

# Non-Negative Matrix Factorization (NMF) for Topic Modeling

---

## Overview

Non-Negative Matrix Factorization (NMF) is a dimensionality reduction technique that decomposes a non-negative matrix into two lower-dimensional non-negative matrices. In the context of Natural Language Processing (NLP), NMF is applied to uncover latent topics within a collection of documents by factorizing the document-term matrix. Each document is represented as a combination of topics, and each topic is characterized by a distribution over terms. [?cite?turn0search5?](#)

---

## Why Use NMF for Topic Modeling?

- **Simplicity and Interpretability:** NMF ensures that the components and coefficients are non-negative, leading to more interpretable results compared to other methods like Singular Value Decomposition (SVD). [?cite?turn0search4?](#)
  - **Sparsity:** The non-negativity constraint often results in sparse representations, making it easier to identify the most significant terms associated with each topic.
  - **Performance:** NMF has been shown to perform well in extracting coherent topics from text data, especially when the data is large and sparse. [?cite?turn0search3?](#)
- 

## Prerequisites

Before running the code, ensure you have the following installed:

- Python 3.x
- Required libraries:
  - `scikit-learn`
  - `numpy`
  - `matplotlib`

You can install the necessary libraries using pip:

```
pip install scikit-learn numpy matplotlib
```

---

## Files Included

- `nmf_topic_modeling.py` : The main script containing the implementation of NMF for topic modeling.
- 

## Code Description

The provided code demonstrates how to perform topic modeling using NMF with a small set of example documents.

1. **Import Necessary Libraries:**

```
from sklearn.decomposition import NMF
from sklearn.feature_extraction.text import TfidfVectorizer
```

- `NMF` from `sklearn.decomposition` is used to perform Non-Negative Matrix Factorization.
- `TfidfVectorizer` from `sklearn.feature_extraction.text` converts the collection of raw documents to a matrix of TF-IDF features.

## 2. Prepare the Document Corpus:

```
documents = [
    "Data science is a multidisciplinary field.",
    "Machine learning provides systems the ability to learn.",
    "Deep learning is a subset of machine learning.",
    "Artificial intelligence encompasses machine learning."
]
```

- A list of sample documents is defined for topic modeling.

## 3. Transform Documents into TF-IDF Matrix:

```
vectorizer = TfidfVectorizer(stop_words='english')
tfidf_matrix = vectorizer.fit_transform(documents)
```

- `TfidfVectorizer` is initialized with English stop words to remove common words that may not be informative.
- The `fit_transform` method is applied to the documents to create the TF-IDF matrix.

## 4. Apply NMF to Extract Topics:

```
nmf = NMF(n_components=2, random_state=42)
nmf.fit(tfidf_matrix)
```

- An NMF model is initialized to extract 2 topics ( `n_components=2` ).
- The model is fitted to the TF-IDF matrix.

## 5. Display the Top Terms for Each Topic:

```
for idx, topic in enumerate(nmf.components_):
    print(f"Topic {idx + 1}:")
    print([vectorizer.get_feature_names_out()[i] for i in topic.argsort()[-5:]])
```

- For each topic, the code retrieves the top 5 terms that contribute the most to that topic.
- `topic.argsort()[-5:]` returns the indices of the top 5 terms for the topic.
- `vectorizer.get_feature_names_out()` provides the mapping from indices to actual terms.

---

## Expected Outputs

When you run the code, you can expect output similar to:

```
Topic 1:
['data', 'field', 'multidisciplinary', 'science']
Topic 2:
['ability', 'provides', 'systems', 'learning', 'machine']
```

This output indicates the top terms associated with each of the two topics identified by the NMF model.

---

## Use Cases

- **Document Clustering:** Grouping similar documents based on underlying topics.
  - **Information Retrieval:** Enhancing search algorithms by understanding the main topics within documents.
  - **Content Recommendation:** Suggesting relevant content to users based on topic similarity.
- 

## Advantages

- **Interpretability:** The non-negativity constraint leads to more interpretable components.
  - **Sparsity:** Results in sparse representations, highlighting the most significant terms for each topic.
  - **Efficiency:** Effective for large, sparse datasets commonly encountered in text mining.
- 

## Future Enhancements

- **Dynamic Topic Number Selection:** Implement methods to automatically determine the optimal number of topics.
  - **Incorporate Additional Features:** Enhance the model by including metadata or other contextual information.
  - **Interactive Visualization:** Develop tools to visualize topics and their relationships interactively.
- 

## References

- [Topic extraction with Non-negative Matrix Factorization and Latent Dirichlet Allocation](#)
- [Topic Modeling Tutorial – How to Use SVD and NMF in Python](#)
- [Topic Modelling using NMF | Guide to Master NLP (Part 14)](<https://www.analyticsvidhya.com/blog/2021/06/part-15-step-by-step-guide-to-master-nlp>)