

Adding affective state to contextonyms

Ovidiu SERBAN⁽¹⁾⁽²⁾, Alexandre PAUCHET⁽¹⁾, Alexandrina ROGOZAN⁽¹⁾,
Horia F. POP⁽²⁾, AND Jean-Pierre PECUCHET⁽¹⁾
email: ovidiu.serban@insa-rouen.fr

1. Contributions

Contextonyms are defined by Ji and Ploux [JPW03] as relevant contextually related words for a target word. By context, they mean choosing a certain number of neighboring words of the target word (from a small-sized window to one or more paragraphs). Unlike synonyms or antonyms, contextonyms are not symmetric or transitive (i.e., when target word W has contextonyms c_1, c_2, \dots, c_n , W is not necessarily a contextonym of c_i ($1 \leq i \leq n$)). In their articles, they provide also a simple and efficient method to filter out irrelevant noise from the context.

Starting from these steps we reduced the method to only the α and β filtering, and we applied the β filtering process recursively, to all the children. Also, because it was unclear how α and β were chosen, we decided to do a grid search in the bi-dimensional space. Each generation of annotated contextonyms is evaluated using a measure described in the next paragraph.

Starting from the annotated contextonyms, we propose two classification methods. The first one is the result of applying the same measure used for building the graph and the second one is a pseudo-LSA decomposition that provides numerical features to a Self Organizing Map classifier.

2. Building the contextonyms

Starting with a given window size, we compute the frequencies of appearance from each word pair in a phrase. After the frequencies are build, a pre-filtering is made in order to eliminate very rare context, with the frequency equals to 1 or 1% of the node frequency. After these steps, we define the α and β filtering as following:

α Filtering

Given the W_i^n word, with the context words: $c_1^i, c_2^i, \dots, c_n^i$ and the parameter α (where $0 < \alpha \leq 1$)

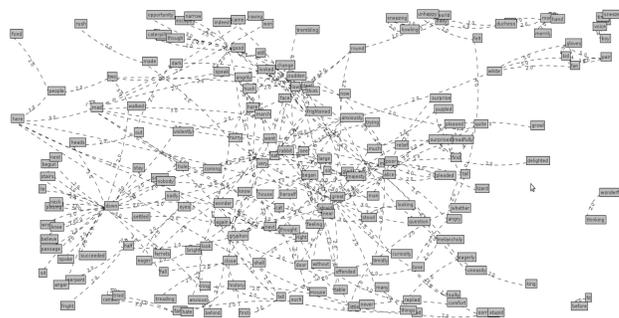
We select the first k words in the W sequence: $k = n * \alpha \rightarrow W_i^n : c_1^i, c_2^i, \dots, c_k^i$

β Filtering

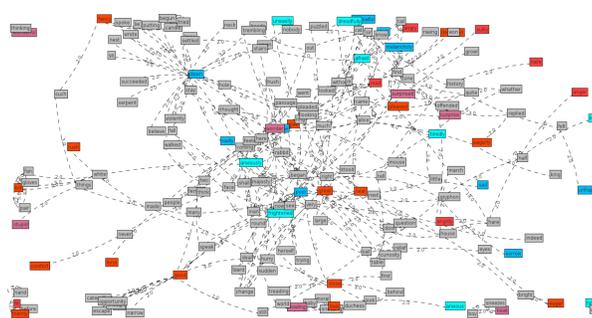
Given parameter β (where $0 < \beta \leq 1$)

We select the first l words for each contextualized word: $l = m * \beta \rightarrow c_j^m : g_1^j, g_2^j, \dots, g_l^j, 1 < j \leq k$

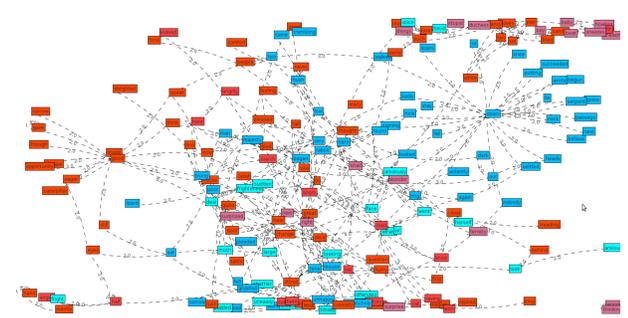
4. Annotation Steps - Example



Step 1: Contextonyms



Step 2: Pre-annotation with strong emotional words



Step 3: The affective contextonyms, the strong emotional words had spread their influence in the graph

5. Classification

The direct application of contextonyms is the measure defined in the previous paragraph. For a certain sentence in a corpus, we can compute the emotional conflicts and the general affective value by evaluating the emotional label on each word, according to (2) equation.

$$(2_c) E(y) = \sum E(x, y), \forall x \in N(y) \wedge x \in \{\text{Contextonyms dictionary}\}$$

$$(2_{cs}) E(\text{sentence}) = \sum E(x), \forall x \in \{\text{sentence}\}$$

where $N(y)$ is the collection of all the neighbors of y

The conflict index is computed the same as (3).

The decision is taken if the conflict index is lower than a given threshold, then the general emotional value of a sentence will be computed by (2_{cs}).

