# Semantic Text Similarity

## Applications

→ Group similar words in semantic contexts

→ Textual entailment

→ Paraphrasing

Wordnet → symantic dictionary → interlinked by semantic relations.

→ organizes information in a hierarchy

Path Similarity → find shortest path between two concepts.

## Lowest Common Subsumer (LCS)

→ find the closest ancestor to both concepts.

Lin Similarity → similarity measure based on the information contained in the LCS of two concepts.

$$LinSim(u, v) = \frac{2 \times \log P\left(LCS(u, v)\right)}{\left(\log P(u) + \log P(v)\right)}$$

$P(u)$ is given by the information content learnt over a large corpus.

## Python

```
from nltk.corpus import wordnet as wn

deer.path - similarity (elk)
```

## Collocations & Distributional Similarity

Two words that frequently appears in similar contexts are more likely to be semantically related.

context → words before, words after; within a small window

→ POS words before, after, within a small window.

→ specific syntactic relation to the target word.

→ frequency of two or more words.

→ frequency of individual words.

→ so, normalisation is important.

Pointwise Mutual Information.

$$PMI(w, c) = \frac{\log P(w, c)}{P(w) \times P(c)}$$

$w \rightarrow$ word

$c \rightarrow$ context.

Python    (NLTK) collocations & associations

     . pmi (      )

## Topic Modeling

### Intuition

Documents are a mixture of topics.

Ex $\begin{cases} \text{genetics} \\ \text{computation} \\ \text{life sciences} \\ \text{anatomy} \end{cases}$

Topic modeling $\rightarrow$ coarse level analysis of what's in a text collection

$\rightarrow$ topics are represented by a word distribution

$\rightarrow$ document is assumed to be a mixture of topics.

What's known → text collection
                number of topics

what's not known → actual topics
                    topic distribution


→ text clustering problem
→ documents + words are clustered
   simultaneously

Topic modeling approaches
───────────────────────────
① PLSA        (1999)

② LDA         (2003)    (better)


Generative Models and LDA
───────────────────────────


Generative Models




is

the harry          Generation        | the movie
potter the  I     ─────────────→     | harey potter
      movie                          | is _ _ _
am,                                  | _ _ _ _ _ _ _

              Inference
          ←───────────────
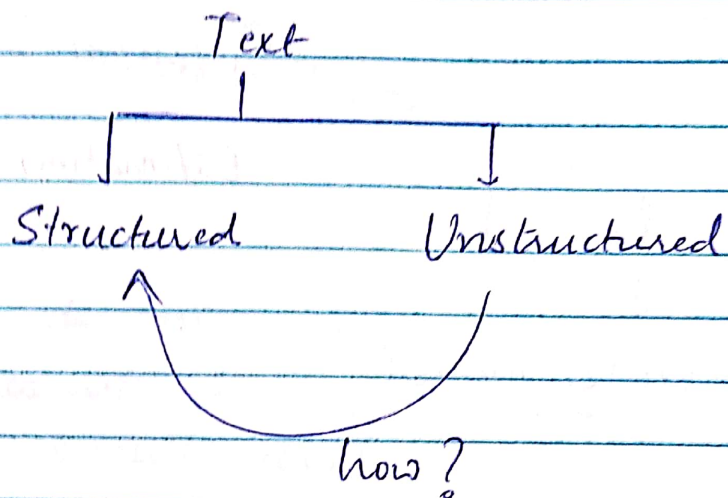$Pr(text/model)$   Estimation.

mixture model.

LDA → generative model of document d
   → mixture of topics.
   → use a topic's multinomial
distribution to output words to fill that
topic's quota.

In practice → How many topics ?
   → Interpreting topics
   → Topics are just word distributions
   → making sense is non-trivial
         and subjective.

Gensim → LDA.

Preprocess → stemming, normalize, stopword
         removal, convert to a
         dtm → document term matrix.

# Information Extraction

Text

```
        Structured        Unstructured
              ↑_____|
                   how?
```

Goal → identify and extracts fields of
interest from free text

headline, author, reviewer, date/
place of publishing, etc.

## Fields of Interest

→ Named Entities. (NEWS)
→ Money, companies (Finance)
→ Diseases, drugs, procedures ( Medicine)

→ Relations (what, who, when, where)

## NER

technique to identify all
mentions.

→ identify the mention/phrase:
boundary detection
→ identify the type: tagging

# Approaches to identify named entities.

→ dates, phone numbers → Regex.
↪ others → ML approach.

## Standard NER Task in NLP

→ four class model
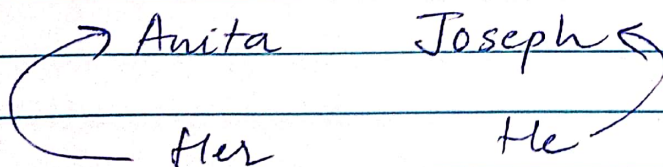
       → PER
       → ORG
       → LOC/GPE
       → other/outside

→ Relation extraction
→ Co-reference resolution.
    • disambiguate mentions and group mentions together.

Anita     Joseph

Her         He

### This is important Q&A's.

Given a question, find the appropriate answer.