama]	63 [romney, about]
19 [romney, or]	42 [obama, an]	64 [obama, happy, won]
20 [romney, in]	43 [on, obama]	\hookrightarrow [glad, obama, won]
21 [romney, is]	44 [obama, it]	65 [yesss, obama]
22 [romney, had]	45 [so, obama]	66 [president, obama, barack]
23 [romney, at]	46 [obama, d]	

Figure 4.3: Itemsets containing ‘barack’ or ‘romeny’ that are closed but not *distinct*

4.3 Strongly Closed Itemsets

To remove the redundancy in *distinct* itemsets we merge similar ones into *strongly closed* itemset clusters. This will also filter out any itemsets that do not belong to any particular topic since they will not be part of a cluster. The similarity of a *distinct* itemset, s_d , and another *distinct* itemset, s_c , is measured as the overlap of the *transactions* containing both of them with the transactions containing s_c . According to equation 4.1 this is the confidence of the association rule that $s_c \rightarrow s_d$. Confidence is not a symmetric measure,

but our proposed clustering makes use of the order in which PPC-extension generates itemsets (as discussed in section 3.1) to determine which itemset should be the antecedent and which should be the consequent.

If PPC-extension uses a total ordering of items that is non-increasing in item frequency, then an itemset generated earlier should be the antecedent, and an itemset generated later should be the consequent. This builds upon the use of variable length N-grams as items, as described in section 3.3, so that the most frequent items are not function words. In that case, following the non-increasing order of item frequency leads to the following:

- Itemsets related to a certain topic are mined close to each other. Recall that PPC-Extension is equivalent to a depth first traversal of a hypothetical prefix tree of closed itemsets. That is, closed supersets of each itemset are generated before moving to an itemset that is not a superset.
- The most discriminative itemset from each topic (according to users' selection) is mined before other itemsets from the same topic, because it is composed of items having the highest support within its topic.
- If an itemset branches out into supersets related to two topics or two opinions within one topic, this itemset is mined before itemsets from either of the two branches. It should not be clustered with either of the branches.
- For an itemset to belong to a topic, it must contain a set of topical words. Since we are processing items in descending order of frequency, then itemsets containing multiple topical words are generated when mining the supersets of the itemset containing only the most frequent of them. Thus, itemsets belonging to a topic will be mined after the topic's most discriminative itemset, even if they are not its superset.

According to the observations above, the *topical* similarity between itemsets can be indicated by the confidence of the association rule that an itemset generated earlier implies an itemset generated later. Using this similarity measure, we cluster itemsets as follows:

1. Itemsets are processed one by one as they are generated by PPC-extension (using a total ordering that is non-decreasing in frequency).

2. If the itemset is not *distinct* it is discarded.
3. A *distinct* itemset is clustered with a previously generated itemset if one exists such that the overlap exceeds a similarity threshold, which we take to be $1 - \kappa$ (the indistinctiveness of s_d from s_c).
4. If more than one satisfies this condition, the *distinct* itemset is clustered with the one having the highest overlap ratio.
5. When a *distinct* itemset is clustered with an itemset that is already part of a cluster, the *distinct* itemset is added to the existing cluster.

Each cluster is aggregated into a *strongly closed* itemset, which is the union of all cluster members, and its support is the size of the union of their postings lists.

This clustering scheme adds an itemset, s_i , to the cluster containing the itemset, s_j , which maximizes the confidence of the rule $conf(s_j \rightarrow s_i)$, with a lower bound on the overlap to maintain distinctiveness. Using \mathcal{D} to denote the set of *distinct* itemsets, we define the desired clustering and the strongly closed itemset represented by each cluster as follows:

$$\begin{aligned}
\mathcal{R} = \{r : r = \bigcup_i (s_i, s_j) \text{ where } s_i \in \mathcal{D} \text{ and } s_j \in \mathcal{D} \\
\text{and } \forall_{(s_i, s_j)} s_j = \mathbf{argmax}_{s_k} conf(s_k \rightarrow s_i) \\
\text{and } conf(s_j \rightarrow s_i) \geq (1 - \kappa) \\
\text{and } (s_j = r.\text{centroid} \text{ or } (s_j, r.\text{centroid}) \in r)\} \\
S_l = \{w : w \in \bigcup_{(s_i, s_j) \in r_l} s_i \text{ where } r_l \in \mathcal{R}\}
\end{aligned} \tag{4.3}$$

A different clustering scheme could be used; the goal is to reduce redundancy, and we devised this scheme because it has several merits that makes it suitable. First, the similarity is measured in the transaction space, so that itemsets are clustered according to *topical* similarity and regardless of the lexically similarity. Jaccard similarity is another possible similarity measure, but it cannot be interpreted as the strength of an implication.

Also, calculating confidence requires calculating the intersection only, so it is easier to terminate the calculation early if the confidence would not be more than the required threshold. Second, this clustering can be implemented efficiently using techniques similar to the ones proposed by Bayardo et al. [8].

The main ideas for an efficient implementation are to limit the comparisons to a few candidates, and to terminate the comparison early if the similarity threshold will not be met. In our case, the postings lists are longer than the itemsets, so we generate candidates for comparison by calculating similarity between itemsets. We calculate the similarity between the postings lists of two itemsets, only if they intersect in one item at least. When calculating the similarity between two postings lists, we can terminate early if the difference exceeds the maximum difference permissible to achieve a similarity of $1 - \kappa$, which can be derived from equation 4.1. Algorithm 2 shows a possible implementation.

To illustrate how the algorithm works, in figure 4.4 we show the clustering hierarchy of itemsets pertaining to the “sham and travesty” tweet. The figure shows a tree structure of how itemsets were added to the cluster. An itemset’s parent in the tree is the itemset that resulted in the rule with the maximum confidence; we call it the best clustering candidate. An itemset is added to a cluster through its best candidate, and if the best candidate does not belong to a cluster a new one is created with the best merge candidate as the centroid. The best merge candidate always has a smaller id, since ids are given sequentially to itemsets and we cluster each itemset with ones that were generated before it. The support of the best candidate is not necessarily higher than the support of the itemset being clustered, but all itemsets within the topic has a support lower than that of the topic’s most discriminative itemset, $n810 = \{\text{sham}\}$. For example, the hierarchy $n1018 \rightarrow n1019 \rightarrow n1075$ has support values 120, 124, 126 respectively, which are all less than 142 (the support of $n810$). Finally, we note that this single *strongly closed itemset* cluster replaces 26 *distinct* itemsets. The overall reduction in the number of itemsets will be discussed in section 5.1.

Another example is given in figure 4.5 to illustrate an itemset that branches into different opinions within one topic. These itemsets are mined from tweets of fans voting for different artists to win the MTV Europe Music Awards (MTVEMA). There are 4 distinct clusters in the figure; one for each artist, and one for itemsets not mentioning any artist.

Input: κ : Minimum distinctiveness threshold

Data: \mathcal{C} : Closed itemsets produced by LCM, in the order they were generated

Result: \mathcal{R} : Strong closed itemset clusters

```

1 for  $i \leftarrow 2$  to  $|\mathcal{C}|$  do
2    $C \leftarrow \{s_c : s_c \in \mathcal{C} \text{ and } c < i \text{ and } |s_c \cap s_i| > 0\};$  // Candidates for clustering
3    $P \leftarrow \{s_p : s_p \in \mathcal{C} \text{ and } p < i \text{ and } s_p \cap s_i = s_p\};$  //  $s_i$ 's subsets (ancestors)
4    $s_p \leftarrow \operatorname{argmax}_{s_p \in P} \frac{|T_{s_i}|}{|T_{s_p}|};$  // Direct parent
5   if  $s_p = \text{null}$  OR  $\frac{|T_{s_i}|}{|T_{s_p}|} < \kappa$  then // Not a distinct itemset
6      $\mathcal{C} \leftarrow \mathcal{C} \setminus s_i;$  // Should not be considered in later iterations
7     continue; // Should not be added to a cluster
8   end
9    $s_m \leftarrow \text{null}, \text{maxConf} \leftarrow 0;$  // Best clustering candidate and its score
10  foreach  $s_c \in C$  do
11     $\Delta \leftarrow (1 - (1 - \kappa))|T_{s_c}|;$  // Maximum difference
12     $\delta \leftarrow \text{difference}(T_{s_c}, T_{s_c \cup s_i}, \Delta);$  // Terminates early if  $\delta > \Delta$ 
13    if  $\delta \leq \Delta$  then
14       $\text{conf} \leftarrow \frac{|T_{s_c}| - \delta}{|T_{s_c}|};$ 
15      if  $\text{conf} > \text{maxConf}$  then
16         $s_m \leftarrow s_c;$  // Best clustering candidate
17         $\text{maxConf} \leftarrow \text{conf};$ 
18      end
19    end
20  end
21  if  $s_m \neq \text{null}$  then
22    if  $\mathcal{R}[s_m] = \text{null}$  then  $\mathcal{R}[s_m] \leftarrow s_m;$  // New cluster with centroid  $s_m$ 
23     $\mathcal{R}[s_i] \leftarrow \mathcal{R}[s_m];$  // Add  $s_i$  to the cluster containing  $s_m$ 
24     $\mathcal{R}[s_m].\text{itemset} \leftarrow \mathcal{R}[s_m].\text{itemset} \cup s_i \cup s_m;$ 
25     $\mathcal{R}[s_m].\text{postingsList} \leftarrow \mathcal{R}[s_m].\text{postingsList} \cup s_i.\text{postingsList} \cup s_m.\text{postingsList};$ 
26  end
27 end
28 return  $\mathcal{R};$ 

