

See discussions, stats, and author profiles for this publication at: <http://www.researchgate.net/publication/265690208>

Atmospheric Environment and Quality of Life Information Extraction from Twitter with the use of Self-Organizing Maps

ARTICLE *in* JOURNAL OF ENVIRONMENTAL INFORMATICS · JULY 2015

Impact Factor: 3.77 · DOI: 10.3808/jei.201500311

5 AUTHORS, INCLUDING:



Marina Riga

The Centre for Research and Technology, Hel...

12 PUBLICATIONS **19** CITATIONS

SEE PROFILE



Markus Stocker

University of Eastern Finland

25 PUBLICATIONS **261** CITATIONS

SEE PROFILE



Kostas D. Karatzas

Aristotle University of Thessaloniki

136 PUBLICATIONS **692** CITATIONS

SEE PROFILE



Mikko Kolehmainen

University of Eastern Finland

89 PUBLICATIONS **1,412** CITATIONS

SEE PROFILE

1 **Authors**

2 Marina Riga^{1,*}, Markus Stocker², Mauno Rönkkö², Kostas Karatzas¹ and Mikko

3 Kolehmainen²

4 ¹ Department of Mechanical Engineering, Aristotle University of Thessaloniki, P.O. Box

5 483, GR-54124, Thessaloniki, Greece

6 ² Department of Environmental Science, University of Eastern Finland, P.O. Box 1627,

7 FI-70211, Kuopio, Finland

8

9 * Corresponding author: Tel: +30 2310 994359; Fax: +30 2310 994176. E-mail address:

10 mriga@isag.meng.auth.gr

- 1 **Title**
- 2 Atmospheric Environment and Quality of Life Information Extraction from Twitter with the
- 3 use of Self-Organizing Maps

DRAFT VERSION

1 **ABSTRACT**

2 The emergence of Web 2.0 technologies has changed dramatically not only the way
3 users perceive the Internet and interact on it but also the way they influence a
4 community and act in real life aspects. With the rapid rise in use and popularity of social
5 media, people tend to share opinions and observations for almost any subject or event in
6 their everyday life. Consequently, microblogging websites have become a rich data
7 source for user-generated information. The leading opportunity is to take advantage of
8 the *wisdom of the crowd* and to benefit from collective intelligence in any applicable
9 domain. Towards this direction, we focus on the problem of mining and extracting
10 knowledge from unstructured textual content, for the atmospheric environment domain
11 and its effect to quality of life. As the main contribution, we propose a combined
12 methodology of unsupervised learning methods for analyzing posts from Twitter and
13 clustering textual data into concepts with semantically similar context. By applying Self-
14 Organizing Maps and k-means clustering, we identify possible inter-relationships and
15 patterns of words used in tweets that can form upper concepts of atmospheric and
16 health related topics of discussion. We achieve to group together tweets, from more
17 generic to more specific description levels of their content, according to the selected
18 number of clusters. Strong clusters with significant semantic relatedness among their
19 content are revealed, and hidden relations between concepts and their related semantics
20 are acquired. The results highlight the potential use of social media text streams as a
21 highly-valued supplement source of environmental information and situation awareness.

22

23 **Keywords**

24 air quality, clustering, computational intelligence, k-means, semantic analysis, self-
25 organizing maps, text mining, twitter

26

1 1. Introduction

2 The advent of Web 2.0 technologies introduced a new way of perceiving the Internet and
3 interacting on it. Within the framework of social media platforms, microblogs and virtual
4 communities, people share different types of content, such as text, images or video,
5 provided either to the general public or towards specific user communities. A novel form
6 of *push-push-pull* communication is established, where direct access, forwarding and
7 further search for information takes places (Kaplan and Haenlein, 2011). Individual users
8 easily communicate news, events, observations and personal opinions, while
9 communities promote their services and information directly to groups of interest. These
10 interaction processes can eventually lead to increased and shared awareness as well as
11 to the discovery of knowledge on a personal level, among the active users.

12 In the current work, we focus on the atmospheric environment domain and its effect to
13 humans' quality of life, where there is a trending involvement and participation of citizens
14 related to environmental and quality of life aspects. With the empowerment of the
15 general public and the provision of technology as a significant step in creating and
16 distributing information, the idea of participatory sensing (Burke et al., 2006) has become
17 more mature and intense nowadays. Ubiquitous data capture can be performed:

- 18 a) either explicitly, by taking advantage of the existing hardware components of
19 commonly used devices in order to record and report directly measurements of a
20 physical entity,
- 21 b) or implicitly, by prompting the user to share its own observations in particular
22 domains of interest.

23 Applications like *PEIR* (Personal Environmental Impact Report) (Mun et al., 2009),
24 *AsTEKa* (Skön et al., 2011) and *Air Quality Egg* (<http://airqualityegg.com/>) fall into the
25 category of automated data capture, while *EnviObserver* (Kotovirta et al., 2012) and
26 *Patient's Hay-fever Diary* (PHD) (<https://www.pollendiary.com/Phd/>) represent cases of

1 direct reporting of personal observation and crowd knowledge from users. In both
2 aforementioned ways of communication, individuals act as *soft sensors* of their
3 surrounding environment with time and location stamp, replacing in this way the
4 traditional, authority-centralized use of *hard sensors* and official monitoring stations (Hall
5 et al., 2008). Sensing the environment through humans' participation can be done near
6 real-time and on a broad scale, even in areas where static monitoring sensors do not
7 operate. Involving citizens in quality of life aspects and making them potential
8 contributors and users of information at the same time, will create a new, bi-directional
9 way of communicating environmental data: from individuals to the community
10 (practitioners) and vice versa. This is expected to transform citizens from passive
11 receivers of information to actors and participants in environmental and quality of life
12 related decision making processes (Karatzas, 2009).

13 By combining Web 2.0 and participatory sensing concepts, our motivation is to benefit
14 from collective intelligence. We aim to mine content from users' activities and social
15 media. Data in any form, either numeric or text, could be analyzed, structured and
16 disseminated properly, according to applications' needs. A lot of effort has been already
17 conducted in the manipulation of numerical air quality measurements (Niska et al., 2004;
18 Voukantsis et al., 2010; Voukantsis et al., 2011) with significant results in mining and
19 forecasting tasks. To the best of our knowledge, there is no published work that deals
20 with the analysis and mining of massive user-generated text data for the atmospheric
21 environment domain, with the aid of computational intelligence. There are authors that
22 have recognized the importance of social media content and its impact to citizens in
23 relation to air pollution (Cairns, 2013). There is the *AirTwitter* project that introduces the
24 idea of potential use of social media to augment air quality event identification
25 (Robinson, 2010). We aim to provide a straightforward methodology for analyzing and
26 clustering text data that are crawled and archived from social media, into classes of

1 similar content. We do not make a real-time analysis; we are interested in analyzing
2 individuals' personal observations expressed in free text and in mining knowledge that
3 can be derived from massive collections of user-generated content, regarding the
4 atmospheric environment and quality of life related issues. The source of data in the
5 current work is Twitter (<http://twitter.com>), one of the most popular microblogging
6 services, which enables its users to post and read short text messages, known as
7 *tweets*, of up to 140 characters. From its launch in 2006, Twitter's user base has been
8 growing exponentially: it counts more than 500 million registered users, where 1/4 of
9 them are considered as extremely active, while 460K accounts are created every day
10 (TechCrunch, 2012). Tweets are publicly visible by default, while users can select to add
11 geo-location information to their posts. Despite the controversial quality and bias of each
12 individual tweet as a unit, the interaction and communication in social media oftenly
13 reflects real-world events and dynamics, especially as the user base of social networks
14 gets wider and more active in producing content about real-world events (Aiello et al.,
15 2013).

16 Our work was motivated by the successful application of computational methods in
17 mining knowledge from social media (Section 2). Still, Twitter is a challenging source of
18 user generated text that differs from a typical document in length and structure. Usually
19 sentences have syntactic or spelling problems, while abbreviations and semantic
20 inconsistencies of words used may increase the complexity of text analysis (Kaufmann
21 and Kalita, 2010). The effective adoption and adaption of text mining methodologies is
22 considered as the main factor for successful results in analyzing the content of our
23 collection.

24 The methodology followed is data-driven. We present an extensive feature selection of
25 words and we create sets-of-words with semantically similar meaning that are used for
26 further encoding of data into vectors. With this implementation we focus on the explicit

1 description of the domain of interest, avoiding at the same time any redundant content
2 that may cause irregularities or misleading results in the clustering process. We
3 demonstrate the use of Self-Organizing Maps (SOMs) and k-means for clustering text
4 into groups with similar patterns of used words and thus similar content and we
5 investigate relationships between formed clusters.

6 The rest of the paper is structured as follows. In Section 2, we review the literature
7 related to our work. In Section 3, we describe the technological background of our
8 methodology as well as the main characteristics of the data (tweets) that have been
9 collected and analyzed. In Section 4, we pinpoint and present the most significant results
10 derived from the clustering process, at different levels of semantic description of clusters.
11 In Section 5, we discuss the results and we highlight the information gained through this
12 process. Finally, in Section 6, we conclude with future aspects and directions of our
13 research.

14

15 **2. Related work**

16 Text mining, or also known as text data mining or text analysis, is the process of
17 extracting non-trivial patterns or high-quality knowledge from unstructured text
18 documents (Tan, 1999). It is an interdisciplinary field that involves different tasks, such
19 as text structure, linguistic preprocessing, pattern recognition, information extraction and
20 visualization.

21 In applied research, a lot of effort has been invested so as to benefit from the potential
22 use and the commercial value of massive collections of textual data. Text mining is
23 applied for business intelligence (Sullivan, 2001) and market trends (Zhang et al., 2011),
24 for event detection (Li et al., 2012; Sadilek et al., 2012), political (Maynard and Funk,
25 2011) or commercial (Mostafa, 2013) opinion mining (Pang and Lee, 2008),
26 security/crisis management (MacEachren et al., 2011) and decision making.

1 Many different methodologies have been developed for text mining tasks, ranging from
2 traditional natural language processing (NLP) techniques and statistics to machine
3 learning methods and their novel extensions. Statistics and probabilistic methods are
4 used in (Hofmann 2001) for latent semantic text analysis of different corpus of
5 documents. The authors map documents to a vector space called latent semantic space
6 that has an order of approximately 130 dimensions. They encode the ‘true’ similarity
7 hidden in the semantics of words with different terms, by forming *sets-of-words* with
8 similar semantic meaning, but they do not address the problem of high dimensionality.
9 Aiello et al. (2013) compare probabilistic, feature-pivot and NLP methods in order to
10 evaluate their effectiveness in detecting events and sensing trending topics in Twitter.
11 Dredze and Paul (2014) employ NLP methods to discover health related issues in social
12 media, while Popescu and Pennacchiotti (2010) use supervised machine learning
13 models for detecting controversial events from Twitter.

