

Do Not Feel The Trolls

Erik Cambria¹, Praphul Chandra²,
Avinash Sharma², and Amir Hussain¹

¹ University of Stirling, Stirling, UK

² HP Labs India, Bangalore, India

{eca, ahu}@cs.stir.ac.uk, {praphul.chandra, sharma}@hp.com

<http://sentic.net>

Abstract. The passage from a read-only to a read-write Web gave people the possibility to freely interact, share and collaborate through social networks, online communities, blogs, wikis and other online collaborative media. The democracy of the Web is what made it so popular in the past decades but such a high degree of freedom of expression also gave birth to negative side effects – the so called ‘dark side’ of the Web. An example of this is trolling, i.e., the exploitation of the anonymity of the Web to post inflammatory and outrageous messages directed to one specific person or community to provoke them into a desired emotional response. Online community masters usually warn users against trolls with messages such as DNFTT (Do Not Feed The Trolls) but so far this has not been enough to stop trolls trolling. The aim of this work is to use sentic computing, a new paradigm for the affective analysis of natural language text, to detect trolls and hence prevent web-users from being emotionally hurt by malicious posts.

Key words: Sentic Computing, AI, Semantic Web, NLP, Opinion Mining and Sentiment Analysis

1 Introduction

In Internet slang, a troll is someone who posts inflammatory, extraneous, or off-topic messages in an online community, such as an online discussion forum, chat room, or blog, with the primary intent of provoking other users into a desired emotional response or of otherwise disrupting normal on-topic discussion [1].

The amount of social data on the Web is on an infinite uphill and online social networking is becoming one of the most prevalent means of expression worldwide. Websites like Twitter, YouTube and Blogger are providing a tunnel to link different parts of the world and also different classes of global society.

The flip-side of the coin, on the other hand, is rather dark, fractious and bizarre. Social web is inherently democratic and user anonymity is gratuitous in this space. Be it real world or virtual social web, existence of malicious faction among inhabitants and users is inevitable.

In social web context, emotional attacks on a person or a group through malicious and vulgar comments in order to provoke response are referred to as ‘trolling’ and the generator is called ‘a troll’. The term was first used in early 1990 and since then a lot of concern has been raised to contain or curb trolls.

This work proposes a technique based on sentic computing [2], a novel paradigm for the affective analysis of natural language text, to automatically detect and check web trolls. We present results that are effective in controlling trolls efficiently. To the best of our knowledge this work has no prior.

The structure of the paper is the following: Section 2 argues about the phenomenon of Internet trolling, Section 3 presents the state of the art of malicious post detection, Section 4 and Section 5 explain in detail the techniques used within this work, Section 6 illustrates the overall process for filtering trolls, Section 7 demonstrates the potential of such process through an evaluation study, and Section 8 comprises concluding remarks and a description of future work.

2 The Internet Trolling Phenomenon

Trolling is a method of fishing where some baited fishing lines are drawn through the water, usually from a slow-moving boat, with the purpose of hooking unwary fish. An online troll does pretty much the same.

The trend of trolling, where anonymous online users bombard victims with offensive messages or abuse, appears to have spread a lot recently and it is alarming most of the biggest social networking sites since, in extreme cases such as abuse, has led some teenagers to commit suicide. These attacks usually address not only individuals but also entire communities. For example, reports have claimed that a growing number of Facebook tribute pages had been targeted, including those in memory of the Cumbria shootings victims and soldiers who died in Afghanistan.

At present users cannot do much rather than manually delete abusive messages. Current anti-trolling methods, in fact, mainly consist in identifying additional accounts that use the same IP address and blocking fake accounts based on name and anomalous site activity e.g. users who send lots of messages to non-friends or whose friend requests are rejected at a high rate.

In July 2010 Facebook launched an application that gives users a direct link to advice, help and the ability to report cyber problems to the Child Exploitation and Online Protection Centre (CEOP) [3]. Reporting trouble through a link or a button, however, is a too slow process since social networking websites usually cannot react instantly to these alarms. A button, moreover, does not stop users from being emotionally hurt by trolls and it is more likely to be pushed by people who actually do not need help rather than, for instance, children who are being sexually groomed and do not realize it.

