

Latent Semantic Analysis

Latent Semantic Analysis is a technique of analyzing relationships between a set of documents and the terms they contain by producing a set of concepts related to the documents & terms.

→ We can generate keywords using LSA

1	MUSIC	FOOD
2	AR 1-85%	AR 3-100%
3		AR 5-73%
4		
5	AR 4-100%	AR 2-100%
6	AR 6-100%	AR 1-15%
	AR 5-27%	
	NEWS	TECH

Bow Model

 w1 w2 w3 ... wn

Doc 1

Doc 2

Doc 3

→ when Generating LSA there will be a huge corpus of document

SVD (Singular Value Decomposition)

$$A_{[m \times n]} = U_{[m \times r]} * S_{[r \times r]} * \left(V_{[n \times r]} \right)^T$$

A: Input Data Matrix

- $m \times n$ matrix (m = number of documents, n = number of words/features)

U = Left Singular matrix

- $m \times r$ matrix (m = number of documents, r = number of concepts)

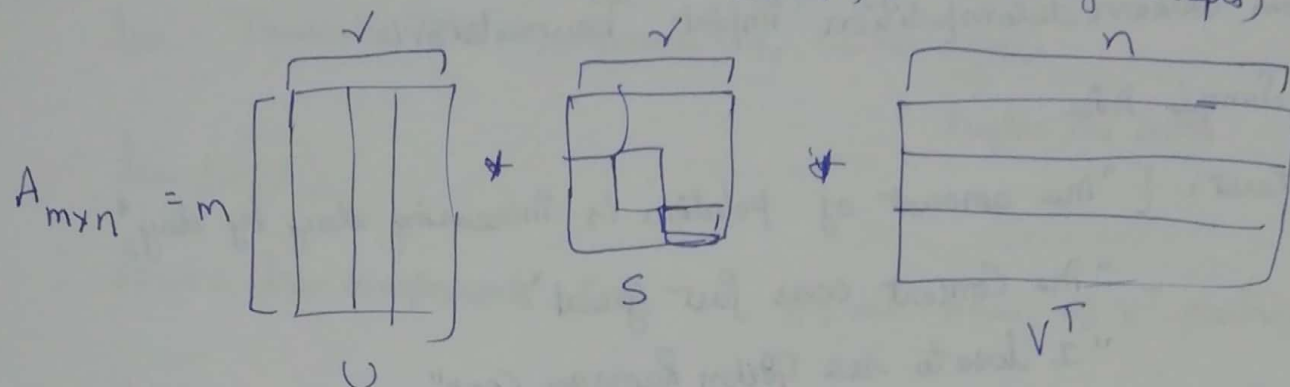
S = Rand matrix

- $r \times r$ matrix (r = rank of A)

→ Diagonal matrix

V = Right Singular matrix

$\rightarrow n \times r$ matrix (n = number of words/features, r = number of concepts)



$V^T = r \times n$ matrix.

Application

- \rightarrow Article Bucketing websites
- \rightarrow Finding relation b/w articles/words
- \rightarrow Page indexing in search engines.

LSA Implementation

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.decomposition import TruncatedSVD

Sample Data

dataset = ["The amount of pollution is increasing day by day,"

"The Concert was just great"

"I love to see Gordon Ramsay Cook"

"Google is introducing new technology"

"AI Robots are Examples of great technology present today"

"All of us were lining in the Concert"

"we have launch campaigns to stop pollution^{and} global warming"]

→ we will try to Create Concepts

→ whether it is able to create a Concept

→ ~~was~~ Important words related to a Concept

→ Identify different documents for a specific Content

→ Preprocessing

→ lower

dataset = [line.lower() for line in dataset]

→ Creating Bow model (Bow/TFIDF/Count Vector)

Vectorizer = TfidfVectorizer()

X = Vectorizer.fit_transform(dataset)

print(X[0])

o/p

(0,34) 0.22786438

(0,2) 0.32114839

(0,24) 0.2278643877

decomposition

we have to convert the matrix X into $U S V^T$

`lsa = TruncatedSVD(n_components=4, n_iter=100)`
↑
higher the better.

`lsa.fit(X)`

`row1 = lsa.components_[0]` // first row of V^T matrix

`term = Vectorizer.get_feature_names()` // gives all the words
`concept_words = {}`

for i, comp in `enumerate(lsa.components_)`:

`ComponentTerms = zip(term, comp)`

`sortedTerms = sorted(ComponentTerms, key = lambda x: x[1], reverse = True)`

`sortedTerms = sortedTerms[:10]`

`print("\nConcept", i, ":")`

for `terms` in `sortedTerms`:

`print(terms)`

`concept_words["Concept " + str(i)] = sortedTerms`

→ we have 7 documents, we classify which document falls into which concept

for `key` in `concept_words.keys()`

`Sentence_score = []`

For `Sentence` is dataset

`words = nltk.word_tokenize(sentence)`

`score = 0`

for `word` in `words`:

for `word_with_score` in `concept_words[key]`:

If `word == word_with_score[0]`:

~~for~~ for

Score += word-wim_score[1]

Sentence_score.append(Score)

print("\n" + key + ":")

for sentence_score in Sentence_scores:

print(Sentence_score)