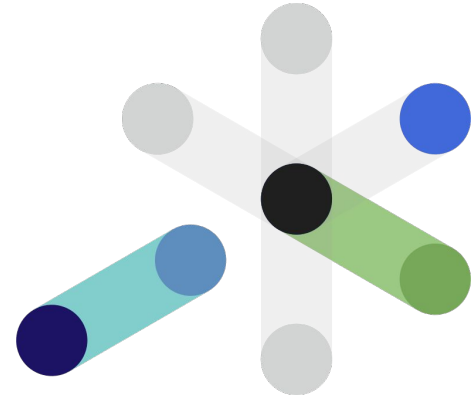


srijan:



Approaching Text Summarization using ML and DNN

RAIT-ACM STTP
25 May 2020



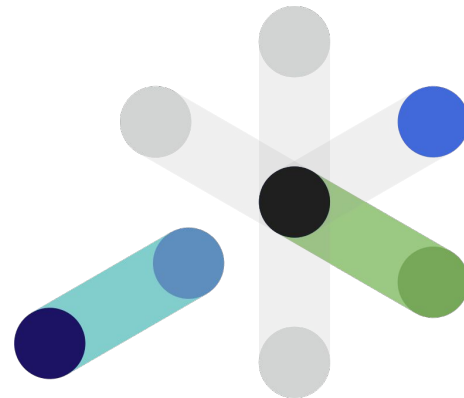
About Me

Mayank Kumar Jha

Data Scientist | Kaggle Competition Expert

Experience Across
Machine Learning, Deep Learning,
Data Ops, Cloud, Algorithms, Optimization

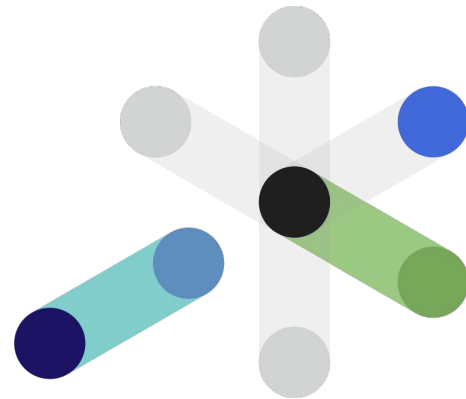
srijan:



Things to cover

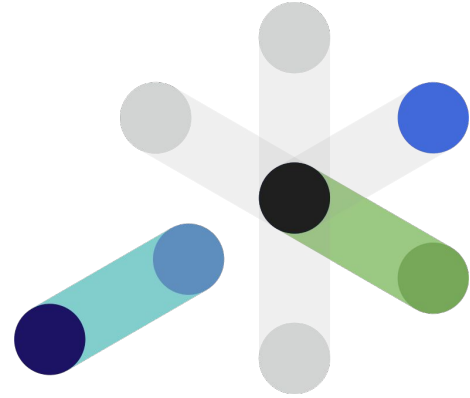
srijan:

- Introduction to Text Summarization
- Various approaches
- Propose possible solutions
- Create a basic solution
- Code Walkthrough
- Query session



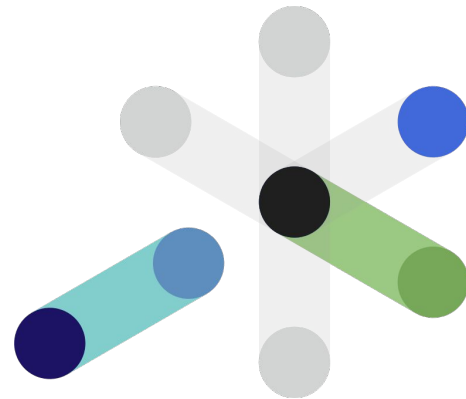
What is Text Summarization?

srijan:



What is Text Summarization?

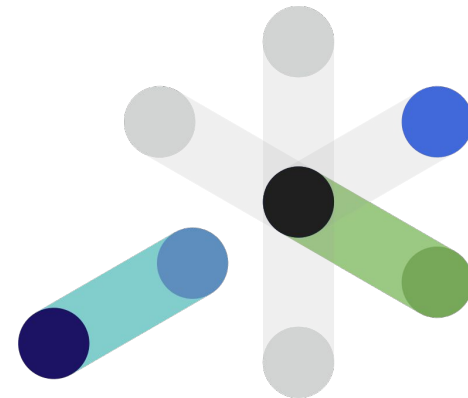
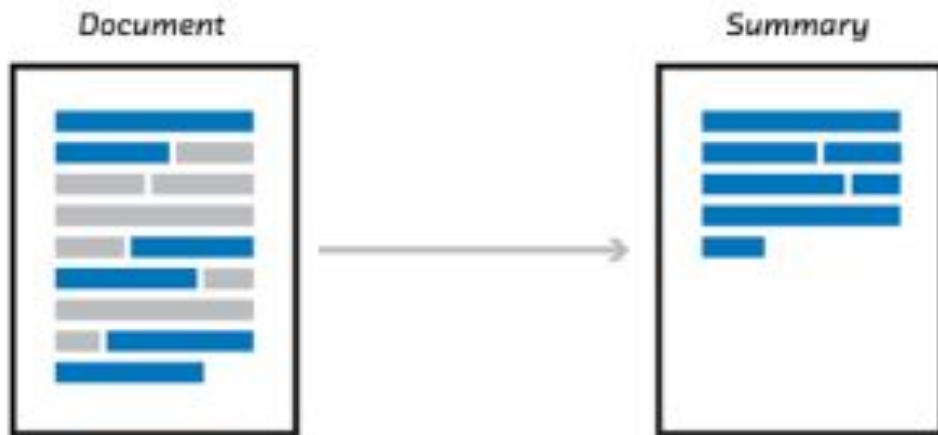
srijan:



What is Text Summarization?

srijan:

- Text summarization is the process of shortening a set of data computationally, to create a subset (a summary) that represents the most important or relevant information within the original content.^{wikipedia}
- Text summarization is the technique for generating a concise and precise summary of voluminous texts while focusing on the sections that convey useful information, and without losing the overall meaning.^{floydhub}

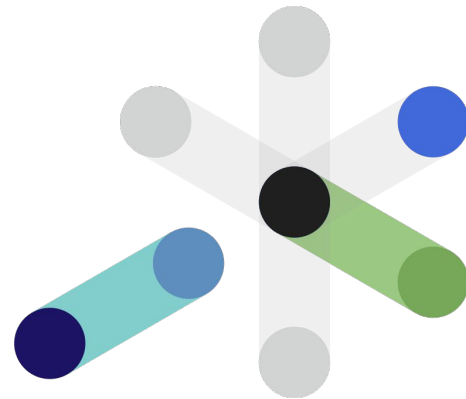


Various approaches for text summarization?

srijan:

- **Extractive Summarization**


- **Abstractive Summarization**



Various approaches for text summarization?

srijan:





- **Extractive Summarization**



Source Text:  Peter and Elizabeth took a taxi to attend the night party in the city.


