Dataset Link :

# 1. Part 1: Text Pre-processing and Embedding Comparison

- Preprocessing steps correctly implemented and explained
- Comparison of word2vec, GloVe, and OpenAI embeddings
- Discussion on the embeddings that provide better semantic understanding

## Libraries Used

In [8]:
```
!pip3 install  scikit-learn
!pip3 install autocorrect
!pip3 install nltk
!pip3 install sentence-transformers
!pip3 install transformers==4.18.0
!pip3 install imblearn
!pip3 install wordcloud
!pip3 install pytorch-pretrained-bert
!pip3 install bertopic
!pip3 install pyspellchecker
!pip3 install numpy
!pip3  install wordcloud
```

Loading [MathJax]/extensions/Safe.js

```
Requirement already satisfied: scikit-learn in c:\users\purus\appdata\roaming\python\python310\site-packages (1.3.0)
Requirement already satisfied: numpy>=1.17.3 in c:\users\purus\appdata\roaming\python\python310\site-packages (from scikit-learn) (1.24.0)
Requirement already satisfied: scipy>=1.5.0 in c:\programdata\anaconda3\lib\site-packages (from scikit-learn) (1.11.1)
Requirement already satisfied: joblib>=1.1.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from scikit-learn) (1.3.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from scikit-learn) (3.1.0)
Requirement already satisfied: autocorrect in c:\programdata\anaconda3\lib\site-packages (2.6.1)
Requirement already satisfied: nltk in c:\programdata\anaconda3\lib\site-packages (3.7)
Requirement already satisfied: click in c:\programdata\anaconda3\lib\site-packages (from nltk) (8.0.4)
Requirement already satisfied: joblib in c:\users\purus\appdata\roaming\python\python310\site-packages (from nltk) (1.3.1)
Requirement already satisfied: regex>=2021.8.3 in c:\programdata\anaconda3\lib\site-packages (from nltk) (2022.7.9)
Requirement already satisfied: tqdm in c:\programdata\anaconda3\lib\site-packages (from nltk) (4.64.1)
Requirement already satisfied: colorama in c:\programdata\anaconda3\lib\site-packages (from click->nltk) (0.4.6)
Requirement already satisfied: sentence-transformers in c:\programdata\anaconda3\lib\site-packages (2.2.2)
Requirement already satisfied: transformers<5.0.0,>=4.6.0 in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (4.18.0)
Requirement already satisfied: tqdm in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (4.64.1)
Requirement already satisfied: torch>=1.6.0 in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (2.0.1)
Requirement already satisfied: torchvision in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (0.15.2)
Requirement already satisfied: numpy in c:\users\purus\appdata\roaming\python\python310\site-packages (from sentence-transformers) (1.24.0)
Requirement already satisfied: scikit-learn in c:\users\purus\appdata\roaming\python\python310\site-packages (from sentence-transformers) (1.3.0)
Requirement already satisfied: scipy in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (1.11.1)
Requirement already satisfied: nltk in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (3.7)
Requirement already satisfied: sentencepiece in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (0.1.99)
Requirement already satisfied: huggingface-hub>=0.4.0 in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers) (0.10.1)
Requirement already satisfied: filelock in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers) (3.9.0)
Requirement already satisfied: requests in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers) (2.28.1)
Requirement already satisfied: pyyaml>=5.1 in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers) (6.0)
Requirement already satisfied: typing-extensions>=3.7.4.3 in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transf
ormers) (4.4.0)
Requirement already satisfied: packaging>=20.9 in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers) (2
2.0)
Requirement already satisfied: sympy in c:\programdata\anaconda3\lib\site-packages (from torch>=1.6.0->sentence-transformers) (1.11.1)
Requirement already satisfied: networkx in c:\programdata\anaconda3\lib\site-packages (from torch>=1.6.0->sentence-transformers) (2.8.4)
Requirement already satisfied: jinja2 in c:\programdata\anaconda3\lib\site-packages (from torch>=1.6.0->sentence-transformers) (3.1.2)
Requirement already satisfied: regex!=2019.12.17 in c:\programdata\anaconda3\lib\site-packages (from transformers<5.0.0,>=4.6.0->sentence-transformer
s) (2022.7.9)
Requirement already satisfied: sacremoses in c:\programdata\anaconda3\lib\site-packages (from transformers<5.0.0,>=4.6.0->sentence-transformers) (0.