The other classifier that we propose is a Self Organizing Map (SOM), used with features extracted by a pseudo-LSA [SM08]. Basically, the pseudo-LSA method is the same the classical LSA approach, but instead of using word to document space of representation, you use word clusters to document space.

6. Conclusion and future work

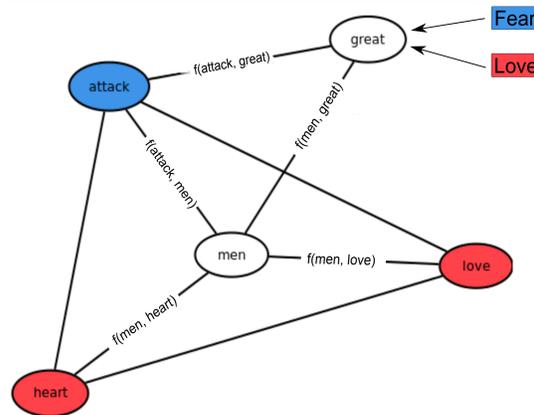
Most of the approaches in text mining are working with large dictionaries in order to detect the emotional valence of a corpus. WordNet Affect [VSS05], ConceptNet [LS04] or SentiWordNet [BES10] were generated for this purpose, but none of them offer a large enough database to classify quickly all the frequent words in English.

Another strong point of our model is that we involve context in the decision process, which is the foundation of our model. This aids the decision especially in the cases of semantic ambiguity or weak emotional presence.

Our goal is to create an annotated contextonym database that will give a more clear image of the English language and also find several classification methods suited for semantic emotion detection.

As for future perspectives in our work, we want to develop a classification engine, that will be able to detect the emotions in real-time and integrate our work in other projects.

3. Annotation steps



$$(1) E(x,y) = f(x,y) * E(x) / f(x), \forall x \in \{\text{"emotional" words}\}, y \in \{\text{"non-emotional" words}\},$$

E - is the emotional state of the word x or an edge (x,y)
 f - is the frequency of appearance of the word x or (x,y) couple

The equation (1) applied on $(men, heart)$ edge is:

$$(1') E(men, heart) = f(men, heart) * E(heart) / f(heart)$$

After each emotional label is computed for each edge, we can compute the emotional label of unlabeled nodes:

$$(2) E(y) = \sum E(x, y), \forall x \in N(y)$$

where $N(y)$ is the collection of all the neighbors of y

$$(2') E(men) = E(men, heart) + E(men, love) + E(men, attack)$$

There are some cases of conflict, which can occur when two opposite emotions (positive or negative) appear on the same label, or in the same clique. In the case of same label, we just consider the dominant emotion and if this is not possible, we consider this a conflictual case.

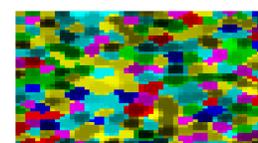
In the case of conflictual cliques, we compute $C(q)$ measure which is the number of the conflicts inside the clique q .

Globally the quality measure for a contextonym graph is defined as following:

$$(3) C(\text{contextonyms}) = \sum C(q), \forall q - \text{clique in the contextonym graph}$$

In [SM08] this method was used with WordNet Affect synsets, but we plan to use it with the cliques in our contextonym graph. Because this decomposition is not so strict as the classical LSA approach, we believe the noise in the data will be reduced by the SOM algorithm.

In order to test if the SOM is suited for the classification of emotional data, we started testing it on the corpus proposed by C. Strapparava and R. Mihalcea at the SemEval 2007 conference, for the task 14 [SM08]. In our first tests we used different decomposition models, because the contextonyms database is not ready yet. Since, the WordNet affect is too short for a proper usage with a pseudo-LSA, we tried also a top 1000 words (as taken from Project Gutenberg).



Dominant emotion visualization using SOM

	Precision	Recall
Anger	18.52%	15.38%
Disgust	8.33%	7.69%
Fear	9.09%	27.67%
Joy	40.49%	64.62%
Sadness	27.08%	19.60%
Surprise	22.50%	4.95%

Dominant emotion classification using SOM, with top 1000 words decomposition

7. References

- [BES10] S. Baccianella, A. Esuli, and F. Sebastiani. Sentiwordnet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In Seventh conference on International Language Resources and Evaluation, Malta. Retrieved May, volume 25, page 2010, 2010.
- [EFJ+98] P. Ekman, W.V. Friesen, JM JENKINS, K. OATLEY, and NL STEIN. Constants across cultures in the face and emotion. Human emotions. A Reader, pages 63-72, 1998.
- [JPW03] Hyungsuk Ji, Sabine Ploux, and Eric Wehri. Lexical knowledge representation with contextonyms. In Proceedings of MT Summit IX, New Orleans, USA, 2003.
- [SM08] C. Strapparava and R. Mihalcea. Learning to identify emotions in text. In Proceedings of the 2008 ACM symposium on Applied computing, pages 1556-1560. ACM, 2008.
- [LS04] H. Liu and P. Singh. ConceptNet practical commonsense reasoning tool-kit. BT technology journal, 22(4):211-226, 2004.
- [VSS05] A. Valitutti, C. Strapparava, and O. Stock. Lexical resources and semantic similarity for affective evaluative expressions generation. Affective Computing and Intelligent Interaction, pages 474-481, 2005.