While in the party, Elizabeth collapsed and was rushed to the hospital.

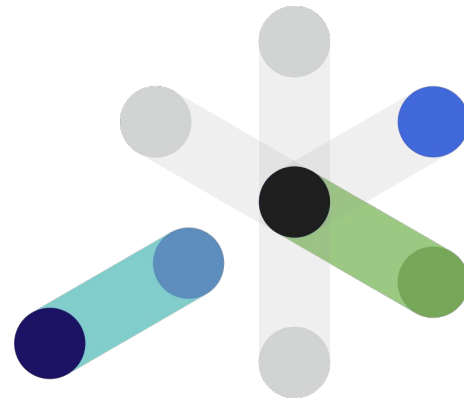
Summary: Peter

- **Abstractive Summarization**

Source Text:  and  took a taxi to  the night  in the city.

While in the party,  collapsed and was rushed to the .

Summary: Elizabeth was hospitalized after attending a party with Peter. 




Various approaches for text summarization?

srijan:

- **Extractive Summarization**

- Extractive summarization means identifying important sections (paragraphs or sentences or even words) of the text and selecting (copy paste) them producing a subset of the text from the original text.





Source Text:  Peter and Elizabeth took a taxi to attend the night party in the city.



While in the party, Elizabeth collapsed and was rushed to the hospital.


Summary: Peter

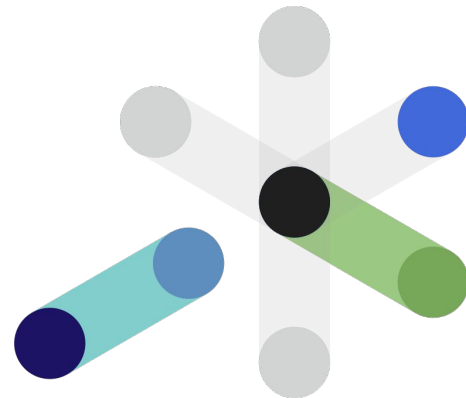
- **Abstractive Summarization**

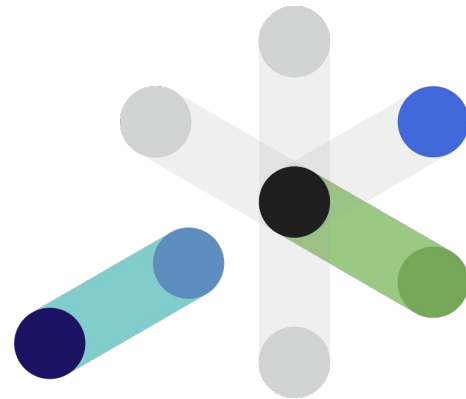
- Abstractive summarization is the technique of generating a summary of a text from its main ideas, not by copying verbatim most salient sentences from text.

Source Text:  and  took a taxi to  the night  in the city.

While in the party,  collapsed and was rushed to the .

Summary: Elizabeth was hospitalized after attending a party with Peter. 

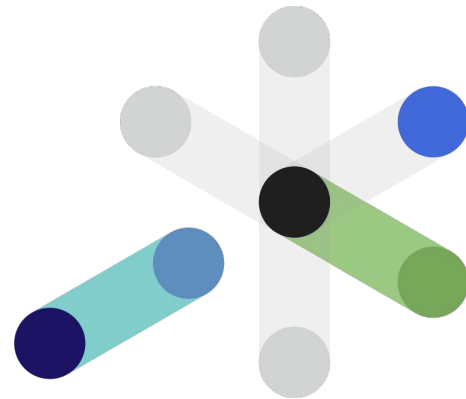




Approaching Extractive Summarization?

srijan:

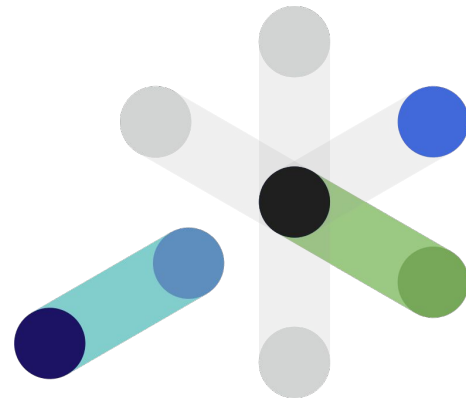
- One approach could be to create a semantic representation of sentences.



Approaching Extractive Summarization?

srijan:

- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer



Approaching Extractive Summarization?

srijan:

- What is TF-IDF Vectorizer:

```
documentA = 'the man went out for a walk'  
documentB = 'the children sat around the fire'
```

Count vectorization

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{i,j}}$$

	walk	for	out	sat	man	children	the	went	around	a	fire
0	0.142857	0.142857	0.142857	0.0	0.142857	0.0	0.142857	0.142857	0.0	0.142857	0.0
1	0.000000	0.000000	0.000000	0.2	0.000000	0.2	0.400000	0.000000	0.2	0.000000	0.2

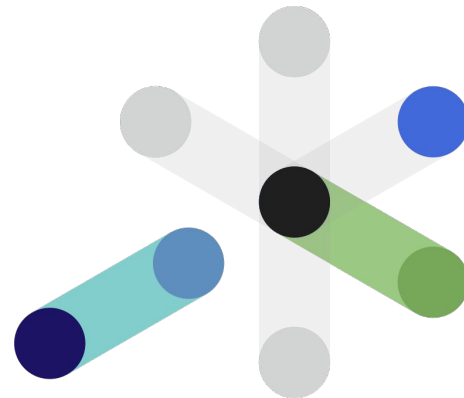
Inverse Document Frequency (IDF)

$$idf(w) = \log\left(\frac{N}{df_t}\right)$$

	walk	for	out	sat	man	children	the	went	around	a	fire
0	0.099021	0.099021	0.099021	0.000000	0.099021	0.000000	0.0	0.099021	0.000000	0.099021	0.000000
1	0.000000	0.000000	0.000000	0.138629	0.000000	0.138629	0.0	0.000000	0.138629	0.000000	0.138629

