

TTCS^E: a Vectorial Resource for Computing Conceptual Similarity

Enrico Mensa

University of Turin, Italy
Dipartimento di Informatica
mensa@di.unito.it

Daniele P. Radicioni

University of Turin, Italy
Dipartimento di Informatica
radicion@di.unito.it

Antonio Lieto

University of Turin, Italy
ICAR-CNR, Palermo, Italy
lieto@di.unito.it

Abstract

In this paper we introduce the TTCS^E, a linguistic resource that relies on BabelNet, NASARI and ConceptNet, that has now been used to compute the conceptual similarity between concept pairs. The conceptual representation herein provides uniform access to concepts based on BabelNet synset IDs, and consists of a vector-based semantic representation which is compliant with the Conceptual Spaces, a geometric framework for common-sense knowledge representation and reasoning. The TTCS^E has been evaluated in a preliminary experimentation on a conceptual similarity task.

1 Introduction

The development of robust and wide-coverage resources to use in different sorts of application (such as text mining, categorization, *etc.*) has known in the last few years a tremendous growth. In this paper we focus on computing conceptual similarity between pairs of concepts, based on a resource that extends and generalizes an attempt carried out in (Lieto et al., 2016a). In particular, the TTCS^E—so named after Text to Conceptual Spaces-Extended—has been acquired by integrating two different sorts of linguistic resources, such as the encyclopedic knowledge available in BabelNet (Navigli and Ponzetto, 2012) and NASARI (Camacho-Collados et al., 2015), and the common-sense grasped by ConceptNet (Speer and Havasi, 2012). The resulting representation enjoys the interesting property of being anchored to both resources, thereby providing a uniform conceptual access grounded on the sense identifiers provided by BabelNet.

Conceptual Spaces (CSs) can be thought of as

a particular class of vector representations where knowledge is represented as a set of limited though cognitively relevant quality dimensions; in this representation a geometrical structure is associated to each quality dimension, mostly based on cognitive accounts. In this setting, concepts correspond to convex regions, and regions with different geometrical properties correspond to different sorts of concepts (Gärdenfors, 2014). The geometrical features of CSs have a direct appeal for common-sense representation and common-sense reasoning, since prototypes (the most relevant representatives of a category from a cognitive point of view, see (Rosch, 1975)) correspond to the geometrical centre of a convex region, the centroid. Also exemplars-based representations can be mapped onto points in a multidimensional space, and their similarity can be computed as the distance intervening between each two points, based on some suitable metrics such as Euclidean or Manhattan distance. *etc.*

The CS framework has been recently used to extend and complement the representational and inferential power allowed by formal ontologies—and in general symbolic representation—, that are not suited for representing defeasible, prototypical knowledge, and for dealing with the corresponding typicality-based conceptual reasoning (Lieto et al., 2017). Also, wide-coverage semantic resources such as DBPedia and ConceptNet, in fact, in different cases fail to represent the sort of common-sense information based on prototypical and default information which is usually required to perform forms of plausible reasoning.¹ In this

¹ Although DBPedia contains information on many sorts of entities, due to its explicit encyclopedic commitment, common-sense information is dispersed among textual descriptions (e.g., in the abstracts) rather than being available in a well-structured formal way. For instance, the *fork* entity can be categorized as an object, whilst there is no structured information about its typical usage. On the other hand, Concept-

paper we explore whether and to what extent a linguistic resource describing concepts by means of *qualitative* and *synthetic* vectorial representation is suited to assess the conceptual similarity between pairs of concepts.

2 Vector representations with the TTCS^E

The TTCS^E has been designed to build resources encoded as general conceptual representations. We presently illustrate how the resource is built, deferring to Section 3 the description of the control strategy designed to use it in the computation of conceptual similarity.

The TTCS^E takes in input a concept c referred through a BabelNet synset ID, and produces as output a vector representation \vec{c} where the input concept is described along some semantic dimensions. In turn, filling each such dimension amounts to finding a set of appropriate concepts: features act like vector space dimensions, and they are based on ConceptNet relationships.² The dimensions are filled with BabelNet synset IDs, so that finally each concept c residing in the linguistic resource can be defined as

$$\vec{c} = \bigcup_{d \in \mathcal{D}} \{ \langle ID_d, \{c_1, \dots, c_k\} \rangle \} \quad (1)$$

where ID_d is the identifier of the d -th dimension, and $\{c_1, \dots, c_k\}$ is the set of values chosen for d .