```

Algorithm 2: Forming strongly closed itemset clusters

Id	Itemset	Support
n810	▼ [sham]	142
n906	▼ [and, sham]	135
n907	▼ [and, sham, travesty]	134
n908	▼ [and, sham, total, travesty]	132
n944	▼ [a, and, is, sham, total, travesty]	131
n945	▼ [a, and, election, is, sham, this, total, travesty]	130
n1016	▼ [@realdonaldtrump, this]	129
n1018	▼ [@realdonaldtrump, a, democracy, not, this]	120
n1019	▼ [@realdonaldtrump, a, and, is, sham, this, total, travesty]	124
n1075	[@realdonaldtrump, election, this]	126
n1076	[@realdonaldtrump, a, and, election, is, sham, this, total, travesty]	125
n1055	▼ [a, and, are, democracy, is, not, sham, total, travesty, we]	125
n1056	[a, and, are, democracy, election, is, not, sham, this, total, travesty, we]	124
n1057	▼ [@realdonaldtrump, a, and, are, is, not, sham, this, total, travesty, we]	120
n1079	[@realdonaldtrump, a, and, are, election, is, not, sham, this, total, travesty, we]	121
n1059	[@realdonaldtrump, a, and, are, democracy, is, not, sham, this, total, travesty, we]	119
n1080	[@realdonaldtrump, a, and, are, democracy, election, is, not, sham, this, total, travesty, we]	120
n963	▼ [a, are, not]	132
n964	▼ [a, are, not, we]	128
n1053	▼ [a, and, are, is, not, sham, total, travesty, we]	127
n1054	[a, and, are, election, is, not, sham, this, total, travesty, we]	126
n965	[a, are, democracy, not]	128
n966	[a, are, democracy, not, we]	127
n1033	[travesty, we]	128
n925	[a, is, total]	134
n943	[a, and, sham, travesty]	133

Figure 4.4: The hierarchy of itemsets in the cluster pertaining to the “sham and travesty” tweet

Id	Itemset	Support
n0	▼ #mtvema	624
n2	▼ #mtvema,be,big,the,will	608
n3	▼ #mtvema,be,big,the,will,winner	607
n4	▼ #mtvema,pick,tweet,winner,your	606
n5	▼ #mtvema,at,pick,tweet,winner,your	605
n6	▶ #mtvema,big,tweet,winner,your	605
n11	▶ URL,at,pick,tweet,your	604
n25	▼ #mtvema,be,bieber,big,justin,the,will,winner	405
n26	▼ #mtvema,at,be,bieber,big,justin,pick,the,tweet,will,winner,your	403
n27	▶ #mtvema,be,bieber,justin,think,will	402
n28	■ #mtvema,URL,at,be,bieber,big,justin,pick,the,tweet,will,winner,your	401
n91	▼ #mtvema,be,big,katy,perry,the,think,will	166
n119	▼ #mtvema,#emawinkaty,be,big,katy,perry,the,think,will	159
n126	▶ #mtvema,#emawinkaty,be,big,i,katy,perry,the,think,will	157
n122	▶ #mtvema,#emawinkaty,be,big,katy,perry,the,think,will,winner	158
n99	▼ #mtvema,be,big,i,katy,perry,the,think,will	164
n104	▶ #mtvema,be,big,i,katy,perry,the,think,will,winner	163
n95	▼ #mtvema,be,big,katy,perry,the,think,will,winner	165
n97	▶ #mtvema,be,big,katy,perry,the,think,tweet,will,winner,your	164
n5207	▼ be,big,gaga,the,will	25
n5275	▼ #mtvema,be,big,gaga,the,tweet,will,winner	24
n5569	▶ #mtvema,at,be,big,gaga,pick,the,tweet,will,winner,your	23
n5562	■ #mtvema,#emawingaga,at,pick,tweet,winner,your	23
n5572	■ #mtvema,be,big,gaga,lady,the,tweet,will,winner	23

Figure 4.5: Several clusters for itemsets pertaining to the MTV Europe Music Awards votes

4.3.1 Bounding Space and Time Requirements

Algorithm 2 requires keeping *distinct* itemsets and their postings lists in memory, so that the similarity between a *distinct* itemset and all the previously generated ones can be efficiently calculated. This is an overhead that the original LCM does not incur, but it is not a large overhead; due to the limited number of *distinct* itemsets, the memory requirements does not exceed 6 GB for mining hour long epochs using a Java implementation that uses a lot of objects. However, it is possible to use less memory and also bound the memory requirement by keeping only a limited number of itemsets in memory.

The number of positions between an itemset and a clustering candidate is 1973.11 on average, with a standard deviation of 1443.64. The large standard deviation indicates that candidates are either found close to the itemset (less than 500 positions behind), or in a position that is several thousands of itemsets behind. According to our observation that itemsets from one topic are generated close to each other, when the total ordering is descending in items frequency, we limit the mining results kept in memory to the last 1000 *distinct* itemsets and their postings lists. This sets a bound on the memory requirement, as well as the runtime which is dominated by the similarity calculations. Without limiting the number of itemsets kept in memory, the runtime for filtering and clustering hour long epochs at $\kappa = 0.25$ is 5.2 seconds on average. This is reduced to less than 1 second if only the last 1000 results from LCM are kept in memory. More detailed discussion of the runtime performance of the algorithm is given in section 5.1. The *quality* of the *strongly closed itemsets* produced with that limit in place is investigated in section 5.2.

4.4 Temporal Ranking

The number of *strongly closed itemsets* from an hour long epoch is 139.1 itemsets on average, which is under 0.25% of the number of closed itemsets returned by LCM without any of our proposed methods (61,505.16 on average). We now propose a ranking function that scores itemsets according to their novelty, such that the ordered list of itemsets can be presented to a user as a summary of what is happening on Twitter in the mined epoch. The ranks can also be used as weights for itemsets when used by other system components,

such as a search component that uses itemsets for query expansion. Novelty is measured compared to a longer period of time leading up to the mined epoch – the background model.

A good indicator of novelty is the pointwise Kullback-Leibler Divergence (KLD) between an itemset’s probability in the current epoch and in the background model. The KLD of the probability of an itemset s_i in the background model Q from the current epoch’s model P can be considered as the Information Gain, IG:

$$\begin{aligned} IG(P(s_i), Q(s_i)) &= -(H(P(s_i)) - H(P(s_i), Q(s_i))) \\ &= \sum P(s_i) \log P(s_i) - \sum P(s_i) \log Q(s_i) \\ &= KLD(P(s_i) || Q(s_i)) \end{aligned} \tag{4.4}$$

This interprets the KLD of an itemset as the number of extra bits in the representation of the itemset used by an encoding based on Q as compared to the number of bits used by an encoding based on P ; hence, the information gain in moving from a prior distribution to a posterior distribution.

To calculate the collective IG of itemsets in a *strongly closed itemset*, we have to take into account that the itemsets of the cluster are not independent. For simplicity we will consider only the pairwise dependence between every itemset and the smallest common subset. The joint probability of the smallest common subset, s_{min} , and any itemset in the cluster, s_j , is basically the probability of the itemset s_j , since $s_{min} \subset s_j$. In other words, the probability of s_{min} conditioned on any itemset s_j in the cluster is 1, and there is no information in the occurrence of s_{min} given that s_j occurred. An alternative interpretation is that there is a shared component in the IG of all itemsets in the cluster, which is the IG of s_{min} . Hence, the IG of an itemset cluster, $r = \{s_{i_1}, \dots, s_{i_m}\}$, can be approximated as the IG of its smallest subset, s_{min} , plus the IG of each itemset in the cluster conditioned on s_{min} :

$$\begin{aligned} IG(P(s_{i_1}, \dots, s_{i_m}), Q(s_{i_1}, \dots, s_{i_m})) &= IG(P(s_{min}), Q(s_{min})) \\ &+ \sum_{j=i_1..i_m} IG(P(s_j | s_{min}), Q(s_j | s_{min})) \end{aligned} \tag{4.5}$$

This formula avoids implicitly adding the information gain from s_{min} multiple times to the total information gain.

The formula above can be used directly for ranking clusters, but it will favour larger clusters. We normalize by the size of the cluster giving our final ranking formula for strongly closed itemsets:

$$score(r) = IG(P(s_{min}), Q(s_{min})) + \frac{\sum_{j=i_1..i_m} IG(P(s_j|s_{min}), Q(s_j|s_{min}))}{m} \quad (4.6)$$

In this thesis, the background model we use for each day is the results of mining the 4 weeks before it, using a *minimum support* value of 1 occurrence. Mining the background model at such a low support increases the number of produced itemsets, which is desirable for a background model. All probability estimates are smoothed by add-one smoothing.

Chapter 5

Evaluation

In this chapter we analyze the performance of the filtering conditions proposed by applying them to the mining results of all one-hour long epochs in the Twitter data. We evaluate their efficiency in reducing the number of mined itemsets in section 5.1, and we assess the effectiveness of the filtering in choosing a high quality synopsis in section 5.2. The experiments were run using a single threaded Java implementation of algorithm 2. The experiments were run on a 1.4GHz processor with 2MB of cache.