14 SOMs have been widely used to perform exploratory analysis of text data. The use of
15 SOM for topic identification of large documents is presented in Yang and Lee (2010).
16 Lagus et al. (2004) developed the WEBSOM method, a software system that extends
17 the SOM principle, by using the so called ‘document maps’ for categorization of massive
18 document collections.

19 Crooks et al. (2013) performed an analysis on Twitter feeds based on hashtags, proving
20 that data streams from social media can be a unique source for rapid geo-located
21 detection of earthquakes. These results motivated our research for analyzing air quality
22 related tweets in order to identify topics of discussion and improve awareness
23 concerning events as they occur, based on observations of users/citizens.

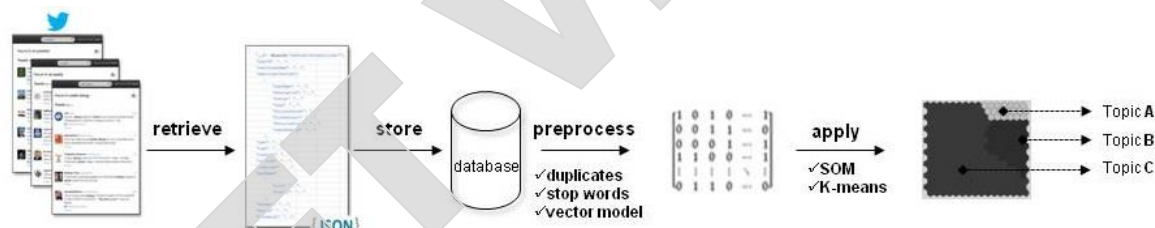
24 Costa et al. (2013) built a model that uses empirical meta-hashtag classes to group
25 together hashtags with similar semantics, for better classification results. This approach

1 is close to our definition of *unified concepts*, words that describe sets-of-words with
 2 similar semantic meaning.

3

4 **3. Materials and Methods**

5 We propose a straightforward methodology for mining text content, retrieved from
 6 Twitter, by combining a set of well-known computational methods that cover the
 7 following main tasks: information retrieval, lexical analysis and preprocessing,
 8 knowledge representation, clustering, information extraction and visualization. The
 9 ultimate aim of this process is to represent text data in an efficient way for further
 10 categorization in semantically similar groups of documents. The overall framework of the
 11 methodology is shown in Figure 1, while each part of it is described separately in the
 12 following subsections (3.1 to 3.3).



13

14

Figure 1: The framework that describes the text mining methodology.

15

16 **3.1 Data retrieval via crawling**

17 With the use of the Twitter Streaming API (<https://dev.twitter.com/docs/streaming-apis>),
 18 we developed a crawler that streams publicly available tweets, targeting at those tweets
 19 that are related to atmospheric environment and health related issues. The distinction of
 20 the topic was feasible by querying tweets that include specific keywords of interest, such
 21 as: air quality, air pollution, pollen allergies, air pollutants, medication, symptoms, etc.
 22 The collection period was from February till middle April 2013, and all tweets included

1 time stamp and geo-location. We only collected tweets that were written in English
 2 language.
 3 The data harvested through Twitter Streaming API are returned in JSON format. In order
 4 to keep the “structured” format of the information derived, we used MongoDB
 5 (<http://www.mongodb.org/>), a well known open source document-oriented database
 6 system. Data were parsed and then stored in a local database for further processing. An
 7 example of the JSON schema is shown in Figure 2.

```

{
  "_id": "511003ee56652adbd1c5d1a",
  "userID": 334068...,
  "userScreenName": "NameXSurnameY",
  "additionalUserInfo":
  {
    "userName": "someUsername",
    "description": "Sport, chocolate, laugh, etc",
    "status": null,
    "lang": "en",
    "followersCount": 220,
    "friendsCount": 216,
    "favoritesCount": 155,
    "isGeoEnabled": true
  },
  "text": "Love running but my asthma doesn't",
  "time": "4.2.2013 18:54:38",
  "location": "Billericay, Essex",
  "latitude": "51.63947377",
  "longitude": "0.4218109",
  "geoNames":
  {
    "city": "Essex",
    "country": "United Kingdom"
  },
  "retweet": 0,
  "streaming": "StatusJSONImpl{createdAt=Mon Feb 04 20:54:38 EET 2013, id=...}"
}

```

8
 9 **Figure 2: An example of the JSON schema used for storing tweets.**

10 A total number of 52,500 tweets were collected from 34,502 unique users. The average
 11 number of words per tweet was approximately 11. The geographical distribution of all
 12 collected tweets is presented in Figure 3. Tweets are concentrated in UK and USA
 13 mostly due to language restrictions in the collection process.



1

2 **Figure 3: Tweets' geographical distribution, via Google Maps (<https://maps.google.com/>).**

3

4 **3.2 Preprocessing data**

5 Tweets are often very peculiar in syntax and use of words, mainly due to the limitation of
6 140 characters and the way that users express themselves in informal text. It is common
7 for Twitter users not to strictly follow the syntactic and grammatical rules of written
8 language; they rather tend to use extensively abbreviations, acronyms, contracted words
9 and emoticons, or even create new words in order to express their opinion in the most
10 compact way.

11 We apply two levels of preprocessing: at first, we clean the text from unnecessary
12 content. Then we turn the text into numerical data in an efficient way, in order to encode
13 the represented information properly and feed it into selected computational methods for
14 further analysis.

15

16 **3.2.1 Removing noisy data**

1 Text analysis refers to the process of deriving high-quality information from text. The
2 process is directly related to the words used in text. For optimized results, it is important
3 to reduce the unnecessary (noisy) content from each tweet, such as:

- 4 a) hyperlinks and re-tweet references (denoted with “http://” and “@” entities,
5 correspondingly),
- 6 b) any punctuation or non-alphabet character,
- 7 c) words with one or two characters, and
- 8 d) English stop-words, meaning words that are frequently used in any text (like
9 articles, pronouns, prepositions, etc.).

10 After the first level of preprocessing, the average number of words per tweet was 9,
11 while the removal of duplicates resulted in a total number of 44,888 unique tweets, being
12 available for further analysis.

13

14 **3.2.2 Representing text documents with vector space model**

15 A common methodology for the analysis and representation of text data in a structured
16 way is to adopt the vector space model, also known as the *bag-of-words* model,
17 introduced by Salton (1975). The general version of this methodology represents each
18 text document d_i as a row vector of t_j features, each of which corresponds to a separate
19 term (word), as shown in Equation (1).

$$20 \quad d_i = [t_1, t_2, \dots, t_n] \quad (1)$$

21 If the term exists in the document, its value in the vector can be either one,
22 corresponding to the occurrence in text, or any non-zero natural number, corresponding
23 to the frequency of its occurrence in text. The whole database will be transformed into an
24 $m \times n$ matrix, of m documents (tweets) and n terms, with zero and non-zero values, for the
25 frequency of occurrence or the non-occurrence of each term, respectively (Berry, 2004).

1 Usually, the collection of words is unordered, thus no information about word position in
2 the document is taken into account.

3 A bag-of-words may include all words used in the total collection of documents being
4 investigated. Such an implementation would lead to a high-dimensional, sparse matrix
5 with few combinations of words per document. The redundant information carried
6 throughout the process, makes the task of extracting similar patterns from data more
7 difficult for clustering algorithms. Dimensionality reduction is of great importance in order
8 to manipulate efficiently massive amounts of text data and improve the solution to the
9 problem at hand. A different implementation of a bag-of-words would include only the n
10 most frequent words in the database, but this representation method would leave out
11 words of high interest that might not appear often in tweets. In our approach, we select
12 to include in the bag-of-words model the most frequent (above a threshold) words used
13 in our collection of tweets, combined with some less frequent words that were empirically
14 considered as important within the context of our domain of interest. We should prompt
15 here that we have already removed any data considered as noisy (Section 3.2.1).

16 For dimensionality reduction reasons, it became evident that we had to take into
17 consideration the meaning (semantics) of words, in the bag-of-words concept. For
18 example, words like *flu*, *cold* and *sick* can be considered as one term, under the label of
19 *medical condition* since these three words describe in general a medical condition which
20 could be related to existing environmental conditions.

21 Hence, by having created the initial bag-of-words, we define sets of similar words that
22 are not necessarily synonyms at the same time, but can be considered as having the
23 same semantical meaning under the context of atmospheric environment and its effect to
24 humans' quality of life. Each set of words S includes a different number of similar words
25 that describe explicitly an upper semantic level of related content (unified concept). All
26 words and sets formed define the *bag-of-sets of words*, as given in detail in Table 1. The

1 selection of words and related unified concept has been done empirically, establishing
 2 thus implicitly the semantic context within which the analysis is going to be performed.
 3 We state that the current form gives a well-defined context of the most representative
 4 aspects that describe atmospheric environment and health related issues in an informal
 5 way. This context can be re-produced in any other combination of words, thus changing
 6 the interpretation of the analysis' results. Any change in the current definition of sets of
 7 words affects the semantic interpretation of results and the potential groups created from
 8 the clustering process.

9 **Table 1: Sets of words, with similar context, forming the bag-of-sets of words for the analysis**

#	Words in set	Unified concept
1	air, atmosphere, atmospheric	air
2	eye, nose, skin, mouth, throat, stomach, head, lungs, face, nasal, heart, respiratory, chest, body	body part/organ
3	pollution, pollute, pollutants	pollution
4	itch, itchiness	itch
5	sneeze, sneezing	sneeze
6	cough, coughing	cough
7	run, running, runny nose*	runny
8	flu, cold, sick, ill, headache, asthma, rhinitis, disease, fever	medical condition
9	city, town, urban, area, region, outdoor, outside, country	outdoor
10	alert, alarm, warning	warning
11	quality, condition	quality
12	breathe, breathing, breath	breathe
13	indoor, house, home, office, inside, school	indoor
14	bad, hard, severe, difficult, sore, suffer, poor, weak, risky	bad/poor
15	problem	problem
16	not bad*, good, great, well	good/well
17	today, now, tonight, morning, evening, night	time/now
18	seasonal, hourly, daily, weekly, monthly, yearly, hour, day, week, month, year	time/period
19	allergy, allergic, sensitive	allergies

20	food, feed, eat	food
21	pollen	pollen
22	hospital, doctor, medical, clinic	hospital
23	medication, medicine, pills	medication
24	high, huge, big, heavy	high
25	levels, standards, index	level
26	citizens, public, people, national, men, women, children, elder	people
27	cars, vehicles, motor, bike, bus, taxi	car
28	water, garbage, waste, soil	physical entity
29	smog, humidity, wind, dust, rain, weather	weather conditions
30	spring, summer, autumn, winter	season
31	pets, dogs, cats, birds	pets
32	particulates, particles, ozone	pollutants
33	horrible, hell, hate, crazy, killing, ugh	bad feelings
34	happy, funny, yeah	good feelings
35	wow, haha, lol	making fun

1 In the preprocessing phase, the case *wordA followed by wordB* was converted to “wordAwordB” (words
2 were collapsed without whitespace) in order to be distinguished from the cases where the same single
3 words exist in free order.