The control strategy implemented by the TTCS^E includes two main steps, *semantic extraction* (composed by the *extraction* and *concept identification* phases) and the *vector injection*.

Net is more suited to structurally represent common-sense information related to typicality. However, in ConceptNet the coverage of this type of knowledge component is sometimes not satisfactory. For similar remarks on such resources, claiming for the need of new resources more suited to represent common-sense information, please also refer to (Basile et al., 2016).

²The full list of the employed properties, which were selected from the most salient properties in ConceptNet, includes: INSTANCEOF, RELATEDTO, ISA, ATLOCATION, DBPEDIA/GENRE, SYNONYM, DERIVEDFROM, CAUSES, USEDFOR, MOTIVATEDBYGOAL, HAS-SUBEVENT, ANTONYM, CAPABLEOF, DESIRES, CAUSESDESIRE, PARTOF, HASPROPERTY, HASPREREQUISITE, MADEOF, COMPOUNDDERIVEDFROM, HASFIRST-SUBEVENT, DBPEDIA/FIELD, DBPEDIA/KNOWNFOR, DBPEDIA/INFLUENCEDBY, DEFINEDAS, HASA, MEMBEROF, RECEIVESACTION, SIMILARTO, DBPEDIA/INFLUENCED, SYMBOLOF, HASCONTEXT, NOTDESIRES, OBSTRUCTEDBY, HASLASTSUBEVENT, NOTUSEDFOR, NOTCAPABLEOF, DESIREOF, NOTHASPROPERTY, CREATEDBY, ATTRIBUTE, ENTAILS, LOCATIONOFACTION, LOCATEDNEAR.

Extraction The TTCS^E takes in input c and builds a bag-of-concepts C including the concepts associated to c through one or more ConceptNet relationships. All ConceptNet nodes related to the input concept c are collected: namely, we take the corresponding ConceptNet node for each term in the WordNet (Miller, 1995) synset of c , $s^c \in \text{WN-syn}_c$. For each such term we extract all terms t linked through d , one of the aforementioned ConceptNet relationships: that is, we collect the terms $s^c \xrightarrow{d} t$ and store them in the set T . Each s^c can be considered as a different lexicalization for the same concept c , so that all t can be grouped in T , that finally contains all terms associated in any way to c .

Since ConceptNet does not provide any direct anchoring mechanism to associate its terms to meaning identifiers, it is necessary to determine which of the terms $t \in T$ are relevant for the concept c . In other words, when we access the ConceptNet page for a certain term, we find not only the association regarding that term with the sense conveyed by c , but also all the associations regarding it in any other meaning. To select only (and possibly all) the associations that concern the sense individuated through the previous phase, we introduce the notion of *relevance*. To give an intuition of this process, the terms found in ConceptNet are considered as relevant (and thus retained) either if they exhibit a heavy weight in the NASARI vector corresponding to the considered concept, or if they share at least some terms with the NASARI vector (further details on a similar approach can be found in (Lieto et al., 2016a)).

Concept identification Once the set of relevant terms has been extracted, we need to lift them to the corresponding concept(s), which will be used as value for the features. We stress, in fact, that dimension fillers are concepts rather than terms (please refer to Eq. 1). In the concept identification step, we exploit NASARI in order to provide each term $t \in T$ with a BabelNet synset ID, thus finally converting it into the bag-of-concepts C .

Given a $t_i \in T$, we distinguish two main cases. If t_i is contained in one or more synsets inside the NASARI vector of c , we obtain c_i (the concept underlying t_i) by directly assigning to t_i the identifier of the heaviest weighted synset that contains it.³ Otherwise, if t_i is not included in any of the

³NASARI *unified* vectors are composed by a head con-

synsets in the NASARI vector associated to c , we need to choose a vector among all possible ones: we first select a list of candidate vectors (that is, those containing t_i in their vector head), and then choose the best one by retaining the vector where c 's ID has highest weight.

