# NICF -TEXT ANALYTICS

## MODULE 8: LINGUISTIC RESOURCES TO IMPROVE CONCEPTUALIZATION

**Dr. Fan Zhenzhen**

**Institute of Systems Science**

**National University of Singapore**

**Email: zhenzhen@nus.edu.sg**

# Objectives

At the end of this module, you can

- Identify common text analytics artifacts or resources

- Develop such artifacts/resources based on domain knowledge

# Outline

- Linguistic/knowledge resources and their roles in text analytics

- Dictionaries
  - General dictionaries
  - Synonym dictionaries
  - WordNet
  - Sentiment/Opinion Lexicon

- Defining patterns using regular expressions

# Linguistic Resources

- ## In machine readable form

  - **Dictionaries** - valid terms, POS information, list of stop words, or words to be filtered

  - **Terminologies** – special domain words and phrases

  - **Patterns/rules** – for information extraction

# Dictionaries

- Text analytics systems may be equipped with dictionaries in different languages for various purposes.
    - Terminology dictionaries for special domains or tasks
        - e.g. Biomedical domain
        - Customer Relation Management
        - IT
        - Market Intelligence
        - Opinions Mining, etc.

- ## A list of valid terms

  - Only terms in the dictionary appear in the document term vector or matrix.

  - It helps to restrict the dimension of the matrix a priori and to focus on specific terms for distinct text mining contexts.

- ## It may include useful information such as POS

# Filter Dictionary

- ## A list of invalid terms

  - stopword list (not complete):

| | | | |
|---|---|---|---|
| *a* | *an* | *because* | *before* |
| *about* | *and* | *been* | *being* |
| *above* | *any* | *before* | *below* |
| *after* | *are* | *being* | *between* |
| *again* | *aren't* | *below* | *both* |
| *against* | *as* | *between* | *but* |
| *all* | *at* | *both* | *by* |
| *am* | *be* | *been* | *...* |

From *http://www.ranks.nl/resources/stopwords.html*

# Synonym Dictionaries

- Typically for known synonyms, user-defined synonyms

> **dislike**, detest

- Can be used as a hard way to deal with inflections if no stemmer is used

> **like**, likes, liked

# Knowledge Resources

- Taxonomy and ontology – a hierarchical conceptual model to map terms to concepts

- Prerequisite for advance text mining, together with terminology lexicon

An ontology identifies and distinguishes concepts and their relationships; it describes content and relationships.
A taxonomy formalizes the hierarchical relationships among concepts and specifies the term to be used to refer to each; it prescribes structure and terminology

- A large lexical database of English

- Created and maintained by the Cognitive Science Laboratory of Princeton University

- *Nouns*, *verbs*, *adjectives* and *adverbs* are grouped into sets of cognitive synonyms (*synsets*), each expressing a distinct concept

**Number of words, synsets, and senses**

| POS | Unique Strings | Synsets | Total Word-Sense Pairs |
|---|---|---|---|
| Noun | 117798 | 82115 | 146312 |
| Verb | 11529 | 13767 | 25047 |
| Adjective | 21479 | 18156 | 30002 |
| Adverb | 4481 | 3621 | 5580 |
| Totals | 155287 | 117659 | 206941 |