53)
Requirement already satisfied: tokenizers!=0.11.3,<0.13,>=0.11.1 in c:\programdata\anaconda3\lib\site-packages (from transformers<5.0.0,>=4.6.0->sent
ence-transformers) (0.11.4)
Requirement already satisfied: colorama in c:\programdata\anaconda3\lib\site-packages (from tqdm->sentence-transformers) (0.4.6)
Requirement already satisfied: click in c:\programdata\anaconda3\lib\site-packages (from nltk->sentence-transformers) (8.0.4)
Requirement already satisfied: joblib in c:\users\purus\appdata\roaming\python\python310\site-packages (from nltk->sentence-transformers) (1.3.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from scikit-learn->sentence-tra
nsformers) (3.1.0)
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in c:\programdata\anaconda3\lib\site-packages (from torchvision->sentence-transformers) (9.4.0)
Requirement already satisfied: MarkupSafe>=2.0 in c:\programdata\anaconda3\lib\site-packages (from jinja2->torch>=1.6.0->sentence-transformers) (2.1.
1)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentenc
e-transformers) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentence-transforme
rs) (3.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentence-t
ransformers) (1.26.14)
Requirement already satisfied: certifi>=2017.4.17 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentence-tran
sformers) (2023.5.7)
Requirement already satisfied: six in c:\programdata\anaconda3\lib\site-packages (from sacremoses->transformers<5.0.0,>=4.6.0->sentence-transformers)
(1.16.0)
Requirement already satisfied: mpmath>=0.19 in c:\programdata\anaconda3\lib\site-packages (from sympy->torch>=1.6.0->sentence-transformers) (1.2.1)
Requirement already satisfied: transformers==4.18.0 in c:\programdata\anaconda3\lib\site-packages (4.18.0)
Requirement already satisfied: filelock in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (3.9.0)
Requirement already satisfied: huggingface-hub<1.0,>=0.1.0 in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (0.10.1)
Requirement already satisfied: numpy>=1.17 in c:\users\purus\appdata\roaming\python\python310\site-packages (from transformers==4.18.0) (1.24.0)
Requirement already satisfied: packaging>=20.0 in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (22.0)
Requirement already satisfied: pyyaml>=5.1 in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (6.0)
Requirement already satisfied: regex!=2019.12.17 in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (2022.7.9)
Requirement already satisfied: requests in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (2.28.1)
Requirement already satisfied: sacremoses in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (0.0.53)
Requirement already satisfied: tokenizers!=0.11.3,<0.13,>=0.11.1 in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (0.11.4)
Requirement already satisfied: tqdm>=4.27 in c:\programdata\anaconda3\lib\site-packages (from transformers==4.18.0) (4.64.1)
Requirement already satisfied: typing-extensions>=3.7.4.3 in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub<1.0,>=0.1.0->transforme
rs==4.18.0) (4.4.0)
Requirement already satisfied: colorama in c:\programdata\anaconda3\lib\site-packages (from tqdm>=4.27->transformers==4.18.0) (0.4.6)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\programdata\anaconda3\lib\site-packages (from requests->transformers==4.18.0) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in c:\programdata\anaconda3\lib\site-packages (from requests->transformers==4.18.0) (3.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\programdata\anaconda3\lib\site-packages (from requests->transformers==4.18.0) (1.26.14)
Requirement already satisfied: certifi>=2017.4.17 in c:\programdata\anaconda3\lib\site-packages (from requests->transformers==4.18.0) (2023.5.7)
Requirement already satisfied: six in c:\programdata\anaconda3\lib\site-packages (from sacremoses->transformers==4.18.0) (1.16.0)
Requirement already satisfied: click in c:\programdata\anaconda3\lib\site-packages (from sacremoses->transformers==4.18.0) (8.0.4)
Requirement already satisfied: joblib in c:\users\purus\appdata\roaming\python\python310\site-packages (from sacremoses->transformers==4.18.0) (1.3.