For example, given in input the concept *bank* intended as a financial institution, we inspect the edges of the ConceptNet node 'bank' and its synonyms. Then, thanks to the relevance notion we get rid of associations such as 'bank ISA flight maneuver' since the term 'flight maneuver' is not present in the vector associated to the concept *bank*. Conversely, we accept sentences such as 'bank HASA branch' (i.e., 'branch' is added to T). Finally, 'branch' goes through the concept identification phase, resulting in a concept c_i and then it is added to C .

Vector injection The bag-of-concepts C is then scanned, and each value is injected in the template for \vec{c} . Each value $\{c_1, \dots, c_n \in C\}$ is still provided with the relationship that linked it to c in ConceptNet, so this value is employed to fill the corresponding feature in \vec{c} . For example, if c_k is extracted from the ConceptNet relation USED FOR (i.e., $c \xrightarrow{\text{USED FOR}} c_k$), the value c_k will be added to the set of entities that are used for c .

2.1 Building the TTCS $^{\mathcal{E}}$ resource

In order to build the set of vectors in the TTCS $^{\mathcal{E}}$ resource, the system took in input 16,782 concepts. Such concepts have been preliminarily computed (Lieto et al., 2016b) by starting from the 10K most frequent nouns present in the Corpus of Contemporary American English (COCA).⁴ Then, for each input concept the TTCS $^{\mathcal{E}}$ scans some 3M ConceptNet nodes to retrieve the terms that appear into the WordNet synset of the input. This step allows to browse over 11M associations available in ConceptNet, and to extract on average 155 ConceptNet nodes for each input concept. Subsequently, the TTCS $^{\mathcal{E}}$ exploits the 2.8M NASARI vectors to decide whether each of the extracted nodes is relevant or not w.r.t. the input concept, and then it tries to associate a NASARI vector to each of them (concept identification step). On av-

cept (represented by its ID in the first position) and a body, that is a list of synsets related to the head concept. Each synset ID is followed by a number that grasps the strength of its correlation with the head concept.

⁴<http://corpus.byu.edu/full-text/>.

erage, 14.90 concepts are used to fill each vector.⁵

3 Computing Conceptual Similarity

One main assumption underlying our approach is that two concepts are similar insofar as they share some values on the same dimension, such as when they are both used for the same ends, they share components, etc.. Consequently, our metrics does not employ WordNet taxonomy and distances between pairs of nodes, such as in (Wu and Palmer, 1994; Leacock et al., 1998; Schwartz and Gomez, 2008), nor it depends on information content accounts either, such as in (Resnik, 1998a; Jiang and Conrath, 1997).

The representation available to the TTCS $^{\mathcal{E}}$ is entirely filled with conceptual identifiers, so to assess the similarity between two such values we check whether both the concept vector \vec{c}_i and the vector \vec{c}_j share the same (concept) value for the same dimension $d \in D$, and our similarity along each dimension basically depends on this simple intuition:

$$\text{sim}(\vec{c}_i, \vec{c}_j) = \frac{1}{|D|} \cdot \sum_{d \in D} |d_i \cap d_j|.$$

The score computed by the TTCS $^{\mathcal{E}}$ system can be justified based on the dimensions actually filled: this explanation can be built automatically, since the similarity between \vec{c}_i and \vec{c}_j is a function of the sum of the number of shared elements in each dimension, so that one can argue that a given score x is due to the fact that along a given dimension d both concepts share some values (e.g., $\text{sim}(\text{table}, \text{chair}) = x$ because each one is a (ISA) 'furniture', both are USED FOR 'eating', 'studying' and 'working'; both can be found AT-LOCATION 'home', 'office'; and each one HASA 'leg').

Ultimately, the TTCS $^{\mathcal{E}}$ collects information along the 44 dimensions listed in footnote 2, so that we are in principle able to assess in how far similar they are along each and every dimension. However, our approach is presently limited by the actual average filling factor, and by the noise that can be possibly collected by an automatic procedure built on top of the BabelNet knowledge base. Since we need to deal with noisy and incomplete information, some adjustments to the above formula have been necessary in order to handle

⁵The final resource is available for download at the URL <http://ttcs.di.unito.it>.