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

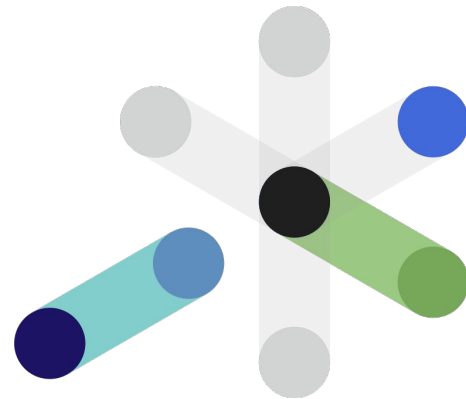
TF-IDF



Approaching Extractive Summarization?

srijan:

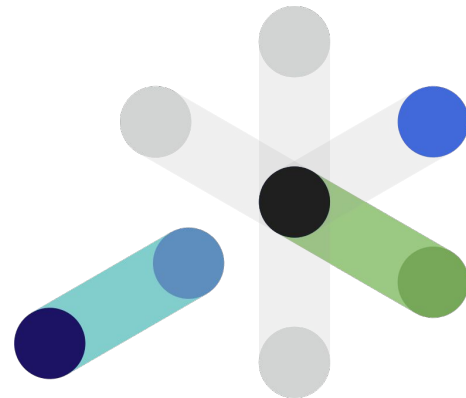
- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer



Approaching Extractive Summarization?

srijan:

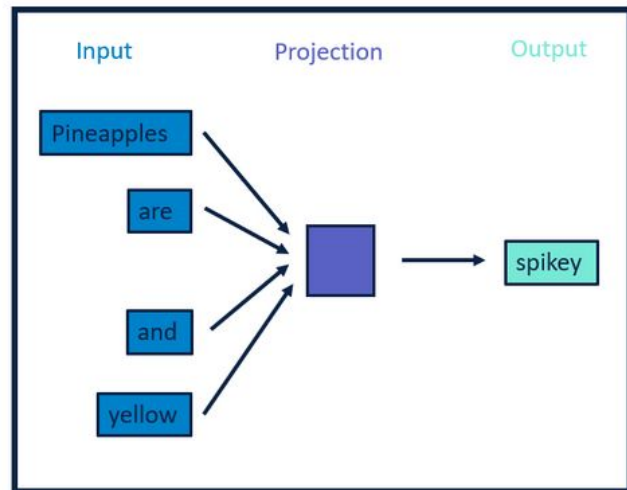
- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove and fastText



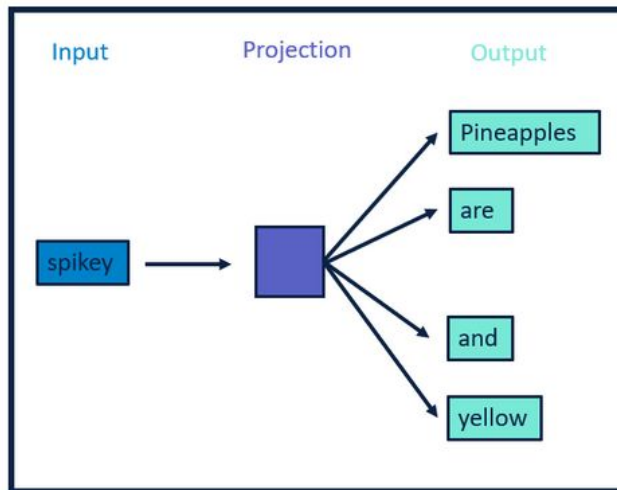
Approaching Extractive Summarization?

srijan:

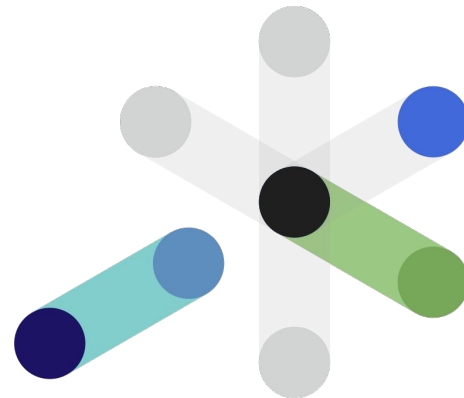
- **What is Word2Vec Vectorizer:**
- Word2vec is a group of related models that are used to produce word embeddings.
[wikipedia]
- Word2vec can utilize either of two model architectures to produce a distributed representation of words: continuous bag-of-words (CBOW) or continuous skip-gram.[wikipedia]
- Both of these techniques learn weights which act as word vector representations.



CBOW



Skip-gram



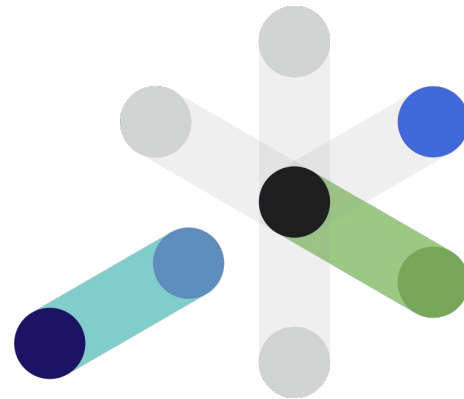
Approaching Extractive Summarization?

srijan:

- **What is CBOW:**
- In the continuous bag-of-words architecture:
 - model predicts the current word from a window of surrounding context words.^[wikipedia]
 - The order of context words does not influence prediction (bag-of-words assumption).^[wikipedia]

Input	Output
Pineapples	Spikey
are	Spikey
and	Spikey
yellow	Spikey

CBOW



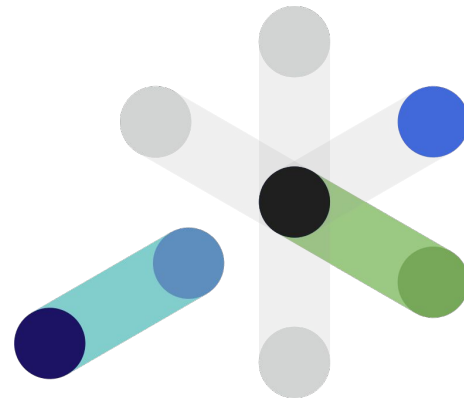
Approaching Extractive Summarization?

srijan:

- **What is SKIP-GRAM:**
- In the continuous skip-gram architecture:
 - model uses the current word to predict the surrounding window of context words. [wikipedia]
 - The skip-gram architecture weighs nearby context words more heavily than more distant context words. [wikipedia]
 - According to the authors' note, CBOW is faster while skip-gram is slower but does a better job for infrequent words. [wikipedia]

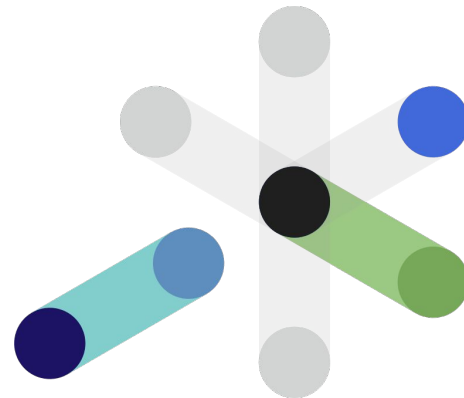
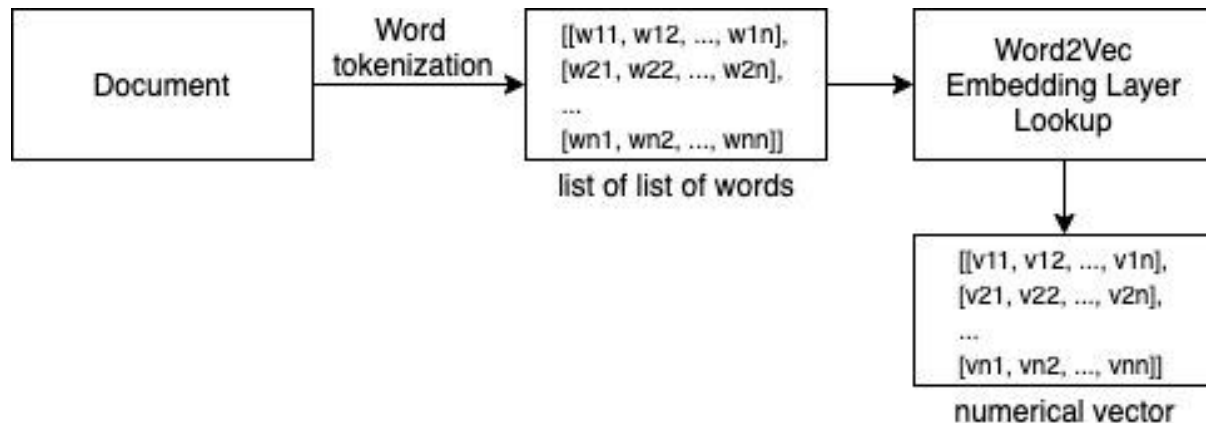
Input	Output
Spikey	Pineapples
Spikey	are
Spikey	and
Spikey	yellow

SKIP-GRAM



- **Using Word2Vec Embedding**

- Embedding is something like key-value pair.
- For each word in our vocabulary, there will be a learned numerical representation for it.
- To handle unknown words, we will treat each of them as OOV (Out of vocabulary) words and will use the learned OOV representation
- Flow would be something like below one:

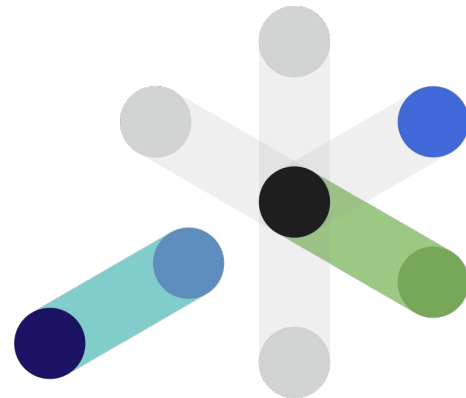


Using Word2Vec to convert text data to numerical semantic vector

Approaching Extractive Summarization?

srijan:

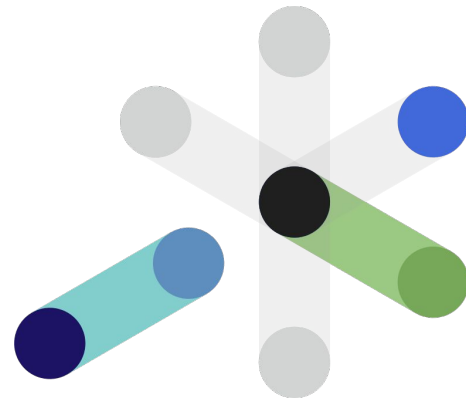
- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove, fastText



Approaching Extractive Summarization?

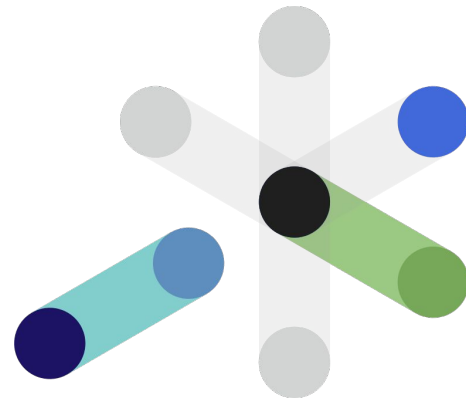
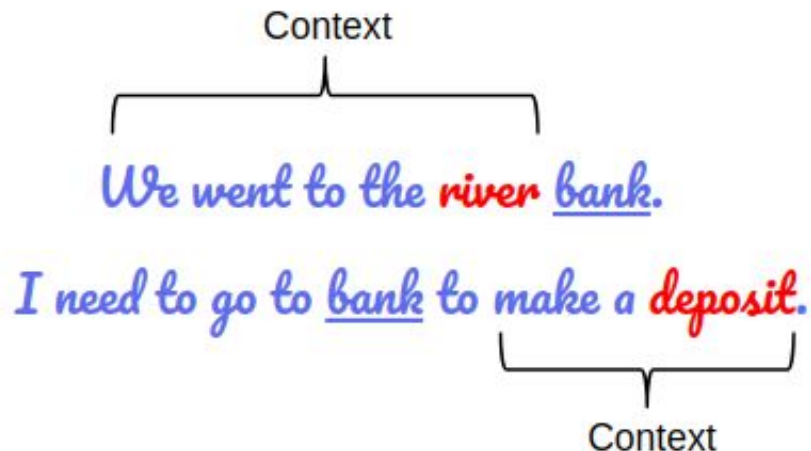
srijan:

- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove, fastText
 - Pretrained SOTA transformers like BERT (or its variants) to better capture context as well.