5.1 Efficiency

We evaluate the efficiency of our proposed filtering and clustering methods in terms of the number of itemsets before and after their application, and the runtime overhead they add on top of the runtime of the LCM algorithm. Using variable length N-grams as items, the average number of itemsets mined from an hour-long epoch is 2439.17 closed itemsets of length 2 or more after flattening; that is, excluding itemsets that are merely a frequent unigram. We start by analyzing the reduction in this number at different values of the distinctiveness threshold, κ .

Figure 5.1 show the effect of varying κ on the mean number of *distinct* and *strongly closed* itemsets. The number of *distinct* itemsets drops as the distinctiveness threshold

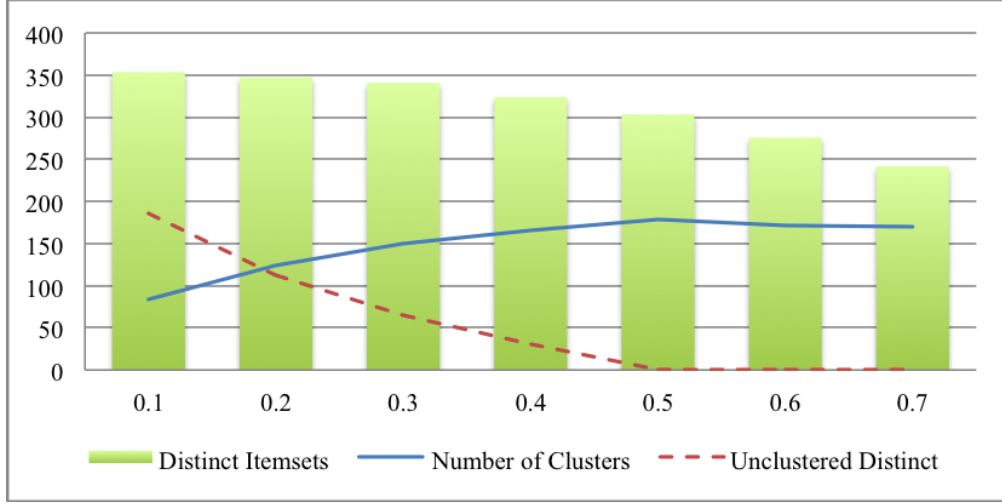


Figure 5.1: Effect of changing κ on mining results

increases. On the other hand, the number of *strongly closed* clusters formed increases as the similarity/indistinctiveness threshold (taken as $1 - \kappa$) decreases. The dashed line shows that the number of unclustered distinct itemsets reaches zero at $\kappa = 0.5$, explaining why the number of clusters changes very slightly after that. We use $\kappa = 0.25$ in this thesis, which is an arbitrary choice based on the definition not the data. The average number of *strongly closed* itemsets mined from one-hour epochs at this value of κ is only 139.1, which is about 5.7% of the number before filtering and clustering.

Figure 5.2 shows the total runtime of the LCM algorithm plus filtering based on the distinct condition and clustering into strong closed itemsets at different epoch spans. The runtime of LCM alone is also plotted for reference. We also plot the performance of another frequent itemset mining algorithm, FP-Zhu [29], which was the runner up at FIMI 2004 [41]. We include it to show that our extensions do not degrade the performance of LCM even in the context of competitions. The y-Axis is in logarithmic scale to keep the scale of the plot suitable for seeing slight differences. The output of LCM is the input to the filtering and clustering step, so it is affected by the number of closed itemsets produced. This explains why it takes slightly longer time for clustering results from the 15-minute epoch and then takes a constant time for epochs of a longer span.

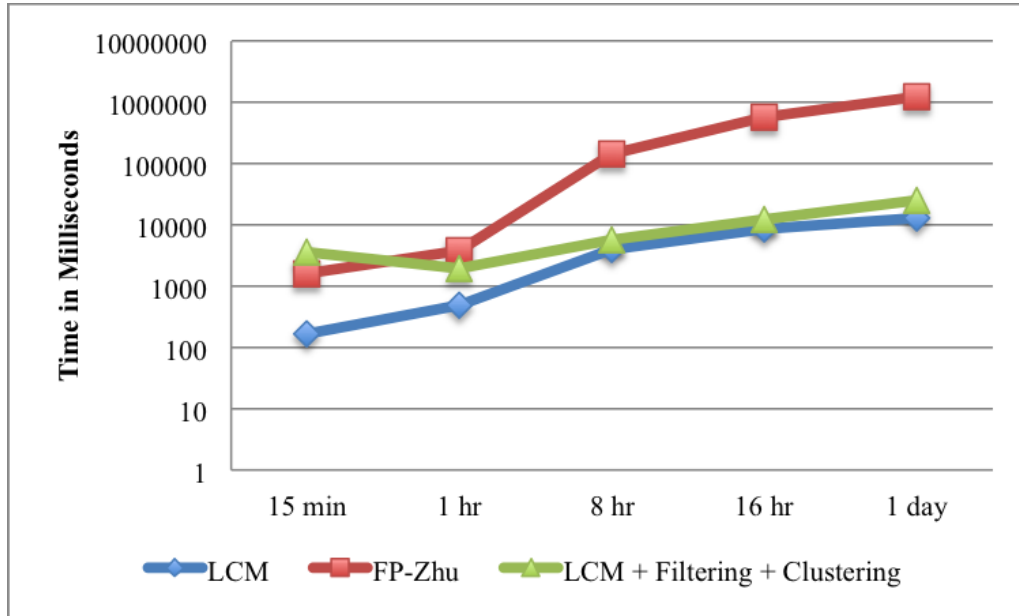


Figure 5.2: Runtime of two itemset mining algorithms alone and of the LCM algorithm with filtering and clustering, at different epoch spans

5.2 Effectiveness

We now show examples of the synopses produced by ranking the *strongly closed* itemsets as discussed in section 4.4. The examples will serve as our assessment of the quality of the synopses. We make an effort to choose epochs from time periods when it is known what people should be talking about. In section 5.2.1, we show how the 2012 US presidential elections day is summarized by our methods, and we compare our summary to the one produced by the MTV algorithm [57]. In section 5.2.2, we investigate the synopses of hours selected by the aid of Google Trends¹, assessing which popular queries were or were not reflected in our synopses.

¹<http://www.google.ca/trends>

5.2.1 US presidential elections day

In this section, we compare the quality of the synopsis created by our proposed methods to the quality of the synopsis created by mining the top itemsets with the MTV algorithm [57] (according to a maximum entropy model). The synopses should be suitable for giving a user an overview of what people were Tweeting about at different times on the 2012 US presidential election day, November 6th.

Table 5.1 shows the top 3 *strongly closed* itemsets from one-hour epochs using a half-hour time step. The itemsets shown are the first appearances of the most interesting itemsets; that is, an hour is shown only if its top 3 feature novel interesting itemsets². The first column is the beginning of the hour, EST time. The second column is the top itemsets. The third column is a commentary to explain the itemsets. The fourth column is the rank of an equivalent itemset if one exists in the top 30 itemsets mined by the MTV algorithm.

We can see how the events of the US presidential elections unwind from “get out and vote” to the projections and debates, all the way to the “acceptance speech”. Early in the day, people were excited and congratulating each for being on the elections day. The excitement continues for one more hour, with some people taking it too far by sending pictures of their ballots. Itemsets about the elections stop occupying the top 3 ranks until the evening. Instead, there were itemsets about UEFA Champions football matches with timely updates of their scores, and about TV shows. We include the top 3 itemsets of 3 PM as an example. In the evening, events about the elections keep occupying the top positions. At 10 PM all the top 5 were relevant and very interesting that we decided to include an extra two itemsets in the example (we do this for MTV as well in the hour of acceptance speech). Shortly after the results of the elections became clear, news that Marijuana and same sex marriage were legalized in some states occupies the top position. This exemplifies the power of social media as a collaborative filter, selecting the news of greatest importance to social media users. The user centric definition of importance is also

²The top 50 itemsets in the hours from 5 AM on the 6th till almost 9 AM on the 7th are available for download at http://plg.uwaterloo.ca/~yaboulna/thesis_results/twitter_synopses/elections-day_nov6-0500_nov7-0850_top50-strongly-closed.txt

evident in attention given to the “lady with a flag in her hair” appearing on TV behind Obama during the acceptance speech.

The top 3 itemset mined by the MTV algorithm [57] are usually itemsets from Tweets sent by services reporting users’ scores in games, or the number of new followers each user got. However, the top 10 itemsets usually contain interesting itemsets⁴. Therefore, in table 5.2 we show itemsets that we find interesting from MTV’s top 10, assuming that the itemsets from Tweets sent by automatic services can be filtered out by another dedicated filter. The actual rank in which the itemset appeared is shown in column 3, and the rank in which it appears among *strongly closed* itemsets is shown in column 4. We show only hours in which the top 10 itemsets contain novel interesting itemsets. The itemsets chosen by MTV tend to be complete Tweets which were retweeted by many users, thus we omit the explanation column for brevity. The input to the algorithm was transactions made up of N-grams up to 5 terms long, which helped the algorithm converge faster since the distribution is flatter. The use of N-grams also overcomes the dominance of language constructs, which are otherwise ranked high in all hours.