—*intra* dimension— the possibly unbalanced number of concepts that characterize the different dimensions; and to prevent —*inter* dimensions— the computation from being biased by more richly defined concepts (i.e., those with more dimensions filled). The computation of the conceptual similarity score is thus based on the following formula:

$$\text{sim}(\vec{c}_i, \vec{c}_j) = \frac{1}{|D^*|} \cdot \sum_{d \in D} \frac{|d_i \cap d_j|}{\beta(\alpha a + (1 - \alpha)b) + |d_i \cap d_j|}$$

where $|d_i \cap d_j|$ counts the number of concepts that are used as fillers for the dimension d in the concept \vec{c}_i and \vec{c}_j , respectively; and a and b are computed as $a = \min(|d_i - d_j|, |d_j - d_i|)$, $b = \max(|d_i - d_j|, |d_j - d_i|)$; and $|D^*|$ counts the dimensions filled with at least one concept in both vectors.

This formula is known as the Symmetrical Tversky’s Ratio Model (Jimenez et al., 2013), which is a symmetrical reformulation for the Tversky’s ratio model (Tversky, 1977). It enjoys the following properties: *i*) it allows grasping the number of common traits between the two vectors (at the numerator); *ii*) it allows tuning the balance between cardinality differences (through the parameter α), and between $|d_i \cap d_j|$ and $|d_i - d_j|, |d_j - d_i|$ (through the parameter β). Interestingly, by setting $\alpha = .5$ and $\beta = 1$ the formula equals the popular Dice’s coefficient. The parameters α and β were set to $.8$ and $.2$ for the experimentation.

4 Evaluation

In the experimentation we addressed the conceptual similarity task at the *sense level*, that is the TTCS^E system has been fed with sense pairs. We considered three datasets,⁶ namely the RG, MC and WS-Sim datasets, first designed in (Rubenstein and Goodenough, 1965; Miller and Charles, 1991) and (Agirre et al., 2009), respectively. Historically, while the first two (RG and MC) datasets were originally conceived for similarity measures, the WS-Sim dataset was developed as a subset of a former dataset built by (Finkelstein et al., 2001) created to test similarity by including pairs of words related through specific relationships, such as synonymy, hyponymy, and unrelated. All senses were mapped onto WordNet 3.0.

The similarity scores computed by the TTCS^E system have been assessed through Pearsons r

⁶Publicly available at the URL <http://www.seas.upenn.edu/~hansens/conceptSim/>.

	ρ	r
RG	0.78	0.85
MC	0.77	0.80
WS-Sim	0.64	0.54

Table 1: Spearman (ρ) and Pearson (r) correlations obtained over the three datasets.

and Spearman’s ρ correlations, that are largely adopted for the conceptual similarity task (for a recent compendium of the approaches please refer to (Pilehvar and Navigli, 2015)). The former measure captures the linear correlation of two variables as their covariance divided by the product of their standard deviations, thus basically allowing to grasp differences in their values, whilst the Spearman correlation is computed as the Pearson correlation between the *rank* values of the considered variables, so it is reputed to be best suited to assess results in a similarity ranking setting where relative scores are relevant (Schwartz and Gomez, 2011; Pilehvar and Navigli, 2015).

Table 1 shows the results obtained by the system in a preliminary experimentation.

Provided that the present task of *sense-level* similarity is slightly different from *word-level* similarity (about this distinction, please refer to (Pilehvar and Navigli, 2015)), and our results can be thus hardly compared to those in literature, the reported figures are still far from the state of the art, where the Spearman correlation ρ reaches 0.92 for the RG dataset (Pilehvar and Navigli, 2015), 0.92 for the MC dataset (Agirre et al., 2009), and 0.81 for the WS-Sim dataset (Halawi et al., 2012; Tau Yih and Qazvinian, 2012).⁷

However, we remark that the TTCS^E employs vectors of a very limited size w.r.t. the standard vector-based resources used in the current models of distributional semantics (as mentioned, each vector is defined, on average, through 14.90 concepts). Moreover, due to the explicit grounding provided by connecting the NASARI feature values to the corresponding properties in ConceptNet, the TTCS^E can be used to provide the scores returned as output with an explanation, based on the shared concepts along some given dimension. At the best of our knowledge, this is a unique feature, that cannot be easily reached by methods

⁷Rich references to state-of-the-art results and works experimenting on the mentioned datasets can be found on the ACL Wiki, at the URL <https://goo.gl/NQ1b6g>.

based on Latent Semantic Analysis (such as those pioneered by (Deerwester et al., 1990)) and can be only partly approached by techniques exploiting taxonomic structures (Resnik, 1998b; Banerjee and Pedersen, 2003). Conversely, few and relevant traits are present in the final linguistic resource, which is thus *synthetic* and more *cognitively plausible* (Gärdenfors, 2014).