- **What is BERT:**

- BERT stands for Bi-directional Encoder Representation from Transformers.
- It is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context.
- BERT is pre-trained on two NLP tasks:
 - Masked Language Modeling
 - Next Sentence Prediction



Approaching Extractive Summarization?

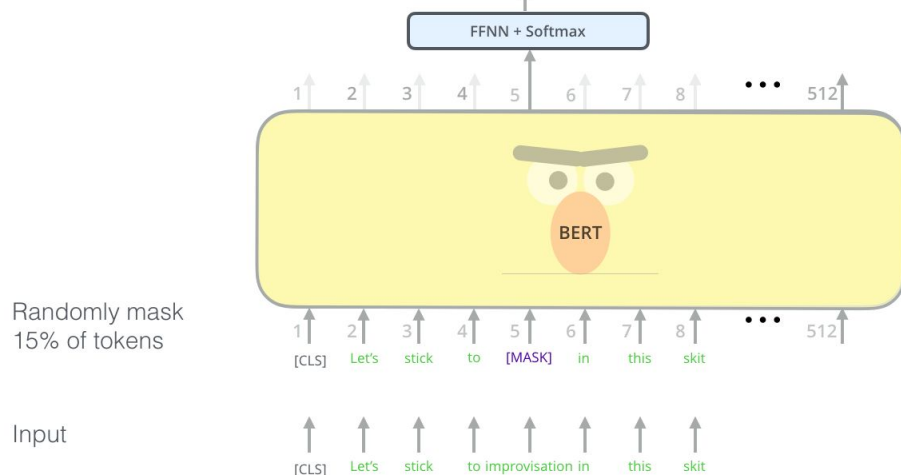
srijan:

- What is Masked Language Modeling:

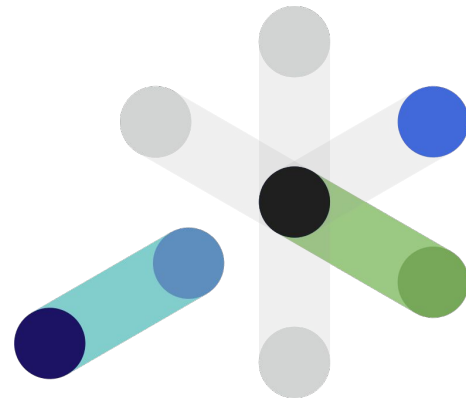
Use the output of the masked word's position to predict the masked word

Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzzzyva



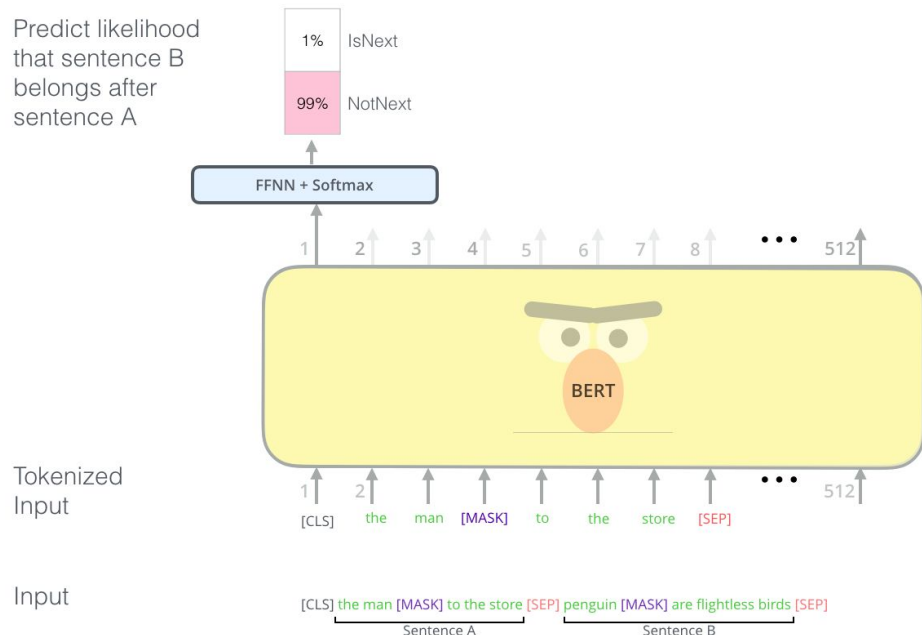
BERT's clever language modeling task masks 15% of words in the input and asks the model to predict the missing word



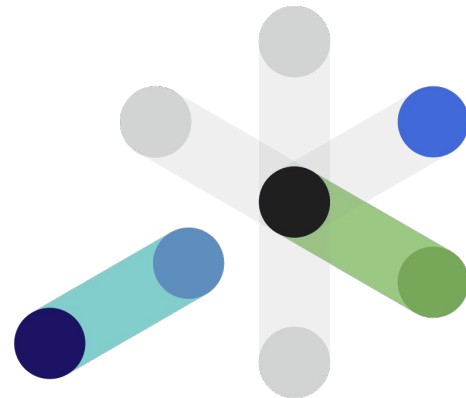
Approaching Extractive Summarization?

srijan:

- What is Next Sentence Prediction:



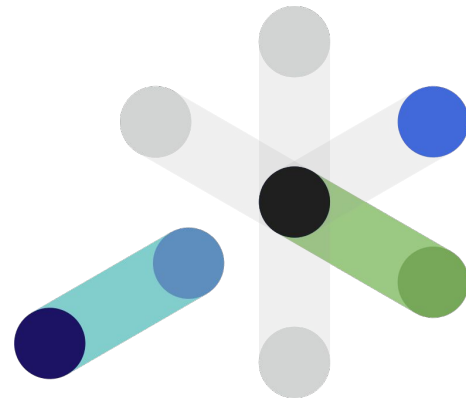
The second task BERT is pre-trained on is a two-sentence classification task. The tokenization is oversimplified in this graphic as BERT actually uses WordPieces as tokens rather than words --- so some words are broken down into smaller chunks.



