

Fundamental difference b/w IR & NLP

IR is focused much more on bag of words model (does not care about the sequence).

A is a student of B & B is a student of A is equivalent for an IR model. (But this is changing now).

We can store the tf-idf even in an evolving corpus without too much overhead because idf remains almost constant (new documents will usually have same distribution & changes are in log scale). \therefore Occasionally recalculating idf is sufficient.

When is phonetic correction required

i) Various words

ii) Queries containing names (& other things) which sound the same but are spelled differently in different languages.

Vector space models

An example of using Euclidean distance is a bad idea:

Consider $A = D_1$ & $B = D_2 D_1$ (D_1 concatenated with D_2)

B will be usually farther to a query than A (considering query to be small); even though their contents are same.

Now we're considering bag of words model, length normalisation is fine because we anyways don't look at the content of words.

For building vector for query, it makes more sense to use $\text{tf} = 1$ than using $\text{tf} = \text{no of items with term}$ cause (cause using the same token multiple times in query should not build a bias towards terms)

Normalising the query vector is not required (because we are only concerned with the order of similarity & not normalising just scales all the similarity by the same value)

Rocchio Algorithm:

Value of q_{top} is obtained by differentiating $(\text{sim}(q_{top}, C_r) - \text{sim}(q_{top}, C_{irr}))$

For IR task: Accuracy is not a good measure of performance.

Reason: Usually, no of non-relevant documents \gg no of relevant documents. Even if our system marks everything as non-relevant, it will still have a very high accuracy even though the model will not be good.

4 services that IR systems provide

→ Search

→ Discovery → identification of documents with desired queries.

→ Summarisation → extract key content from document for user.

→ Recommendation → algorithmic or rule-based recommendation.

Implicit feedback:

Used in a search session.

Relevance feedback: Uses documents to improve the query.

Query expansion: Uses query log to improve query.

Precision can increase as recall increases. But usually that is not the case.

Recall is not easy to measure for an unannotated database.

Why precision @ k is not a good measure:

Consider queries that do not have large no. of relevant documents (unpopular queries). Even if system retrieves all these docs, precision @ k will be very low.

NDCG

Initially we want the documents to be presented in the decreasing order of relevance. We take this retrieval (docs in decreasing relevance order) as the ideal retrieval to calculate $Z_{k,j}$.

Link prediction problem

Given a set of nodes we have to predict which pairs of nodes have an edge b/w them.

Given a graph & a node, find predict the set of nodes this node will have an edge to.

NDCG can be used to measure the accuracy of link prediction algorithm.

(Rank the links according to the probability the algorithm predicts & use the ground truth relevance to calculate NDCG).

Kappa measure for inter-judge agreement

Measure how much the human judges agree.

(Used to eliminate human bias with respect to rating)

$$P(A) = 0$$

$$P(A) = 1$$

Kappa value ranges from $\frac{-P(E)}{1-P(E)}$ to 1

★ ★ Note $P(A)$ can be less than $P(E)$ (They're independent)

We don't have a fair measure for marginal relevance. Usually measured first by using click stream data & ranked.

Kappa measure for inter judge agreement:

$P(E) \rightarrow$ Dependent on the density of the various outcomes & not dependent on the individual behaviour.

i.e. $P(E)$ depends on the ground truth number of relevant & non-relevant documents. But because we don't have the ground truth values, we use the annotations as a proxy. If we had ground truth values, $P(E)$ would have been

$$\frac{(\# \text{relevant docs})^2 + (\# \text{irrelevant docs})^2}{\# \text{total docs}}$$

VB encoding

Average waste = $4 \text{ bits} + \frac{\alpha K \text{ bits}}{1 \text{ gib}} \text{ per document}$

gamma encoding

Seldom used because not byte aligned.

#

Language model

P(STOP) regulates the length of sentences.

#

Parameter estimation

$$\hat{P}(b|M_d) = \frac{1}{M_d} \sum_{i=1}^{M_d} \delta(b_i)$$

This is an estimate because this document is not an actual language model but just a sample from a master distribution (i.e. this document does not represent a language but a sample from the prob distribution of the language). We are just estimating master distribution from documents distribution.

#

Which is better measure of info amount: collection frequency or document frequency?

Document frequency because there might be some documents which might contain a term large no. of times irrationally increasing the collection freq.
thus,

- # Advantages of vector space / probabilistic retrieval model over classification (learning) & generation models:
- Vector space / probabilistic models are computationally cheap.
 - generation models are usually closer to info much (but the information might be hallucinated).

Item doc frequency matrix: Has t/f instead of 1/0 values.

Time complexity of

$$\text{OR/ AND} \rightarrow O(|x| + |y|)$$

$$\text{NOT (just NOT)} \rightarrow O(N)$$

$$\text{XOR AND NOT } y \rightarrow O(|x| + |y|)$$

$$\text{XOR NOT } y \rightarrow O(N)$$

F score

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(P^2 + 1) PR}{P^2 + R}$$

New encoding

Kappa measure

$$= \frac{P(A) - P(E)}{1 - P(E)}$$

#

Probabilistic IR

$$P(R=1 | d, q) = P(\vec{x}^d | R=1, q) P(R=1 | q)$$

$$P(\vec{x}^d | q).$$

$$\therefore O(R | d, q) = P(\vec{x}^d | R=1, q) P(R=1 | q)$$

$$\cdot \frac{P(\vec{x}^d | R=0, q) P(R=0 | q)}{P(\vec{x}^d | R=0, q)}$$

$$= O(R | q) \frac{P(\vec{x}^d | R=1, q)}{P(\vec{x}^d | R=0, q)}$$

$$P(\vec{x}^d | R=0, q).$$

* separate into term products, where $\mu_t = p_t = 0.5$ for all terms not in query, you get:

$$O(R | d, q) = O(R | q) \prod_{t: x_t = q_t = 1} \frac{p_t}{1-p_t} \frac{1-\mu_t}{\mu_t}$$

$$RSV_d = \sum_{t: x_t = q_t = 1} \log \frac{p_t}{1-p_t} + \log \frac{\mu_t}{1-\mu_t} \log \left(\frac{1-\mu_t}{\mu_t} \right)$$

Estimate $\mu_t = \frac{\# d/t}{N}$, $p_t = \text{constant} = 0.5$

(can also be estimated as
no. of docs containing d_t that
have t)

$$RSV_d = \sum_{t: x_t = q_t = 1} \log \frac{\# d/t}{N}$$

Chapt / BM 25

$$RSV_d = \sum_{t: \alpha_t = \phi_t = 1} \frac{\log \frac{N}{df_t}}{(k_i + 1) t f_{t,d}} - \frac{k_i ((1-b) + b \times \frac{L-d}{\text{avg}}) + t f_{t,d}}{k_i ((1-b) + b \times \frac{L-d}{\text{avg}}) + t f_{t,d}}$$

Language models for IR.

$$P(b|d) = \lambda P(b|M_d) + (1-\lambda) P(b|M_c).$$

longer queries \rightarrow prefer lower λ shorter queries \rightarrow prefer larger λ .

NDCG

$$= \sum_{k=1}^K \frac{R(i, m)}{\log_2 (1 + i)}$$

#

Mine Bayes assumption in IR:

Positional independence: Probability of token t occurring in document d given class c is independent of its position.

Although usually violated, does not matter much because the order usually remains same (order of docs)

#

Usually: Micro average has higher value than macro average:

Reason: Precision is usually low for the tail classes. In micro macro average, all classes receive the same weight (& tail classes are usually high in number)

#

Mutual importance:

For this part we look at terms other than the query terms that are important to answer the query

#

Applications of retrieval systems:

#

Summarization

For a document of length L , a summary of $\frac{L}{2}$ length was considered acceptable, but with increase of docs on web, we want summaries to be orders of magnitude less than L . (Nothing greater than $\frac{L}{10}$ is considered acceptable).

Types of summaries

i) Inductive v/s informative

↓
like headlines

Very precise & short

↳ Some content drawn from the main doc

ii) Extractive v/s abstractive (generative)

↓
Extract parts from the main doc to... ↳ Rewrite the info in the
build the summary doc in its own "sufficiently
abstractive" version.

iii) Generic v/s query oriented

↓
Provides the author's view

It provides the summary as per the
user query / demand.

iv) Background v/s just in the msg

↳
Most relevant
to news
articles

Assume that the reader
has no prior knowledge

Assume that the reader is well
aware of the situation

v) Single v/s multi document

↓
↳
Summary from single doc
Multiple docs are collated to generate
summary.

Some important parts of documents that help in summarization

- Abstracts & conclusions
- Section headings in scientific papers
- Many subtle linguistic cues hidden in general documents
Ex: Phrases like "in summary" are direct sources of summary.
- Font differences (bold, italics, larger font size) are usually cues to important info.

Graph based summary algorithms

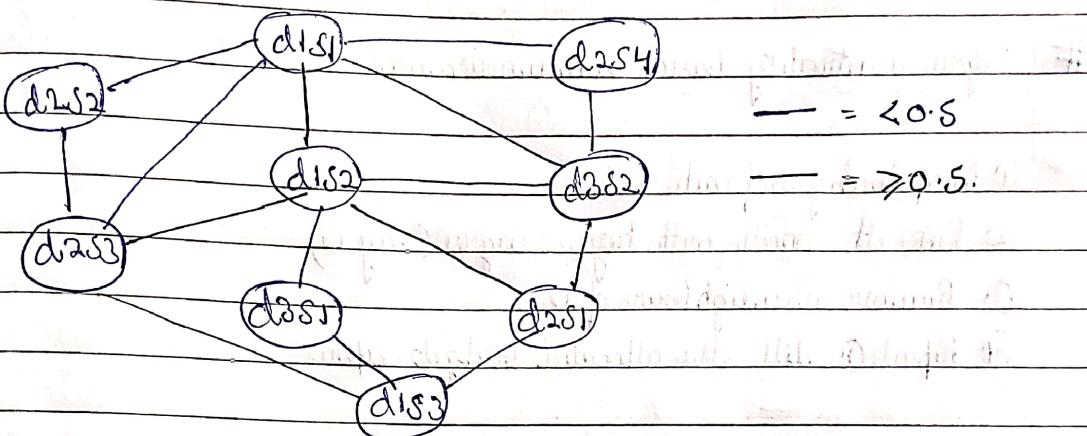
- Build a graph with nodes & edges that somehow represent the document.
- Measure certain properties of this graph to build the summary.
- The algorithm anchors upon the concept of "Prestige" in social networks.
- Imagine you're a graph where: nodes are sentences & edges are similarity b/w sentence pairs.

- Nodes constructed using tf-idf vectors of sentences.
- Edges constructed using cosine similarity b/w the tf-idf vectors.

Each sentence represented with an ID : d_{XSY}

~~d_{2S3} : 3rd sentence of 2nd doc.~~ doc X, sentence Y.

Consider the graph



In the above graph, it is apparent that d_{1s2} must be present in the summary of texts: high degree: similar to most other sentences & most of its edges have high weight? This idea is called degree based summarization approach.

noisy links:

Links that cause spurious similarities. Must be removed from the graph.

Removed by setting a threshold T to an similarity score below which we'll remove the edges.

High $J \rightarrow$ extremely sparse graph
low $J \rightarrow$ too much redundant info (noise)

Value of J chosen on validation set

Binarization of graph

After removal of edges below certain threshold, forget the weights (convert weighted graph to unweighted graph).

Degree centrality based summarization

- ① Compute degree of each node
- ② Pick the node with highest degree (say v)
- ③ Remove all neighbours of v
- ④ Repeat ① till the allocated budget expires.

PageRank

Recursive idea of recursive importance: a node's importance is determined by how important its neighbours are.

The graph remains same as before but rather than computing degree, use pagerank algorithm.

① Run pagerank

②

2 possible algorithms:

- I) ① Compute pagerank of nodes
② Select node with highest pagerank
③ Remove neighbours
④ Go to ②

- II) ① Compute pagerank of nodes
② Select node with highest pagerank
③ Remove neighbours
④ Recompute pagerank, go to ②.

Called the terrank / lerank algorithm.

Linguistic summarisation algorithms

Based on the idea of lexical chains.

lexical chains summarisation:

Collecting information from correlated chain of words

Build chains:

- ① Pick a set of candidate words (generally, the words are nouns / noun compounds)
- ② For each candidate word, build a chain based on its relatedness with other words.
- ③ If the candidate word is already related to words in a particular chain, then fuse it with that chain rather than building a new chain for it.

Measure of relatedness

Usually rely on lexical resources like WordNet.

Choosing one of the chains as a document representation# Choose the strong chains

↳ Determined by a scoring mechanism.

Heuristics for scoring:

i) length: No. of occurrences of a particular word of the chain.

ii) Homogeneity Index:

$$= 1 - \frac{(\text{No. of distinct occurrences of a word of a chain})}{\text{length}}$$

$\Delta_{\text{CH}}(\text{chain})$: length \times Homogeneity index.

Ex:

word # of occurrences

Chain: Microsoft (10) | Concord (1) | Company (6) | entertainment (1)

Enterprise (1) | MIT (1)

$$\text{length} = 20 \quad (10 + 1 + 6 + 1 + 1 + 1)$$

$$\text{HII} = 1 - \frac{6}{20} = \frac{14}{20}$$

$$\Delta_{\text{CH}} = 14.$$

Strong chain: chains with score higher than Average(score) + $2 \times$ standard deviation(score).

If there are 3 strong chains (say 1, 2 & 3): Use these strong chains to build summary.

Possible summary strategies:

- i) Choose the sentence that has the first occurrence of word from the chain (i.e. the first sentence that contains any word from sentence). Stop doing till budget runs out.
- ii) Choose the first sentence that contains the occurrence of the key word (candidate word that started the chain) of the chain.
- iii) Choose the sentences that has the first occurrence of a high density (choose a threshold) of words from the chain.

For practical purposes, ii) method works best (because its works based on the key-word & other words in the chain are related to this word).

#

Wordnet based summarization

Idea:

- ① Subsentences based on the semantic context:
↳ meaning
- ② The context is given relative importance w.r.t the semantics of the whole text.
- ③ Reduce the text
- ④ Reduce the pieces of text corresponding to the same semantic context: Reduction of redundancy.

#

Key steps:

- i) Preprocessing
- ii) Subgraph construction from wordnet
- iii) Hyperm graph ranking
- iv) Sentence selection and their connection with words
- v) PCA
- vi) Final pruning.

#

Preprocessing:

- a) split the text into sentences (except for languages like Mandarin)
- b) POS tagging (every word in the sentence is tagged with its most relevant POS): Helps to detect the correct sense of the word.

② Identifying collocations, words that typically appear together in a sentence.

Ex: all idiomatic phrases.

d) Remove stopwords

④ Note: The sequence of the preprocessing steps is very important. Due to phrases like "take off"

Subgraph construction:

a) Mark all the words & collocations that appear in the text to be summarized in the wordnet

b) Traverse the generalisation edges upto a fixed depth & mark the synsets you visit. → hyponymy edges of wordnets

↳ groupings of synonymous words that express the same concepts (parts of wordnets)

Ex: book (n. 02) → has [book - n. 01, collocation - n. 02, expression - n. 06,
magazine - n. 01...]
Each node also has a gloss (the meaning) & a synset

c) Construct a graph containing only the marked synsets as nodes & the generalisation relationships as edges → Synset subgraph.

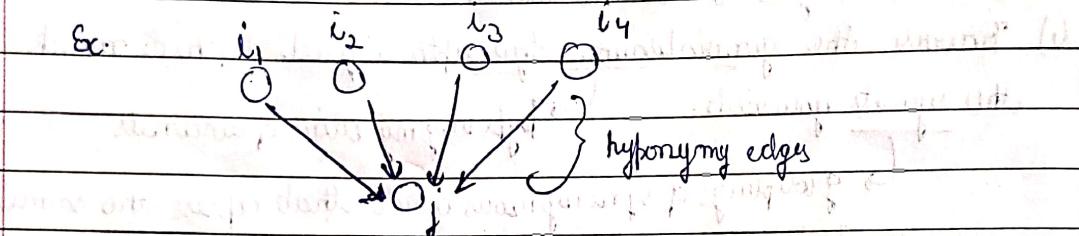
Symbol ranking: Rank the symbols based on their relevance in the text (to be summarized)

Hints

- construct a rank vector R corresponding to each node of the graph. (R has 1 entry for each node)
each entry = $\frac{1}{\sqrt{n}}$ n = no of nodes/symbols in graph.

ii) Authority matrix

$$A(i,j) = \begin{cases} 1 & \text{if } j \text{ is child of } i \\ \text{num of predecessor}(j) & \rightarrow \text{predicessor in wordnet} \end{cases}$$



Authority matrix determines how many nodes does j derive its meaning from.

iii) Update rank vector

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

$|R_{\text{old}} * A|$

until R_{new} changes less than a predefined threshold:

Each time we compute R_{new} , we add a small bias to all its entries to prevent underflow.

Sentence selection

(a) Construct a matrix M with m rows & n columns \rightarrow no of sentences in the text to be summarised
 \hookrightarrow no of words in subgraph

(b) For each sentence S_i , again traverse the subgraph following the generalisation edges with the words present in S_i .

Find all the symbols that are reachable $\cdot S_{Y_i}$

For each $S_{Y_j} \in S_{Y_i}$, $\cup M[S_i][S_{Y_j}] = R[S_{Y_j}]$

PCA

\rightarrow Compute the principal components of M .

\hookrightarrow Eigen vectors of M with largest eigen values

These eigen vectors will capture the important meaning directions.

\rightarrow Sort the eigen vectors based on eigen values, take the top few eigen vectors & compute projection on each sentence.

$$Pr(\vec{v})_{S_i} = \vec{v} \cdot \vec{s}_i / \|\vec{s}_i\|$$

\rightarrow How many sentences to project on?

Project each \vec{v} eigen vector on K sentences, where

$$K \propto \frac{1}{\sum_j \lambda_j}$$

$$K = \left(\frac{\lambda_i}{\sum_j \lambda_j} \right) N \rightarrow N = \text{no of sentences in target budget}$$

→ A project on all sentences & choose the top k sentences for each eigen vector.

Q8

Final pruning

Removal of undefined references:

(a) Remove sentences that start with pronouns.

(b) Remove sentences within quotes.