All the topically relevant itemsets chosen by MTV are present in the top 50 *strongly closed* itemsets – most of them are in the top 30. We do not know if all the top ranked *strongly closed* itemsets would also be picked by MTV if we use it to mine more itemsets, since we did not have time to get its results for all epochs at $k > 30$, because the algorithm’s space and time requirements of itemsets requested is increased. At $k > 100$ the algorithm fails because of a memory error, even though it was being run on a machine with 256 GB of RAM, with no limit on its resource consumption. The runtime performance at $k \leq 100$ is plotted in figure 5.3. The small value of k that has to be used makes MTV inadequate for mining social media, where the high volume of posts and the variety of topics users talk about is likely to result in more itemsets than what MTV would pick as the top k . On the other hand, *strongly closed* itemsets are efficiently mined and are not restricted to a predefined number.

⁴The top 30 itemsets in the hours from 5 AM on the 6th till 5 AM on the 7th are available for download at http://plg.uwaterloo.ca/~yaboulna/thesis_results/twitter_synopses/elections-day_nov6-0500_nov7-0500_top30-MTV.txt

Time	Itemset	Explanation	MTV
08:30	<i>get out and vote</i>	Encouraging participation	17
	<i>happy elections day</i>	Congratulations	27
	<i>it's elections day</i>	Excitement	-
09:30	<i>if romney wins</i>	People dreading the thought	-
	<i>of your ballot</i>	Telling people to stop sending pictures of ballots	-
	<i>in line to vote</i>	Tweeting while standing in line	28
15:00	<i>Borussia Dortmund Real Madrid 2</i>	Score update of UEFA match	17
	<i>incroyable talent</i>	France's version of the Got Talent TV show	-
	<i>City Ajax 2</i>	Score of Manchester City vs Ajax	-
19:30	<i>Linda McMahon</i>	Linda McMahon loses CT senate race	-
	<i>Florida is so ...</i>	tight or close; the number of votes	-
	<i>Romney is winning</i>	Speculations about the results	-
20:00	<i>New Hampshire</i>	Obama won in NH	9
	<i>Obama to win / is winning</i>	Two itemsets speculating	-
	<i>too close to call</i>	The results are still not clear	22
21:00	<i>Ohio and Florida</i>	Obama won 2 more states	-
	<i>Sarah Palin</i>	Sarah Palin appears on Fox news	-
	<i>concession speech</i>	Todd Akin's concession speech	-
22:00	<i>The electoral college</i>	Republicans complaining about the voting system	9
	<i>President of the unites states</i>	Barack Obama is the 44th president of the USA	8
	<i>who voted for</i>	People turning against each other	-
	<i>once you go black you never go [back]</i>	A popular cultural reference	15
	<i>my president is black</i>	Americans proud about the equality of opportunities	-
22:30	<i>@realdonaldtrump this elections (see table 5.2 for the rest of the tweet)</i>	One of Donald Trump's famous melt down tweets. This cluster merges 26 <i>distinct</i> itemsets.	1
	<i>concede to</i>	The losing candidate has to concede	-
	<i>Obama won. I ...</i>	People's reactions to the results	-
23:00	<i>same sex marriage</i>	Some states legalized same sex marriage	27
	<i>for recreational use</i>	Colorado legalized Marijuana for recreational use	21
	<i>acceptance speech</i>	People anticipating the speech	26
00:30	<i>Delivered, sealed, signed</i>	Stevie Wonder's song used in Obama's campaigns	10
	<i>The best is yet to come</i>	Quote from Obama's speech	2
	<i>with the flag in her [hair]</i>	During the speech, attention on social media goes to a woman appearing behind Obama on TV!	14
01:00	<i>Four more years</i>	Barack Obama's tweet; the most retweeted in 2012 ³	20
	<i>@paulmyers47: spots</i>	Pyramidal marketing; people retweet to make money	12
	<i>President Barack Obama</i>	News headlines about the speech uses the formal title	-

Table 5.1: Top 3 (or 5) *strongly closed* itemsets for hours in US presidential elections day

Time	Itemset	Ranks	
19:00	<i>A partir de que idade você considera alguém velho?</i>	1	1
	<i>A qué edad consideras que alguien es viejo?</i>	4	10
	<i>We'll find out if it's true. America went black, will it go back?</i>	9	24
20:30	<i>Obama rhymes with Ohana. Ohana means family & family means nobody gets left behind. Mitt rhymes with shit.</i>	3	3
	<i>New Hampshire</i>	9	1
	<i>#IamSickOf</i>	10	-
21:00	<i>We're all in this together. That's how we campaigned, and that's who we are. Thank you. -bo</i>	1	22
	<i>@barackobama This happened because of you, thank you.</i>	4	20
	<i>@barackobama Four more years URL</i>	6	6
22:00	<i>@realDonaldTrump: We can't let this happen. We should march on Washington and stop this travesty. Our nation is totally divided!</i>	5	44
	<i>President of the United States</i>	8	2
	<i>@realDonaldTrump: The electoral college is a disaster for a democracy.</i>	9	1
22:30	<i>@realDonaldTrump: This election is a total sham and a travesty. We are not a democracy!</i>	1	1
	<i>@realDonaldTrump: Our country is now in serious and unprecedented trouble...like never before.</i>	3	24
	<i>@realDonaldTrump: He lost the popular vote by a lot and won the election. We should have a revolution in this country!</i>	5	31
23:00	<i>@realDonaldTrump: Lets fight like hell and stop this great & disgusting injustice! The world is laughing at us.</i>	4	48
	<i>@realDonaldTrump: House of Representatives shouldn't give anything to Obama unless he terminates Obamacare.</i>	6	15
	<i>@adamlevine: That's what happens when you fu** with Sesame Street</i>	8	18
00:30	<i>The best is yet to come</i>	2	2
	<i>Michelle, I have never loved you more</i>	5	12
	<i>We are an American Family, and we rise ...</i>	8	20
	<i>Delivered, sealed, signed, I'm yours</i>	10	1

Table 5.2: The most interesting itemsets out of the top 10 picked by the MTV algorithm for hours in US presidential elections day. The actual rank given to each itemset by MTV is in column 3 and the rank of its equivalent *strongly closed* itemset is in column 4.

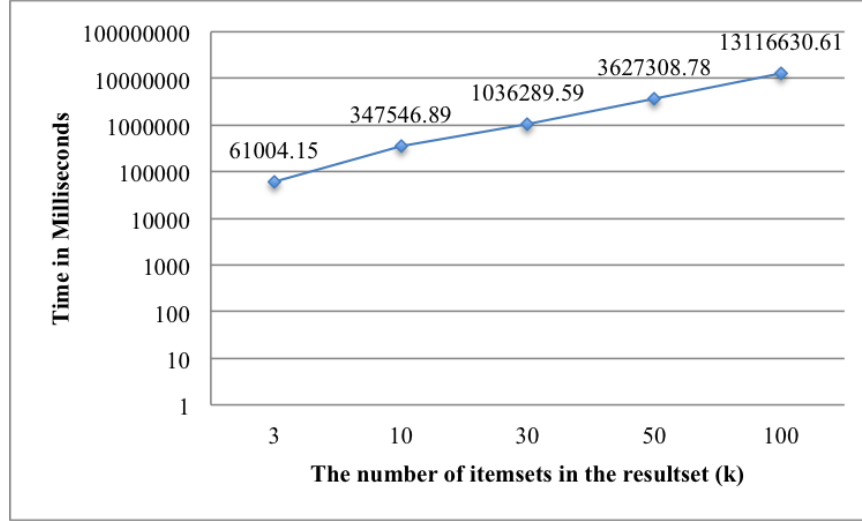


Figure 5.3: Runtime of the MTV algorithm [57] at different resultset sizes (top k)

5.2.2 Google Trends

Google Trends is a freely available service provided by Google. It analyzes queries submitted to the Google search engine by millions of users worldwide, and it can be used to know the most popular queries for different geographies⁵. It can also be used to compute trends of user specified queries over time⁶. Queries that are named entities according to Google’s Knowledge Graph⁷ are automatically tracked and aggregated by category into Top Charts⁸ – lists of real-world people, places and things ranked in order of search volume during user specified intervals of time⁹. We use the Top Charts of November 2012 for the people category¹⁰ (only available for queries made from the USA) as an aid for finding out what social media users are expected to be talking about. We observed that named entities that make it to the top 10 chart in the people category are divided into finer groups. The queries about people in each group have similar volumes over time, with peaks at the