For these reasons, we need systems able to automatically analyze semantics and sentics, i.e., cognitive and affective information, associated to natural language in order to filter out inopportune messages and, hence, stop users from ‘feeling’ the trolls.

3 Related Work

A prior analysis of the trustworthiness of statements published on the Web has been presented by Rowe and Butters [4]. Their approach adopts a contextual trust value determined for the person who asserted a statement as the trustworthiness of the statement itself. This study, however, does not focus on the problem of trolling but rather on defining a contextual accountability for the detection of web, email and opinion spam.

Existing approaches in these fields, in particular, can be grouped into three main categories: keyword spotting [5][6], in which text is classified according to the presence of fairly unambiguous spam words, lexical affinity [7][8], which assigns arbitrary words a probabilistic affinity for spam content, and statistical methods [9][10], which consist in calculating the valence of keywords, punctuation and word co-occurrence frequencies on the base of a large training corpus.

The problem with these approaches is that they mainly rely on parts of text in which web, email and opinion spam is explicitly expressed through spam links, commercial terms or abusive words. But, more generally, spam manifests implicitly through context and domain dependent concepts, which makes keyword-based approaches extremely ineffective.

To overcome this problem we need to use natural language processing (NLP) techniques that rely on semantics rather than syntactics. Within this work, in particular, we exploit two sentic-computing tools to extract semantics and sentics from web posts and, eventually, process the results in order to detect and filter trolls.

4 Sentic Computing

Sentic computing is a new opinion mining and sentiment analysis paradigm which exploits AI and Semantic Web techniques to better recognize, interpret and process opinions and sentiments in natural language text.

The term sentic computing derives from the Latin ‘sentire’ (the root of words such as sentiment and sensation) and ‘sense’ (intended as common sense) and concerns a kind of computing that relates to, arises from and influences opinions and sentiments in natural language text.

In sentic computing, the analysis of text is not based on statistical learning models but rather on common sense reasoning tools [11] and domain-specific ontologies [12]. Differently from statistical classification, which generally requires large inputs and thus cannot appraise texts with satisfactory granularity, sentic computing enables the analysis of documents not only on the page or paragraph-level but also on the sentence-level.

Within this work, in particular, we exploit the combination of two sentic-computing tools for the extraction of semantics and sentics from web posts, i.e., a multi-dimensional vector space of common sense and affective knowledge (Section 4.1) coupled with a novel emotion categorization model born from the idea that our mind consists of four independent emotional spheres, whose different levels of activation make up the total emotional state of the mind (Section 4.2).

4.1 AffectiveSpace

AffectiveSpace [13] is a language visualization system which transforms natural language from a linguistic form into a multi-dimensional space. AffectiveSpace is built by blending ConceptNet [14], a semantic network of common sense knowledge, and WordNet-Affect [15], a linguistic resource for the lexical representation of emotions. This alignment operation yields *AffectNet*: a new dataset in which common sense and affective knowledge coexist, i.e., a matrix $14,301 \times 117,365$ whose rows are concepts (e.g. ‘dog’ or ‘bake cake’), whose columns are either common sense and affective features (e.g. ‘isA-pet’ or ‘hasEmotion-joy’), and whose values indicate truth values of assertions.

Therefore, in *AffectNet*, each concept is represented by a vector in the space of possible features whose values are positive for features that produce an assertion of positive valence (e.g. ‘a penguin is a bird’), negative for features that produce an assertion of negative valence (e.g. ‘a penguin cannot fly’) and zero when nothing is known about the assertion. The degree of similarity between two concepts, then, is the dot product between their rows in *AffectNet*. The value of such a dot product increases whenever two concepts are described with the same feature and decreases when they are described by features that are negations of each other. When performed on *AffectNet*, however, these dot products have very high dimensionality (as many dimensions as there are features) and are difficult to work with. In order to approximate these dot products in a useful way, we project all of the concepts from the space of features into a space with many fewer dimensions, i.e., we reduce the dimensionality of *AffectNet* by means of principal component analysis (PCA). In particular, we perform truncated singular value decomposition (TSVD) [16] on *AffectNet* and obtain a new matrix, *AffectNet**, which forms a low-rank approximation of the original data. This estimation is based on minimizing the Frobenius norm of the difference between *AffectNet* and *AffectNet** under the constraint $\text{rank}(\text{AffectNet}^*) = k$ and it represents the best approximation of *AffectNet* in the least-square sense.