4

5 Given the specific bag-of-sets of words, we can examine the occurrence (one) or non-
6 occurrence (zero) of each unified concept in the text. A x^{th} tweet (T_x) can be represented
7 as a vector of the form in Equation (2):

$$8 \quad T_x = [Value_{S_1} \quad Value_{S_2} \quad Value_{S_3} \quad \dots \quad Value_{S_n}] \quad (2)$$

9 where S_k is the k^{th} set of words that includes v distinct words and the corresponding
10 value for this set will be: $Value_{S_k} = 0$, if none or $Value_{S_k} = 1$, if at least one of the v words
11 in S_k exists/occurs in the examined tweet.

12 The overall occurrence matrix M (i.e the matrix including all information on the
13 occurrence of the words in sets, as in Table 1) will be of the form:

$$M = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & 0 & 1 & 1 & \dots & 0 \end{bmatrix}$$

1 with m number of documents (tweets) per row and n number of sets of words per
2 column.

3 By following the aforementioned methodology, we first analyze each tweet in its
4 corresponding words and then we transform its textual representation into a binary
5 vector with zero/one values according to whether a word from a predefined set of words
6 (Table 1) belongs or not to the tweet. All tweets in the database are converted to array
7 vectors with zero/one values and can thus be considered as vectors in a
8 multidimensional space. Therefore, vectors that are close to each other in a topological
9 manner (according to distance metrics applied in the learning phase) can be regarded as
10 representatives of documents with similar word patterns. *Semantic similarity* can also be
11 assigned to documents based on vector operations, due to the fact that each set of
12 words has a general semantic interpretation. With the transformation of text into binary
13 data and the use of a similarity metric, we move from lower levels of co-occurrence of
14 words to upper levels of a realistic semantic interpretation of documents that have similar
15 patterns of occurrence of words and thus analogous semantic content and similarity.

17 **3.3 Learning, Mining and Evaluation Methods**

18 The way text documents are encoded is of great importance for the performance of any
19 learning and mining method. In the current paper, the key idea relies on a two-level
20 approach, based on SOMs and k-means clustering. At the first level we apply SOM for
21 its advantages of dimensionality reduction, data compression and topological relation
22 properties. At the second level we feed the new vector space (weights of SOMs) to k-
23 means clustering in order to produce clusters with similar characteristics. Clusters are

1 formed by learning at the same time the structure of data from SOM weights and its
2 segmentation using both distance and density information (k-means and Euclidean
3 distance). The number of clusters formed is based on the Davies-Bouldin index, while
4 additional methods, like Sammon mapping and silhouette metric, are taken into account
5 to evaluate and validate the efficiency of the clustering results. The methodological
6 background is described in the following subsections.

7

8 **3.3.1 Self-Organizing Map**

9 Kohonen's Self-Organizing Map (SOM) is a competitive, unsupervised learning method
10 that maps high dimensional data into a low dimensional space by preserving their spatial
11 correlation (Kohonen, 1990). SOMs consist of neurons, each of which is a set of weights
12 (or prototype vectors) that correspond to the projection of the instances of a high-
13 dimensional dataset onto a usually 2-dimensional grid. The neuron whose weight vector
14 is most similar to the input is called the best matching unit (BMU) and this can be
15 determined based on criteria like the Euclidean distance.

16 Since its first introduction, a lot of research has been conducted concerning SOMs and
17 many different implementations have been proposed, for various application domains.

18 Our research has been inspired by Ritter and Kohonen (1989) as well as Honkela
19 (1997). The authors introduce a methodological background for SOMs that use word
20 categories (so called self-organizing semantic maps) in order to preserve semantic
21 relationships reflected in the data according to their relative distances in the map.

22 The value of SOMs lies in easily interpretable results through the investigation of
23 topologically similar areas so as to discover relations between examined parameters. A
24 matrix of occurrence or co-occurrence of words in documents can be visualized by
25 SOMs as a topological grid that represents every correlation between sets of words in
26 the same topological space. We take advantage of SOMs' potentials for dimensionality

1 reduction and easily interpretable results, in order to form clusters of similar patterns of
2 words' occurrence and extract topic categories from the underlying information.

3

4 **3.3.2 K-means clustering**

5 K-means is a well known non-hierarchical clustering technique which partitions a set of
6 observations (data points) into k numbers of clusters (Hartigan and Wong, 1979). In the
7 initial step of the algorithm, k "means" are randomly defined as centroids of clusters
8 within the data domain. Then, an iterative refinement process is performed, where
9 clusters are formed by associating each data point with the closest centroid (nearest
10 mean based on Euclidean Distance) and centroids are updated by calculating the new
11 mean values of data points belonging to each cluster. The iterative process continues
12 until a convergence criterion (minimal decrease of total sum of squared errors, minimal
13 reassignment of data points into different clusters, maximum number of iterations, etc.) is
14 met.

15 The efficiency of k-means is highly depended on the determination of a suitable number
16 (k) of clusters. This number should be somehow close to the number of "natural" clusters
17 that are present in the data, otherwise:

- 18 a) if k is higher than needed, we might get "garbage cans" (clusters with no
19 significant distinction between them) or empty clusters, due to the force of
20 splitting, or
- 21 b) if k is lower, we will get more generic divisions.

22 Since there is no prior knowledge of what number to select, we adopt an internal
23 evaluation criterion of clustering results: we make use of the Davies-Bouldin (DB) index
24 for different numbers of clusters, which is calculated according to Equation (3) (Davies
25 and Bouldin, 1979):

$$DB = \frac{1}{k} \sum_{i=1}^k \max_{j,j \neq i} \frac{S_i + S_j}{d_{ij}} \quad (3)$$

where k is the number of clusters. The within (S_i) and between (d_{ij}) cluster distances are calculated using the cluster centroids as follows:

$$S_i = \frac{1}{|C_i|} \sum_{x \in C_i} \|x - m_i\| \quad (4)$$

$$d_{ij} = \|m_i - m_j\| \quad (5)$$

where m_i is the centroid of cluster C_i , with $|C_i|$ being the number of data points that belong to the cluster. According to Equation (3), it becomes evident that the objective of DB-index is to minimize its value and, consequently, to have the minimum “within-cluster” dispersion and the maximum “between-clusters” separation (Davies and Bouldin, 1979). Here, k-means and DB-index is combined with SOM as a two-step approach for clustering, in order to define distinct areas of compact clusters in map representation.

3.3.3 Sammon mapping

Whereas the visualization of high dimensional data with the use of SOMs can reveal topological relations among data points and clusters, it doesn't give any interpretation of their actual closeness and inner distance between the neurons of the map. Usually, additional coloring schemes, like the unified distance matrix (U-matrix) (Ultsch and Siemon, 1990), can be used in the analysis to show the boundaries between the formed, and depict low or high distance estimation between them.

In order for a map to depict the data structure, the distance property must be retained along with the topology. Sammon mapping (Sammon, 1969), is a non-linear methodology that can preserve the topology of the input data as well as the distances between inter-points within the map. This projection of data can reveal cluster structure and tendency, showing how “strong” or “loose” the connections are between data points in each cluster and between clusters.

1 Sammon mapping is an iterative process with random initialization of the d-space
 2 configuration that performs pseudo-Newton minimization of the error function given in
 3 Equation (6), also known as Sammon's stress, by calculating the inter-point distances d_{ij}^*
 4 between the i^{th} and the j^{th} object in the original space, and the distance d_{ij} between their
 5 projections (Sammon, 1969):

$$6 \quad E = \frac{1}{\sum_{i < j} d_{ij}^*} \sum_{i < j}^N \frac{(d_{ij}^* - d_{ij})^2}{d_{ij}^*} \quad (6)$$

7 We apply Sammon mapping on the actual clusters derived from SOM and k-means, by
 8 projecting the weights of neurons of the SOMs and visualizing their Euclidean distance in
 9 the 2-dimensional space.

10

11 3.3.4 Silhouette metric

12 Silhouette values can act as a metric of validity of the already assigned clusters, with
 13 respect to factors like: (a) cohesion (how compact each cluster is), and (b) separation
 14 (how clearly clusters are separated from each other). A silhouette value $s(i)$ for an object
 15 i is calculated according to the Equation (7):

$$16 \quad s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}, \quad -1 \leq s(i) \leq 1 \quad (7)$$

17 where $a(i)$ is the average dissimilarity of object i to all the other objects within the cluster
 18 in which i falls, and $b(i)$ is the minimum average dissimilarity of object i to all objects of
 19 each assigned cluster, thus defining the neighboring cluster of i (Rousseeuw, 1987).

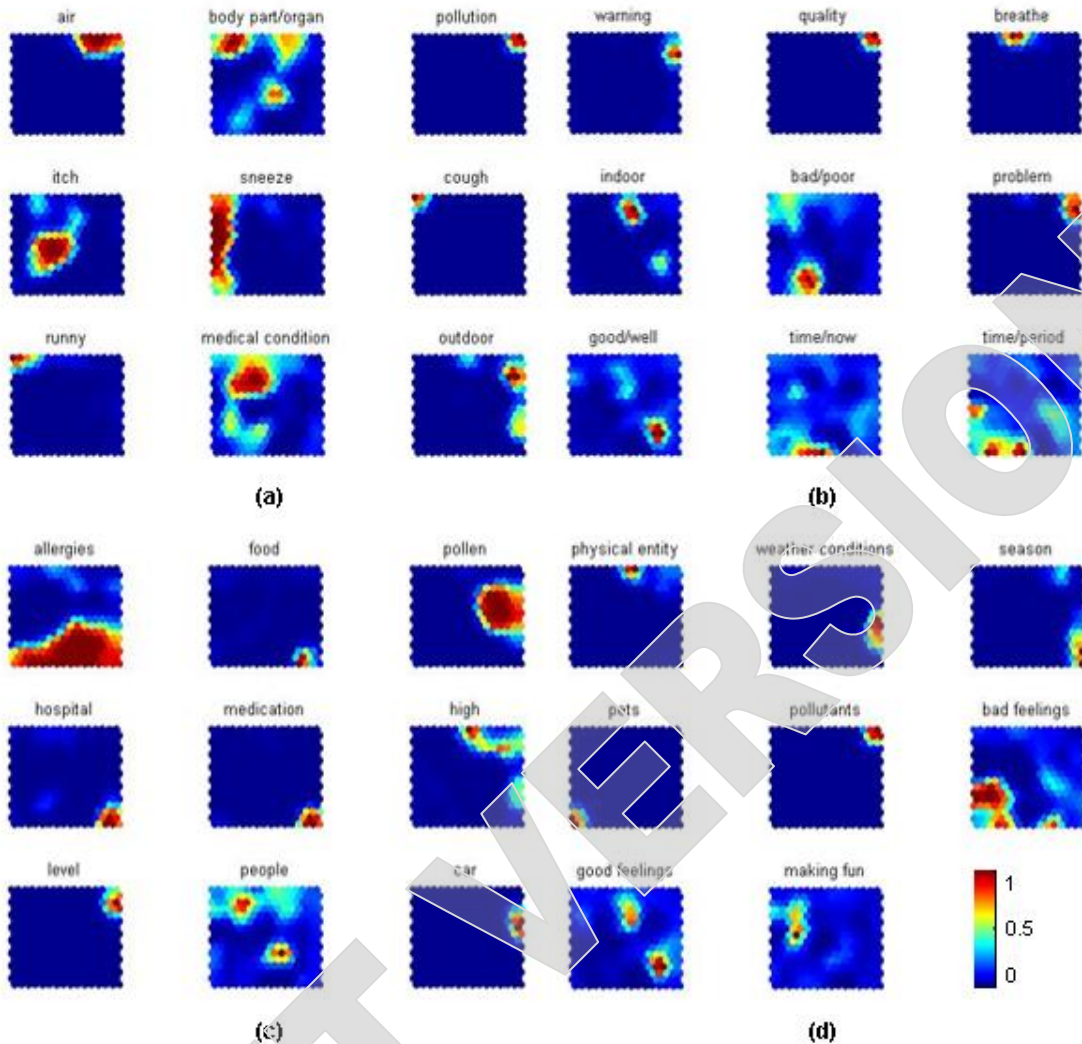
20 Silhouette values can be used as an indicator of how well objects have been classified in
 21 a cluster, or which objects are misclassified. By plotting the silhouette values of objects
 22 per cluster, we can distinguish compact and clearly separated clusters from non-tight
 23 ones. If a cluster contains many objects with low silhouette values, or even negative
 24 ones, this indicates that the cluster is not well-separated from other clusters. On the
 25 contrary, the more close to 1 the mean of silhouette values of all objects in a cluster, the

1 better separation from its neighboring clusters has been achieved. Similarly, we can
2 follow the mean silhouette value per cluster and conclude to corresponding results.

3

4 **4. Results**

5 In the preprocessing phase, tweets are encoded as a high dimensional matrix of
6 44888x35 zero/one values in relation to the non-occurrence/occurrence of 35 different
7 sets of words (Table 1) in posts. We apply the SOM algorithm by using the SOM Toolbox
8 2.0 (<http://www.cis.hut.fi/somtoolbox/>) for Matlab. SOM neurons were positioned in a
9 16x16 orthogonal grid arranged on a hexagonal lattice of a 2-dimensional map. With
10 SOM representation, we manage to reduce the dimensionality of the initial encoded data
11 up to 99.98%. The overall results of the mapping process are shown in Figure 4.



1
2 **Figure 4: Four sets (a, b, c, d) of SOMs for 35 sets of words in total, derived from the analysis of**
3 **feature vector values of tweets with the use of SOM algorithm.**

4
5 Each map has 256 (16x16) colored neurons (cells) that correspond to the weight values,
6 calculated by SOM, for the corresponding examined variable. According to the given
7 scale in Figure 4, dark blue color corresponds to the lowest value (zero) of the variable,
8 meaning the non-occurrence of a set of words in text. On the other hand, dark red color
9 corresponds to the highest value (one) of the variable of interest, meaning the
10 occurrence of a set of words in text. Other colors apart from the aforementioned two
11 correspond to intermediate weight values between the range [0,1]. The interpretation
12 behind intermediate values is that such neurons do not clearly belong to any of the two

1 defined cases that zero/one values represent. Instead, the closer the values are to
2 zero/one, the most certainly they belong to the corresponding cases.

3 The visual inspection of SOMs in Figure 4 and the topology characteristics of each map
4 can reveal qualitative information about the use of each variable and the inter-correlation
5 between them. The latter can be examined by the co-occurrence (values equal to one) of
6 two or more sets-of-words in posts. Thus, similar areas per variable should be
7 investigated per case. For example, we can imply that *allergies* are related to *pets* within
8 the tweets being analyzed. This can be supported by inspecting the bottom-left corner of
9 the corresponding maps in Figure 4. This area for both variables, *allergies* and *pets*, has
10 red colored cells, meaning that there are a number of tweets that use words from both
11 aforementioned sets in the same sentence. On the other hand, we can spot negative
12 relation, like for example between *pets* and *pollutants*. According to Figure 4, there are
13 no data with high values in both variables' visualizations. As expected, concepts like *pets*
14 and *pollutants* are not related to each other.

15 Through this process, we can state similar possible correlations (positive/negative)
16 among words and unified concepts, by investigating topologically related areas of SOM
17 for different sets of words. Though, it is hard to form manually groups of data (clusters)
18 with similar behavior from these distinct maps. The U-matrix in Figure 7a can give a
19 rough representation of the "natural" clusters that exist in the dataset. It represents the
20 Euclidean distance between the high dimensional data and the SOM weight vectors. The
21 colored scale next to it corresponds to the distance values. Cells with blue color in the U-
22 matrix can be considered as clusters (minimum distance between neurons), while cells
23 with colors related to higher values of the scale show the boundaries between clusters
24 (maximum distance between them). According to the formed U-matrix we can stand out
25 approximately 5-7 clearly defined clusters but there are also some additional clusters
26 characterized by fuzzy districts.

1 As we do not have clear knowledge on the number of clusters within the data, it is
2 important to experiment with a range of values for k . In order to select the number of
3 clusters that sufficiently and discretely divide the area of SOMs with the use of k-means,
4 we made use of the DB-index within a big range of k values (from 1 to 50), as presented
5 in Figure 5. It is acknowledged that the lower the DB-index value the better the clustering
6 configuration, in terms of compactness and separation (Kovács et al., 2006). Here, the
7 global minimum of DB-index is for $k = 22$ but there are three local minima that appear
8 before it, for $k = 3, 10$ and 16 , which are of special interest. The DB-index is not highly
9 varying from $k = 10$ to $k = 27$, while within the range that these two values define, the
10 global minimum occurs. After $k = 27$, sudden changes of DB-index appear, which can be
11 used as a criterion that the situation from this point is unstable. The clustering algorithm
12 is forced to do splitting of already well-defined clusters, thus producing “artificial” clusters
13 with loose connections between data points. Thus, the global minimum can be
14 considered as the upper limit of k that provide well defined clusters in the examined
15 case. In the following subsections, we summarize clustering results for the
16 aforementioned local minima and the global one, and we describe the evaluation
17 process.
18

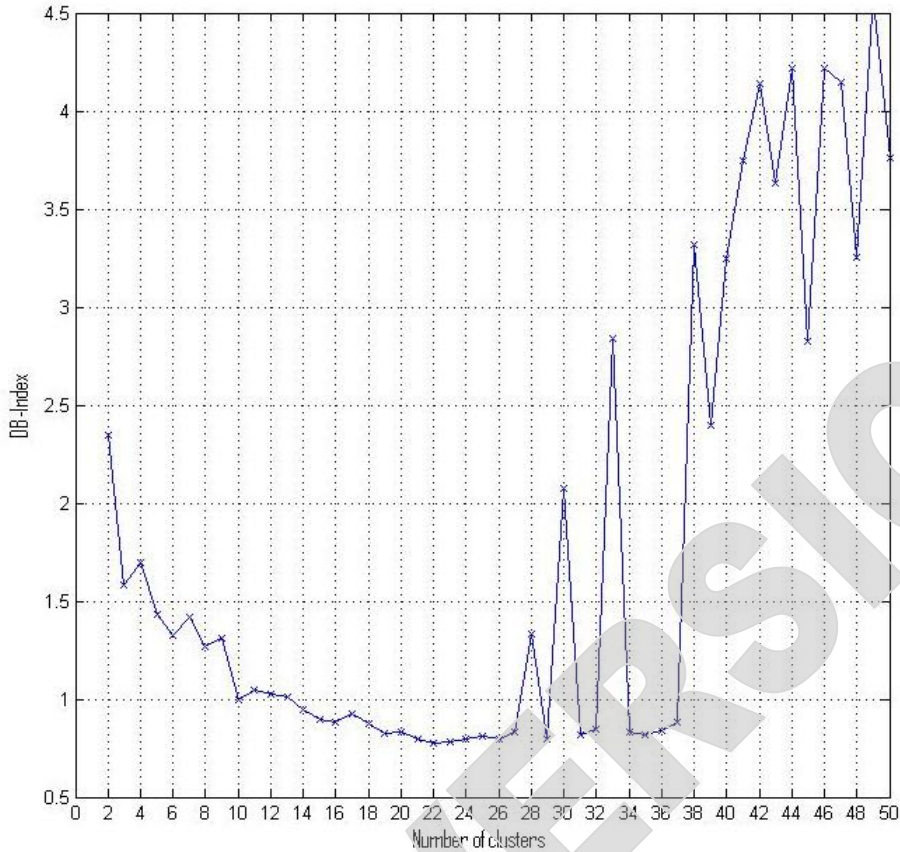


Figure 5: DB-index values for different numbers of k clusters.