⁵<http://www.google.ca/trends/>

⁶<https://support.google.com/trends/answer/87276?hl=en>

⁷<http://www.google.com/insidesearch/features/search/knowledge.html>

⁸<http://www.google.ca/trends/topcharts>

⁹<https://support.google.com/trends/answer/3076011?hl=en>

¹⁰<http://www.google.ca/trends/topcharts#vm=chart&cid=people&geo=US&date=201211>

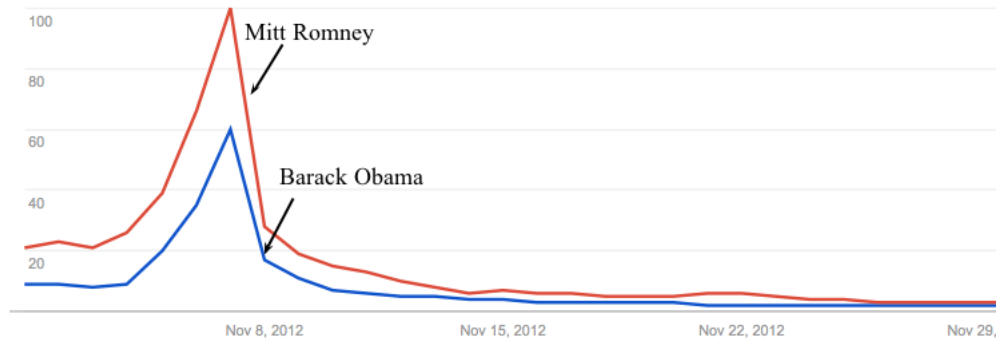


Figure 5.4: Google Trends' volume with time plot for the queries “Barack Obama” and “Mitt Romney”, the top 2 in the Top Chart for the *people* category in November 2012

times of real world events in which those people in the group took part. For example, the top two most searched for people are Barack Obama and Mitt Romney. Figure 5.2.2 is the plot showing the volume of queries about either of them, as displayed by Google trends ¹¹. The two curves in the plot are similar; both have peaks at the elections day and are otherwise flat. Note that the y-Axis in the figure is scaled¹² so that the maximum volume per day for either of the queries becomes 100 (to our surprise the curve for “Barack Obama” is much lower than that of “Mitt Romney”, even though they occupy positions 1 and 2 respectively in the chart). Since entities can be grouped together according to the events which made users search for them, we use the Top Charts to know about real world events which were of interest to internet users. We evaluate the use of *strongly closed* itemsets as temporal synopses by investigating the mining results from the epochs when these events were happening. We look for itemsets about those entities and events in the top 30 *strongly closed* itemset, and assess which of them were or were not reflected in our synopses. The number of itemsets from each hour is limited to 30 because this is a number that can be conveniently displayed as a summary for a human user.

¹¹<http://www.google.ca/trends/explore#q=Barack%20Obama%2C%20Mitt%20Romney&geo=US&date=11%2F2012%201m&cmpt=q>

¹²https://support.google.com/trends/answer/87282?hl=en&ref_topic=13975

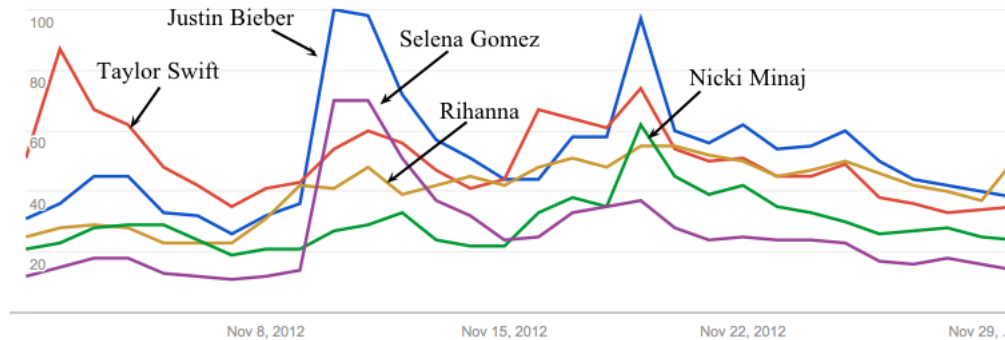


Figure 5.5: Google Trends' volume with time plot for the queries about pop artists in the Top Chart for the *people* category in November 2012

Pop artists and music awards events

Pop artists occupy most of the top 10 chart for the *people* category. Five out of the top 10 are pop artists, and they occupy ranks as high as 3 (Justin Bieber) and 4 (Taylor Swift). This is due to music awards events, specially the MTV Europe Music Awards on November 11th¹³ and the American Music Awards on November 18th¹⁴. Figure 5.2.2 shows the query volumes for the artists in the top 10 chart. The curves for all of the artists show relatively high volume when those events were happening.

The top 30 *strongly closed* itemsets mined from hours leading to the MTV Europe Music Awards contain many itemsets about the events and the artists, especially that this event used social media to collect audience votes. Voting for the award winner is a good example of a topic where people have strongly different opinions. We have shown earlier, in figure 4.5, that such different opinions are all reported as separate *strongly closed* itemsets, because clustering using the postings lists avoid forming incohesive clusters. We also note that different *strongly closed* itemsets from one such topic can all occupy high ranks. For example, the top 10 itemsets for the hour starting 9:30 EST on November 11th are shown in table 5.2.2. The number in parenthesis next to each item is the number of

¹³http://en.wikipedia.org/wiki/2012_MTV_Europe_Music_Awards

¹⁴http://en.wikipedia.org/wiki/American_Music_Awards_of_2012

1	#emawinbieber(43), #mtvema(92), URL(31), at(47), be(77), bieber(62), big(81), i(24), justin(62), pick(50), the(62), think(48), tweet(74), will(77), winner(67), your(65)
2	#jonasbrothers(3), #emawinjoej(4), #joejonas(3)
3	#emawinbieber(1), #justinbieber(2)
4	#emawinbieber(2), #justinbieber(2), #emavotebieber(4)
5	chelsea(2), liverpool(2), vs(3)
6	#staycurrent(3), you(1), @buysocialclout(2), can(1), follow(1), today(1)
7	#emavoteonedirection(1), URL(2), emavoteonedirection(1)
8	#emawinbieber(1), #justinbieber(2), @justinbieber(2)
9	remembrance(2), day(1)
10	veteran(2), day(1), s(2)

Table 5.3: Top 10 itemsets in an hour in the morning of the day of MTVEMA

itemsets in the cluster which contain this item. Items in the itemset are sorted by their frequency in this hour, in descending order. The first itemset in the top 10 reflect the winner, Justin Bieber. Itemsets about other topics also appear in the top 10; there was a football game between Chelsea and Liverpool, and November 11th is the day World War I ended (observed as Remembrance day in the United Kingdom and the Commonwealth of Nations, and Veteran’s day in the USA). Several other hours have itemsets supporting different artists in the top 30, which are all variations of the itemset: {#emawinARTIST, ARTIST, NAME, #mtvema, i, think, will, be, the, big, winner, tweet, your, pick, at, URL¹⁵}.

The American Music Awards (AMAs) did not use social media to collect votes, so there was not as much posts about it as MTVEMA. However, itemsets about it still occupied top positions. There are several hours where the itemset {american, music, awards} occupied the top position, but the hour we choose to show as an example does not. We chose this hour because it contains an itemset that offers more information than what is contained in other itemsets from other hours; itemsets during this event were mainly names of artists, or the name of the event. The itemsets from this hour, shown in table 5.2.2, contains an itemset (at position 26) about how Justin Bieber appeared with his mother Pattie

¹⁵<http://ema-twittertracker.mtv.com/live/predict.html>

Mallette¹⁶ instead of his usual escort Selena Gomez, his ex-girlfriend. This is interesting news, specially following Justin and Selena’s break up that warranted both of their names a surge in their popularity as search queries on November 9th and 10th. The news of the break up also appeared in the top 30 itemsets. Tweets from the postings lists of the itemset {justin, bieber, and, selena, gomez} at that time period include shares of news articles such as: “*Justin Bieber and Selena Gomez Break Up* <http://eonli.ne/SByVnK>” *WHAT?!?!?! Is this truee? @justinbieber @selenagomez*, sad comments such as: *Justin and Selena broke up? O_o the awkward moment when he wrote a song about her and another with her name in it...*, and some hopeful comments: *apparently there is a rumor that justin bieber and selena gomez broke up? GO @jessleigh1122 GO!!!!*.

The itemsets in table 5.2.2 also include the name of another artists that is not in the Top Chart. Carly Rae Jepsen also performed and won at the American Music Awards¹⁷.

People appearing in Top Charts but not in the synopses

Out of the five pop artists appearing in the Top Chart, as shown in figure 5.2.2, the names of only four appeared repeatedly in the top 30 itemsets of hours in the days of their respective higher search volume. The name of Rihanna never appeared in any *strongly closed* itemset. Actually it rarely appeared in any itemset of length 2 or more. It is not unusual for an itemset about an artist to be merely his or her name, and several artists, like Rihanna, are known by their first name only. It is impossible for an itemset of length 1 to be *distinct*, since this condition requires that the itemsets be implied by one of its subsets with high confidence. Therefore *distinct* and *strongly closed* itemsets would never capture information in an itemset of length 1. The same happens with one more of the names in the Top Chart, Abraham Lincoln. During the month of November people were talking on Twitter about Lincoln the movie and this was reflected in mining the term as a closed 1-itemset from several hours throughout the month. However, there were not any closed itemsets longer than 1 item, and thus there was no mention of Lincoln the movie in

¹⁶<http://hollywoodlife.com/2012/11/18/justin-bieber-brings-mom-to-amas-without-selena-gomez/#!1/nokia-jenny-mccarthy/>

¹⁷<http://blog.muchmusic.com/carly-rae-jepsen-wins-at-the-american-music-awards-puts-on-fun-and-flir>