In particular, we choose to discard all but the first 100 principal components and hence obtain AffectiveSpace (Fig. 1), a 100-dimensional space in which different vectors represent different ways of making binary distinctions among concepts and emotions. In AffectiveSpace common sense and affective knowledge are in fact combined, not just concomitant, i.e., everyday life concepts like ‘have breakfast’, ‘meet people’ or ‘watch tv’ are linked to a hierarchy of affective domain labels. By exploiting the information sharing property of TSVD, concepts with the same affective valence are likely to have similar features, i.e., concepts concerning the same opinion tend to fall near each other in the vector space. Concepts and emotions are represented by vectors of 100 coordinates: these coordinates can be seen as describing concepts in terms of ‘eigenmoods’ that form the axes of AffectiveSpace i.e. the basis e_0, \dots, e_{99} of the vector space. For example, the most significant eigenmood, e_0 , represents concepts with positive affective valence. That is, the larger a concept’s component in the e_0 direction is, the more affectively positive it is likely to be. Consequently concepts with negative e_0 components have negative affective valence.

3. the user is comfortable with the interface (Sensitivity)
4. the user is disposed to use the application (Aptitude)

Each affective dimension is characterized by six levels of activation, called ‘sentic levels’, which determine the intensity of the expressed/perceived emotion as a float $\in [-3,3]$. These levels are also labeled as a set of 24 basic emotions (six for each of the affective dimensions) in a way that the model can specify the affective information associated to text both in a dimensional and in a discrete form. The dimensional form, in particular, is called ‘sentic vector’ and it is a four dimensional vector that can potentially express any human emotion in terms of Pleasantness, Attention, Sensitivity and Aptitude. Some particular sets of sentic vectors have special names as they specify well-known compound emotions. For example the set of sentic vectors with a level of Pleasantness $\in (1,2]$ (‘joy’), a null Attention, a null Sensitivity and a level of Aptitude $\in (1,2]$ (‘trust’) are called ‘love sentic vectors’ since they specify the compound emotion of ‘love’.

