

## Lecture 14: Distributed Word Representations for Information Retrieval

### Section 9.2.2 Query Expansion: How can we robustly match a user's search intent?

- *Synonymy*: In most collections, the same concept may be referred to using different words.
  - Has an impact on the *recall* of most IR systems
  - Users often attempt to manually refine their queries
  - How could an IR system help with query refinement?
  - We want to understand a query, rather than simply matching keywords.  
We want to better understand *when* query and documents match
- *Query expansion*: Users give *additional* input on query words or phrases, possibly suggesting additional query terms
  - The users opting to use one of the alternative query suggestions
- How to generate alternative or expanded queries for the user?
  - Global analysis: For each term in a query, automatically expand the query using synonyms and related words from thesaurus
  - Local analysis: Analyze the documents in the current results set
    - \* Feedback on documents or on query terms
- How to build a thesaurus?
  - Use of a controlled vocabulary maintained by human editors: Canonical terms for each concepts
  - Manual thesaurus: Synonymous names for concepts, without designating a canonical term.
  - Automatically derived thesaurus
    - \* Using word co-occurrence statistics: words co-occurring in a document or paragraph are likely to be in some sense *similar or related in meaning*
    - \* Exploiting grammatical relations or dependencies: less robust than co-occurrence statistics, but more accurate
    - \* Quality of the resulting terms often not so good
    - \* Since those terms are highly correlated in documents anyway, this method may not retrieve that many additional documents.
  - Query reformulations based on query log mining: Exploit the *manual query reformulations* of other users
- Use of query expansion generally increases recall
  - A domain-specific thesaurus is required.
  - May significantly decrease precision, particularly when the query contains ambiguous terms.

### How can we represent term relations?

- Under the standard symbolic encoding of terms, different terms have no direct way of representing their similarities.

- Basic IR is scoring on  $q^T d$ . Can we learn parameters  $W$  to rank via  $q^T W d$ ?
  - Berger and Lafferty 1999, Query translation model
  - $W$  is huge ( $> 10^{10}$ ): Sparsity is the problem
- We could learn a dense low-dimensional representation of a word in  $\mathcal{R}^d$ , such that dot products  $u^T v$  express word similarity.
- Supervised Semantic Indexing shows successful use of learning  $W$
- This lecture will however consider *direct similarity*