In some cases —27 concept pairs out of the overall 190 pairs— the system was not able to retrieve an ID for one of the concepts in the pair: such pairs were excluded from the computation of the final accuracy. Missing concepts may be lacking in (at least one of) the resources upon which the TTCS^E is built: including further resources may thus be helpful to overcome this limitation. Also, difficulties stemmed from insufficient information for the concepts at stake: this phenomenon was observed, e.g., when both concepts have been found, but no common dimension has been filled. This sort of difficulty shows that the coverage of the resource still needs to be enhanced by improving the extraction phase, so to add further concepts *per* dimension, and to fill more dimensions.

5 Conclusions

In this paper we have introduced a novel resource, the TTCS^E, which is compatible with the Conceptual Spaces framework and aims at putting together encyclopedic and common-sense knowledge. The resource has been employed to compute the conceptual similarity between concept pairs. Thanks to its representational features it allows implementing a simple though effective heuristics to assess similarity: that is, concepts are similar insofar as they share some values along the same dimension. However, further heuristics will be investigated in the next future, as well.

A preliminary experimentation has been run, employing three different datasets. Provided that we consider the obtained results as encouraging, the experimentation clearly points out that there is room for improvement along two main axes: dimensions must be filled with further information, and also the quality of the extracted information should be improved. Both aspects will be the object of our future efforts.

References

- [Agirre et al.2009] Eneko Agirre, Enrique Alfonseca, Keith Hall, Jana Kravalova, Marius Paşca, and Aitor Soroa. 2009. A study on similarity and relatedness using distributional and wordnet-based approaches. In *Procs of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 19–27. ACL.
- [Banerjee and Pedersen2003] S. Banerjee and T. Pedersen. 2003. Extended gloss overlaps as a measure of semantic relatedness. In *Proceedings of the Eighteenth International Joint Conference on Artificial Intelligence*, pages 805–810.
- [Basile et al.2016] Valerio Basile, Soufian Jebbara, Elena Cabrio, and Philipp Cimiano. 2016. Populating a knowledge base with object-location relations using distributional semantics. In Eva Blomqvist, Paolo Ciancarini, Francesco Poggi, and Fabio Vitali, editors, *EKAW*, volume 10024 of *Lecture Notes in Computer Science*, pages 34–50.
- [Camacho-Collados et al.2015] José Camacho-Collados, Mohammad Taher Pilehvar, and Roberto Navigli. 2015. NASARI: a novel approach to a semantically-aware representation of items. In *Proceedings of NAACL*, pages 567–577.
- [Deerwester et al.1990] S. Deerwester, S. Dumais, G. Furnas, T. Landauer, and R. Harshman. 1990. Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41:391–407.
- [Finkelstein et al.2001] Lev Finkelstein, Evgeniy Gabrilovich, Yossi Matias, Ehud Rivlin, Zach Solan, Gadi Wolfman, and Eytan Ruppín. 2001. Placing search in context: The concept revisited. In *Proceedings of the 10th international conference on World Wide Web*, pages 406–414. ACM.
- [Gärdenfors2014] Peter Gärdenfors. 2014. *The geometry of meaning: Semantics based on conceptual spaces*. MIT Press.
- [Halawi et al.2012] Guy Halawi, Gideon Dror, Evgeniy Gabrilovich, and Yehuda Koren. 2012. Large-scale learning of word relatedness with constraints. In Qiang Yang, Deepak Agarwal, and Jian Pei, editors, *KDD*, pages 1406–1414. ACM.
- [Jiang and Conrath1997] Jay J Jiang and David W Conrath. 1997. Semantic similarity based on corpus statistics and lexical taxonomy. *arXiv preprint cmp-lg/9709008*.
- [Jimenez et al.2013] Sergio Jimenez, Claudia Becerra, Alexander Gelbukh, Av Juan Dios Bátiz, and Av Mendizábal. 2013. Softcardinality-core: Improving text overlap with distributional measures for semantic textual similarity. In *Second Joint Conference on Lexical and Computational Semantics*, volume 1, pages 194–201.
- [Leacock et al.1998] Claudia Leacock, George A Miller, and Martin Chodorow. 1998. Using corpus statistics and wordnet relations for sense identification. *Computational Linguistics*, 24(1):147–165.

