



di.unito.it

DIPARTIMENTO
DI INFORMATICA

TLN-LAB

Automatic summarisation with NASARI

Daniele Radicioni

NASARI



credits

- the following slides have been built based on

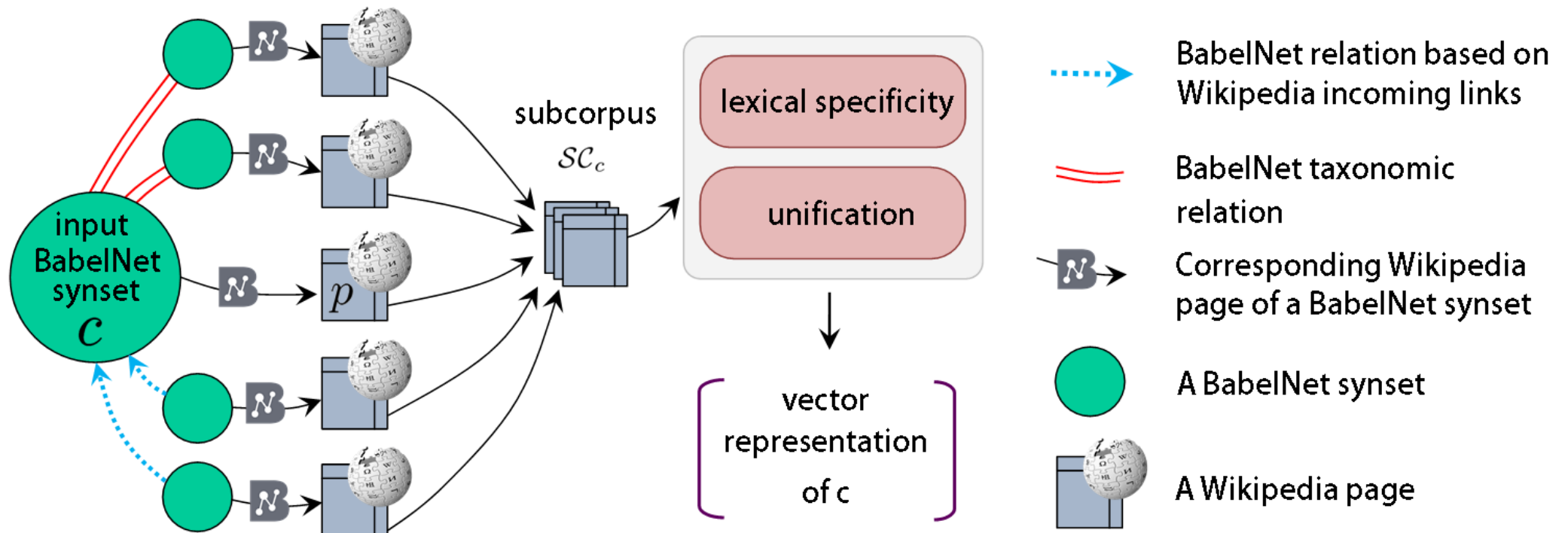
José Camacho-Collados, Mohammad Taher Pilehvar and Roberto Navigli. NASARI: a Novel Approach to a Semantically-Aware Representation of Items. In *Proceedings of the North American Chapter of the Association of Computational Linguistics (NAACL 2015)*, Denver, USA, pp. 567-577, 2015

José Camacho-Collados, Mohammad Taher Pilehvar and Roberto Navigli. A Unified Multilingual Semantic Representation of Concepts. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics (ACL 2015)*, Beijing, China, July 27-29, pp. 741-751, 2015

NASARI

- A Novel Approach to a Semantically-Aware Representation of Items
 - The prevailing methods for the computation of a **vector space representation** are **based on distributional semantics**.
 - These approaches are **unable to model individual word senses** or concepts, since they **conflate different meanings of a word into a single vectorial representation**.
- Chen et al. (2014) addressed this issue and obtained vectors for individual word senses by leveraging WordNet glosses.

MUFFIN (Multilingual, UniFied and Flexible INterpretation)



- For the given synset the contextual information is gathered from Wikipedia by exploiting knowledge from the BabelNet semantic network.
- Then, by analysing the corresponding contextual information and comparing and contrasting it with the whole Wikipedia corpus, a vectorial representation of the given synset is obtained.

two sorts of vectors

- On the basis of lexical specificity two types of representations are built: lexical and unified.
 - The *lexical vector representation* lex_c of a concept c has lemmas as its individual dimensions.
 - The *unified representation* has concepts as individual dimensions.

