# Wordnet based summarization."

~~Date~~ Idea:

① Select sentences based on their /Semantic content
   └→ meaning.

② Given relative importance w.r.t ~~in conjunction~~ to the semantics of the whole text.

③ Reduce the text, pieces of text ~~corr~~ corresponding to the same semantic content.
   └→ reduction of redundancy.

## Key steps.

① Pre processing

② Subgraph construction form the WordNet

③ Synset Ranking

④ Sentence selection

⑤ PCA

⑥ Final pruning.

① **Preprocessing.**

(a) Split the text into sentences.

(b) POS tagging. (every word in the sentence is tagged with its most relevant POS) → NLTK.

⤷ detect the correct sense of the word.

pant ⟨ → noun (clothing).

⟩ → verb (fast breathing)

(c) Identifying collocations.
words that typically appear together in a sentence ___ 4miles per hour

⤷ all idiomatic phrases.

(d) Remove stop words.
it, of etc.

⟶ the sequence is very important. "take off"

## ② Sub-graph Construction.

ⓐ Mask all the words and collocations that appear in the text (to be summarized) in the WordNet hypernymy.

ⓑ Traverse the **generalization edges** upto a fixed depth & mask the (Synsets) you visit

⤷ groupings of synonymous words that express the same concept.

book.(n)(02) → [ book.n.01, collocation.n.02, impression.n.06, magazine.n.04, volume.n.04 ]

noun.

⤷ physical objects of a number of pages bound together.

© Construct a graph containing only the marked synsets as nodes & the generalization relationships as edges —— synset sub-graph.

---

③ Synset Ranking:

Rank the synsets ~~that~~ based on their relevance in the text (to be summarized)

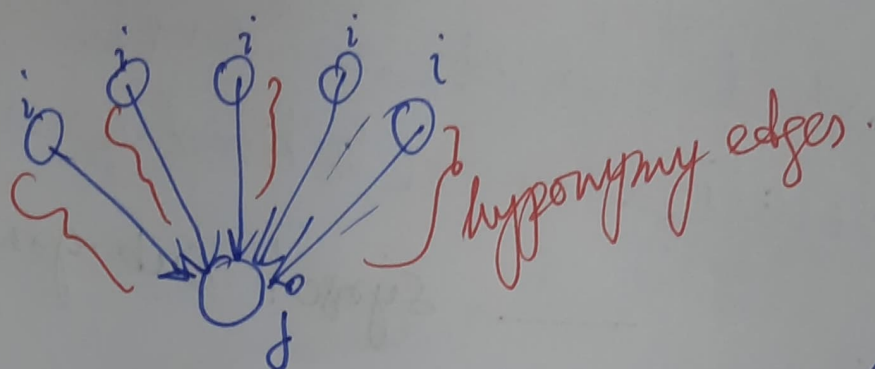ⓐ Construct a rank vector R corresponding to each node of the graph. R is of dimension $n$. $n$ is the no. of nodes/~~synset~~ synsets in the graph.

Each entry $= \frac{1}{\sqrt{n}}$

ⓑ Authority matrix

$$A(i,j) = \frac{1}{(number\text{-}of\text{-}predecessors(j))} \quad \text{if } j \text{ is a child of } i$$

$$= 0 \quad otherwise.$$

*hyponymy edges.*

↳ how many nodes does $j$ ~~log~~ draw its __meaning__ from.

ⓒ Update Rank vector:

$$R_{new} = \frac{R_{old} * A}{| R_{old} * A |}$$  } *Analogous to PageRank.*
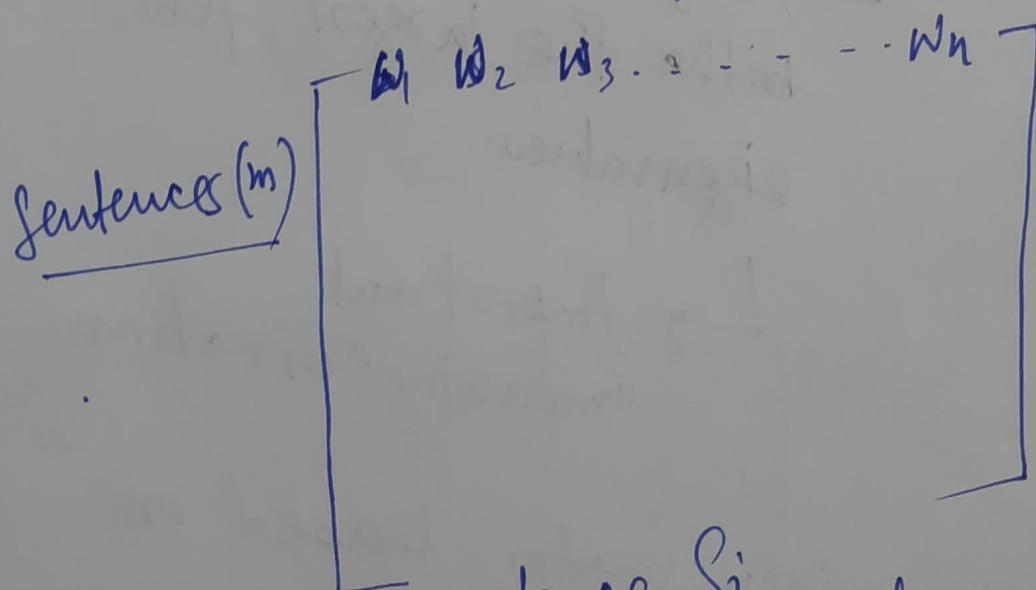
unit normal in that direction.

till $|R_{new}|$ changes less than a predefined threshold.

Higher values correspond to better ranked nodes.

(A) Sentence selection.

(a) Contruct a matrix $M$ with $m$ rows
    $\downarrow$ no. of sentences in the text to be summarized.
& $n$ columns
    $\downarrow$ no. of nodes in the subgraph
    words/synsets ($n$)

Sentences ($m$)
$$\begin{bmatrix} W_1 & W_2 & W_3 & \cdots & \cdots & W_n \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix}$$

(b) For each sentence $S_i$
Traverse the subgraph following the gen.
edges with the words present in $S_i$
Find all the reachable synsets $SY_i$

For each $sy_j \in SY_i$ set $M[S_i][sy_j] = R[sy_j]$.

(5) PCA.

~~The primer~~

(a) Compute the principal components of M.
        ↳ eigenvectors of M with the largest few eigen values.

               ↳ Important "meaning" directions.

(b) Sort the eigen vector based on their eigen values.
take the top few eigen vectors
& compute projection on each sentence.

$$Pr(\vec{e})_{\vec{s_i}} = \frac{\vec{e} \cdot \vec{s_i}}{|\vec{s_i}|}$$

Max

$(\lambda_1)$ $e_1$ ← m how many sentences?

Second max

$(\lambda_2)$ $e_2$ ← m how many sentence?

↳ K sentences.

$$K = \alpha \frac{\lambda_{ik}}{\sum_j \lambda_j}$$

through the eigenvector i.

$$K \simeq \left(\frac{\lambda_i}{\sum \lambda_i}\right) N$$ ← N is no. of sentences in the target summary (budget).

⑥ Ⓡ <u>Final pruning.</u>

<u>Removal of undefined references.</u>

ⓐ Remove sentences that start with pronouns he/she/it.

ⓑ Remove sentences within quotes.