- [Lieto et al.2016a] Antonio Lieto, Enrico Mensa, and Daniele P. Radicioni. 2016a. A Resource-Driven Approach for Anchoring Linguistic Resources to Conceptual Spaces. In *Procs of the XV International Conference of the Italian Association for Artificial Intelligence*, volume 10037 of *LNAI*, pages 435–449. Springer.
- [Lieto et al.2016b] Antonio Lieto, Enrico Mensa, and Daniele P. Radicioni. 2016b. Taming sense sparsity: a common-sense approach. In *Proceedings of Third Italian Conference on Computational Linguistics (CLiC-it 2016) & Fifth Evaluation Campaign of Natural Language Processing and Speech Tools for Italian*.
- [Lieto et al.2017] Antonio Lieto, Daniele P. Radicioni, and Valentina Rho. 2017. Dual PECCS: A Cognitive System for Conceptual Representation and Categorization. *Journal of Experimental & Theoretical Artificial Intelligence*, 29(2):433–452.
- [Miller and Charles1991] George A Miller and Walter G Charles. 1991. Contextual correlates of semantic similarity. *Language and cognitive processes*, 6(1):1–28.
- [Miller1995] George A Miller. 1995. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41.
- [Navigli and Ponzetto2012] Roberto Navigli and Simone Paolo Ponzetto. 2012. BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network. *Artificial Intelligence*, 193:217–250.
- [Pilehvar and Navigli2015] Mohammad Taher Pilehvar and Roberto Navigli. 2015. From senses to texts: An all-in-one graph-based approach for measuring semantic similarity. *Artif. Intell.*, 228:95–128.
- [Resnik1998a] Philip Resnik. 1998a. Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language. *Journal of Artificial Intelligence Research*, 11(1).
- [Resnik1998b] Philip Resnik. 1998b. Semantic similarity in a taxonomy: An information-based measure and its application to problems of ambiguity in natural language. *Journal of Artificial Intelligence Research*, 11(1).
- [Rosch1975] Eleanor Rosch. 1975. Cognitive Representations of Semantic Categories. *Journal of experimental psychology: General*, 104(3):192–233.
- [Rubenstein and Goodenough1965] Herbert Rubenstein and John B Goodenough. 1965. Contextual correlates of synonymy. *Communications of the ACM*, 8(10):627–633.
- [Schwartz and Gomez2008] Hansen A Schwartz and Fernando Gomez. 2008. Acquiring knowledge from the web to be used as selectors for noun sense disambiguation. In *Procs of the Twelfth Conference on Computational Natural Language Learning*, pages 105–112. ACL.
- [Schwartz and Gomez2011] Hansen A Schwartz and Fernando Gomez. 2011. Evaluating semantic metrics on tasks of concept similarity. In *Proc. Int. Florida Artif. Intell. Res. Soc. Conf.(FLAIRS)*, page 324.
- [Speer and Havasi2012] Robert Speer and Catherine Havasi. 2012. Representing General Relational Knowledge in ConceptNet 5. In *LREC*, pages 3679–3686.
- [Tau Yih and Qazvinian2012] Wen Tau Yih and Vahed Qazvinian. 2012. Measuring word relatedness using heterogeneous vector space models. In *HLT-NAACL*, pages 616–620. The Association for Computational Linguistics.
- [Tversky1977] Amos Tversky. 1977. Features of similarity. *Psychological review*, 84(4):327.
- [Wu and Palmer1994] Zhibiao Wu and Martha Palmer. 1994. Verbs semantics and lexical selection. In *Proceedings of the 32nd annual meeting on Association for Computational Linguistics*, pages 133–138. ACL.