Statistics from WordNet website
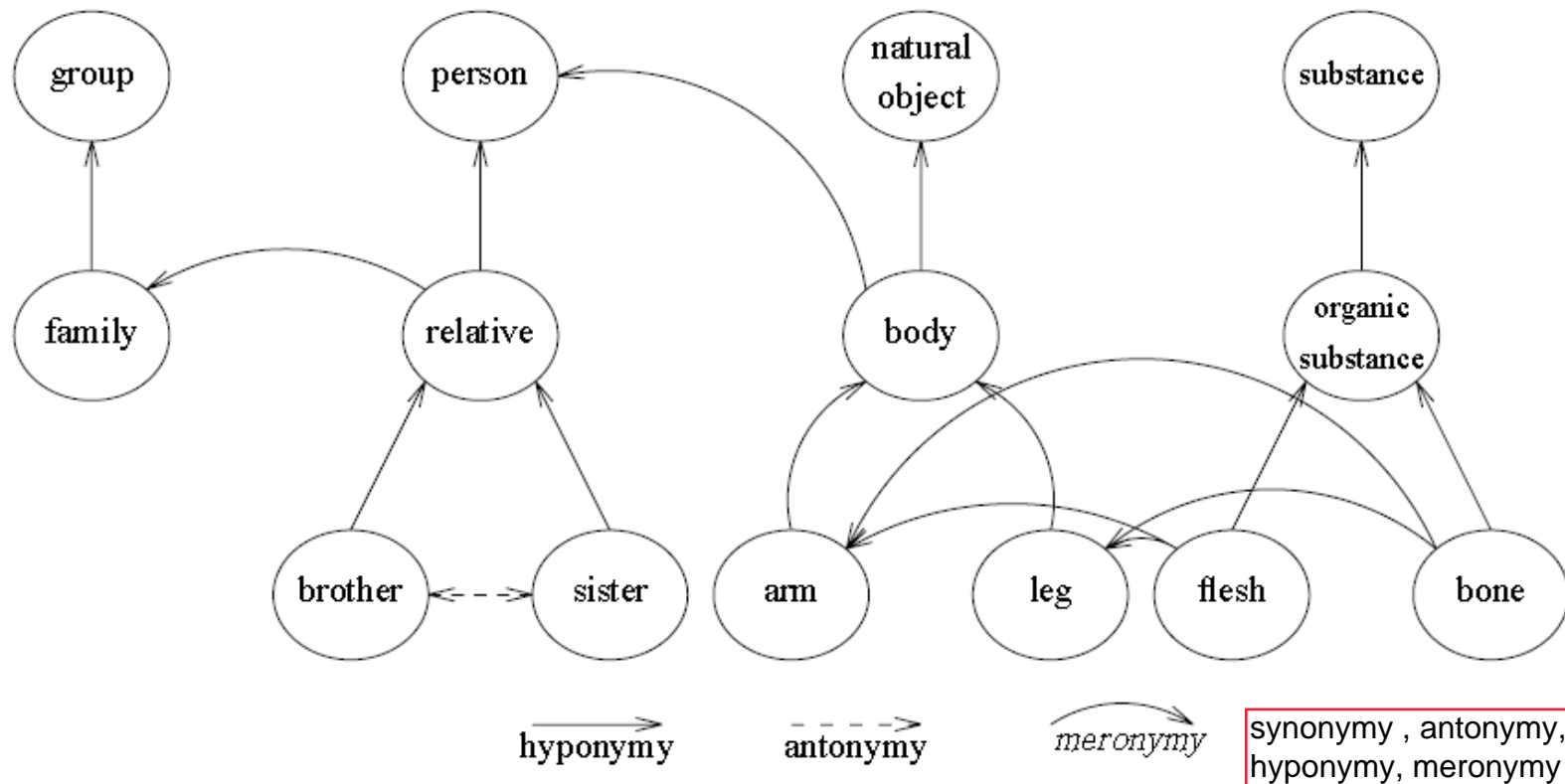http://wordnet.princeton.edu/wordnet/man/wnstats.7WN.html

# WordNet

- Synsets are linked by conceptual-semantic and lexical relations
  - Lexical relations
    - Synonymy – e.g. *shut* and *close, happy* and *joyful*
    - Antonymy – e.g. *wet* and *dry*, *young* and *old*, *happy* and *sad*
    - Morphological relations
  - Semantic relations
    - Hyponymy (or ISA relation, super-subordinate relation) – e.g. *apple* and *fruit*, *bed* and *furniture*, *communicate* and *talk* and *whisper*
    - Meronymy (part-whole relation) – e.g. *leg* and *chair*
  - And more…

Semantics is a branch of linguistics that looks at the meanings of words and language, including the symbolic use of language. It also refers to the multiple meanings of words as well.
Examples of Semantics: A toy block could be called a block, a cube, a toy.