1)
Requirement already satisfied: imblearn in c:\programdata\anaconda3\lib\site-packages (0.0)
Requirement already satisfied: imbalanced-learn in c:\users\purus\appdata\roaming\python\python310\site-packages (from imblearn) (0.11.0)
Requirement already satisfied: numpy>=1.17.3 in c:\users\purus\appdata\roaming\python\python310\site-packages (from imbalanced-learn->imblearn) (1.2
4.0)
Requirement already satisfied: scipy>=1.5.0 in c:\programdata\anaconda3\lib\site-packages (from imbalanced-learn->imblearn) (1.11.1)
Requirement already satisfied: scikit-learn>=1.0.2 in c:\users\purus\appdata\roaming\python\python310\site-packages (from imbalanced-learn->imblearn)
(1.3.0)
Requirement already satisfied: joblib>=1.1.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from imbalanced-learn->imblearn) (1.3.
1)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from imbalanced-learn->imblear
n) (3.1.0)
Requirement already satisfied: wordcloud in c:\users\purus\appdata\roaming\python\python310\site-packages (1.9.2)
Requirement already satisfied: numpy>=1.6.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from wordcloud) (1.24.0)
Requirement already satisfied: pillow in c:\programdata\anaconda3\lib\site-packages (from wordcloud) (9.4.0)
Requirement already satisfied: matplotlib in c:\programdata\anaconda3\lib\site-packages (from wordcloud) (3.7.0)
Requirement already satisfied: contourpy>=1.0.1 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.0.5)
Requirement already satisfied: cycler>=0.10 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (0.11.0)
Requirement already satisfied: fonttools>=4.22.0 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (4.25.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.4.4)
Requirement already satisfied: packaging>=20.0 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (22.0)
Requirement already satisfied: pyparsing>=2.3.1 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (3.0.9)
Requirement already satisfied: python-dateutil>=2.7 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.8.2)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil>=2.7->matplotlib->wordcloud) (1.16.0)
Requirement already satisfied: pytorch-pretrained-bert in c:\users\purus\appdata\roaming\python\python310\site-packages (0.6.2)
Requirement already satisfied: torch>=0.4.1 in c:\programdata\anaconda3\lib\site-packages (from pytorch-pretrained-bert) (2.0.1)
Requirement already satisfied: numpy in c:\users\purus\appdata\roaming\python\python310\site-packages (from pytorch-pretrained-bert) (1.24.0)
Requirement already satisfied: boto3 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pytorch-pretrained-bert) (1.28.1)
```

Loading [MathJax]/extensions/Safe.js

```
Requirement already satisfied: requests in c:\programdata\anaconda3\lib\site-packages (from pytorch-pretrained-bert) (2.28.1)
Requirement already satisfied: tqdm in c:\programdata\anaconda3\lib\site-packages (from pytorch-pretrained-bert) (4.64.1)
Requirement already satisfied: regex in c:\programdata\anaconda3\lib\site-packages (from pytorch-pretrained-bert) (2022.7.9)
Requirement already satisfied: filelock in c:\programdata\anaconda3\lib\site-packages (from torch>=0.4.1->pytorch-pretrained-bert) (3.9.0)
Requirement already satisfied: typing-extensions in c:\programdata\anaconda3\lib\site-packages (from torch>=0.4.1->pytorch-pretrained-bert) (4.4.0)
Requirement already satisfied: sympy in c:\programdata\anaconda3\lib\site-packages (from torch>=0.4.1->pytorch-pretrained-bert) (1.11.1)
Requirement already satisfied: networkx in c:\programdata\anaconda3\lib\site-packages (from torch>=0.4.1->pytorch-pretrained-bert) (2.8.4)
Requirement already satisfied: jinja2 in c:\programdata\anaconda3\lib\site-packages (from torch>=0.4.1->pytorch-pretrained-bert) (3.1.2)
Requirement already satisfied: botocore<1.32.0,>=1.31.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from boto3->pytorch-pretrai
ned-bert) (1.31.1)
Requirement already satisfied: jmespath<2.0.0,>=0.7.1 in c:\programdata\anaconda3\lib\site-packages (from boto3->pytorch-pretrained-bert) (0.10.0)
Requirement already satisfied: s3transfer<0.7.0,>=0.6.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from boto3->pytorch-pretrai
ned-bert) (0.6.1)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\programdata\anaconda3\lib\site-packages (from requests->pytorch-pretrained-bert) (2.0.