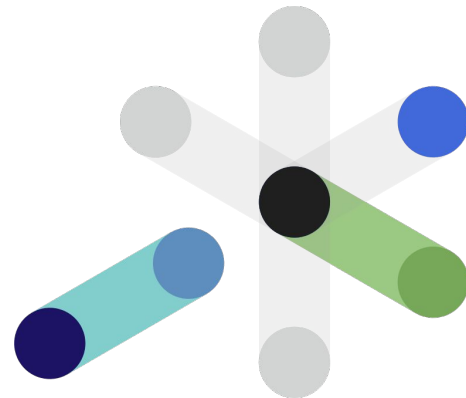
Approaching Extractive Summarization?

srijan:

- **Using BERT for feature extraction**
 - Feed the text to BERT Wordpiece Tokenizer

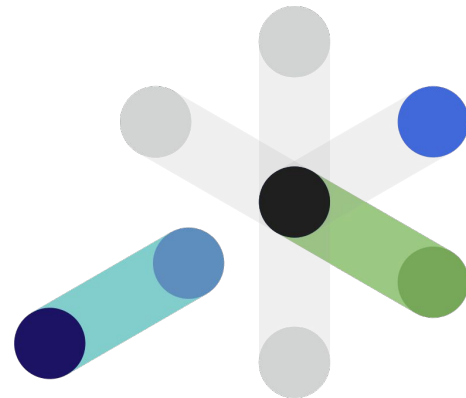


- **Using BERT for feature extraction**
 - Feed the text to BERT Wordpiece Tokenizer
 - BERT tokenizer will split the sentence as well as words incase needed to match vocabulary.



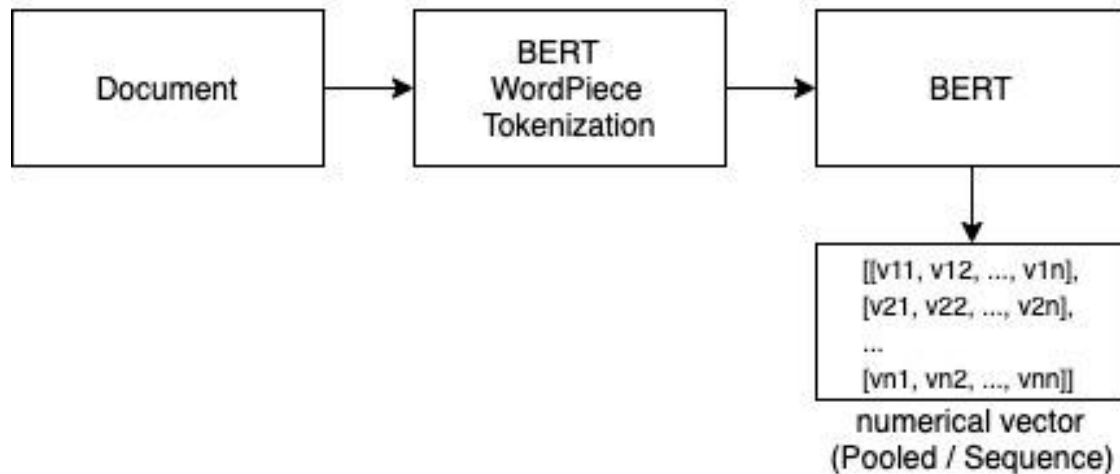
- **Using BERT for feature extraction**

- Feed the text to BERT Wordpiece Tokenizer
- BERT tokenizer will split the sentence as well as words incase needed to match vocabulary.
- Words are splitted on the basis of their probability of occurrences in that context which gives BERT advantages to handle contraction as well as spelling errors

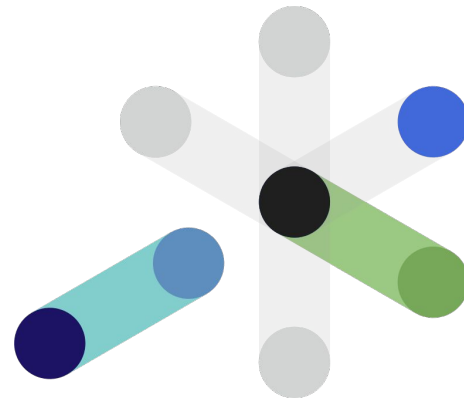


- **Using BERT for feature extraction**

- Feed the text to BERT Wordpiece Tokenizer
- BERT tokenizer will split the sentence as well as words incase needed to match vocabulary.
- Words are splitted on the basis of their probability of occurrences in that context which gives BERT advantages to handle contraction as well as spelling errors
- Flow would be something like below one:



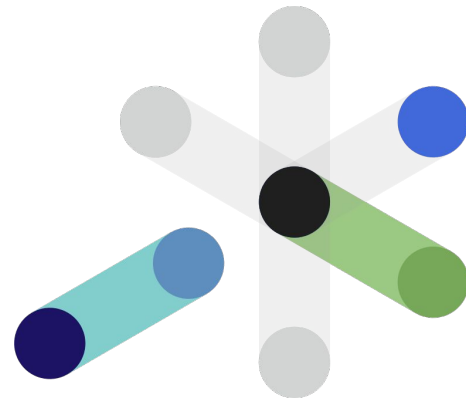
Using BERT to convert text data to numerical vector



Approaching Extractive Summarization?

srijan:

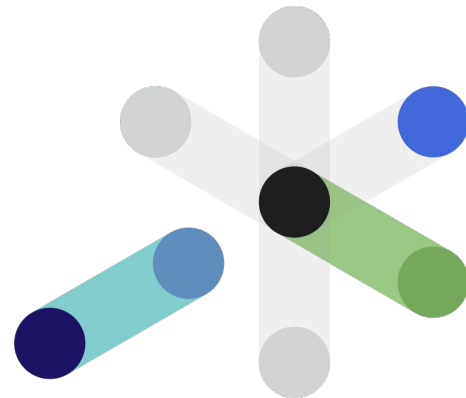
- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove, fastText
 - Pretrained SOTA transformers like BERT (or its variants) to better capture context as well.