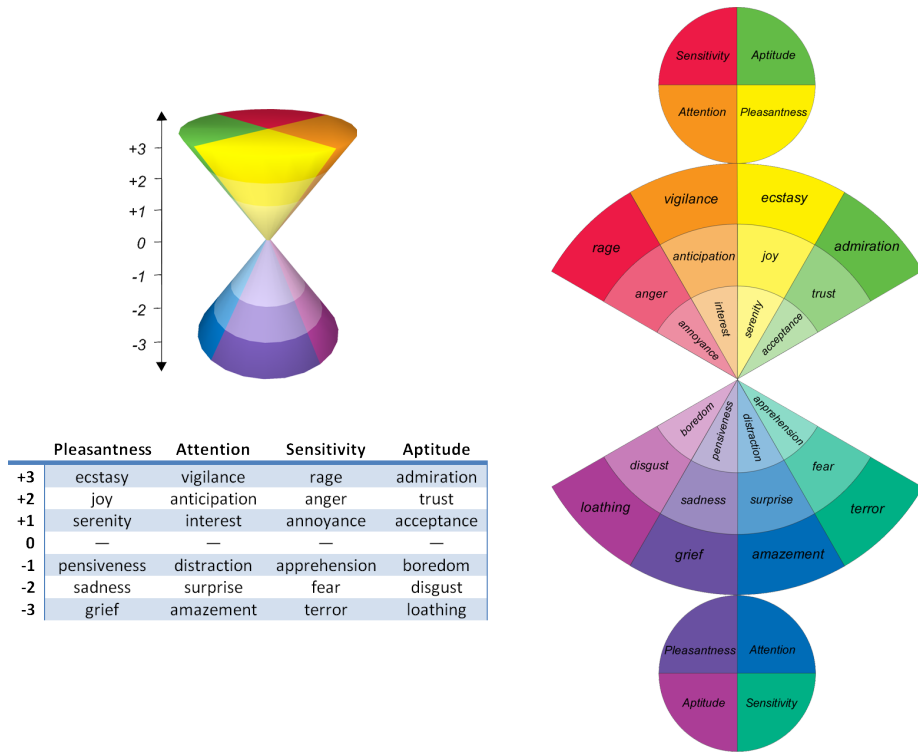


Fig. 2. The Hourglass of Emotions

5 Troll Detector

The main aim of the Troll Detector is to identify malicious contents in natural language text with a certain confidence level. To train the detector, we first identify the concepts most commonly used by trolls (Section 5.1) and then expand the resulting knowledge base with semantically related concepts (Section 5.2). We finally define a method to calculate *trollness* i.e. the probability for a post to be edited by a troll (Section 5.3).

5.1 CF-IOF Weighting

The technique we use to identify the concepts commonly used by trolls is called CF-IOF [19] (concept frequency – inverse opinion frequency) and it is an approach similar to TF-IDF weighting which evaluates how important a concept is to a set of opinions concerning the same topic.

We first calculate the frequency of a concept c_i for a given topic j by counting the occurrences of the concept c_i in the set of available j -tagged opinions and divide the result by the sum of occurrences of the same concept in the whole set of opinions concerning j . We then multiply this frequency by the logarithm of the total number of opinions divided by the number of opinions containing the concept c_i , that is:

$$(CF-IOF)_i = \sum_j \frac{n_{i,j}}{\sum_k n_{k,j}} \log \frac{|O|}{|\{o : c_i \in o\}|}$$

where $n_{i,j}$ is the number of occurrences of the considered concept c_i in the opinions tagged with the topic j , $|\{o : c_i \in o\}|$ the number of opinions where c_i appears and $|O|$ the total number of opinions.

A high weight in CF-IOF is reached by a high concept frequency (in the given opinions) and a low opinion frequency of the concept in the whole collection of opinions. Therefore, thanks to CF-IOF weights, we manage to filter out common concepts and detect relevant concepts that are usually used by trolls to emotionally attack unaware users.

5.2 Spectral Association

In order to expand the set of concepts previously obtained by applying CF-IOF, we use a technique called spectral association [20] that involves assigning values, or activations, to ‘seed concepts’ and applying an operation that spreads their values across the ConceptNet graph.

This operation, an approximation of many steps of spreading activation, transfers the most activation to concepts that are connected to the key concepts by short paths or many different paths in common sense knowledge. In particular, we build a matrix C that relates concepts to other concepts, instead of their features, and add up the scores over all relations that relate one concept to another, disregarding direction.

Applying C to a vector containing a single concept spreads that concept’s value to its connected concepts. Applying C^2 spreads that value to concepts connected by two links (including back to the concept itself). But what we’d really like is to spread the activation through any number of links, with diminishing returns, so perhaps the operator we want is:

$$1 + C + \frac{C^2}{2!} + \frac{C^3}{3!} + \dots = e^C$$

We can calculate this odd operator, e^C , because we can factor C . C is already symmetric, so instead of applying Lanczos’ method to CC^T and getting the SVD, we can apply it directly to C and get the spectral decomposition $C = V\Lambda V^T$. As before, we can raise this expression to any power and cancel everything but the power of Λ . Therefore, $e^C = Ve^{\Lambda}V^T$. This simple twist on the SVD lets us calculate spreading activation over the whole matrix instantly.

As with the SVD, we can truncate these matrices to k axes and therefore save space while generalizing from similar concepts. We can also rescale the matrix so that activation values have a maximum of 1 and do not tend to collect in highly-connected concepts such as ‘person’, by normalizing the truncated rows of $Ve^{\Lambda/2}$ to unit vectors, and multiplying that matrix by its transpose to get a rescaled version of $Ve^{\Lambda}V^T$.