1	@vinoalan(1), #vinofollowspree(2)
2	vote(2), #peopleschoice(2), voted(1), retweet(2), @peopleschoice(4), URL(3), for(1), i(1), just(1), to(2), via(3)
3	@simoncowell(2), follow(2), me(2), please(2), to(1)
4	grinch(2), stole(2), all(3), christmas(2), for(2), haters(3), how(1), is(1), the(5), this(1)
5	you(1), @real.liam.payne(2), thank(2)
6	grinch(3), stole(2), how(1), the(3)
7	all(3), amas(1), for(3), is(2), the(4), this(2), watching(2)
8	grinch(1), stole(2), christmas(2)
9	jepsen(3), carly(1), rae(2)
10	dumb(1), dumber(2)
11	3(2), @real.liam.payne(2), URL(1), live(1), on(1)
12	hi(2), say(1), to(2)
13	oz(1), wizard(2), of(2), the(2)
14	liam(1), @real.liam.payne(1), URL(1), ama(1), live(1), on(2), the(1)
15	american(2), awards(2), music(3)
16	ama(3), award(3), first(1), night(1), of(3), s(2), the(6), watching(1)
17	@simoncowell(1), @real.liam.payne(1), call(2), follow(2), me(3), so(2)
18	pone(3), triste(2), qu(1), te(3)
19	is(1), on(1), red(1), the(2), wanted(1)
20	@simoncowell(1), simon(3), call(1), follow(3), me(2), please(3)
21	american(2), awards(1), music(2), the(2)
22	carpet(2), red(1)
23	oz(1), wizard(3), of(2)
24	imessage(2), is(1)
25	series(1), survivor(2)
26	and(3), justin(2), pattie(2)
27	you(4), @real.liam.payne(3), URL(1), found(2), i(2), live(1), love(1), on(1)
28	#getglue(2), URL(1)
29	#gameinsight(1), #android(2), #androidgames(1), URL(3), android(1)
30	#gameinsight(1), #android(2), #androidgames(3)

Table 5.4: Top 30 itemsets from 18:30 EST on the day of American Music Awards (AMAs)

the synopses. This shortcoming can be overcome by setting a lower support threshold so that longer itemsets are mined.

Another two celebrities who appear in the top 10 chart but not the synopses are Kim Kardashian and Michael Jordan. The volume of queries about Kim Kardashian is different from other artists, since she is not a pop music artist and does not participate in music awards events. For these two celebrities, the volume of queries does not have high peaks, and the curve is relatively flat. Itemsets mined from Twitter rarely mention either of them. Of all itemsets mined from hour long epochs throughout the month of November, only 37 closed itemsets contain “Kim Kardashian” or “@kimkardashian”, and none contain “Michael Jordan” or “@jumpman23”. This means that the number of occurrences of their names is lower than the support threshold in all or almost all of the 1440 overlapping hour epochs in the month. We do not know why the interest in Kim Kardashian and Michael Jordan shown by Google search users was not reflected by Twitter users, or maybe this is an artifact of the algorithms used to enumerate the Top Charts.

Chapter 6

Conclusion and Future Work

We have proposed a method for efficiently creating temporal synposes of social media streams, based on a frequent itemset mining algorithm that is suitable for sparse data, LCM. The synposes can be used to explore social media, and they can also reflect what is going on in the real world. Several studies have shown that social media is a good source of information about real world events, offering timely information from a variety of sources. Actually some of the information on social media is not available in traditional news media, or might be available on social media first then get picked up by news media. Social media allows citizen journalists and ordinary people to make posts that can get disseminated widely based on the quality of the post, and regardless of the fame of the poster. Our methods focus on the content of posts, without explicitly giving different weights to posts made by different types of account owners. Thus the results reflect what social media users find interesting at any point in time, giving voice to ordinary people and allowing the dynamics of social media to surface content that users find interesting.

6.1 Contributions

Our method summarizes an hour of Twitter data (119,035.49 tweets on average) into 139.1 itemsets in 1,945.68 milliseconds on average, and scales for longer epochs of data. The

direct application of LCM on one-hour epochs of Twitter data results in an average of 61,505.16 closed itemsets and takes 2,506.58 milliseconds on average. The improvement is due to the following contributions:

- Using variable length term N-grams as items, to mitigate the effect of the skewness of the frequency distribution of unigrams. Frequently occurring sequential multi-word expressions are concatenated into one item during tokenization, resulting in transactions made up of items coming from a relatively flat distribution. We proposed a simple method to do this; tokenizing documents (Tweets) into unigrams, then repeatedly replacing any token whose probability exceeds a threshold by two tokens formed by concatenating to it the unigrams before and after it. This simple preprocessing step results in dropping the number of itemsets mined from an hour long epoch from more than 60,000 to just above 6,000 itemsets. The reduction comes from avoiding the creation of itemsets that are combinations of function words. This preprocessing step was successfully used to improve results from the MTV algorithm [57] as well as our algorithm.
- The use of the *distinct* condition to select only itemsets which are distinctly different from their closed subsets. The condition strengthens the closure property such that an itemset is included in the mining results only if its support is different from the support of its longest subset by a substantial proportion. The closed property is easily satisfied by itemsets that are trivially different from their subsets, which results in many uninteresting itemsets specially short ones – itemsets made up of a commonly used term along with another term that does not have a strong association with it. Such itemsets are filtered out by the *distinct* condition, given that transactions are made up variable length N-grams as items.
- *Distinct* itemsets are clustered into *strongly closed* itemsets for a further reduction in the number of itemsets. We propose an efficient clustering scheme that exploits the order in which itemsets are generated by LCM. The resulting clusters group together itemsets belonging to a certain fine grained topic, or an opinion within a topic, and avoids forming non-coherent clusters. The union of items in each cluster is a *strongly*

closed itemset, which groups together several *distinct* itemsets, reducing redundancy in the mining results.

Another important contribution is a formula for ranking *strongly closed* itemsets based on their temporal novelty, using the collective Information Gain of itemsets in the cluster, taking into account the dependence between itemsets in a cluster.

6.2 Results

We have shown the effectiveness of our methods by applying them to tweets posted in the time frames of various real world events. We use a minimum support threshold and an epoch length that we set according to certain characteristics of the data, which can be easily found out if the dataset is changed. However, we set the parameter that controls the distinctiveness between two itemsets, κ , to an arbitrary value according to its definition.

We have shown that using our proposed methods and ranking formula the top *strongly closed* itemsets from different hours are relevant to the real world events known to be taking place at those hours. We have shown how the top 3 itemsets from different hours in the day the US presidential elections can be used to compose a summary of how the day unfolded. We also compared our synopses for that day to the ones created by the MTV algorithm [57]. We have shown that our algorithm efficiently creates synopses that cover the information in the synopses created by MTV – itemsets that are picked by the MTV algorithm are almost all among the top 30 *strongly closed* itemsets. We have also assessed the quality of the synopses against data from Google Trends. Google Trends was used to learn about real world events and entities in which Google search users were particularly interested at different periods of time. The synopses were then examined to see if they also reflect similar user interest. We have found that there is a strong similarity in the majority of cases, and we analyzed the cases where our synopses did not reflect the same user interest as Google Trends.

6.2.1 The Complete Results of Mining the Twitter Dataset

The complete results of mining all one hour epochs (with half hour time steps) of Twitter data collected in the last 3 months of 2012 is available for download¹. The mining results at different values of κ , for *minimum support* 10, are available as GZipped Tar balls. The result of each epoch is a tab separated file including the itemsets, their postings lists, and values of several ranking formulae. The itemsets are sorted in the descending order of the ranking formula used in this thesis (column 14).

6.3 Future Work

This work provides a foundation for the use of social in temporal information retrieval. The mining results captures people’s interest at different times, and can be used to create a temporal profile for a query. This profile can then be used for weighting documents according to the time in which they were published.

The synopses can be also directly used for temporal query expansion. Terms from itemsets relevant to a query can be used for query expansion, thus acting as precomputed results of pseudo-relevance feedback. It is also possible to use terms from documents where relevant itemsets occur as the resultset to be used as input for pseudo-relevance feedback. The ranks of itemsets can be used to weight expansion terms, if temporal novelty is desired.

The efficiency of the proposed algorithms makes it possible to use the mining results as input to other real-time systems. For example, association rules created from the itemsets such that the consequent is a hashtag can be used to help a user discover hashtags he can add to his post (this can also be done as automatic document expansion).

Another possible future direction is to use itemsets that appear as a sequence for building extractive coherent summaries of the social media stream.

We also wish to explore ways to make use of the temporal signal during mining, such as when calculating similarity during clustering.

¹http://plg.uwaterloo.ca/~yaboulina/thesis_results/twitter_synopses/1hr+30min_ngram5-relsupp10_oct-nov-dec/

APPENDICES

Appendix A

Unigram Tokenization

The data is tokenized into unigrams using white space and punctuation marks as separators. The characters '@', '#' and '_' are allowed within a unigram, while the rest of the punctuation marks are treated as delimiters. Only Latin characters (Unicode code points less than 'u024F') and numbers are allowed within a unigram. Any other characters that are not delimiters are skipped. If a unigram is a number then the dot and the comma characters are allowed within it. The tokenizer also does the following:

- All URLs (from 'http[s]:' to the next whitespace) are replaced by the unigram "URL".
- Runs of the same character is reduced to only 3 repetitions (for example, "coool" is replaced by "cool").
- Hashtags are stored twice, with and without the '#' character. The tag without the '#' character is left at tag's positions, while the other is appended at the end of the tweet. This handles cases where the hashtag is used in place of a word, such as "president #obama...".
- The apostrophe used in contractions (can't, don't, ..etc) is removed from the unigram, but does not act as a delimiter.