# WordNet

Lexical representations, or rather more technically, lexical concepts, represent the semantic pole of linguistic units, and are the mentally- instantiated abstractions which language users derive from conceptions and the specific semantic contribution perceived to be associated with particular forms.

Figure 2. Network representation of three semantic relations among an illustrative variety of lexical concepts



hyponymy    antonymy    meronymy    synonymy , antonymy, hyponymy, meronymy

同义词，反义词，下位词，代名词

From *Nouns in WordNet: A Lexical Inheritance System*

- Example information in Wordnet for "happy":

**Adjective**

- (37)<u>S:</u> (adj) **happy#1** (enjoying or showing or marked by joy or pleasure)
- (2)<u>S:</u> (adj) <u>felicitous#2</u>, **happy#2** (marked by good fortune)
- <u>S:</u> (adj) <u>glad#2</u>, **happy#3** (eagerly disposed to act or to be of service)
- <u>S:</u> (adj) **happy#4**, <u>well-chosen#1</u> (well expressed and to the point)

- Expanded view:

- (37)<u>S:</u> (adj) **happy#1** (enjoying or showing or marked by joy or pleasure)
  - *see also*
  - *similar to*
    - <u>S:</u> (adj) <u>blessed#6</u> (characterized by happiness and good fortune)
    - <u>S:</u> (adj) <u>blissful#1</u> (completely happy and contented)
    - <u>S:</u> (adj) <u>bright#9</u> (characterized by happiness or gladness)
    - <u>S:</u> (adj) <u>golden#2</u>, <u>halcyon#2</u>, <u>prosperous#3</u> (marked by peace and prosperity)
    - <u>S:</u> (adj) <u>laughing#1</u>, <u>riant#1</u> (showing or feeling mirth or pleasure or happiness)
  - *attribute*
  - *antonym*
    - <u>W:</u> (adj) <u>unhappy#1</u> [Opposed to: <u>happy</u>] (experiencing or marked by or causing sadness or sorrow or discontent)

# WordNet

- Free and open source

- Proved useful for a wide range of Natural Language Processing applications

  - Word sense disambiguation

  - Word semantic distance measuring

  - Mono- and cross-lingual Information retrieval,

  - Question-answering systems

  - Machine translation

  - Document structuring and categorisation

- *Sentiment words, opinion words, polar words,* or *opinion-bearing words*.

- Lexicons or dictionaries of words or phrases that convey *positive* or *negative* sentiments, for example:

> *beautiful, wonderful, amazing…*
> *bad, poor, awful…*

# Defining Patterns using Regular Expressions

# Defining patterns/rules

- With regular expression, we can extract strings containing certain characters, or not containing certain characters, or strings with pre-specified patterns of letters or numbers.

- Such patterns can be defined in a very compact way

  - E.g. regular expression for email addresses

  [A-Z0-9._-]+@([A-Z0-9.-]+\.)+[A-Z]{2,4}

  - Strings matching this expression can then be extracted

    - E.g. *zhenzhen@nus.edu.sg*

Regular expressions are very useful in extracting concepts expressed in a certain way, e.g. *currency*, *dates*, *e-mail addresses, phone numbers*, etc.

# Common Operators

- Special characters (operators) are used to define character patterns

| Operator | Purpose |
|---|---|
| . (period) | Match any single character<br>E.g. .in matches both **Win**dows, and **Lin**ux |
| ^ | Match the empty string that occurs at the beginning of a line or string<br>E.g. ^tre will not match **stre**tch |
| $ | Match the empty string that occurs at the end of a line |
| \d | Match any single digit |
| \D | Match any single non-digit character |
| \w | Match any single alphanumeric character |