Unified representation

- The algorithm first clusters together those words that have a sense sharing the same hypernym according to the BabelNet taxonomy
- Next, the specificity is computed for the set of all the hyponyms, even those that do not appear in the sub-corpus SC_c ;
- The binding of a set of sibling words into a single cluster represented by their common hypernym
 - transforms the representations to a unified semantic space;
 - allows to see the clustering as an implicit disambiguation process.

Crane (bird)			Crane (machine)		
English	French	German	English	French	German
shore_bird _n ¹	‡famille_des_oiseaux _n ¹	‡vogel-familie _n ¹	*lifting device _n ¹	*dispositif de levage _n ¹	*hebevorrichtung _n ¹
bird _n ¹	*limicole _n ¹	*charadrii _n ¹	‡construction _n ⁴	navire _n ¹	radfahrzeug _n ¹
*wading_bird _n ¹	oiseau_aquatique _n ²	†vogel_gattung _n ¹	platform _n ¹	limicole _n ¹	†lenkfahrzeug _n ¹
oscine_bird _n ¹	tollé _n ²	wirbeltiere _n ²	warship _n ¹	◇vaisseau _n ²	regler _n ³
†bird_genus _n ¹	gallinacé _n ¹	fleisch _n ¹	electric circuit _n ¹	spationef _n ¹	reisebus _n ¹
‡bird_family _n ¹	◇classe _n ¹	tier um _n ¹	◇vessel _n ²	‡construction _n ²	charadrii _n ¹
◇taxonomic_group _n ¹	occurence _n ¹	reiherr _n ¹	boat _n ¹	†véhicule _n ³	güterwagen _n ²

$word^p_n$ is the p^{th} sense of the word with part of speech n .

Word senses marked with the same symbol across languages correspond to the same BabelNet synset.

Set of concepts associated to words

- Given these representations for individual word senses, the goal is **to associate the set of concepts**, i.e., BabelNet synsets, $C_w = \{c_1, \dots, c_n\}$ **with a given word w** .
 - If w exists in the BabelNet dictionary, the set of associated senses of the word can be obtained as defined in the **BabelNet sense inventory**.
 - Use of *piped links*. Piped link is a hyperlink appearing in the body of a Wikipedia article, providing a link to another Wikipedia article, such as `[[dockside_crane|Crane_(machine)]]` is a **hyperlink that appears as *dockside_crane* in the text, but takes the user to the Wikipedia page titled *Crane_(machine)***.
- In so doing, a set of concepts for the words not covered by BabelNet can be obtained.

Application: Semantic Similarity

- Once we have the set C_w of concepts associated with each word w , we first retrieve the set of corresponding unified vector representations.
- Then, the square-rooted Weighted Overlap (Pilehvar et al., 2013) as vector comparison method can be used,

$$WO(v_1, v_2) = \frac{\sum_{q \in O} \left(rank(q, v_1) + rank(q, v_2) \right)^{-1}}{\sum_{i=1}^{|O|} (2i)^{-1}}$$

where O is the set of overlapping dimensions between the two vectors and $rank(q, v_i)$ is the rank of dimension q in the vector v_i . The similarity between words w_1 and w_2 is computed as the similarity of their closest senses

$$sim(w_1, w_2) = \max_{v_1 \in C_{w_1}, v_2 \in C_{w_2}} \sqrt{WO(v_1, v_2)}$$

NASARI taste & see

Each line in the vectors corresponds to a BabelNet synset. In the cases where the BabelNet synset is associated with a Wikipedia page, the Wikipedia page title is written in the second column. Otherwise it is written -NA-.

1. [Lexical vectors] the dimensions correspond to lemmas.

Files: NASARI_lexical_*.txt

Format (TAB separated):

BabelSynsetId WikipediaPageTitle lemma1_weight1 lemma2_weight2

...

2. [Unified vectors] the dimensions correspond to the BabelNet synsets.

Files: NASARI_unified_*.txt

Format (TAB separated):

BabelSynsetId WikipediaPageTitle synset1_weight1

synset2_weight2 ...

The dimensions of lexical and unified vectors are separated from the weights using an underscore. Also, the vectors are truncated to the non-zero dimensions only and sorted according to the weights of their dimensions.

NASARI taste & see

3. [Embed vectors] are embedded vector representations of 300 dimensions:

Files: NASARI_embed_*.txt

Format (SPACE separated):

```
BabelSynsetId  WikipediaPageTitle  dimension1  
dimension2 ... dimension300
```

Continuous vector representations (NASARI_embed) of BabelNet synsets constructed by combining lexical vectors and the pre-trained models of Word2Vec (300 dimensions).

I vettori di NASARI sono disponibili all'URL
<http://lcl.uniroma1.it/nasari/>

automatic summarization



credits

- E. Hovy, Chapter *Text Summarization*, in R. Mitkov (Ed.), *The Oxford handbook of computational linguistics*, Oxford University Press, 2005
- D. Jurafsky and J. H. Martin, *SPEECH and LANGUAGE PROCESSING, An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*, Prentice Hall, 2009
- Eduard Hovy and Daniel Marcu, *ACL Tutorial on Text Summarization*, ACL 1998, Université de Montréal Montréal, Québec, Canada

a definition

- The goal of text summarization is to produce an **abridged version of a text** which contains the important or **relevant** information.
 - an **abstract** of a scientific article, a **summary** of email threads, a **headline** for a news article, or the short **snippets** returned by web search engines to describe each retrieved document.

goals

- **Indicative**: give an idea of what is there, provides a reference function for selecting documents for more in-depth reading
- **Informative**: a substitute for the entire document, covers all the salient information in the source at some level of detail
- **Critical**: evaluates the subject matter of the source, expressing the abstractor's view on the quality of the work of the author

kinds of automatic summarization

- **Extracts** are summaries created by reusing portions (words, sentences, etc.) of the input text verbatim, while
- **Abstracts** are created by re-generating the extracted content
 - Paraphrase, generation

kinds of automatic summarization

- Output: **User-focused** (or topic-focused or query focused): summaries that are tailored to the requirements of a particular user or group of users
- Background: Does the reader have the needed **prior knowledge**?
 - Expert reader vs. Novice reader
- General: summaries aimed at a particular –usually broad –**readership community**

Summarisation approaches

- Shallow approaches
 - Syntactic level at most
 - Typically produce extracts
 - Extract salient parts of the source text and then arrange and present them in some effective manner
- Deeper approaches
 - Sentential semantic level
 - Produce abstracts and the synthesis phase involves natural language generation.
 - Knowledge-intensive, may require some domain specific coding

single document versus multiple document summarisation

- In [single document summarisation](#) we are given a single document and produce a summary.
 - Single document summarisation is thus used in situations like [producing a headline or an outline](#), where the final goal is to characterise the content of a single document.
- In [multiple document summarisation](#), the input is a group of documents, and our goal is to produce a condensation of the content of the entire group.
 - We might use multiple document summarisation when we are summarising [a series of news stories on the same event](#), or whenever we have web content on the same topic that we'd like to synthesise and condense.

parameters

- **Compression rate** (summary length/source length)
- **Audience** (user-focused vs. generic)
- Relation to source (extract vs. abstract)
- **Function** (indicative vs. informative vs. critical)
- **Coherence**: the way the parts of the text gather together to form an integrated whole
 - Coherent vs. incoherent
 - Incoherent: unresolved **anaphors**, **gaps in the reasoning**, sentences which repeat the same or similar meaning (redundancy) a lack of organisation

approaches comparison

- NLP/IE:

- Approach: try to 'understand' text—re-represent content using 'deeper' notation; then manipulate that.
- Need: rules for text analysis and manipulation, at all levels.
- Strengths: higher quality; supports abstracting.
- Weaknesses: speed; still needs to scale up to robust open-domain summarisation.

- IR/Statistics:

- Approach: operate at lexical level—use word frequency, collocation counts, etc.
- Need: large amounts of text.
- Strengths: robust; good for query-oriented summaries.
- Weaknesses: lower quality; inability to manipulate information at abstract levels.

relevance criteria



Position in the text

- Important sentences occur in specific positions
 - “*lead-based*” summary (just take first sentence(s)!)
 - Important information occurs in specific sections of the document (introduction/conclusion)
 - Experiments:
 - In 85% of 200 individual paragraphs the topic sentences occurred in initial position and in 7% in final position

Title method

- Title of document indicates its content
 - Not true for novels usually
 - What about blogs ...?
- Words in title help find relevant content
 - Create a [list of title words](#), remove “stop words”
 - [Use those as keywords](#) in order to find important sentences

Optimum Position Policy (OPP)

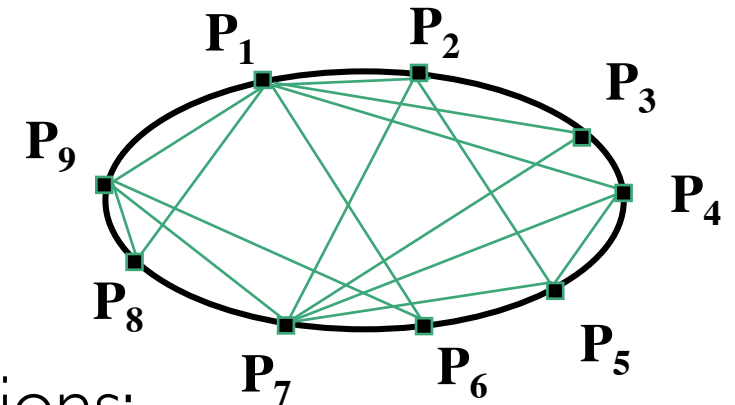
- Relevant sentences are located at positions that are genre-dependent; these positions can be either known or determined automatically through training
 - Step 1: For each article, determine the overlap between sentences and the index terms (e.g., title terms)
 - Step 2: Determine a partial ordering over the locations where sentences containing important words occur: Optimal Position Policy (OPP)

Cue phrases method

- Important sentences contain cue words/indicative phrases,
 - “The main aim of the present paper is to describe...”
 - “The purpose of this article is to review...”
 - “In this report, we outline...”
 - “Our investigation has shown that...”
- Some words are considered **bonus** others **stigma**
 - bonus: comparatives, superlatives, conclusive expressions, etc.
 - stigma: negatives, pronouns, etc. non-important sentences contain ‘stigma phrases’ such as hardly and impossible.
- These phrases can be detected automatically
- Method: Add to sentence score if it contains a bonus phrase, penalise if it contains a stigma phrase.

Cohesion-based methods

- Important sentences/paragraphs are the highest connected entities in more or less elaborate semantic structures.
- Classes of approaches
 - word co-occurrences;
 - local salience and grammatical relations;
 - co-reference;
 - lexical similarity (WordNet, lexical chains);
 - combinations of the above.



Cohesion: word co-occurrence

- Apply IR methods at the document level: texts are collections of paragraphs
 - Use a traditional, IR-based, word similarity measure to determine for each paragraph P_i the set S_i of paragraphs that P_i is related to.
- Method:
 - determine relatedness score S_i for each paragraph,
 - extract paragraphs with largest S_i scores.

3 (to 1) steps

- Text summarisation systems are generally described by their solutions to the following three problems:
 - *Content Selection*: What information to select from the document(s) we are summarising. We usually make the simplifying assumption that the granularity of extraction is the sentence or clause. Content selection thus mainly consists of choosing which sentences or clauses to extract into the summary.
 - *Information Ordering*: How to order and structure the extracted units.
 - *Sentence Realisation*: What kind of clean up to perform on the extracted units so they are fluent in their new context.

unsupervised algorithm

- The simplest unsupervised algorithm is to select sentences that have more salient or informative words.
 - Sentences that contain more informative words tend to be more extract-worthy.
- *Saliency* is usually defined by computing the topic signature, a set of salient or signature terms, each of whose saliency scores is greater than some threshold θ .
 - Saliency could be measured in terms of simple word frequency, but frequency has the problem that a word might have a high probability in English in general but not be particularly topical to a particular document.
- *Lexical specificity* can thus be adopted in order to individuate the most salient terms, and to score the sentences where they appear.

a simple *extractive* algorithm

- reduce the document size of e.g., 10%, 20%, 30%
- 1. **individuate the topic** of the text being summarised; the topic can be referred to as a (set of) NASARI vector(s):
$$v_{t1} = \{term_1_score, term_2_score, \dots, term_{10_score} \}$$
$$v_{t2} = \{term_1_score, term_2_score, \dots, term_{10_score} \}$$
$$\dots$$
- 2. **create the context**, by collecting the vectors of terms herein (this step can be repeated, by dumping the contribution of the associated terms at each round);
- 3. **retain paragraphs whose sentences contain the most salient terms**, based on the Weighted Overlap, $WO(v_1, v_2)$
- rerank paragraphs weight by applying at least one of the mentioned approaches (*title, cue, phrase, cohesion*).

NASARI (lexical) subset

- two distribution files are provided for NASARI, that require different resources allocation.
 - [*dd-nasari.txt*](#). a subset of NASARI (obtained by truncating vectors at 10 features). 3,587,754 vectors, ~600MB;
<https://goo.gl/85BubW>
 - [*dd-small-nasari-15.txt*](#). a subset of NASARI. same filtering as above, with 15 features + intersection with 60K lemmas in the Corpus of Contemporary American English: 13,084 vectors, 2MB storage (many entities removed here...).
- the second one has been extracted for starting our experimentation; the second one is intended to explore the resource in a richer (though reduced) flavour.

documents for summarisation

- text documents are provided for summarisation purposes:
 - *Andy-Warhol.txt*
 - *Ebola-virus-disease.txt*
 - *Life-indoors.txt*
 - *Napoleon-wiki.txt*
- do experiment with different compression rates: 10%, 20% and 30%.