References

- [1] Younos Abounaga and Charles L. A. Clarke. Frequent itemset mining for query expansion in microblog ad-hoc search. In *Proceedings of the 21st TREC Conference*, Text Retrieval Evaluation Conference (TREC), Gaithersburg, MD, USA, 2012.
- [2] Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. Mining association rules between sets of items in large databases. In *ACM SIGMOD Record*, volume 22, pages 207–216. ACM, 1993.
- [3] Rakesh Agrawal, Ramakrishnan Srikant, et al. Fast algorithms for mining association rules. In *Proc. 20th Int. Conf. Very Large Data Bases, VLDB*, volume 1215, pages 487–499, 1994.
- [4] Mubarak Albathan, Yuefeng Li, and Abdulmohsen Algarni. Using patterns co-occurrence matrix for cleaning closed sequential patterns for text mining. In *2012 IEEE/WIC/ACM International Conference on Web Intelligence*, pages 201–205. IEEE, 2012.
- [5] James Allan. *Topic detection and tracking: event-based information organization*, volume 12. Kluwer Academic Pub, 2002.
- [6] James Allan, Jaime G Carbonell, George Doddington, Jonathan Yamron, and Yiming Yang. Topic detection and tracking pilot study final report. 1998.
- [7] James Allan, Victor Lavrenko, Daniella Malin, and Russell Swan. Detections, bounds, and timelines: Umass and tdt-3. In *Proceedings of Topic Detection and Tracking Workshop (TDT-3)*, pages 167–174. Vienna, VA, 2000.

- [8] Roberto J Bayardo, Yiming Ma, and Ramakrishnan Srikant. Scaling up all pairs similarity search. In *Proceedings of the 16th international conference on World Wide Web*, pages 131–140. Citeseer, 2007.
- [9] Roberto J Bayardo Jr and Rakesh Agrawal. Mining the most interesting rules. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 145–154. ACM, 1999.
- [10] Hila Becker, Mor Naaman, and Luis Gravano. Beyond trending topics: Real-world event identification on twitter. In *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media (ICWSM’11)*, 2011.
- [11] Albert Bifet and Ricard Gavaldà. Learning from time-changing data with adaptive windowing. In *SIAM International Conference on Data Mining*, pages 443–448, 2007.
- [12] Jean-Francois Boulicaut, Artur Bykowski, and Christophe Rigotti. Approximation of frequency queries by means of free-sets. In *Principles of Data Mining and Knowledge Discovery*, pages 75–85. Springer, 2000.
- [13] Jean-François Boulicaut, Artur Bykowski, and Christophe Rigotti. Free-sets: a condensed representation of boolean data for the approximation of frequency queries. *Data Mining and Knowledge Discovery*, 7(1):5–22, 2003.
- [14] Gerlof Bouma. Normalized (pointwise) mutual information in collocation extraction. *Proceedings of GSCL*, pages 31–40, 2009.
- [15] Danah Boyd, Scott Golder, and Gilad Lotan. Tweet, tweet, retweet: Conversational aspects of retweeting on twitter. In *System Sciences (HICSS), 2010 43rd Hawaii International Conference on*, pages 1–10. IEEE, 2010.
- [16] Sergey Brin, Rajeev Motwani, and Craig Silverstein. Beyond market baskets: generalizing association rules to correlations. In *ACM SIGMOD Record*, volume 26, pages 265–276. ACM, 1997.
- [17] Toon Calders and Bart Goethals. Mining all non-derivable frequent itemsets. In *Principles of Data Mining and Knowledge Discovery*, pages 74–86. Springer, 2002.

- [18] Toon Calders and Bart Goethals. Non-derivable itemset mining. *Data Mining and Knowledge Discovery*, 14(1):171–206, 2007.
- [19] Deepayan Chakrabarti and Kunal Punera. Event summarization using tweets. In *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*, pages 66–73, 2011.
- [20] Jaeho Choi and W Bruce Croft. Temporal models for microblogs. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, pages 2491–2494. ACM, 2012.
- [21] Edith Cohen, Mayur Datar, Shinji Fujiwara, Aristides Gionis, Piotr Indyk, Rajeev Motwani, Jeffrey D. Ullman, and Cheng Yang. Finding interesting associations without support pruning. *Knowledge and Data Engineering, IEEE Transactions on*, 13(1):64–78, 2001.
- [22] Jacques JF Commandeur and Siem Jan Koopman. *An introduction to state space time series analysis*. OUP Oxford, 2007.
- [23] John N Darroch and Douglas Ratcliff. Generalized iterative scaling for log-linear models. *The annals of mathematical statistics*, 43(5):1470–1480, 1972.
- [24] Miles Efron, Peter Organisciak, and Katrina Fenlon. Improving retrieval of short texts through document expansion. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, pages 911–920. ACM, 2012.
- [25] Venkatesh Ganti, Johannes Gehrke, and Raghu Ramakrishnan. Mining data streams under block evolution. *ACM SIGKDD Explorations Newsletter*, 3(2):1–10, 2002.
- [26] Liqiang Geng and Howard J Hamilton. Interestingness measures for data mining: A survey. *ACM Computing Surveys (CSUR)*, 38(3):9, 2006.
- [27] Chris Giannella, Jiawei Han, Jian Pei, Xifeng Yan, and Philip S Yu. Mining frequent patterns in data streams at multiple time granularities. *Next generation data mining*, 212:191–212, 2003.

- [28] Anna C Gilbert, Yannis Kotidis, S Muthukrishnan, and Martin J Strauss. Surfing wavelets on streams: One-pass summaries for approximate aggregate queries. In *Proceedings of the International Conference on Very Large Data Bases*, pages 79–88, 2001.
- [29] Gösta Grahne and Jianfei Zhu. Reducing the main memory consumptions of fpmax* and fpclose. In *Proc. Workshop Frequent Item Set Mining Implementations (FIMI 2004, Brighton, UK), Aachen, Germany*. Citeseer, 2004.
- [30] Jiawei Han, Jian Pei, and Yiwen Yin. Mining frequent patterns without candidate generation. In *ACM SIGMOD Record*, volume 29, pages 1–12. ACM, 2000.
- [31] Jiawei Han, Jian Pei, Yiwen Yin, and Runying Mao. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data mining and knowledge discovery*, 8(1):53–87, 2004.
- [32] Matthew Hoffman, David M Blei, and Francis Bach. Online learning for latent dirichlet allocation. *Advances in Neural Information Processing Systems*, 23:856–864, 2010.
- [33] Liangjie Hong and Brian D Davison. Empirical study of topic modeling in twitter. In *Proceedings of the First Workshop on Social Media Analytics*, pages 80–88. ACM, 2010.
- [34] Yifan Hu. Efficient, high-quality force-directed graph drawing. *Mathematica Journal*, 10(1):37–71, 2005.
- [35] Jonathan Hurlock and Max L Wilson. Searching twitter: Separating the tweet from the chaff. *Proc. ICWSM 2011*, 2011.
- [36] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613. ACM, 1998.
- [37] Szymon Jaroszewicz and Dan A Simovici. A general measure of rule interestingness. In *Principles of Data Mining and Knowledge Discovery*, pages 253–265. Springer, 2001.

- [38] Szymon Jaroszewicz and Dan A Simovici. Interestingness of frequent itemsets using bayesian networks as background knowledge. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 178–186. ACM, 2004.
- [39] Akshay Java, Xiaodan Song, Tim Finin, and Belle Tseng. Why we twitter: understanding microblogging usage and communities. In *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis*, pages 56–65. ACM, 2007.
- [40] Rosie Jones and Fernando Diaz. Temporal profiles of queries. *ACM Transactions on Information Systems (TOIS)*, 25(3):14, 2007.
- [41] Roberto J. Bayardo Jr., Bart Goethals, and Mohammed Javeed Zaki, editors. *FIMI '04, Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations, Brighton, UK, November 1, 2004*, volume 126 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2004.
- [42] Won-Young Kim, Young-Koo Lee, and Jiawei Han. Ccmine: Efficient mining of confidence-closed correlated patterns. In *Advances in Knowledge Discovery and Data Mining*, pages 569–579. Springer, 2004.
- [43] Jon Kleinberg. Bursty and hierarchical structure in streams. *Data Mining and Knowledge Discovery*, 7(4):373–397, 2003.
- [44] Marzena Kryszkiewicz. Concise representation of frequent patterns based on disjunction-free generators. In *Data Mining, 2001. ICDM 2001, Proceedings IEEE International Conference on*, pages 305–312. IEEE, 2001.
- [45] Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. What is twitter, a social network or a news media? In *Proceedings of the 19th international conference on World wide web*, pages 591–600. ACM, 2010.
- [46] Cher Han Lau, YueFeng Li, and Dian Tjondronegoro. Microblog retrieval using topical features and query expansion. In *Proceedings of the 20th TREC Conference, Text Retrieval Evaluation Conference (TREC)*, Gaithersburg, MD, USA, 2011.

- [47] Young-Koo Lee, Won-Young Kim, Y Dora Cai, and Jiawei Han. Comine: Efficient mining of correlated patterns. In *Proceedings of the Third IEEE International Conference on Data Mining*, page 581, 2003.
- [48] Janette Lehmann, Bruno Gonçalves, José J Ramasco, and Ciro Cattuto. Dynamical classes of collective attention in twitter. In *Proceedings of the 21st international conference on World Wide Web*, pages 251–260. ACM, 2012.
- [49] Jure Leskovec, Lars Backstrom, and Jon Kleinberg. Meme-tracking and the dynamics of the news cycle. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 497–506. ACM, 2009.
- [50] David D Lewis. An evaluation of phrasal and clustered representations on a text categorization task. In *Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 37–50. ACM, 1992.
- [51] Yuefeng Li, Abdulmohsen Algarni, and Ning Zhong. Mining positive and negative patterns for relevance feature discovery. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 753–762. ACM, 2010.
- [52] Zhenmin Li, Zhifeng Chen, Sudarshan M Srinivasan, and Yuanyuan Zhou. C-miner: Mining block correlations in storage systems. In *Proceedings of the 3rd USENIX Conference on File and Storage Technologies*, volume 186. USENIX Association, 2004.
- [53] Jimmy Lin and Gilad Mishne. A study of “churn” in tweets and real-time search queries. In *Sixth International AAAI Conference on Weblogs and Social Media*, 2012.
- [54] Bing Liu, Wynne Hsu, and Yiming Ma. Mining association rules with multiple minimum supports. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 337–341. ACM, 1999.
- [55] Guimei Liu, Haojun Zhang, and Limsoon Wong. Finding minimum representative pattern sets. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 51–59. ACM, 2012.

- [56] GILAD LOTAN, ERHARDT GRAEFF, MIKE ANANNY, DEVIN GAFFNEY, IAN PEARCE, and DANAH BOYD. The revolutions were tweeted: Information flows during the 2011 tunisian and egyptian revolutions. 2011.
- [57] Michael Mampaey, Nikolaj Tatti, and Jilles Vreeken. Tell me what i need to know: succinctly summarizing data with itemsets. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 573–581. ACM, 2011.
- [58] Michael Mathioudakis and Nick Koudas. Twittermonitor: trend detection over the twitter stream. In *Proceedings of the 2010 international conference on Management of data*, pages 1155–1158. ACM, 2010.
- [59] Ahmed Metwally, Divyakant Agrawal, and Amr El Abbadi. Duplicate detection in click streams. In *Proceedings of the 14th international conference on World Wide Web*, pages 12–21. ACM, 2005.
- [60] Mor Naaman, Hila Becker, and Luis Gravano. Hip and trendy: Characterizing emerging trends on twitter. *Journal of the American Society for Information Science and Technology*, 62(5):902–918, 2011.
- [61] Mor Naaman, Jeffrey Boase, and Chih-Hui Lai. Is it really about me?: message content in social awareness streams. In *Proceedings of the 2010 ACM conference on Computer supported cooperative work*, pages 189–192. ACM, 2010.
- [62] Brendan O’Connor, Michel Krieger, and David Ahn. Tweetmotif: Exploratory search and topic summarization for twitter. *Proceedings of ICWSM*, pages 2–3, 2010.
- [63] Edward R Omiecinski. Alternative interest measures for mining associations in databases. *Knowledge and Data Engineering, IEEE Transactions on*, 15(1):57–69, 2003.
- [64] Ruchi Parikh and Kamalakar Karlapalem. Et: events from tweets. In *Proceedings of the 22nd international conference on World Wide Web companion*, pages 613–620. International World Wide Web Conferences Steering Committee, 2013.

- [65] Nicolas Pasquier, Yves Bastide, Rafik Taouil, and Lotfi Lakhal. Discovering frequent closed itemsets for association rules. In *Database Theory—ICDT’99*, pages 398–416. Springer, 1999.
- [66] Nicolas Pasquier, Yves Bastide, Rafik Taouil, and Lotfi Lakhal. Efficient mining of association rules using closed itemset lattices. *Information systems*, 24(1):25–46, 1999.
- [67] Jian Pei, Jiawei Han, Runying Mao, et al. Closet: An efficient algorithm for mining frequent closed itemsets. In *ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery*, volume 4, pages 21–30, 2000.
- [68] Saša Petrović, Miles Osborne, and Victor Lavrenko. Streaming first story detection with application to twitter. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 181–189. Association for Computational Linguistics, 2010.
- [69] Ana-Maria Popescu, Marco Pennacchiotti, and Deepa Paranjpe. Extracting events and event descriptions from twitter. In *Proceedings of the 20th international conference companion on World wide web*, pages 105–106. ACM, 2011.
- [70] Bruno Pôssas, Nivio Ziviani, Wagner Meira Jr, and Berthier Ribeiro-Neto. Set-based model: a new approach for information retrieval. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 230–237. ACM, 2002.
- [71] Carlos Ramisch, Vitor De Araujo, and Aline Villavicencio. A broad evaluation of techniques for automatic acquisition of multiword expressions. In *Proceedings of ACL 2012 Student Research Workshop*, pages 1–6. Association for Computational Linguistics, 2012.
- [72] Alan Ritter, Oren Etzioni, Sam Clark, et al. Open domain event extraction from twitter. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1104–1112. ACM, 2012.

- [73] Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, and Mike Gatford. Okapi at trec-3. *NIST SPECIAL PUBLICATION SP*, pages 109–109, 1995.
- [74] Masakazu Seno and George Karypis. Finding frequent patterns using length-decreasing support constraints. *Data Mining and Knowledge Discovery*, 10(3):197–228, 2005.
- [75] Abraham Silberschatz and Alexander Tuzhilin. What makes patterns interesting in knowledge discovery systems. *Knowledge and Data Engineering, IEEE Transactions on*, 8(6):970–974, 1996.
- [76] Avi Silberschatz and Alexander Tuzhilin. On subjective measures of interestingness in knowledge discovery. In *Proceedings of KDD-95: First International Conference on Knowledge Discovery and Data Mining*, pages 275–281, 1995.
- [77] Pang-Ning Tan, Vipin Kumar, and Jaideep Srivastava. Selecting the right interestingness measure for association patterns. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 32–41. ACM, 2002.
- [78] Nikolaj Tatti. Maximum entropy based significance of itemsets. *Knowledge and Information Systems*, 17(1):57–77, 2008.
- [79] Jaime Teevan, Daniel Ramage, and Meredith Ringel Morris. # twittersearch: a comparison of microblog search and web search. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 35–44. ACM, 2011.
- [80] Fernando Tusell. Kalman filtering in r. *Journal of Statistical Software*, 39(2):1–27, 2011.
- [81] Takeaki Uno, Masashi Kiyomi, and Hiroki Arimura. Lcm ver. 2: Efficient mining algorithms for frequent/closed/maximal itemsets. In *Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI 04)*, 2004.

- [82] Jilles Vreeken, Matthijs Van Leeuwen, and Arno Siebes. Krimp: mining itemsets that compress. *Data Mining and Knowledge Discovery*, 23(1):169–214, 2011.
- [83] Chao Wang and Srinivasan Parthasarathy. Summarizing itemset patterns using probabilistic models. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 730–735. ACM, 2006.
- [84] Ke Wang, Yu He, and Jiawei Han. Mining frequent itemsets using support constraints. In *Proceedings of the 26th International Conference on Very Large Data Bases*, pages 43–52, 2000.
- [85] Ke Wang, Liu Tang, Jiawei Han, and Junqiang Liu. *Top down fp-growth for association rule mining*. Springer, 2002.
- [86] Jianshu Weng and Bu-Sung Lee. Event detection in twitter. In *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media*, volume 3, 2011.
- [87] Sheng-Tang Wu, Yuefeng Li, and Yue Xu. Deploying approaches for pattern refinement in text mining. In *Data Mining, 2006. ICDM’06. Sixth International Conference on*, pages 1157–1161. IEEE, 2006.
- [88] Sheng-Tang Wu, Yuefeng Li, Yue Xu, Binh Pham, and Phoebe Chen. Automatic pattern-taxonomy extraction for web mining. In *Web Intelligence, 2004. WI 2004. Proceedings. IEEE/WIC/ACM International Conference on*, pages 242–248. IEEE, 2004.
- [89] Dong Xin, Hong Cheng, Xifeng Yan, and Jiawei Han. Extracting redundancy-aware top-k patterns. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 444–453. ACM, 2006.
- [90] Dong Xin, Jiawei Han, Xifeng Yan, and Hong Cheng. Mining compressed frequent-pattern sets. In *Proceedings of the 31st international conference on Very large data bases*, pages 709–720. VLDB Endowment, 2005.

- [91] Hui Xiong, P-N Tan, and Vipin Kumar. Mining strong affinity association patterns in data sets with skewed support distribution. In *Data Mining, 2003. ICDM 2003. Third IEEE International Conference on*, pages 387–394. IEEE, 2003.
- [92] Xifeng Yan, Hong Cheng, Jiawei Han, and Dong Xin. Summarizing itemset patterns: a profile-based approach. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 314–323. ACM, 2005.
- [93] Xintian Yang, Amol Ghoting, Yiye Ruan, and Srinivasan Parthasarathy. A framework for summarizing and analyzing twitter feeds. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 370–378. ACM, 2012.
- [94] Yiming Yang, Tom Pierce, and Jaime Carbonell. A study of retrospective and on-line event detection. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 28–36. ACM, 1998.
- [95] Limin Yao, David Mimno, and Andrew McCallum. Efficient methods for topic model inference on streaming document collections. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 937–946. ACM, 2009.
- [96] Mohammed J Zaki and Ching-Jui Hsiao. Charm: An efficient algorithm for closed itemset mining. In *2nd SIAM international conference on data mining*, volume 15, pages 457–73, 2002.
- [97] Mohammed J Zaki, Srinivasan Parthasarathy, Mitsunori Ogihara, Wei Li, et al. New algorithms for fast discovery of association rules. In *3rd Intl. Conf. on Knowledge Discovery and Data Mining*, volume 20, pages 283–286, 1997.
- [98] Yunyue Zhu and Dennis Shasha. Efficient elastic burst detection in data streams. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 336–345. ACM, 2003.