5.3 Calculating Trollness

In order to calculate the probability for a post to be edited by a troll, we exploit both the semantics and the sentsics associated to it.

For each concept contained in the post, the Troll Detector checks if this belongs to the set of ‘troll concepts’ calculated through spectral association and exploits its relative sentic vector to check if it carries malicious affective charge. By analyzing a set of 1000 offensive phrases extracted from Wordnik [21], in fact, we found that, statistically, a post is likely to be edited by a troll when its average sentic vector has a high absolute value of Sensitivity and a very low polarity. Hence we defined the *trollness* t_i associated to a concept c_i as a float $\in [0, 1]$ such that:

$$t_i(c_i) = \frac{s_i(c_i) + |Snsit(c_i)| - p_i(c_i)}{5}$$

where s_i (float $\in [0, 1]$) is the semantic similarity of c_i wrt any of the CF-IOF seed concepts, p_i (float $\in [-1, 1]$) is the polarity associated to the concept c_i and 5 is the normalization factor (the maximum value of the numerator in fact is given by a similarity of 1, a Sensitivity of 3 or -3 and a polarity equal to -1). In particular, p_i is defined [22] as:

$$p_i(c_i) = \frac{Plsnt(c_i) + |Attnt(c_i)| - |Snsit(c_i)| + Aptit(c_i)}{9}$$

where 9 is the normalization factor (since the numerator’s maximum value is given by the sentic vectors $[3, \pm 3, 0, 3]$ and the minimum by $[-3, 0, \pm 3, -3]$).

In the formula, Attention and Sensitivity are taken in absolute value since, from the point of view of polarity rather than affection, all of their sentic values represent positive and negative values respectively (e.g. ‘anger’ is positive in the sense of level of activation of Sensitivity but negative in terms of polarity and ‘surprise’ is negative in the sense of lack of Attention but positive from a polarity point of view).

Hence, the total *trollness* of a post containing N concepts is defined as:

$$t = \frac{5}{9} \sum_{i=1}^N \frac{9 s_i(c_i) + 10 |Snsit(c_i)| - Plsnt(c_i) - |Attnt(c_i)| - Aptit(c_i)}{N}$$

This information is stored, together with post type and content plus sender and receiver ID, in an interaction database that keeps trace of all the messages and comments interchanged between users within the same social network.

Posts with a high level of *trollness* (current threshold has been set, using a trial and error approach, to 60%) are labeled as troll posts and, whenever a specific user addresses more than two troll posts to the same person or community, his/her sender ID is labeled as troll for that particular receiver ID.

All the past troll posts sent to that particular receiver ID by that specific sender ID are then automatically deleted from the website (but kept in the database with the possibility for the receiver to either visualize them in an apposite *troll folder* and, in case, restore them). Moreover, any new post with a high level of *trollness* edited by a user labeled as troll for that specific receiver is automatically blocked i.e. saved in the interaction database but never displayed in the social networking website.

6 Troll Filtering Process

The process for filtering trolls (illustrated in Fig. 3) comprises four main components: a NLP module, which performs a first skim of the document, a Semantic Parser, whose aim is to extract concepts from the lemmatized text, AffectiveSpace, for the extraction of sentics from the given concepts, and the Troll Detector, whose aim is to detect and eventually block the troll.

The NLP module interprets all the affective valence indicators usually contained in text such as special punctuation, complete upper-case words, onomatopoeic repetitions, exclamation words, negations, degree adverbs and emoticons, and eventually lemmatizes text.

The Semantic Parser then deconstructs text into concepts and provides, for each of them, the relative frequency, valence and status i.e. the concept’s occurrence in the text, its positive or negative connotation, and the degree of intensity with which the concept is expressed.

The AffectiveSpace module projects the retrieved concepts into the vector space, clustered wrt the Hourglass model, and it infers the affective valence of these, in terms of Pleasantness, Attention, Sensitivity and Aptitude, according to the positions they occupy in the space.

