

Wordnet based summarization.

~~Basic~~ 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.
pieces of
- ③ Reduce the text ~~can~~ corresponding to the same semantic content.
↳ reduction of redundancy.

Key steps.

- ① Pre processing
- ② Subgraph construction from the WordNet
- ③ Synset Ranking
- ④ Sentence selection
- ⑤ PCA
- ⑥ Final pruning.

① Preprocessing.

① a Split the text into sentences.

② 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).
pant → verb (fast breathing)

③ c Identifying collocations.

words that typically appear together in a sentence — 4 miles per hour

↳ all idiomatic phrases.

④ d Remove stopwords.
it, of etc.

→ The sequence is very important. "take off"

② Sub-graph Construction.

(a) Mark all the words and collocations that appear in the text (to be summarized) in the WordNet ^{hypernymy}.

(b) Traverse the generalization edges upto a fixed depth & mark 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.01, volume. n.04]
↑
noun.

↳ physical objects of a number of pages bound together.

② Construct a graph containing only the marked synsets as nodes & the generalization relationships as edges
—— synset sub-graph.

③ Synset Ranking:

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

(a) Construct a rank vector R corresponding to each node of the graph. R is of dimension n .

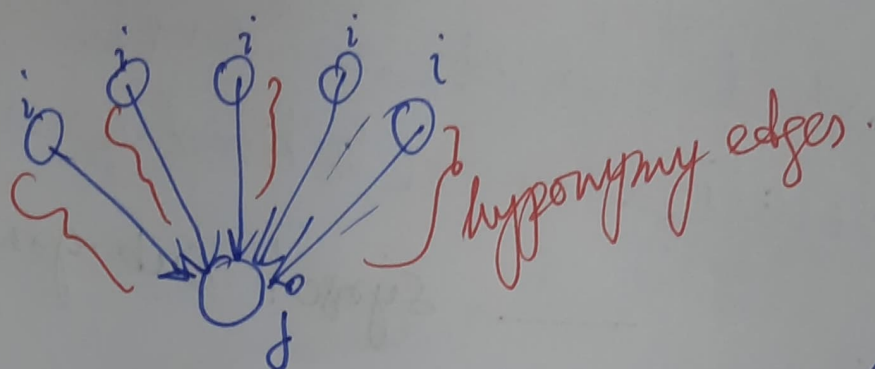
$$\text{Each entry} = \frac{1}{\sqrt{n}}$$

n is the no. of nodes / ~~synsets~~ in the graph.

(b) Authority matrix

$$A(i, j) = \frac{1}{(\text{number of predecessors}(j))}$$
$$= 0 \quad \text{otherwise.}$$

If j is a child of i



how many nodes
does j ~~draw~~ draw
its meaning from.

③ Update Rank vector:

$$R_{\text{new}} = \frac{R_{\text{old}} * A}{|R_{\text{old}} * A|}$$

Analogous to
PageRank.

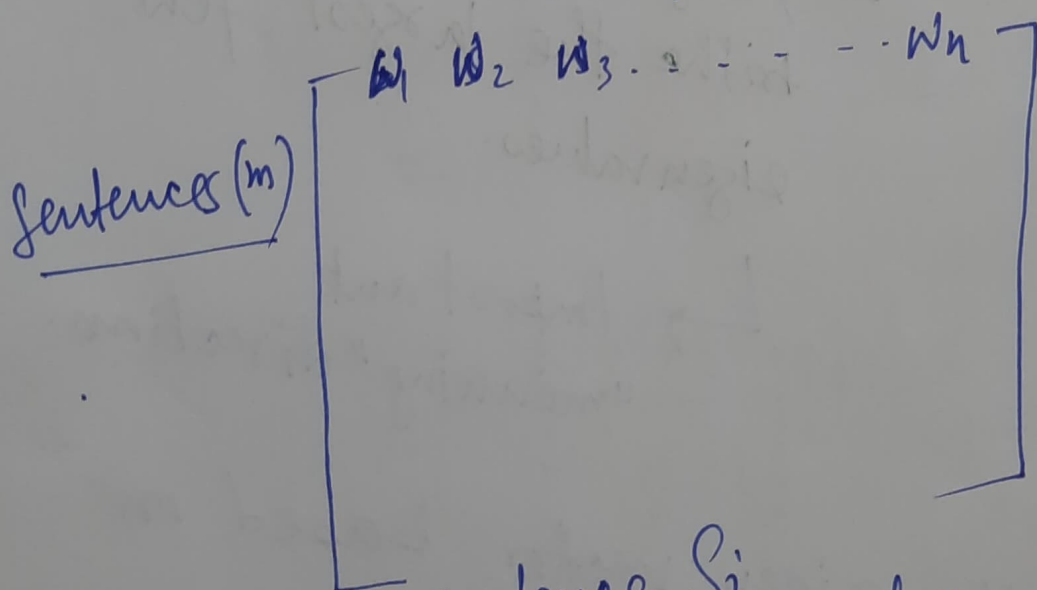
unit
normal
in that
direction.

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

Higher values correspond to better
ranked nodes.

Sentence selection -

- 4) Sentence summarization
- a) Construct a matrix M with m rows
& n columns
↳ no. of nodes in the subgraph
words/synsets (n)
↳ no. of sentences in the text to be summarized.



- (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 prince~~

① Compute the principal components of M .

↳ eigenvectors of M with the largest few eigenvalues.

↳ Important "meaning" directions.

② Sort the eigen vectors based on their eigen values.
take the top few eigenvectors

& compute projection on each

sentence .

$$Pr(\vec{e}) \vec{s}_i = \frac{\vec{e} \cdot \vec{s}_i}{|\vec{s}_i|}$$

Max λ_1 \leftarrow on how many sentences?
 Second max λ_2 \leftarrow on how many sentence?
 λ_1 λ_2
 } K sentences.

$$K \propto \frac{\lambda_i}{\sum_j \lambda_j}$$

through the eigenvector i .

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

⑥ Final pruning.

Removal of undefined references.

① Remove sentences that start with
pronouns he/she/it.

② Remove sentences within quotes.