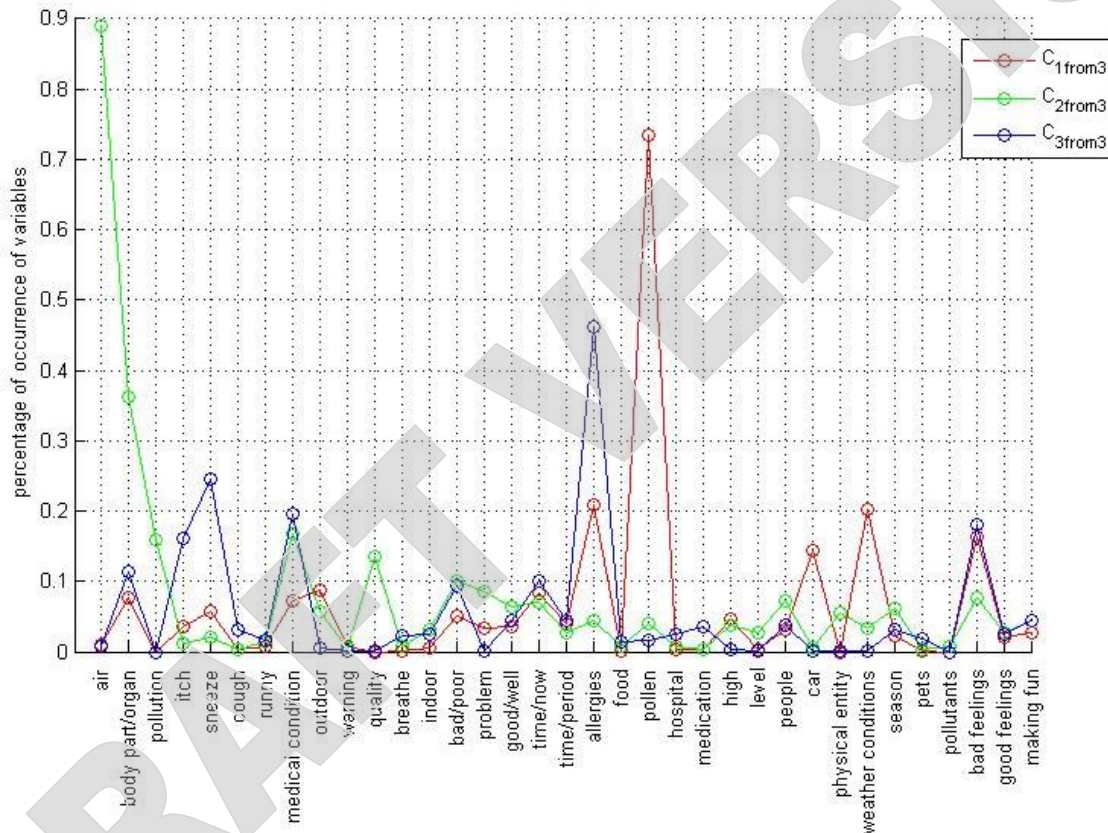
4.1 Clustering results for different numbers of clusters

By applying k-means clustering to the weight vectors of the SOM we can form k number of clusters, based on distance metrics. Each neuron of the map will correspond to a specific cluster. Since neurons are linked to high dimensional data, we can easily define the actual data behind the clustering results.

For each $C_{i\text{from}k}$, meaning the i^{th} cluster of k in total, we can identify the instances (tweets) it includes. We can thus calculate the total number of instances per cluster, as well as the percentage value of occurrence of each set of words per cluster. The latter values are depicted in Figure 6.

According to Figure 6, each cluster has a peak value for a specific set of words. For example, almost 90% of tweets that are included in cluster $C_{2\text{from}3}$ include in their text

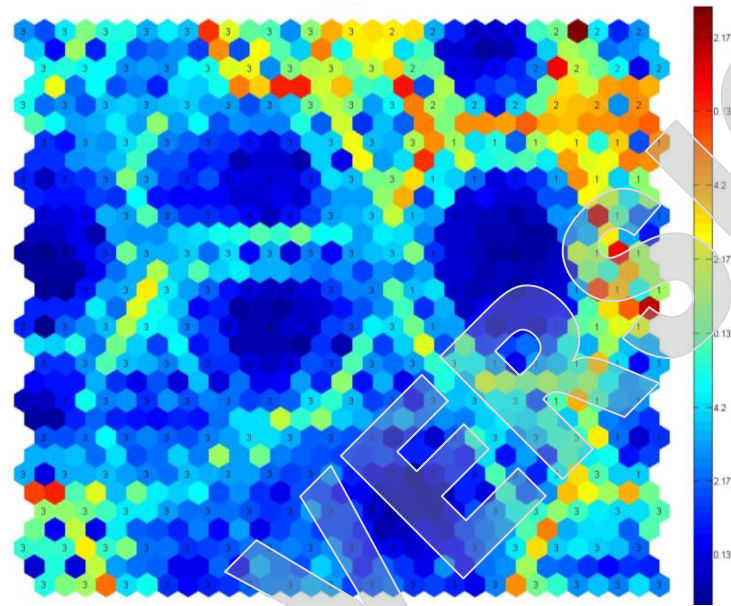
1 one or more words from the set of words that corresponds to the unified concept *air*.
 2 Clusters C_{1from3} and C_{3from3} include tweets with words related to *pollen* and *allergies*
 3 respectively. Their percentage of occurrence of corresponding words is lower than the
 4 one calculated for cluster C_{2from3} , but still of significant level (73% and 46%
 5 approximately). The name of the dominant concept can clearly define the label of the
 6 formed clusters of tweets, and can thus provide with the general topic of discussion per
 7 case. In other words, through this clustering process, we map text data into concepts
 8 and their semantics.



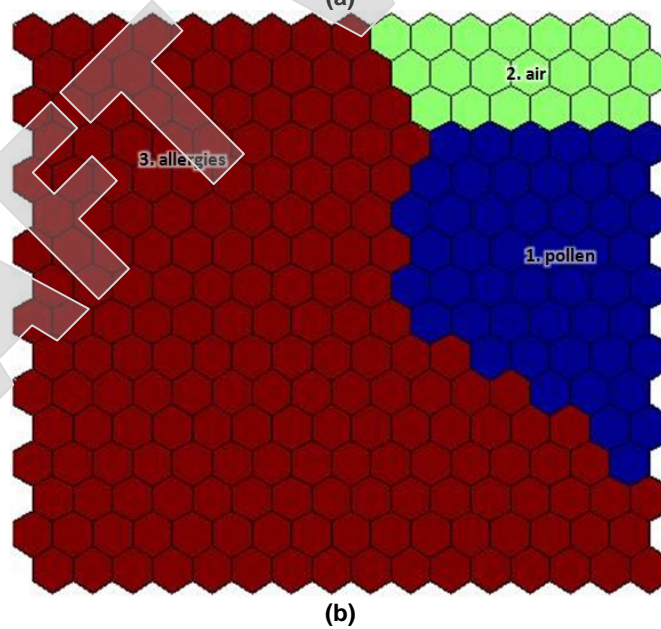
9
 10 **Figure 6: Percentage of occurrence of each concept (set of words) for 3 formed clusters.**

11 By using the U-Matrix along with the results from k-means clustering, we can draw the
 12 identification number of each cluster per node in the map, as shown in Figure 7a. An
 13 integrated area per cluster is formed and we can get the visualization of clusters' content
 14 and coverage by adding the corresponding labels of description on the map (Figure 7b).

1 The cluster named *allergies* covers the biggest area in the map, as it includes 80.35% of
 2 tweets, while *pollen* is the second biggest cluster with 14.19% of tweets.



3
4

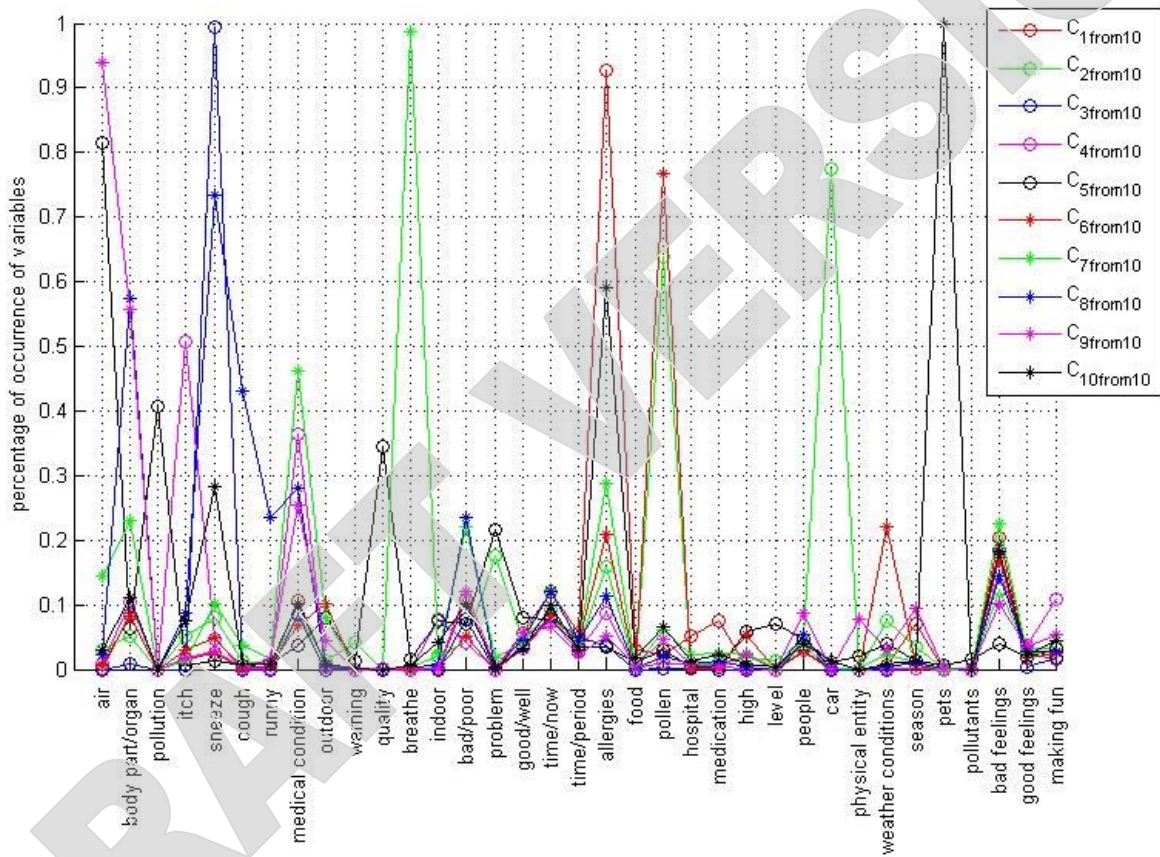


5
6
7
8

Figure 7: (a) U-Matrix and clusters per neuron, (b) areas and labels of 3 formed clusters on 2-dimensional map.

