

## Concept level ILP.

X text unit.

X sentence level

✓ concept level.

The important concepts that a ~~summary~~ summary should cover.

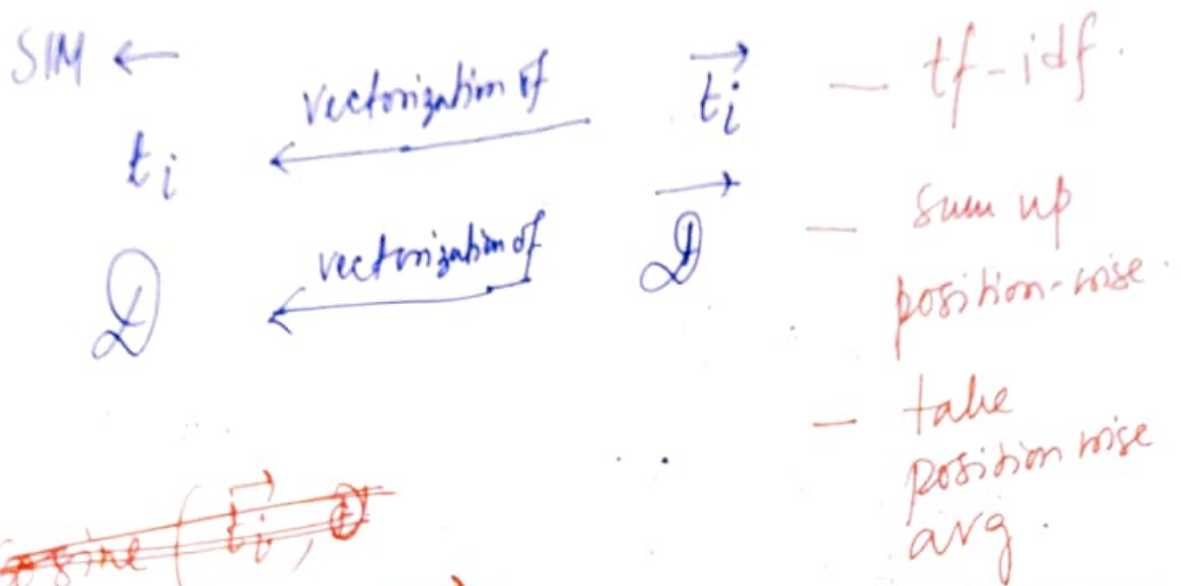
① Summary benefits by including a particular concept only once!

② That means redundancy is implicitly captured.

$c_i$  ← ~~idea~~ indicator variable for the <sup>presence of</sup> concept <sub>1</sub>  $i$  in the target summary.

$w_i$  ← weight of this concept  $i$

ConceptNet. { Concept → "word bigrams"  
Weight → no. of occurrences of the word bigrams in the input documents.



~~$\text{Cosine}(\vec{t_i}, \vec{Q})$~~

$SIM = \text{Cosine}(\vec{t_i}, \vec{Q})$

$Rel(i) = SIM(t_i, Q) \rightarrow \text{Cosine}(\vec{t_i}, \vec{Q})$

Query focused summarization:

$\rightarrow$  given a query  $Q \rightarrow$  relevant summary.

$Rel(i) = SIM(t_i, Q) + SIM(t_i, Q)$

## Sentence level ILP formulation:

Optimization function:

$$\text{maximize } \sum_i \alpha_i \text{Rel}(i) - \sum_{i < j} \alpha_{ij} \text{Red}(i, j).$$

Constraints.  $(\forall i, \forall j)$

$$\alpha_i, \alpha_{ij} \in \{0, 1\} \quad (1)$$

← indicator variables.

~~$\sum_i \alpha_i \text{Rel}(i) \leq K$~~

$$\sum_i \alpha_i \text{Rel}(i) \leq K \quad (2)$$

← budget constraint.

whether or not a textual unit or a pair are included in the target summary or not.

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

$$\alpha_{ij} - \alpha_j \leq 0$$

← (3)/(4)

if  $t_i$  &  $t_j$  are both included then they must be individually included.

$$\alpha_i + \alpha_j - \alpha_{ij} \leq 1 \quad (5)$$

← inverse of (3) & (4).

How to operationalize  $Rel(\cdot)$  &  $Red(\cdot, \cdot)$ .

Domain for summarization.

# News summarization.

~~( $Rel(\cdot)$ )~~

$\mathcal{D} \leftarrow$  Document collection.

$D$  is a single document.

$$Rel(i) = \frac{POS(t_i, D)^{-1}}{(t_i \in D, D \in \mathcal{D})} + \frac{SIM(t_i, \mathcal{D})}{(t_i \in D, D \in \mathcal{D})}$$

$POS(i, D) \leftarrow$  position of the text unit  $t_i$  in the Document  $D$ .

$SIM(t_i, \mathcal{D}) \leftarrow$  similarity of text unit  $t_i$  with the overall collection of documents

→ Empirical obs: Initial sentences in a news doc are usually very important (Headlines)



# Optimization based summarization.

## Global inferencing method.

$D$  is a document

$t_n$  ~~number~~ textual units.

$$D = t_1, t_2, \dots, t_n.$$

[textual  
units

↓  
individual sentences].

→  $Rel(i)$  : the relevance of  $t_i$  to the target summary.

→  $Red(i, j)$  : Redundancy ~~of~~ between  $t_i$  &  $t_j$ .

→  $l(i)$  is the length of  $t_i$  (in terms of the no. of words).

## Inferencing:

Problem is to select a subset  $S$  of textual units from  $D$  such that the summary score  $\phi$ , i.e.,  $\phi(S)$  is maximized.

$$\phi(S) = \max_{S \subseteq D} \left[ \sum_{t_i \in S} \text{Rel}(i) - \sum_{t_i, t_j \in S, i < j} \text{Red}(i, j) \right]$$

↑  
increase the relevance of  $t_i$  to the target summary  $S$

↘  
reduce the redundancy ~~between~~ among the choice of text ~~units~~ units.

→ Greedy ~~or~~ approaches give us approximate results.

→ Better solutions.

Integer Linear Programming (ILP).

→ GNU ← ILP solver.

---

Recast the problem into a

constraint optimization formulation.

Rel(i) , Red(i,i) ←

↑  
presence  
of some  
special  
entities.  
etc.

↑  
enforces  
diverse  
information.

### Greedy Solution.

1. Sort D so that  $Rel(i) > Rel(i+1) \forall i$
2.  $S = \{t_1\}$  ← most relevant text unit
3. while  $\sum_{t_i \in S} \underline{l(i)} < K$  ← budget
4.  $t_j = \arg \max_{t_j \in D - S} f(S \cup \{t_j\})$
5.  $S = S \cup \{t_j\}$
6. return S.



## Evaluation of summaries.

ROUGE score.

RUT

(Recall oriented Understudy for  
Gisting Evaluation)

How much ← "coverage"

Precision ← Safety-critical system.  
Medical diagnostics  
(Decision support  
systems).

## Revised ILP.

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

subject to:

$$\sum_j l_j s_j \leq K \quad (1)$$

length of  
the  
target summary  
(budget)

$c_i \leftarrow$  indicator of  
concept  $i$  in the target  
summary

$s_i \leftarrow$  indicator of  
sentence  $i$  in the  
target summary.

$Occ_{ij} \leftarrow$  occurrence of  
concept  $i$  in sentence  
 $j$

$$s_j Occ_{ij} \leq c_i \quad \forall i, j \quad (2)$$

$\hookrightarrow$  If you select a sentence then it  
mandates that all concepts in that  
sentence should be selected.

$$\sum_j s_j Occ_{ij} \geq c_i \quad (3)$$

$$c_i \in \{0, 1\} \quad \forall i \quad (5)$$

$$s_j \in \{0, 1\} \quad \forall j \quad (6)$$

If a concept is ever selected then it  
mandates that the sentence containing  
it must be selected.

## ROUGE-1

### ROUGE-1 precision:

$$\frac{3}{5} = 0.6$$

### ROUGE-1 recall:

$$\frac{3}{6} = 0.5$$

### ROUGE-2 precision:

$$\frac{1}{4} = 0.25$$

### ROUGE-2 recall:

$$\frac{1}{5} = 0.2$$

Example:

R: The cat is on the mat  
C: The cat and the dog.

(the, cat, the)

R: {the cat, cat is, is on, on the, the mat}

C: {the cat, cat and, and the, the dog}

# ROUGE - N.

$N = 1, 2, \dots$

$N \rightarrow N\text{-gram}$

M/c generated summary (Candidate summary)

C

Reference summary  $\rightarrow$  Human written

summary  $\rightarrow$  gold standard summary - R.

R: The cat is on the mat

C: The cat and the dog.

ROUGE-N precision

Ratio of the number of  $N$ -grams in C that also appear in R over the number of  $N$ -grams in C.

ROUGE-N recall (More important)

Ratio of the number of  $N$ -grams in C that also appears in R over the no. of  $N$ -grams in R.

$$\text{ROUGE F1} = \frac{2 \times \text{precision} \times \text{recall}}{(\text{precision} + \text{recall})}$$

$$= 0.22.$$

ROUGE-L

*L: Longest common subsequence.*

longest sequence — not necessarily consecutive.

R: The cat is on the mat } ← nuggets  
 C: the cat and the dog } ← nuggets.

LCS = the cat the.

$$|LCS| = 3.$$

ROUGE-L precision

ROUGE-L recall

$$= \frac{3/5 = 0.6}{3/6 = 0.5}$$

Numerator is the LCS of R, C.  
 Denominator is the unigram count.



ROUGE-S.

(Skip-Grams)

R: The cat is on the mat

C: The ~~gony~~ cat and the dog.

• ROUGE-2. (Skip = 1)  
↳ unigram shipping)

The cat  
The gony cat } match.



ILP summary.

Selected some sentences / text units.

Ordering.

News summarization:

Chronological ordering.

Coherence.

→ Choose orderings that make neighboring sentences / text units similar (cosine).

→ Multi entity summary:

Choose orderings that bring similar/same entities closer in the summary.

Topicality.

Make the ~~sum~~ summary topically coherent

↓ finding the topics.



ILP summary.

Selected some sentences / text units.

Ordering.

News summarization:

Chronological ordering.

Coherence.

→ Choose orderings that make neighboring sentences / text units similar (cosine).

→ Multi entity summary:

Choose orderings that bring similar/same entities closer in the summary.

Topicality.

Make the ~~sum~~ summary topically coherent

↓ finding the topics.

## Simplifying sentences.

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

Initial adverbials: ~~that~~ "On the other hand", "as a matter of fact".

PPs without named entities:

↳ prepositional phrases.

E.g. The commercial fishing restrictions in Washington will not be lifted ~~to~~ unless the salmon population increases to a substantial number. (PP removed).

Attribution clauses:

E.g. Rebels agreed to talks with govt officials, (international observers said Tuesday). <sup>attribution</sup>

Appositives:

Rajan, 28, ~~an artist who was living at~~  
~~the time in Philadelphia,~~ found the  
inspiration in the back of city magazines.