↳ Because usually represents opinions & not facts

Optimisation based summarisation

→ A global pruning method.

→ Notation:

D is a document

→ D has t_n textual units (sentences in this case)

→ $D = t_1, t_2, \dots, t_n$

→ $Ru(i)$: The relevance of t_i to the target summary

→ $Rd(i, j)$: Redundancy b/w t_i & t_j

→ $L(i)$: Length of text unit t_i (in terms of no of words)

→ The pruning is modelled as follows:

• Problem is to select a subset S of the document's textual units

from D such that the summary score $S(S)$ is maximised

$$\delta(S) = \sum_{i \in S} \text{Rel}(i) - \sum_{i, j \in S, i < j} \text{Red}(i, j)$$

$$\therefore A(S) = \underset{S \subseteq D}{\operatorname{argmax}} \left[\sum_{i \in S} \text{Rel}(i) - \sum_{i, j \in S, i < j} \text{Red}(i, j) \right]$$

↑ ↑

Increase the relevance of text units in S .
reduce the redundancy among the choice of text units.

$\rightarrow \text{Rel}(i), \text{Red}(i, j) \rightarrow$ Should be something that enforces diversity.

would be something like presence of special entities

Solving this optimisation problem:

i) Greedy solution

\rightarrow Take i such that $\text{Rel}(i) > \text{Rel}(i+1) \forall i$.

$\rightarrow S = \{t_i\} \rightarrow$ most relevant to last word

\rightarrow While $|S| \leq k \quad \sum_{i \in S} l(i) \leq k \leftarrow$ budget
 $t_j = \underset{t_j \in D-S}{\operatorname{argmax}} \delta(S \cup \{t_j\})$

Output S

* Greedy solution provides approximate results. We need better solution

* We'll ~~solve~~^{state} this optimisation problem as an Integer Linear Programming problem

has many solving methods

Recasting this problem as unconstrained optimisation function

Anterior level ILP formulation:

Optimisation function.

$$\text{maximise } \sum_i \alpha_i \text{Rel}(i) - \sum_{i,j} \alpha_{ij} \text{Rel}(i, j)$$

constraints: (t_i, t_j)

$\alpha_i, \alpha_{ij} \in \{0, 1\}$ ← Indicator variables: tells us whether a textual unit or textual pair are indicated in the target summary or not

$$\sum_i \alpha_i \text{ll}(i) \leq K \leftarrow \text{Budget constraint}$$

$$\boxed{\alpha_{ij} - \alpha_i \leq 0}$$

$\alpha_{ij} - \alpha_j \leq 0$ ← If $b_i \& b_j$ pair is included then they must be individually included.

$$\boxed{\alpha_i + \alpha_j - \alpha_{ij} \leq 1}$$

$\alpha_i + \alpha_j - \alpha_{ij} \leq 1$ ← If $b_i \& b_j$ are individually included then the pair of $b_i \& b_j$ must be included

How to operationalise $\text{Rel}(i)$ & $\text{Rel}(i, j)$:

Depends on a lot of things such as:

i) Domain of the summarisation

Ex: For news summarisation, we could use

$$\text{Rel}(i) = \text{Pos}(t_i, D)^{-1} + \text{SIM}(t_i, D)$$

where D = document collection

D = documents being summarised

$\text{Pos}(t_i, D)$ = position of textual unit t_i in documents D .

$\text{SIM}(t_i, D)$ = similarity of text unit with overall collection of documents.

Empirical observation: initial sentence in a news document is very important.

$\text{SIM}(t_i, D) \leftarrow \vec{t}_i \rightarrow$ vectorisation of document t_i

$\vec{D} \rightarrow$ vectorisation of $D \rightarrow$ one way is to sum up the vectors of all individual docs.

$$\text{SIM}(t_i, D) = \text{cosine sim}(\vec{t}_i, \vec{D})$$

Ex: Query based summarisation:

↳ from query Q : vector

$$\text{III}^{\text{by}} \text{ Rel}(i, j) = \text{SIM}(t_i, b_j) = \text{cosine sim}(\vec{t}_i, \vec{b}_j)$$

b) Sc: Query focused summarization

Given a query $Q \rightarrow$ we want a relevant document summary.

Now position of t_j in doc is not much imp but presence of query terms in b_i is imp.

$$\text{Rel}(i) = \text{SIM}(t_i, Q) + \text{SIM}(b_i, Q)$$

$$\text{Rel}(i, j) = \text{SIM}(t_i, t_j)$$

Alternative problem formulation: concepts level ILP:

No longer at word units or sentence level but at the concept level.

Tells all the important concepts that the summary should cover.

- ① Summary benefits by including a particular concept only once
- ② This means redundancy is implicitly captured.

$b_i c_i \leftarrow$ indicator variable for the concept i : (present or not)

$w_i \leftarrow$ weight of this concept

Revised ILP

$$\text{maximise} \sum_i w_i c_i$$

$$\text{subject to: } \sum_i b_i s_i \leq K$$

$$s_i c_i \leq c_i \quad \forall i, j$$

\hookrightarrow If you take a sentence: its coverage = coverage of concept i

mandates all concepts in that sentence should be chosen.

$s_i =$ indicators of sentence

j is summary

$c_i =$ coverage of concept i

in sent j



$$\sum_j s_j c_{ij} \geq c_i$$

↳ If a concept is ever selected, its mandatory that at least one sentence containing that concept must be selected.

$$c_i \in \{0, 1\} + i$$

$$s_j \in \{0, 1\}, + j$$

• One simple definition of concepts: "word bigrams" & weight = no of occurrences of word bigrams in the input documents.

Evaluating summaries

i) ROUGE score: (Recall oriented Understudy of existing Evaluation)

→ Recall is important in summarisation because we are more interested in "coverage" than precision (because we want to cover more in limited budget).

Precision is more important in safety critical system.

a) ROUGE- N ($N = 1, 2, \dots$)

→ Given a machine generated summary (candidate summary): C & a reference summary / human written summary / gold standard summary: R.

Ex: R: the cat is on the mat

C: the cat is and the dog

Rouge- N precision

Ratio of the number of Ngrams in C that also appear in R by the number of Ngrams in C

Rouge N recall (more imp)

Ratio of number of Ngrams in C that also appears in R over no of Ngrams in R.

For the given example:

ROUGE-1 =

$$\text{precision} = \frac{3}{5}$$

intersection unigrams: {the, cat, the}.
unigrams in C: {The, cat, and, the, dog}.

$$\text{recall} = \frac{3}{6}$$

unigrams in R: {the, cat, is, on, the, mat}.

$$\text{Rouge-1 precision} = \frac{1}{4}$$

unigrams in C: {the, cat, cat and, and
the, the dog}.

$$\text{Rouge-1 recall} = \frac{1}{5}$$

unigrams in R: {this cat, cat is, is on, on the,
the mat}.

$$\text{Rouge-N F1} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{precision} + \text{recall}$$

b) ROUGE-L: longest common subsequence

longest common subsequence \rightarrow not necessarily consecutive.

\Rightarrow LCS = the cat the

$$|\text{LCS}| = 3$$

Denominator = no of unigrams of C (for precision) & R (for recall)

$$\text{Rouge-L precision} = \frac{3}{5}$$

$$\text{Rouge-L recall} = \frac{3}{6}$$

If R & C contain multiple sentences: compute value of each sentence pair of R & C & average.

c) Rough - S \rightarrow skip grams

R: The cat is on the mat.

C: The gray cat ~~is~~ with and the dog.

Skip = S: we can skip S unigrams & if the n-grams still match, we count it as match.

Ex: For the above example: $S=0 \rightarrow 0$ match

$S=1 \rightarrow 1$ match.

suppose we have generated the TLP summary; i.e. we have selected some sentences.

Further issues:

i) Ordering of selected sentences

Domain dependent.

Ex: In news summarization: order in the same order as they appear in doc.

\rightarrow If coherence is imp: choose ordering that makes neighbouring sentences similar (closely)

\rightarrow In multi entity summary (summary that tells abt 2 or more things): choose ordering that brings similar entities closer.

\rightarrow Topicality: if summary covers multiple topics, make the sentences relevant as well as topically coherent.

ii) Simplifying sentences:

Parse the output summary & decide based on some rules which parts to skip off.

Some basic rules that can be followed:

i) Initial adverbials: Ex: "on the other hand", "as a matter of fact" etc.

ii) Prepositional phrases with named entities:

Ex: The commercial fishing restrictions in Washington will not be lifted until salmon population increases to a substantial number. X

iii) Attribution clauses:

attributes the sentence to this layer.

Ex: Rebels agreed to talks with officials local observers said on Tuesday X

iv) Affixatives: parts in b/w commas

Rajam, 28, an artist who was at that time living in Philly.

~~found inspiration in Bauhaus~~

Information extraction

① We want to obtain structured text from unstructured text

Objectives of an IE system:

- Find & understand limited relevant parts of text
- Gather info from many different texts
- Produce a structured representation of the relevant info gathered (may be a schema, knowledge base etc)

End goals:

- Make this data "useful" for users
- Organize the info in a semantically precise manner so that it can be used by algorithms (computers)

Input: unstructured / semi-structured data $\xrightarrow{\text{IE}}$ Structured data.

→ concerned with exact & factual information. We do not want opinions.

Roughly: An IE system would want to answer:

Who did What to When & With.

Ex: In 1998, Larry Page & Sergey Brin founded Google Inc.

Info: Founders (Google Inc., Larry Page)

Founders (Google Inc., Sergey Brin)

FoundedIn (Google Inc., 1998)

{ abstraction }

{ Entity1 }

{ Entity2 }

Applications

Biomedical domain:

- Start with large no. of scientific publications
- We want to have a table lookup that allows us to discover easily innovations related to i) particular genus. ii) particular bacteria etc.
- Biotech entities have many ambiguous names & lots of synonymous terms.
- We want to automatically identify the biomedical entities & link them (store as a graph) in knowledge base.

We want to convert the running text into tables like

| subject | relation | object |
|-----------|--------------|------------|
| b53 | is a | protein |
| b53 | has function | aff�tosis |
| aff�tosis | involved in | cell death |

New articles:

Relations here are slightly more well defined:

| Relations | Examples | Types |
|------------------|-----------------------------|---------------|
| affiliation | Married to, mother of | PER → PER |
| → personal | | |
| → organisational | president of, Spokesman for | PER → ORG |
| → Artifactual | owns, invented | PER/ORG → ART |

| | | |
|---------------|----------|-----------|
| spatial | near to | LOC → LOC |
| → proximity | | |
| → directional | south of | LOC → LOC |

Part of → Organisations
 ↳ Political a unit of / part of
 annexed / acquired ORG → ORG
 GPE → GPE
 ↳ geopolitical entity

Some trivial methods

- Handwritten rules (suffix based)
- Bootstrapping methods
- Supervised methods.
- Disturb supervision
- Unsupervised

REGEX based IE

Use finite automaton for noun groups

(proper noun)

PN → 0

ART

→ 1

N → 2

adj (adjective)

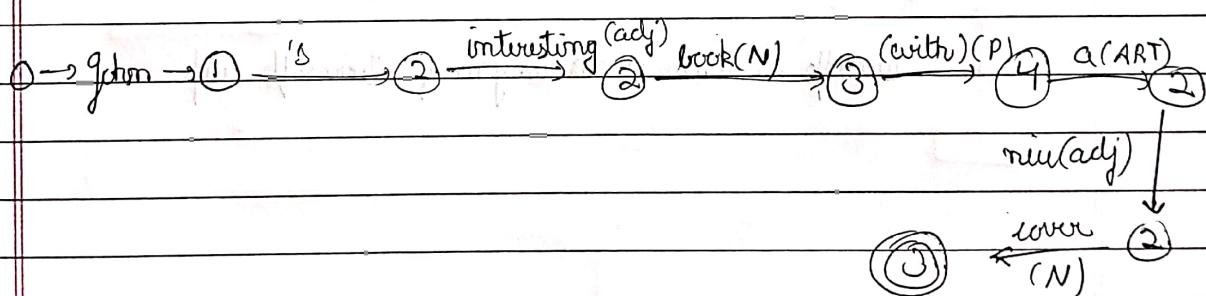
PN → 3

P

→ 4

N → @ non proper noun

Ex: John's interesting book with a nice cover.



Many such sugars to extract info.

Ex: To extract which person holds what position in what organisation, we would use sugars like

[person], [position] of [org]

Ex: Brnjkovic, leader of the Serbian Radical movement -
or [person] (approxid), named to [position] prep [offic]

#

Hyponymy as a key to IE (Specialised sugars)

Rules: (some representative rules)

Y such as X ($, X^*$) (and/or) X

such Y as X: such author as Shakespeare

X or other Y: bruises, wound or other injuries

X & other Y: temples, monasteries & other civic buildings

Y including X: common law countries including Canada & England

Y, especially X

X ← hyponym

Also called Heads patterns

Microsoft research is a part of Microsoft corporation: neither hyponymy nor hyponymy.

These are called ~~also~~ Meronymy & Holonymy.

Microsoft Research is meronym of Microsoft corp

Meronymy, hyponymy based Hearst patterns

Ex: ~~bark~~ [NN-PL] of [PREP] {^{the/a} [DET]} [JJ|NN] ^{*} whole [NN].

Ex: basement of a building

Ex: whole (s) [POS-variety] bark
[NN]/[NN-PL] (CNN-PL)

→ Think which pos is main fib.

Disadvantage of Hearst patterns

- Hard to write rules for every situation
- Hard to maintain rules.
- Expands very fast
- Changes with domain of documents (news / biomedical etc)

Generally hard to only get accuracy of 66% on Hyponymy & abt 55% on meronymy (not acceptable)

B

#

Bootstrapping approach

↳ we need something to bootstrap on.

- Not too much data is available for full fledged supervised training.
- But we have a small seed set of annotated data.
- We follow semi-supervised approach.

Ex: Target relation: Burial Place

Seed set: [Mark Twain, Elmira]

→ Expanding the knowledge using seed set:

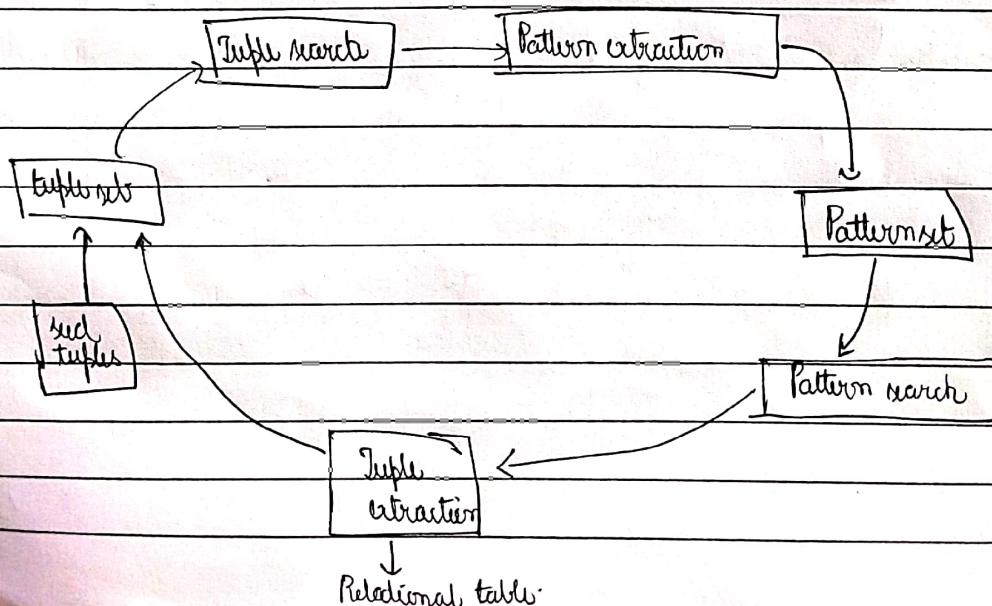
Google search "Mark Twain" & "Elmira".

We get results like "Mark Twain is buried in Elmira"

"The grave of Mark Twain is in Elmira" & "Elmira is Mark Twain's final resting place"

→ Build Harst pattern's using the obtained results. e.g. X is Y's final resting place.

→ Search Google for other patterns to obtain more tuples to expand the seed set & repeat.



Disadvantages

- Extremely sensitive to the seed set.
- Large number of hyperparameters (when to stop, how many search results to take, how many features to add)
- Probabilistic notion is missing (Not sure how confident we are with the results).

Supervised relation (information) extraction

- Choose a set of relations that you want to discover.
- Find & label data. (Expensive)
 - * choose a representative corpus
 - * ~~label entities~~ label entities
 - * label relations (expensive because needs expert intervention)
- Divide into train, dev & test
- Train a supervised ML model & predict ~~the~~ on the test set.

We model supervised IE as a classification task

We break it into a two step process

① Given a pair of entities, we predict if they're related or not

② If a pair is related, predict the relation type.

↳ multi-label classifier / multi-class classifier

Feature engineering for supervised learning.

→ consider the entity as bag of words

Consider the example: American Airlines, a unit of AMR, followed with the same move, spokesperson Tim Wagner said.

→ Feature 1: Bag of words

W M1: {American, Airlines}.

W M2: {Tim, Wagner}.

→ Feature 2: Head words.

H M1: American H M2: Wagner, H M12: Airlines + Wagner

→ Feature 3: left context / right context

M2-1: spokesperson M2+1: said.

→ Feature 4: Span of bag of words b/w the entities (models long range dependency)

{a, AMR, of, unit, followed, with, move, said, spokesperson}

→ Feature 5: Category of the named entity (gives type of the relation)

M1 → org M2 → PER.

M12 → ORG + PER.

→ Feature 6: Name / Nominal / Pronoun.

M1: EL → name

M2: EL → name

Needed to correctly link the pronouns to their named entities

→ Feature 7: Dependency parse (see slides for tree)

Can we the parse in the dependency tree (give dependency info)

Ex: HIDW1: followed: airlines

HIDW2: said: Wagner.

also: path: {airlines, followed, said, Wagner}.

→ Feature 8 : constituency parse features (See slides for this)

→ ~~Find~~ Feature 9: Triggers words which are possible sources of relationship.

Ex: parents, father, wife, husband

Can be recognized using Wordnet & some gazetteers like
country name list for geopolitical relations.

Use these features on standard ML models like SVM.