Approaching Extractive Summarization?

srijan:

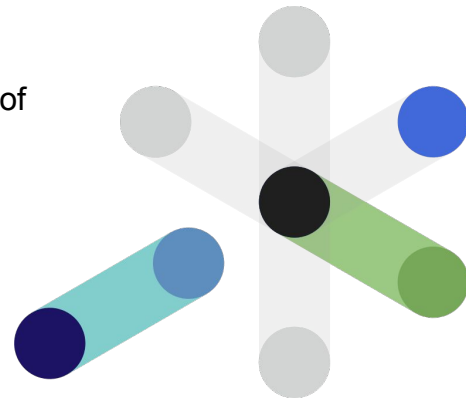
- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove, fastText
 - Pretrained SOTA transformers like BERT (or its variants) to better capture context as well.
 - Train your own (quite cumbersome :()



Approaching Extractive Summarization?

srijan:

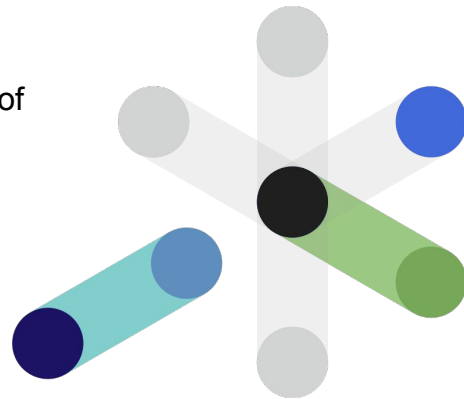
- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove, fastText
 - Pretrained SOTA transformers like BERT (or its variants) to better capture context as well.
 - Train your own (quite cumbersome :()
- Now comparison can be done between each sentences as they are no more a sequence of characters and words but a numerical vector now.



Approaching Extractive Summarization?

srijan:

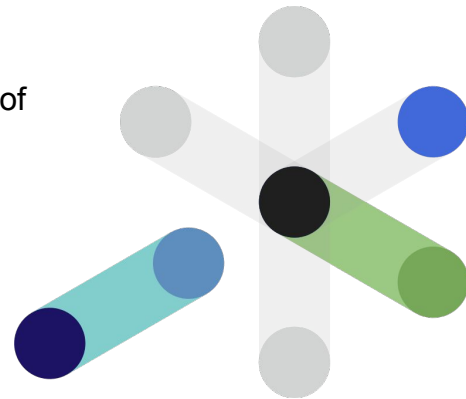
- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove, fastText
 - Pretrained SOTA transformers like BERT (or its variants) to better capture context as well.
 - Train your own (quite cumbersome :()
- Now comparison can be done between each sentences as they are no more a sequence of characters and words but a numerical vector now.
- Use the comparisons to score each sentences and pick the top scored ones as your summary.



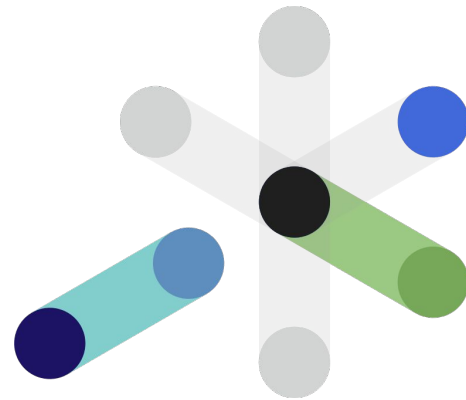
Approaching Extractive Summarization?

srijan:

- One approach could be to create a semantic representation of sentences.
- Following can be used to create a semantic representation for texts:
 - Count-based techniques like CountVectorizer, Tf-Idf Vectorizer
 - Pretrained word embeddings based techniques like Word2Vec, Glove, fastText
 - Pretrained SOTA transformers like BERT (or its variants) to better capture context as well.
 - Train your own (quite cumbersome :()
- Now comparison can be done between each sentences as they are no more a sequence of characters and words but a numerical vector now.
- Use the comparisons to score each sentences and pick the top scored ones as your summary.
- Scoring technique needs to be intelligent enough to properly evaluate what are the information content of a sentence and thus score it accordingly.

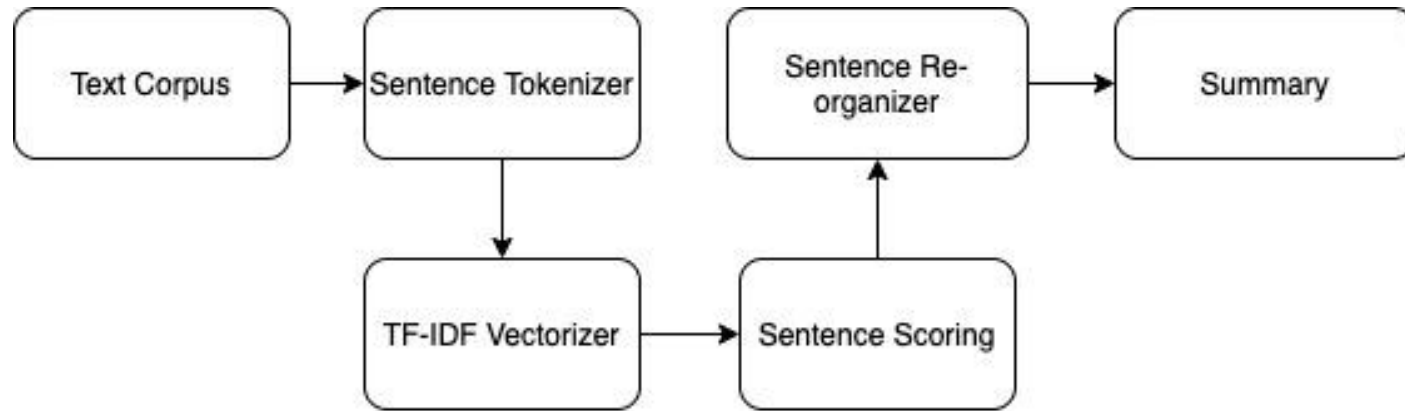


**Let's build our
basic Solution**

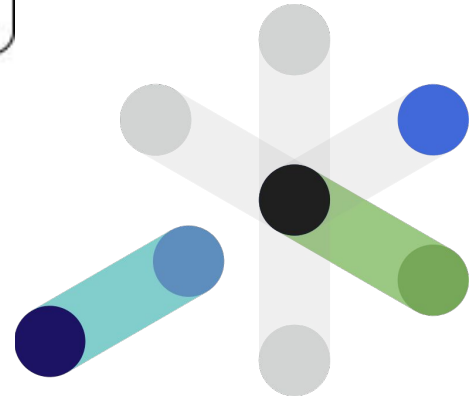


Approaching Extractive Summarization?