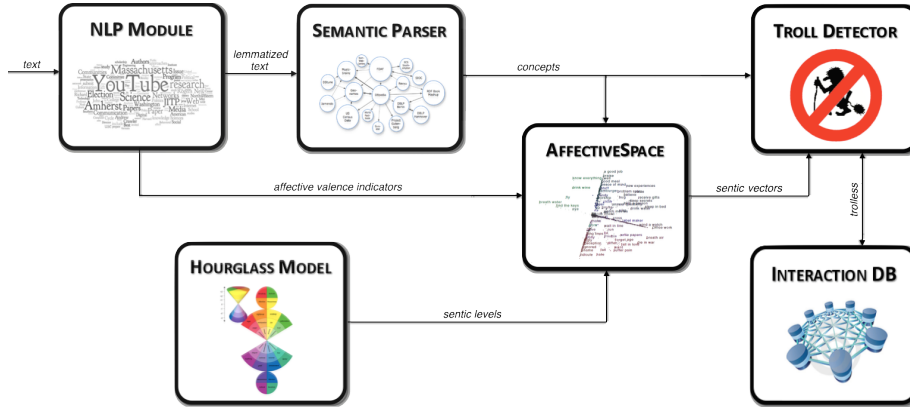


Fig. 3. Troll Filtering Process

This information, encoded as a sentic vector, is given as input to the Troll Detector which exploits it, together with the semantic information coming directly from the Semantic Parser, to calculate the post’s *trollness* and, eventually, to detect and block the troll (according to the information stored in the interaction database). As an example of Troll Filtering Process output, we can consider a troll post recently addressed to the Indian author Chetan Bhagat: “You can’t write, you illiterate douchebag, so quit trying, I say!!!”. In this case we have a very high level of Sensitivity (corresponding sentic level ‘rage’) and a negative polarity, which give a high percentage of *trollness*, as shown below:

< Concept: !‘write’>
 < Concept: ‘illiterate’>
 < Concept: ‘douchebag’>
 < Concept: ‘quit try’>
 < Concept: ‘say’>

Semantics: 0.69
 Sentics: [0.0, 0.48, 2.7, -1.22]
 Polarity: -0.38
 Trollness: 0.75

7 Evaluation

In order to perform a first evaluation of our system, we considered a set of 500 tweets (most of which fetched from Wordnik) manually annotated as troll and non-troll posts. We considered true positives those posts with both a positive troll-flag and a *trollness* $\in [0.6, 1]$ and those with both a negative troll-flag and a *trollness* $\in [0, 0.6)$. The threshold has been set to 60% based on trial and error over a separate dataset of 50 tweets.

Results show that, by using the Troll Filtering Process, inflammatory and outrageous messages can be identified with good precision (82%) and decorous recall rate (75%). In particular, the F-measure value (78%) is significantly high compared to the corresponding F-measure rates of the baseline methods (53% for keyword spotting, 59% for lexical affinity, 66% for statistical methods).

However, we expect to obtain much better results by evaluating the process at interaction-level rather than just at post-level. In the next future, in fact, we plan to evaluate the Troll Filtering Process by monitoring not just single posts but also users' holistic behavior within the same social network (i.e. contents and recipients of their interaction) and submit further results elsewhere for publication.

8 Conclusion and Future Efforts

As the Web plays a more and more significant role in people's social lives, it contains more and more information concerning their opinions and feelings. After the explosion of Web 2.0, a lot of users have been exploiting this trend, together with the anonymity of the Web, to attack specific people or communities with inflammatory and outrageous messages and, hence, provoke them into a desired emotional response.

For their fiendish nature, these users have been labeled as trolls. Online community masters have desperately tried to warn users against these mischievous people with messages such as DNFTT (Do Not Feed The Trolls) but so far this has not been enough to stop trolls trolling.

Within this work we exploited sentic computing, a new paradigm for the affective analysis of natural language text, to design a process capable to extract semantics and sentics from web-posts and infer from these the truthfulness of user interaction.

The main aim of the Troll Filtering Process, in fact, is to exploit the cognitive and affective information associated to natural language text to define a level of *trollness* of each post and, according to this, classify users and prevent the malicious ones from emotionally hurting other people or communities within the same social network.