# Common Operators

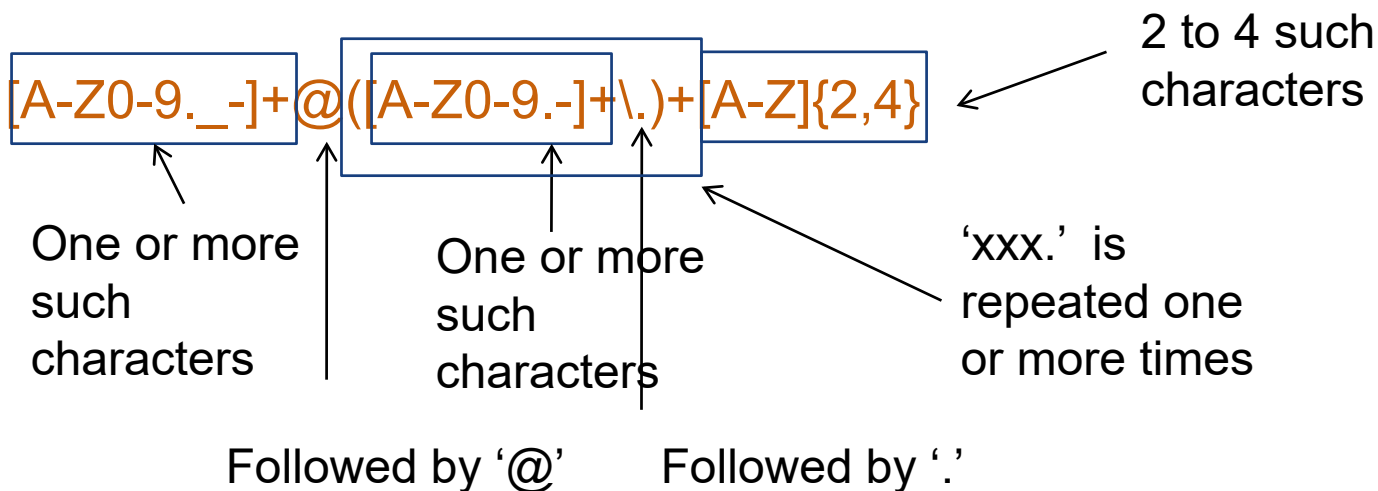| Operator | Purpose |
|---|---|
| ? | Match the preceding character 0 or 1 time<br>E.g. colou?r matches **color** *(0)* and **colour** *(1)* |
| * | Zero or more of the preceding character<br>E.g. tre* matches **tree** *(2)*, **tre**ad *(1)*, and **tr**ough *(0)* |
| + | Match the preceding character 1 or more times<br>E.g. tre+ matches **tree**, and **tre**ad |
| [...] | Match anything inside the square brackets for one character position once<br>E.g. [0-9] matches any character in the range 0-9<br>     [abc] matches **a**, **b**, or **c** |
| [^...] | Match any character excluding those in the square brackets<br>E.g. [^A-M]in matches **Windows**, but not **Lin**ux |

# Common Operators

| Operator | Purpose |
|---|---|
| {n} | Match the preceding character, or character range, n times<br>E.g. [0-9]{3}-[0-9]{4} matches local phone number like *123-4567* |
| {n,m} | Match the preceding character at least n times but not more than m times<br>E.g. [A-Z]{2,4} matches *com*, *sg*, but not *abcde* |
| () | Group parts of search expression together |
| \| | Separate two alternative values<br>E.g. gr(a\|e)y matches both *gray* and *grey* |
| \b | Match empty string, frequently used to indicate a word boundary<br>E.g. \bhis\b matches *his* only, not *this* or *history* |

- Take a look at our email pattern regex again:

[A-Z0-9._-]+@([A-Z0-9.-]+\.)+[A-Z]{2,4}

2 to 4 such characters

One or more such characters

One or more such characters

'xxx.' is repeated one or more times

Followed by '@'

Followed by '.'

- GA Miller. WordNet: A Lexical Database for English, *Communications of the ACM*, 1995

- GA Miller. Nouns in WordNet: A Lexical Inheritance System, *International Journal of lexicography*, Oxford University Press, 1990

- Morato, Marzal, Llorens and Moreiro. WordNet Applications, in Proceedings of Global WordNet Conference, pp. 270-278, 2004.

- B. Liu. *Sentiment Analysis and Opinion Mining*, Morgan & Claypool, 2012.

- Regular Expression Tutorial:

  http://www.zytrax.com/tech/web/regex.htm