9 Up to this point, the three distinct concepts named *allergies*, *pollen* and *air* describe at a
 10 general level the content of the database. We thus gained a conceptual level from
 11 unstructured data, with semantic interpretation as well. By dividing the areas into more

1 clusters, as the DB-index suggests, we can obtain a more detailed overview of the
 2 concepts that are discussed in tweets.
 3 When applying k-means clustering of SOMs for $k = 10$, different clusters with different
 4 content are formed, which is depicted in form of percentage values of occurrences of
 5 sets of words per cluster, as shown in Figure 8. Concepts per cluster are also
 6 summarized in Table 2.



7
 8 **Figure 8: Percentage of occurrence of each concept (set of words) for 10 formed clusters.**

1 Table 2: Defining concepts that describe 10 clusters, based on the percentage of occurrence of sets
 2 of words in each cluster

Cluster ID	Instances (%)	Set of words based on their occurrence in clusters (% values)			
		(90, 100]	(70, 90]	(50, 70]	(30, 50]
$C_{1from10}$	35.34%	allergies	-	-	-
$C_{2from10}$	2.63%	-	car	pollen	-
$C_{3from10}$	13.36%	sneeze	-	-	-
$C_{4from10}$	23.20%	-	-	itch	medical condition
$C_{5from10}$	2.15%	-	air	-	pollution
$C_{6from10}$	11.38%	-	pollen	-	-
$C_{7from10}$	1.87%	breathe	-	-	medical condition
$C_{8from10}$	5.33%	-	sneeze	body part/organ	cough
$C_{9from10}$	3.31%	air	-	body part/organ	-
$C_{10from10}$	1.44%	pets	-	allergies	-

3 Following a similar way in the interpretation of results, we can conclude that there are
 4 clusters with more than one dominant concept that cover at least 50% of the
 5 corresponding tweets. In other words, clusters' content can be defined by a combined
 6 label of both dominant concepts, where applicable. With the use of the combined label
 7 we get a more detailed semantic interpretation of the content and the tweets included in
 8 each cluster.

9 Through this refinement, from 3 to 10 numbers of clusters, cluster named *allergies*
 10 ($C_{1from10}$) remains the biggest one, including approximately 35% of total tweets, while it
 11 was reduced compared to the results obtained for $k = 3$. From its division, new clusters
 12 arose, namely *sneeze* ($C_{3from10}$), *itch & medical condition* ($C_{4from10}$), *breathe & medical*
 13 *condition* ($C_{7from10}$), *sneeze & body part/organ* ($C_{8from10}$), and *pets & allergies* ($C_{10from10}$),
 14 forming clusters with more detailed description of the content. Similarly, cluster *air* from
 15 the previous clustering results was split into two sub-clusters: *air & pollution* ($C_{5from10}$) and
 16 *air & body part/organ* ($C_{9from10}$). Finally, from cluster *pollen* a new sub-cluster was
 17 created, named *car & pollen* ($C_{2from10}$). Labels and areas of clusters in the topological
 18 space are shown in Figure 9.



1

2

Figure 9: Areas and labels of 10 formed clusters on 2-dimensional map.

3

From the so far analysis, it becomes evident that by increasing the number of clusters,

4

new sub-clusters arise with lower number of instances. In addition, labels derived give

5

more detailed description of clusters' content. Based on the dominant sets of words

6

presented in groups of tweets, we arrive to more detailed results for 16 and 22 numbers

7

of clusters, as shown in Figures 10 and 11 respectively.



1
2
3

Figure 10: Areas and labels of 16 formed clusters on 2-dimensional SOM.

DRAFT VERSION



Figure 11: Areas and labels of 22 formed clusters on 2-dimensional SOM.

1

2

3

4.2 Evaluation of clustering results

As a next step, we quantitatively evaluate the clustering results, by using the silhouette metric. We can investigate how strong or loose is the relationship between the instances that form each cluster, by following the mean as well as the minimum silhouette value per cluster. The optimum value of the metric equals to 1. All calculated values for cases with 3 and 10 numbers of clusters are summarized in Table 3.

10

1 Table 3: Mean and minimum silhouette value per cluster, for different values of k clusters

(a) SOM and k-means for $k = 3$			
Cluster ID	Name/ concept	Mean silhouette value	Min silhouette value
C_{1from3}	pollen	0.3680	0.0349
C_{2from3}	air	0.4563	0.1154
C_{3from3}	allergies	0.3559	-0.0085
(b) SOM and k-means for $k = 10$			
Cluster ID	Name/ concept	Mean silhouette value	Min silhouette value
$C_{1from10}$	allergies	0.3384	-0.2542
$C_{2from10}$	car & pollen	0.4403	-0.1371
$C_{3from10}$	sneeze	0.6402	0.2009
$C_{4from10}$	itch & medical condition	0.3143	-0.1522
$C_{5from10}$	air & pollution	0.7499	0.5789
$C_{6from10}$	pollen	0.5532	0.1742
$C_{7from10}$	breathe & medical condition	0.5922	-0.0280
$C_{8from10}$	sneeze & body part/organ	0.4144	0.1427
$C_{9from10}$	air & body part/organ	0.7803	0.4646
$C_{10from10}$	pets & allergies	0.5941	-0.1097

2 In the first clustering analysis (where $k = 3$), all three clusters have low mean silhouette
3 values (Table 3a). Also, their minimum silhouette value is close to zero. These results
4 denote high distances between neurons within clusters and the connections between
5 instances of clusters tend to be loose. Such clusters are unstable and more likely to split
6 in a next clustering process, with increased number of clusters.

7 When we increase the number of clusters ($k = 10$), the situation is different (Table 3b):
8 we still get some weakly connected clusters, like $C_{4from10}$ (*itch & medical condition*) which
9 has the lowest mean silhouette value (0.3143) and it is more likely to be split in a next
10 division. On the other hand, there are clusters, like $C_{5from10}$ (*air & pollution*) or $C_{9from10}$ (*air
11 & body part/organ*) that have high mean silhouette value (closer to 1) and thus these
12 clusters are considered to be tight and their instances are strongly connected. We prove
13 that such clusters remain stable in the next clustering process.

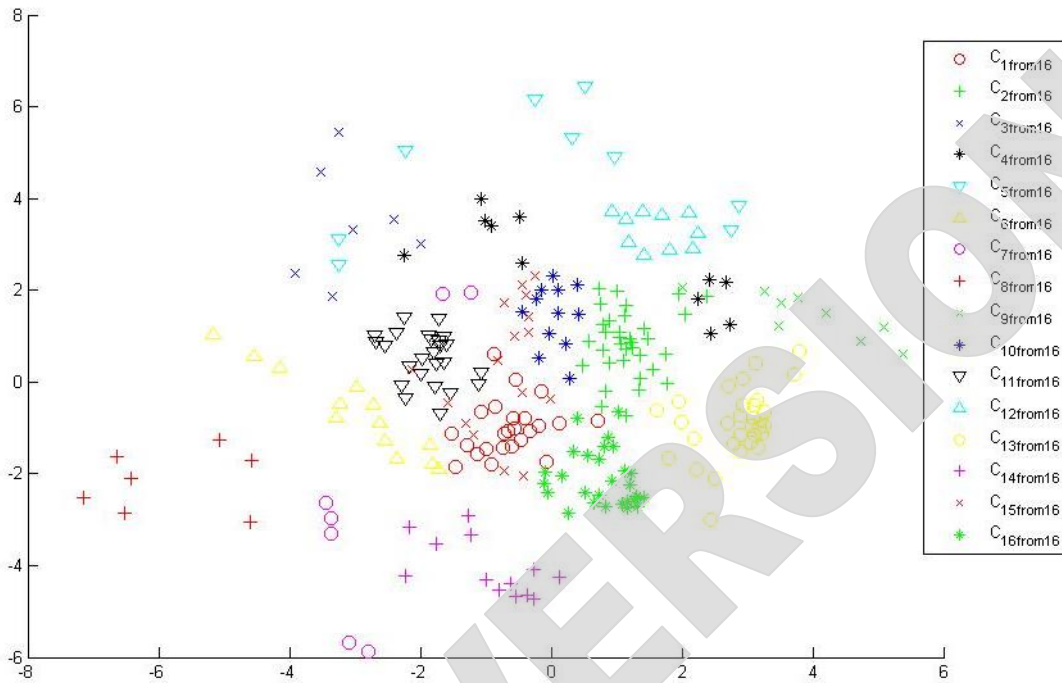
1 By increasing once more the number of clusters from 10 to 16, we get new clustering
2 results. Clusters with low mean silhouette value in the previous step were split or
3 transformed. For example, cluster *itch & medical condition* ($C_{4\text{from}10}$) became two
4 separate clusters with one dominant concept per each: *itch* ($C_{16\text{from}16}$) and *medical*
5 *condition* ($C_{1\text{from}16}$). On the contrary, clusters with high mean silhouette value, like *air &*
6 *pollution* ($C_{5\text{from}10}$) and *air & body part/organ* ($C_{9\text{from}10}$) remain intact (now as $C_{8\text{from}16}$ and
7 $C_{14\text{from}16}$ correspondingly) from the current clustering process, including almost the same
8 number of instances in both divisions.

9 By plotting clusters in two dimensions using Sammon mapping, we obtain the distribution
10 of data points per cluster and we can visualize the distance properties among them. We
11 demonstrate Sammon mapping for the case where $k = 16$ clusters (results are presented
12 in Figure 12). Following this visualization, we can distinguish three main categories of
13 clusters:

- 14 a) clusters that are clearly separated, like $C_{8\text{from}16}$ (*air & pollution*), $C_{13\text{from}16}$ (*pollen*)
15 and $C_{14\text{from}16}$ (*air & body part/organ*),
- 16 b) clusters that are overlapping, like $C_{15\text{from}16}$ (*bad feelings & medical condition &*
17 *allergies*), and
- 18 c) clusters whose data points are spread over the map, even though they belong to
19 the same cluster, like $C_{2\text{from}16}$ (*allergies*) and $C_{4\text{from}16}$ (*indoor & allergies*).

20 Clusters of categories (b) and (c) are more likely to be split in the next clustering phase.

21 This fact is also confirmed by the corresponding low silhouette values of these clusters.



1

2

Figure 12: Sammon mapping of SOM weights of 16 clusters in 2-dimensional space.

3

Similarly, by incrementing the number k for k -means from 16 to 22, those clusters that

4

were more likely to split, generated new, smaller and more detailed ones in the next

5

level. Overall, the Silhouette Coefficient (SC), calculated as the mean of all silhouette

6

values per clustering level, increases (see Table 4), meaning that we achieve to get

7

more stable clusters in each clustering step.

8

Table 4: Silhouette coefficient values per different total number of clusters

Total number of clusters	Silhouette coefficient
3	0.3661
10	0.4512
16	0.5448
22	0.5669

9

1

2 **5. Discussion**

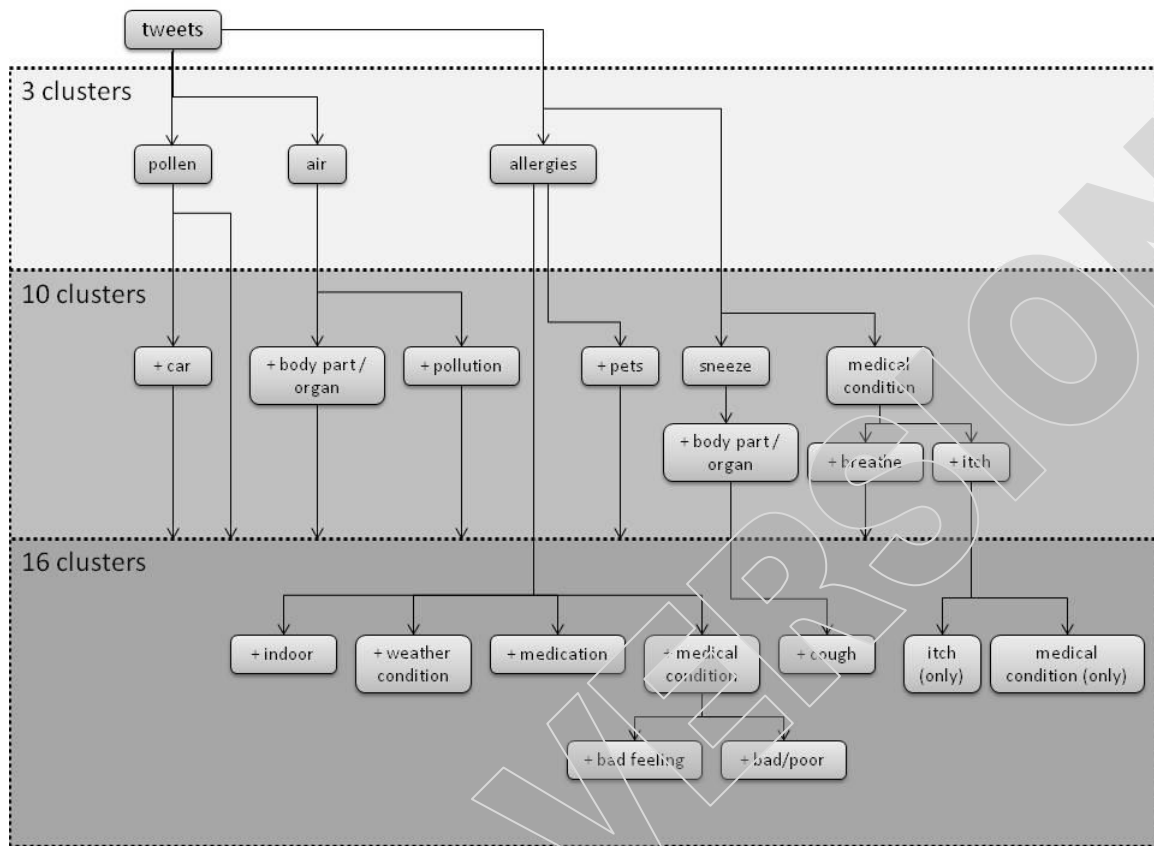
3 Our results demonstrate that SOM is an efficient and convenient method to structure,
4 analyze and interpret user generated text from social media. The representation of text
5 data with vectors, based on the occurrence of specified sets-of-words in text, can
6 transform original text into a binary matrix. Sets-of-words are assigned with a unified
7 concept that maps them into a general semantic meaning. Clusters may appear with one
8 or more dominant concepts that describe their content. We assign topics of discussion in
9 groups of tweets and we enrich formed clusters with conceptual semantics, based on
10 their dominant concepts.

11 Through the performed analysis, we derive to a clear visual representation of the actual
12 content of the text database (see Figure 11). Clusters' labeling corresponds to the real
13 content of the text data included and to the semantic knowledge implicitly assigned to
14 each concept. The size of each cluster reflects the number of instances included on it.
15 We identified big clusters, like *pollen*, *allergies*, *itch* and *sneeze*, where general concepts
16 uniquely describe their content. We also identified smaller, more descriptive clusters,
17 such as *pets & allergies*, *car & pollen*, *air & pollution* and *body part/organ & sneeze* that
18 combine more than one concept in order to characterize their content. Such clusters
19 acquire labels that state clearly and in more detail their content (topic of discussion).
20 The potential of determining the relationships between words as well as correlations
21 among concepts is a highly-valued prospect of this method. Not only logically apparent
22 concepts, but also new, non-typical relations can be discovered. A good example of
23 unexpected interrelation learned is the cluster *car & pollen*. Its name defines two
24 dominant concepts that describe the content of the cluster. However, such a correlation
25 does not seem to make "sense". A thorough investigation of the actual data reveals that
26 pollen is related to cars in a functional way, i.e. people are expressing their

1 disappointment or frustration because their vehicle is covered with pollen during the
2 pollen season, and requires cleaning. We should also note here that both concepts were
3 introduced into the bag-of-sets-of-words due to the increased frequency of use in the
4 collection of documents, even though *car* seems to be unrelated to the context of our
5 research.

6 Our results also highlight the ability of SOM and k-means to organize practically related
7 topics of discussion into neighboring clusters. As an example, the main cluster *allergies*
8 adjoins with concepts of similar content, like *allergies & body part/organ*, *allergies &*
9 *bad/poor & medical condition*, *food & allergies*, while *air* related concepts are in a
10 different area of the map.

11 The number k of clusters that divide adequately and efficiently the actual data is based
12 on the DB-index value. By investigating different numbers for k and not only the one that
13 corresponds to the global minimum of DB-index, we derive different levels of description
14 of clusters' content. It is characteristic that the lower the number of clusters the more
15 generic the resulting concepts. With increasing number of clusters, these generic
16 clusters are subdivided into more concrete ones. Hence, given the content of a
17 database, we can build a *hierarchy of concepts*, as shown in Figure 13. The derived
18 hierarchy can be considered as a tree diagram, where upper nodes are the generic
19 description concepts and leaf nodes are the detailed concepts. Concepts of low levels
20 describe the content by combining concept nodes from upper levels of the tree. Here,
21 the hierarchy is given only for the three out of four levels of the clustering process that
22 we performed, due to clarity reasons.



1
2 **Figure 13: The hierarchy of concepts derived after three levels of clustering.**

3 With Sammon mapping we can visualize the relationship between individual neurons
4 and defined clusters, in terms of their relative distance. Clusters with low distances
5 between their individual points are more tight and concentrated in an area of 2-
6 dimensional space, while disperse clusters with high distances among their instances
7 are more likely to be split in an expanded clustering process. Results through the
8 evaluation phase confirm the aforementioned statement.

9 The efficiency of the methodology is evaluated by using the silhouette metric, which
10 gives quantitative information for how strong or loose the connections between the
11 neurons within a cluster are. We showed that, by performing the optimum clustering as
12 the DB-index denotes, clusters derived are more stable, having at the same time higher
13 mean silhouette values than the ones of less extensive clustering processes.

1 On the other hand, the qualitative efficiency of the methodology can be evaluated
2 subjectively by: (a) the clustering results derived, (b) the accuracy of the actual data
3 selected per cluster, and (c) the correlation between the content and the semantic
4 annotation. Based on our extended investigation, all three factors are achieved in the
5 current analysis.

6 The results of the proposed analysis are of interest for several reasons. Clustering
7 results can be used to perform efficient classification of newly retrieved tweets into
8 different description groups of content. They can also be used as an awareness or
9 decision factor. The semantic interpretation of clusters can play a key role in different
10 applications, such as recommendation systems and event detection services. By
11 exploiting the massive content from social media and mining their text streams,
12 practitioners can monitor events like dispersion of allergies or bad quality conditions,
13 even in areas where there is no other means of monitoring air quality.

14

15 **6. Conclusions**

16 We focused on the task of social media analysis and we investigated the use of Twitter
17 as a source of user-generated information concerning the atmospheric environment
18 domain and its effect to humans' quality of life. We proposed a straightforward
19 methodology for encoding and analyzing tweets based on the occurrence or non-
20 occurrence of selected words in text. The ultimate goal that we achieved was to cluster
21 text into groups with similar content, giving also a semantic prospect to the interpretation
22 of the results. We demonstrated the use of SOM and k-means, along with the DB-index
23 to perform the clustering process efficiently. We were able to extract well-defined
24 clusters and derive topics of discussion from the massive collection of posts.
25 Relationships among concepts that describe the content can further form hierarchical
26 structures. Both interpretation of visualizations and evaluation of clustering results can

1 underscore the efficiency and the consistency of the methodology. By structuring the
2 knowledge derived in an easily interpretable way, we manage to provide a rich content
3 for individuals and practitioners for quality of life issues.
4 In future work, we aim to overcome the limitation of language used in text collection and
5 add spatiotemporal dimensions to the related extracted concepts. Profiling concepts over
6 a longer time interval, i.e. within a calendar year, or analyzing the content near real-time
7 will create a reference of event detection or supplementary serve as an
8 observation/alarm service (in our example of air quality issues), even in areas where
9 monitoring data for air quality are not available. The hierarchy of concepts as well as
10 additional semantic description can be used to form a data-driven ontology for the air
11 quality domain. Finally, we are interested to research the applicability of the methodology
12 in other domains of interest.

13

14 **Acknowledgements**

15 We would like to thank the Centre for International Mobility CIMO (<http://www.cimo.fi/>) for
16 the CIMO Fellowship funding which enabled the work described in this article.