srijan:



A high level solution for frequency based extractive summarizer



Approaching Extractive Summarization?

srijan:

- Sentence scoring using TF-IDF:

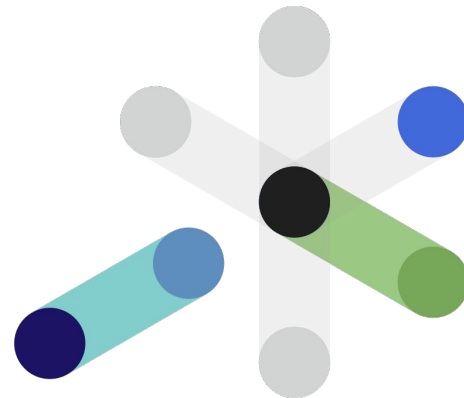
```
documentA = 'the man went out for a walk'  
documentB = 'the children sat around the fire'
```

TF-IDF

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

	walk	for	out	sat	man	children	the	went	around	a	fire
0	0.099021	0.099021	0.099021	0.000000	0.099021	0.000000	0.0	0.099021	0.000000	0.099021	0.000000
1	0.000000	0.000000	0.000000	0.138629	0.000000	0.138629	0.0	0.000000	0.138629	0.000000	0.138629

TF-IDF



Approaching Extractive Summarization?

srijan:

- Sentence scoring using TF-IDF:

```
documentA = 'the man went out for a walk'  
documentB = 'the children sat around the fire'
```

TF-IDF

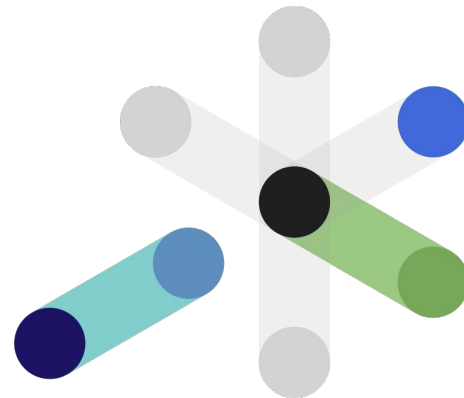
$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

	walk	for	out	sat	man	children	the	went	around	a	fire
0	0.099021	0.099021	0.099021	0.000000	0.099021	0.000000	0.0	0.099021	0.000000	0.099021	0.000000
1	0.000000	0.000000	0.000000	0.138629	0.000000	0.138629	0.0	0.000000	0.138629	0.000000	0.138629

Adding TF-IDF weights to
score each sentence

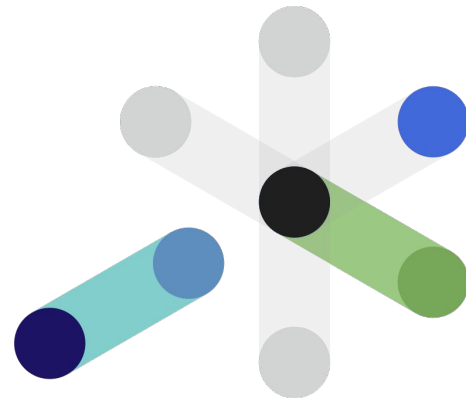
	walk	for	out	sat	man	children	the	went	around	a	fire	score
0	0.099021	0.099021	0.099021	0.000000	0.099021	0.000000	0.0	0.099021	0.000000	0.099021	0.000000	0.594126
1	0.000000	0.000000	0.000000	0.138629	0.000000	0.138629	0.0	0.000000	0.138629	0.000000	0.138629	0.554518

TF-IDF



- **What is Sentence Reorganizer:**

Text : Stop Words are words which do not contain important significance to be used in Search Queries. Usually, these words are filtered out from search queries because they return a vast amount of unnecessary information. Each programming language will give its own list of stop words to use.

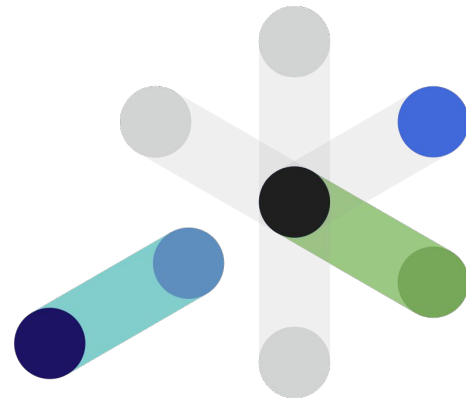


- **What is Sentence Reorganizer:**

Text : Stop Words are words which do not contain important significance to be used in Search Queries. Usually, these words are filtered out from search queries because they return a vast amount of unnecessary information. Each programming language will give its own list of stop words to use.



Sentence Order Indexing : { “Stop Words are words which do not contain important significance to be used in Search Queries”: 0,
“Usually, these words are filtered out from search queries because they return a vast amount of unnecessary information”: 1 ,
“Each programming language will give its own list of stop words to use”: 2 }

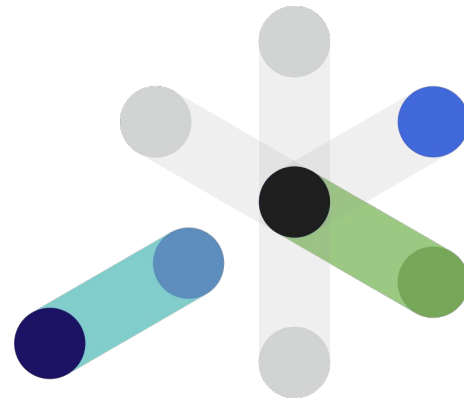


- What is Sentence Reorganizer:

Sentence Order Indexing : { “Stop Words are words which do not contain important significance to be used in Search Queries”: 0,
“Usually, these words are filtered out from search queries because they return a vast amount of unnecessary information”: 1 ,
“Each programming language will give its own list of stop words to use”: 2 }

Scored Sentences : { “Each programming language will give its own list of stop words to use”: 1.32
,
“Stop Words are words which do not contain important significance to be used in Search Queries”: 0.79,
“Usually, these words are filtered out from search queries because they return a vast amount of unnecessary information”: 0.32 }

Result : Stop Words are words which do not contain important significance to be used in Search Queries. Each programming language will give its own list of stop words to use.



Fire up your Notebooks



Any Questions ?

Thank You



srijan:

