# Data Mining and Machine Learning

## Lecture 5
## Query Ex

Peter Jančovič

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Objectives

- To understand how the use of semantic relationships between words can improve the performance of a text IR system

  - Query expa https://eduassistpro.github.io/
  - Generalisat
  - Synonyms, hypernyms & h
  - WordNet

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Query Processing

- Remember how we previously processed a query:
- Example:
  - "I need information on distance running" <span style="color:red">Assignment Project Exam Help</span>
- Stop word re<span style="color:red">https://eduassistpro.github.io/</span>
  - information
- Stemming <span style="color:red">Add WeChat edu_assist_pro</span>
  - information, distance, run
- But what about:
  - "The London marathon will take place…"

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Query Expansion

- Add terms to the query to increase the overlap between it and potentially relevant documents...

Assignment Project Exam Help

- ...but not irrelevant documents

- Two approa https://eduassistpro.github.io/

  – User feedback

  Add WeChat edu_assist_pro

  – Linguistic knowledge

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Feedback-based Query Expansion

- User provides feedback on the results of retrieval
  - Which of the returned documents are particularly relevant
  - Which are irrelevant

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Query reformulation

- Revise the query in response to the user feedback
  - Query expansion: Find terms in the 'relevant' documents that are not in the query. Add them to the query (of maybe just those w
  - Term rewei                                    t of query terms in relevant docume                     ght of query terms in irrelevant docume                     mple

  $$w_{td} = \boxed{\lambda} \times f_{td} \times IDF(t)$$

  - Various methods for determining $\lambda$ have been proposed

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Knowledge-Based Query Expansion

- Recall:
  - $q$ = "I need information on distance running"
  - $d$ = "The London marathon will take place..."
- We know th                                    tween
  - run, distanc
- Words with the same mea                         onyms
- If a $q$ contains $w_1$ and $w_2$ is a synonym of $w_1$, then add $w_2$ to $q$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Thesaurus

- A thesaurus is a 'dictionary' of synonyms and semantically related words and phrases

- E.G: Roget's Thesaurus

- Example: physician

```
syn: || croaker,               tor, MD,
medical, mediciner, medico ||
rel: medic, general practitioner,
surgeon
```

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Peter Mark Roget 1779 –1869

- Born London 1779

- Founder of the Royal Society of Medicine

- Invented the log-log slide rule

- Professor of P                                    titution, 1834

- Retired 1840

- Roget's *Thesaurus of English                   hrases Classified and Arranged so as to Facilitate the Expression of Ideas and Assist in Literary Composition* appeared in 1852.

- Died 1869. Buried St James' Church, West Malvern, Worcestershire.

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Hyponyms

- Not only synonyms are useful for query expansion

- Query $q$ = "Tell me about England"

- Document $d$ should be on everyone's it

- 'London' is a hyponym of

- Hyponym ~ subordinate ~ subset

- If a query $q$ contains a word $w_1$ and $w_2$ is a hyponym of $w_1$, then $w_2$ should be added to $q$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Hypernyms

- Hypernyms are also useful for query expansion

- Query $q$ = "Tell me about England"

- Document $d$ <span style="color:blue">e British Isles"</span>

- 'British Isles<span style="color:red">https://eduassistpro.github.io/</span>gland

- Hypernym ~ generalisation <span style="color:red">Add WeChat edu_assist_pro</span>

- If a query $q$ contains a word $w_1$ and $w_2$ is a hypernym of $w_1$, then $w_2$ should be added to $q$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# WordNet

- Online lexical database for the English Language

- http://www.cogsci.princeton.edu/~wn

Assignment Project Exam Help

| Category | | Meanings (syn sets) |
|---|---|---|
| Nouns | 57,00 | 48,800 |
| Adjectives | 19,500 | 10,000 |
| Verbs | 21,000 | 8,400 |

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

See Belew, chapter 6

UNIVERSITY OF
BIRMINGHAM

# WordNet

- Organised as a set of hierarchical trees

- For example, 25 trees for nouns

- 'Children' of hyponyms

- Words become more specific as you move deeper into the tree

British Isles    (Hypernym)

London    Birmingham    (Hyponym)

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

| Noun Categories | |
|---|---|
| act, action, activity | natural object |
| animal, fauna | natural phenomenon |
| artefact | person, human being |
| attribute, property | plant, flora |
| body, corpus | possession |
| cognition, | s |
| comm | n
| event, happening | |
| feeling, emotion | |
| food | state, condition |
| group, collection | substance |
| location, place | time |
| motive | |

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Query-document scoring

- A query *q* is expanded to include hyponyms and synonyms
- Recall that for a document *d*

$$w_{td} = f_{td} \cdot IDF(t)$$

$$Sim(q,d) = \frac{\sum\limits_{t \in q \cap d} w_{td} \cdot w}{\|d\| \cdot \|q\|}$$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Query expansion

- Suppose:
  - $t$ is the original term in the query,
  - $t'$ is a synonym or hyponym of $t$ which occurs in $d$

- Then we cou<inline>https://eduassistpro.github.io/</inline>

$$w_{t'd} = \lambda_{tt'} \times f_{t'd} \times IDF(t) \quad 0 \leq \lambda_{tt'}$$

- Where $\lambda_{tt'}$ is a weighting depending on how 'far' $t$ and $t'$ are apart according to WordNet ($\lambda_{tt}=1$)

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example

- Query *q* is:
  - *Is the Dark Knight on at the town cinema?*
  - *q* becomes: *dark knight town cinema*

- Document *d*
  - *The latest Batman movie pl          ed crusader in a dark urban environment*
  - *d* becomes: *late batman move cape crusade dark urban environment*

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example (continued)

- In the similarity calculation, $q \cap d = \{dark\}$

- But:

  - *move* and *cinema* are synonyms (compare "go to the cinema" wi
  - *crusader* is $k$
  - *urban* is a hypernym of *tow*

- Therefore, after query expansion,

  $q \cap d = \{dark, move (syn(cinema)), crusade(hypo(knight)), urban(hyper(town))\}$

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example (continued)

- So, if $\lambda = 1$, $\lambda_{syn} = 0.8$, $\lambda_{hypo} = 0.5$ and $\lambda_{hyper} = 0.3$, then the numerator in the calculation of *sim(q,d)* becomes <span style="color:red">Assignment Project Exam Help</span>

$$w$$ <span style="color:red">https://eduassistpro.github.io/</span>

$$+ \ 0.8*w_{movie,d} \ * \ w_{cinema}$$ <span style="color:red">Add WeChat edu_assist_pro</span>

$$+ \ 0.5*w_{crusader,d} \ * \ w_{knight,q}$$

$$+ \ 0.3*w_{urban,d} \ * \ w_{town,q}$$

Note: this is just a 'made up' example. I haven't consulted WordNet for synonym, hyponym or hypernym information and the weights $\lambda$ are just for illustration

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example (continued)

- The drawback of query expansion is that as well as increasing the overlap between a query *q* and a *relevant* document *d*, it may also increase the overlap with an *irrelevant* doc

- Consider:

- *The crusades were a dark peri〔...〕tory when knights moved from across Europe to j〔...〕s to the holy land*

- This becomes: *crusade dark period history knight move europe crusade holy land*

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Example (continued)

- In this case

    $q \cap d = \{dark, knight, move\ (syn(cinema)),$

    $2 \times crusade(hypo(knight)),$

    $urban(hyper(town)), land(hyper(town))\}$

- This document imilarity than the previous one

- So, the challenge is:

    - Expand queries *enough* to promote overlap with relevant documents...

    - ...but not so much that they overlap with irrelevant documents

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Summary

- Query expansion
  - Feedback-based
  - Knowledge-based: Synonyms, hyponyms and hypernyms
- Goal is to in https://eduassistpro.github.io/ relevant documents
- WordNet
- Generalization
- Example "toy" calculation

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM