# Concept level ILP.

X text unit
X sentence level
✓ concept level.

The important concepts that a ~~summary~~
summary should cover.
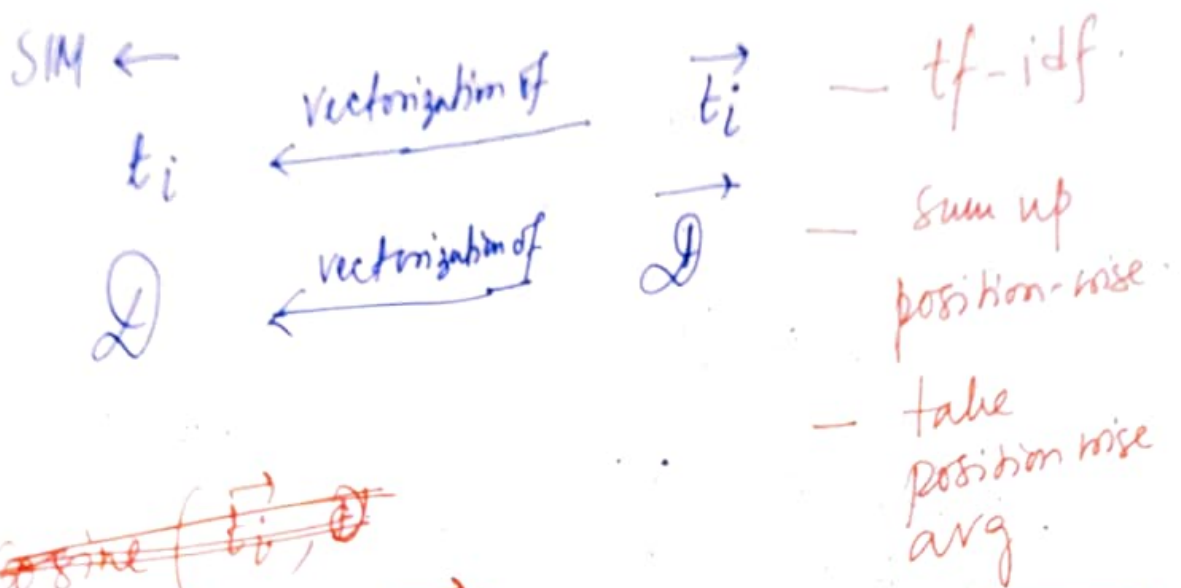
① Summary benefits by including a
particular concept only once!

② That means ___redundancy is implicitly___
___captured.___

$c_i$ ← ~~thing~~ indicator variable for the concept
i in the target summary.       presence of
the concept

$w_i$ ← weight of this concept i

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

$$\text{SIM} \leftarrow$$

$$t_i \xleftarrow{\text{vectorization of}} \vec{t_i} \quad - \text{ tf-idf}$$

$$\mathcal{D} \xleftarrow{\text{vectorization of}} \vec{\mathcal{D}} \quad - \begin{array}{l} \text{sum up} \\ \text{position-wise} \end{array}$$

$$\quad - \begin{array}{l} \text{take} \\ \text{position wise} \\ \text{avg} \end{array}$$

~~cosine $(\vec{t_i}, \vec{\mathcal{D}})$~~

$$\text{SIM} = \text{cosine}\left(\vec{t_i}, \vec{\mathcal{D}}\right)$$

$$\boxed{\text{Red}(ij) = \text{SIM}(t_i, t_j') \longrightarrow \text{cosine}\left(\vec{t_i}, \vec{t_j}\right)}$$

## Query focused summarization

$$\hookrightarrow \text{given a query } \mathcal{Q} \rightarrow \begin{array}{l} \text{relevant} \\ \text{summary} \end{array}$$

$$\text{Rel}(i) = \text{SIM}(t_i, \mathcal{Q}) + \text{SIM}(t_i, \vec{\mathcal{D}})$$

# Sentence level ILP formulation:

## Optimation function:

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

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

$$\alpha_i, \alpha_{ij} \in \{0, 1\} \quad (1) \longleftarrow \text{indicator variables.}$$

~~$\sum_i \alpha_i \, l(i)$~~

$$\sum_i \alpha_i \, l(i) \leq K \quad (2) \nearrow \text{budget constraint.}$$

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

$$\boxed{\begin{array}{l} \alpha_{ij} - \alpha_i \leq 0 \\ \alpha_{ij} - \alpha_j \leq 0 \end{array}} \quad (3),(4) \nwarrow \text{if } t_i \, \& \, t_j \text{ are both included then they must be individually included.}$$

$$\alpha_i + \alpha_j - \alpha_{ij} \leq 1 \quad (5) \leftarrow \text{inverse of } (3) \& (4)$$

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

Domain for summarization.

\# Newsummarization.

~~(Rel())~~

$\mathcal{D}$ ← Document collection.

D is a single document.

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

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

$SIM(t_i, \mathcal{D})$ ← 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

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

[textual units ↓ individual sentences].

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

→ $Red(i,j)$ : Redundancy 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 $\vartheta$, i.e., $\vartheta(S)$ is maximized.

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

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

reduce the redundancy among the choice of text units.

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

→ Better solutions.

↑

Integer Linear Programming (ILP).

→ . GNU ← ILP solver.

---

Recast the problem into a

↑

constraint optimization formulation.

$Rel(i)$ , $Red(i,j)$ ←

↑
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 $\sum_{t_i \in S} \underline{\ell(i)} < K$ ← —— budget

4. $t_j = \arg\max_{t_j \in D - S} \phi(S \cup \{t_j\})$

5. $S = S \cup \{t_j\}$

6. return $S$.

# Evaluation of summaries.

**ROUGE** Score

(**Recall** oriented Understudy for Gisting Evaluation)

ROUGE → RUJ

How much ← "Coverage"

Precision ← Safety - Critical System.
Medical diagnostics
(Decision support systems).

# Revised ILP.

maximize $\sum_i w_i c_i$

subject to:

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

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

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

$\swarrow$ length of the target summary (budget)

$Occ_{ij} \leftarrow$ 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 mandate 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

$N = 1, 2, \ldots$

$N \rightarrow$ N-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 | ROUGE-N recall (More important) |
|---|---|
| Ratio of the number of N-grams in C that also appears in R over the number of N-grams in C. | 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

$$|LCS| = 3.$$

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

ROUGE-L precision $= \left( \dfrac{3/5 = 0.6.}{3/6 = 0.5.} \right.$

ROUGE-L recall

# ROUGE-S.

R: The cat is on the mat

C: The gray cat and the dog.

· ROUGE-2.  (skip = 1)

           ↳ unigram shipping).

The cat  } match.

The (gray) cat

## 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).

↓ finding the topics.

## Topicality.
Make the ~~coh~~ summary topically coherent

## 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).

↓ finding the topics.

## Topicality.

Make the ~~some~~ summary topically coherent

# Simplifying sentences.

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

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

PPs <u>without named entities</u>:
↳ prepositional phrases.

E.g. The commercial fishing restrictions in Washington will not be lifted ⊗ 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)~~ attribution.

Appositives:

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