

RAKE  $\Rightarrow$  Rapid Automatic Keyphrase extraction

SAARC is consist of South Asian countries, India is supporting SAARC.

## ① candidate generation

Splits data  $\rightarrow$  Stopword  
punctuation

List of candidates

[ SAARC, consist, South Asian countries, India,  
supporting SAARC ]

## ② candidate scoring

$$\text{Score} = \frac{\text{Degree}}{\text{Frequency}}$$

unique words	Frequency	Degree	Score
SAARC	2	2	$2/2 = 1$
consist	1	1	$1/1 = 1$
South	1	0	$0/1 = 0$
Asian	1	1	1
countries	1	1	1
India	1	1	1
supporting	1	0	0

	SAARC	consist	South	Asian	countries	India	supporting
SAARC	1	0	0	0	0	0	0
consist	0	1	0	0	0	0	0
South	0	0	0	1	0	0	0
Asian	0	0	0	0	1	0	0
countries	0	0	0	0	0	0	0
India	0	0	0	0	0	1	0
supporting	1	0	0	0	0	0	0
Degree	2	1	0	1	1	1	0

[SAARC, consist, South Asian countries, India, supporting SAARC]

"SAARC SAARC"

"SAARC consist"

"SAARC"

"consist consist"

"consist"

"South South"

"South"

Score of Keyphrases:

$$\begin{aligned}
 \text{Score (South Asian countries)} &= \text{Score of South} + \text{Score of Asian} + \text{Score of countries} \\
 &= 0 + 1 + 1 \\
 &= 2
 \end{aligned}$$

$$\text{Score (Supporting SAARC)} = \text{Score of Supporting} + \text{Score of SAARC}$$

$$= 0 + 1$$

$$= 1$$

RAKE  $\Rightarrow$  Keyphrase Score  $\uparrow \Rightarrow$  Keyphrase important  $\uparrow$

YAKE  $\rightarrow$  Yet Another Keyphrase extractor.

SAARC is <sup>x</sup>consist of South Asian countries,  
<sup>x</sup>India is supporting SAARC.

① candidate generation

Splits data  $\rightarrow$  <sup>x</sup>stopword <sup>x</sup>\_\_\_\_\_ <sup>x</sup>stopword  
<sup>x</sup>punctuation <sup>x</sup>\_\_\_\_\_ <sup>x</sup>stopword

candidates  $\rightarrow$  SAARC, South Asian countries,  
 supporting SAARC

② candidate scoring

YAKE  $\rightarrow$  For every word 5 scores

$$\begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix}$$

a - score

$$\textcircled{a} \text{ casing}(w) = \frac{\max(\text{count}(w \text{ is capital}), \text{count}(w \text{ is acronym}))}{1 + \log(\text{count}(w))}$$

Acronym = SAARC, NGO

Capital = Title case

$$\text{count}(w) = 4 + 3 + 2 = 9$$

SAARC = 4 times

Saarc = 3 times

saarc = 2 times

$$\text{casing}(\text{SAARC}) = \frac{\text{maxc}(3, 4)}{1 + \log(9)} = \frac{4}{1 + \log(9)}$$

b - score

⑥ word position =  $\log(\log(3 + \text{median}(\text{sen}(w))))$

$\text{sen}(w)$  = list of positions of word

1	SAARC	
<hr/>		
10	SAARC	SAARC <sup>19</sup>
<hr/>		
25	SAARC	SAARC <sup>36</sup>
<hr/>		

$$\text{sen}(\text{SAARC}) = [1, 10, 19, 25, 36]$$

word position (SAARC)

$$= \log(\log(3 + \text{median}([1, 10, 19, 25, 36])))$$

$$= \log(\log(3 + 19))$$

c - score

⑦ word frequency score = count of word

$$\text{mean}(\text{count}) + \text{std.dev}(\text{count})$$

$$\text{word frequency (SAARC)} = \frac{\text{count of SAARC}}{\text{mean}([c_1, \dots, c_{100}]) + \text{std.dev}([c_1, \dots, c_{100}])}$$

$$1 \text{ document} = 100 \text{ words} = w_1, w_2, w_3, \dots, w_{100}$$

$$\text{count of word} = \begin{matrix} \downarrow & \downarrow & \downarrow & & \downarrow \\ c_1 & c_2 & c_3 & \dots & c_{100} \end{matrix}$$

$$\text{count} = [c_1, c_2, c_3, \dots, c_{100}]$$

d-score

(d) Word Relatedness

$$= 1 + (WR + WL) \times \frac{\text{count}(w) + PL + PR}{\text{max count}}$$

$$WR = \frac{\text{No. of unique words on Right}}{\text{Total words on Right}}$$

$$WL = \frac{\text{No. of unique words on left}}{\text{Total words on left}}$$

Max count = Total words in the document

$$PR = \frac{\text{Total words on Right}}{\text{Max count}}$$

$$PL = \frac{\text{total words on left}}{\text{Max count}}$$

We are learning NLP and NLP is consist  
of both Machine learning and Deep learning.

$$\text{word Relatedness (consist)} = 1 + \left( \frac{6}{7} + \frac{6}{7} \right) \times \frac{1}{15} + \frac{7}{15} + \frac{7}{15}$$

$$WR = \frac{6}{7}, \quad WL = \frac{6}{7}$$

$$\text{Max count} = 15$$

$$PR = \frac{7}{15}, \quad PL = \frac{7}{15}$$

e-score

$$\textcircled{c} \text{ word different score} = \frac{\text{No. of sentences containing that term}}{\text{total no. of sentences}}$$

$$\text{word different (SAARC)} = \frac{10}{100}$$

$$\text{Score of } (w) = \frac{d \times b}{a + \left( \frac{c}{d} \right) + \left( \frac{e}{d} \right)}$$

Here,  $a$  = casing score

$b$  = Position score

$c$  = frequency score  
 $d$  = word relatedness  
 $e$  = word different score

$$\text{Keyphrase score} = \frac{\text{Product}(\text{score}(A), \text{score}(B))}{1 + (\text{sum of score of } A, B) \times \text{count}(\text{Keyphrase})}$$

$(A, B)$

Supporting SAARC

↳ "SS"

\_\_\_\_\_  
 "SS"  
 \_\_\_\_\_  
 "SS"  
 \_\_\_\_\_

↓  
 3

YAKE  $\Rightarrow$  Keyphrase score  $\downarrow \Rightarrow$  Keyphrase important  $\uparrow$

RAKE  $\Rightarrow$  Keyphrase score  $\uparrow \Rightarrow$  Keyphrase important  $\uparrow$