In the next future, we plan to improve the process by using a much bigger dataset for training the Troll Detector and also to perform an evaluation of the system at interaction-level rather than just at post-level, in order to better understand, and hence prevent, trolls' behavior.

Eventually, we plan to enhance the system by making most of its functionalities available as web-services in a way that the Troll Filtering Process could be easily embedded in any social networking website and, hence, change the meaning of the popular acronym often displayed in these websites, DNFTT, from a shadowy and often ineffective suggestion to a reassuring and deterrent slogan – Do Not Feel The Trolls.

References

1. [http://en.wikipedia.org/wiki/Troll_\(Internet\)](http://en.wikipedia.org/wiki/Troll_(Internet)) – Wikipedia
2. Cambria, E., Speer, R., Havasi, C., Hussain, A.: SenticNet: A publicly available semantic resource for opinion mining. In: AAAI CSK, pp. 14–18, Arlington (2010)
3. <http://fw.to/W9zFwtW> – The Telegraph (2010)
4. Rowe, M., Butters, J.: Assessing Trust: Contextual Accountability. In: SPOT at ESWC, Heraklion (2009)
5. Dave, K., Lawrence, S. Pennock, D.: Mining the Peanut Gallery: Opinion Extraction and Semantic Classification of Product Reviews. In: WWW, Budapest (2003)
6. Chandrasekaran, M., Karayanan, K., Upadhyaya, S.: Towards Phising E-Mail Detection Based on Their Structural Properties. In: SCSS, New York (2006)
7. Hu, M., Liu, B.: Mining and Summarizing Customer Reviews. In: KDD, Seattle (2004)
8. Jindal, N., Liu, B.: Analyzing and Detecting Review Spam. In: ICDM, Omaha (2007)
9. Li, W., Zhong, N., Liu, C.: Combining Multiple Email Filters Based on Multivariate Statistical Analysis. In: ISMIS, Bari (2006)
10. Jindal, N., Liu, B.: Opinion Spam and Analysis. In: WSDM, Palo Alto (2008)
11. Cambria, E., Hussain, A., Havasi, C., Eckl, C.: Common Sense Computing: From the Society of Mind to Digital Intuition and Beyond. LNCS, vol. 5707, pp. 252–259. Springer-Verlag, Berlin Heidelberg (2009)
12. Cambria, E., Grassi, M., Hussain, A., Havasi, C.: Sentic Computing for Social Media Marketing. Multimedia Tools and Applications, DOI 10.1007/s11042-011-0815-0
13. Cambria, E., Hussain, A.: Sentic Computing: Techniques, Tools, and Applications. Dordrecht, Netherlands: Springer, DOI 10.1007/978-94-007-5070-8
14. Havasi, C., Speer, R., Alonso, J.: ConceptNet 3: a Flexible, Multilingual Semantic Network for Common Sense Knowledge. In: RANLP, Borovets (2007)
15. Strapparava, C., Valitutti, A.: WordNet-Affect: an Affective Extension of WordNet. In: LREC, Lisbon (2004)
16. Wall, M., Rechtsteiner, A., Rocha, L.: Singular Value Decomposition and Principal Component Analysis. In: Berrar, D. et al. (eds.) A Practical Approach to Microarray Data Analysis. pp. 91–109. Kluwer, Norwell (2003)
17. Plutchik, R.: The Nature of Emotions. *American Scientist* 89(4), 344–350 (2001)
18. Minsky, M.: The Emotion Machine. Simon and Schuster, New York (2006)
19. Cambria, E., Hussain, A., Durrani, T., Havasi, C., Eckl, C., Munro, J.: Sentic Computing for Patient Centered Applications. In: IEEE ICSP10, pp. 1279–1282, Beijing (2010)
20. Havasi, C., Speer, R., Holmgren, J.: Automated Color Selection Using Semantic Knowledge. In: AAAI CSK10, Arlington (2010)
21. <http://wordnik.com> – Wordnik
22. Cambria, E., Hussain, A., Havasi, C., Eckl, C., Munro, J.: Towards Crowd Validation of the UK National Health Service. In: WebSci10, Raleigh (2010)