1 **References**

- 2 Aiello, L., Petkos, G., Martin, C., Corney, D., Papadopoulos, S., Skraba, R., Goker, A.,
3 Kompatsiaris, I. and Jaimes, A. (2013). Sensing trending topics in Twitter. IEEE
4 Transactions on Multimedia. 15 (6), 1268-1282. doi: 10.1109/TMM.2013.2265080
- 5 Berry, M.W. (ed.) (2004). Survey of Text Mining: Clustering, Classification and Retrieval.
6 Springer Science + Business Media, Inc.
- 7 Burke, J., Estrin, D., Hansen, M., Parker, A., Ramanathan, N., Reddy, S. and Srivastava,
8 M.B. (2006). Participatory Sensing, Proc. of the Workshop on World-Sensor-Web
9 (WSW'06): Mobile Device Centric Sensor Networks and Applications, Colorado, USA,
10 117-134.
- 11 Cairns, C. (2013). Air Pollution, Social Media, and Responsive Authoritarianism in China,
12 UCLA COMPASS Conference in Comparative Politics.
- 13 Costa, J., Silva, C., Antunes, M. and Ribeiro, B. (2013). Defining Semantic Meta-
14 hashtags for Twitter Classification. LNCS Adaptive and Natural Computing Algorithms.
15 7824, 226-235. doi: 10.1007/978-3-642-37213-1_24.
- 16 Crooks, A., Croitoru, A., Stefanidis, A. and Radzikowski, J. (2013). #Earthquake: Twitter
17 as a Distributed Sensor System. Transactions in GIS. 17 (1). 124-147. doi:
18 10.1111/j.1467-9671.2012.01359.x.
- 19 Davies, D. and Bouldin, D.W. (1979). A Cluster Separation Measure. IEEE Transactions
20 on Pattern Analysis and Machine Intelligence. PAMI-1 (2), 224-227. doi:
21 10.1109/TPAMI.1979.4766909.
- 22 Dredze, M. and Paul M.J. (2014). Natural Language Processing for Health and Social
23 Media, 64-67.
- 24 Hall, D.L., Llinas, J., McNeese, M. and Mullen, T. (2008). A Framework for Dynamic
25 Hard/Soft Fusion, Proc. of the Eleventh International Conference on Information Fusion,
26 Cologne, Germany.

- 1 Hartigan, J.A. and Wong, M.A. (1979). Algorithm AS 136: A K-Means Clustering
2 Algorithm. *Applied Statistics*. 28(1), 100-108. doi: 10.2307/2346830.
- 3 Hofmann, T. (2001). *Unsupervised Learning by Probabilistic Latent Semantic Analysis*.
4 *Machine Learning*. 42, 177-196. doi: 10.1023/A:1007617005950.
- 5 Honkela, T. (1997). *Self-Organizing Maps of words for natural language processing*
6 *applications*. Ph.D. Dissertation, Neural Networks Research Center, Helsinki University
7 of Technology, Espoo, Finland.
- 8 Kaplan, M. A. and Haenlein, M. (2011). The early bird catches the news: Nine things you
9 should know about micro-blogging. *Business Horizons*. 54(2), 105-113. doi:
10 10.1016/j.bushor.2010.09.004.
- 11 Karatzas, K. (2009). Informing the public about atmospheric quality: air pollution and
12 pollen. *Allergo Journal*. 18 (3/09), 212-217.
- 13 Kaufmann, M. and Kalita, J. (2010). Syntactic normalization of twitter messages, *Proc. of*
14 *the International Conference on Natural Language Processing*, Kharagpur, India.
- 15 Kohonen, T. (1990). The self-organizing map, *Proc. of the IEEE*, 78 (9), 1464-1480. doi:
16 10.1109/5.58325.
- 17 Kotovirta, V., Toivanen, T., Tergujeff, R. and Huttunen, M. (2012). Participatory Sensing
18 in Environmental Monitoring – Experiences, *Proc. of the Sixth Conference on Innovative*
19 *Mobile and Internet Services in Ubiquitous Computing*, Palermo, 155-162. doi:
20 10.1109/IMIS.2012.70.
- 21 Kovács, F., Legány, C. and Babos A. (2006). Cluster Validity Measurement Techniques,
22 *Proc. of the Fifth WSEAS International Conference on Artificial Intelligence, Knowledge*
23 *Engineering and Data Bases*, 388-393.
- 24 Lagus, K., Kaski, S. and Kohonen, T. (2004). Mining massive documents collections by
25 the WEBSOM method. *Journal of Information Sciences – Special issue: Soft computing*
26 *data mining*. 163 (1-3), 135-156. doi: 10.1016/j.ins.2003.03.017.

- 1 Li, R., Lei, H. K., Khadiwala, R. and Chen-Chuan Chang, K. (2012). TEDAS: a Twitter-
2 based Event detection and Analysis System, Proc. of the IEEE 28th International
3 Conference on Data Engineering (ICDE), Arlington, Virginia USA, 1273-1276. doi:
4 10.1109/ICDE.2012.125.
- 5 MacEachren, A. M., Robinson, A. C., Jaiswal, A., Pezanowski, S., Savelyev, A.,
6 Blanford, J. and Mitra, P. (2011). Geo-Twitter Analytics: Applications in Crisis
7 Management, Proc. of the 25th International Cartographic Conference, Paris, France, 1-
8 8.
- 9 Maynard, D. and Funk, A. (2011). Automatic detection of political opinions in tweets,
10 Proc. of the Eighth International Conference on The Semantic Web: ESWC 2011,
11 Heraklion, Greece, LNCS 7117, 88-99.
- 12 Mostafa, M. M. (2013). More than words: Social networks' text mining for consumer
13 brand sentiments. *Expert Systems and Applications*. 40 (10), 4241-4251. doi:
14 10.1016/j.eswa.2013.01.019.
- 15 Mun, M., Reddy, S., Shilton, K., Yau, N., Burke, J., Estrin, D., Hansen, M., Howard, E.,
16 West, R. and Boda, P. (2009). PEIR, the personal environmental impact report, as a
17 platform for participatory sensing systems research, Proc. of the Seventh International
18 Conference on Mobile systems, applications and services, Kraków, Poland, 55-68. doi:
19 10.1145/1555816.1555823.
- 20 Niska, H., Hiltunen, T., Karppinen, A., Ruuskanen, J. and Kolehmainen, M. (2004).
21 Evolving the neural network model for forecasting air pollution time series. *Engineering*
22 *Applications of Artificial Intelligence*. 17 (2), 159-167. doi:
23 10.1016/j.engappai.2004.02.002.
- 24 Pang, B. and Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and*
25 *Trends in Information Retrieval*. 2 (1-2), 1-135. doi: 10.1561/1500000011.

- 1 Popescu, A. M. and Pennacchiotti, M. (2010). Detecting controversial events from
2 Twitter, Proc. of the 19th ACM International Conference on Information and Knowledge
3 Management, 1873-1876. doi: 10.1145/1871437.1871751.
- 4 Ritter, H. and Kohonen, T. (1989). Self-Organizing Semantic Maps. *Bioclinical*
5 *Cybernetics*. 61 (4), 241-254. doi: 10.1007/BF00203171.
- 6 Robinson, E. (2010). Integration of multi-sensory earth observations for characterization
7 of air quality events. Ph.D. Dissertation, School of Engineering and Applied Science,
8 Department of Energy, Environmental and Chemical Engineering, Washington University
9 in St. Louis, USA.
- 10 Rousseeuw, P.J. (1987). Silhouettes: a graphical aid to the interpretation and validation
11 of cluster analysis. *Journal of Computational and Applied Mathematics*. 20, 53-65. doi:
12 10.1016/0377-0427(87)90125-7.
- 13 Sadilek, A., Kautz, H. and Silenzio, V. (2012). Predicting Disease Transmission from
14 Geo-Tagged Micro-Blog Data, Proc. of the 26th AAAI Conference on Artificial
15 Intelligence, Toronto, Ontario, Canada, AAAI Press 2012.
- 16 Salton, G., Wong, A. and Yang, C.S. (1975). A vector space model for automatic
17 indexing. *Communications of the ACM*. 18(11), 613-620. doi: 10.1145/361219.361220.
- 18 Sammon, J.W. (1969). A Nonlinear Mapping for Data Structure Analysis. *IEEE*
19 *Transactions on Computers*. C-18 (5), 401-409. doi: 10.1109/T-C.1969.222678.
- 20 Skön, J.P., Kauhanen, O. and Kolehmainen M. (2011). Energy Consumption and Air
21 Quality Monitoring System, Proc. of the Seventh International Conference on Intelligent
22 Sensors, Sensor Networks and Information Processing, Adelaide, Australia, 163-167.
23 doi: 10.1109/ISSNIP.2011.6146606.
- 24 Sullivan, D. (2001). Document Warehousing and Text Mining: Techniques for Improving
25 Business Operations, Marketing, and Sales. John Wiley & Sons.

- 1 Tan, A.H. (1999). Text Mining: The state of the art and the challenges, Proc. of the
2 PAKDD Workshop on Knowledge Discovery from Advanced Databases, Beijing, China,
3 65-70.
- 4 TechCrunch (2012). Twitter May have 500M+ Users but only 170M Are Active, 75% on
5 Twitter's Own Clients. [http://techcrunch.com/2012/07/31/twitter-may-have-500m-users-
6 but-only-170m-are-active-75-on-twitters-own-clients/](http://techcrunch.com/2012/07/31/twitter-may-have-500m-users-but-only-170m-are-active-75-on-twitters-own-clients/) (accessed Apr 09, 2013).
- 7 Ultsch, A. and Siemon, H.P. (1990). Kohonen's Self Organizing Maps for Exploratory
8 Data Analysis, Proc. of the International Neural Networks Conference (INNC), Paris,
9 France, 305-308.
- 10 Voukantsis, D., Karatzas, K., Kukkonen, J., Räsänen, T., Karppinen, A. and
11 Kolehmainen, M. (2011). Intercomparison of air quality data using principal component
12 analysis, and forecasting of PM₁₀ and PM_{2.5} concentrations using artificial neural
13 networks, in Thessaloniki and Helsinki. Science of the Total Environment. 409 (7), 1266-
14 1276. doi: 10.1016/j.scitotenv.2010.12.039.
- 15 Voukantsis, D., Niska, H., Karatzas, K., Riga, M., Damialis, A. and Vokou, D. (2010).
16 Forecasting daily pollen concentrations using data-driven modeling methods in
17 Thessaloniki, Greece. Atmospheric Environment. 44 (39), 5101-5111. doi:
18 10.1016/j.atmosenv.2010.09.006.
- 19 Yang, H.C. and Lee, C.H. (2010). A novel self-organizing map algorithm for text mining,
20 Proc. of International Conference on System Science and Engineering, 417-420. doi:
21 10.1109/ICSSE.2010.5551734.
- 22 Zhang, X., Fuehres, H. and Gloor, A. P. (2011). Predicting Stock Market Indicators
23 Through Twitter "I hope it is not bad as I fear". Procedia - Social and Behavioral
24 Sciences. 26, 55-62. doi: 10.1016/j.sbspro.2011.10.562.