4)
Requirement already satisfied: idna<4,>=2.5 in c:\programdata\anaconda3\lib\site-packages (from requests->pytorch-pretrained-bert) (3.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\programdata\anaconda3\lib\site-packages (from requests->pytorch-pretrained-bert) (1.26.14)
Requirement already satisfied: certifi>=2017.4.17 in c:\programdata\anaconda3\lib\site-packages (from requests->pytorch-pretrained-bert) (2023.5.7)
Requirement already satisfied: colorama in c:\programdata\anaconda3\lib\site-packages (from tqdm->pytorch-pretrained-bert) (0.4.6)
Requirement already satisfied: python-dateutil<3.0.0,>=2.1 in c:\programdata\anaconda3\lib\site-packages (from botocore<1.32.0,>=1.31.1->boto3->pytor
ch-pretrained-bert) (2.8.2)
Requirement already satisfied: MarkupSafe>=2.0 in c:\programdata\anaconda3\lib\site-packages (from jinja2->torch>=0.4.1->pytorch-pretrained-bert) (2.
1.1)
Requirement already satisfied: mpmath>=0.19 in c:\programdata\anaconda3\lib\site-packages (from sympy->torch>=0.4.1->pytorch-pretrained-bert) (1.2.1)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil<3.0.0,>=2.1->botocore<1.32.0,>=1.31.1->bo
to3->pytorch-pretrained-bert) (1.16.0)
Requirement already satisfied: bertopic in c:\users\purus\appdata\roaming\python\python310\site-packages (0.15.0)
Requirement already satisfied: numpy>=1.20.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from bertopic) (1.24.0)
Requirement already satisfied: hdbscan>=0.8.29 in c:\users\purus\appdata\roaming\python\python310\site-packages (from bertopic) (0.8.30)
Requirement already satisfied: umap-learn>=0.5.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from bertopic) (0.5.3)
Requirement already satisfied: pandas>=1.1.5 in c:\users\purus\appdata\roaming\python\python310\site-packages (from bertopic) (2.0.3)
Requirement already satisfied: scikit-learn>=0.22.2.post1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from bertopic) (1.3.0)
Requirement already satisfied: tqdm>=4.41.1 in c:\programdata\anaconda3\lib\site-packages (from bertopic) (4.64.1)
Requirement already satisfied: sentence-transformers>=0.4.1 in c:\programdata\anaconda3\lib\site-packages (from bertopic) (2.2.2)
Requirement already satisfied: plotly>=4.7.0 in c:\programdata\anaconda3\lib\site-packages (from bertopic) (5.9.0)
Requirement already satisfied: cython>=0.27 in c:\users\purus\appdata\roaming\python\python310\site-packages (from hdbscan>=0.8.29->bertopic) (0.29.3
6)
Requirement already satisfied: scipy>=1.0 in c:\programdata\anaconda3\lib\site-packages (from hdbscan>=0.8.29->bertopic) (1.11.1)
Requirement already satisfied: joblib>=1.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from hdbscan>=0.8.29->bertopic) (1.3.1)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\programdata\anaconda3\lib\site-packages (from pandas>=1.1.5->bertopic) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in c:\programdata\anaconda3\lib\site-packages (from pandas>=1.1.5->bertopic) (2022.7)
Requirement already satisfied: tzdata>=2022.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pandas>=1.1.5->bertopic) (2023.
3)
Requirement already satisfied: tenacity>=6.2.0 in c:\programdata\anaconda3\lib\site-packages (from plotly>=4.7.0->bertopic) (8.0.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from scikit-learn>=0.22.2.post1
->bertopic) (3.1.0)
Requirement already satisfied: transformers<5.0.0,>=4.6.0 in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers>=0.4.1->bertopic)
(4.18.0)
Requirement already satisfied: torch>=1.6.0 in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers>=0.4.1->bertopic) (2.0.1)
Requirement already satisfied: torchvision in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers>=0.4.1->bertopic) (0.15.2)
Requirement already satisfied: nltk in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers>=0.4.1->bertopic) (3.7)
Requirement already satisfied: sentencepiece in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers>=0.4.1->bertopic) (0.1.99)
Requirement already satisfied: huggingface-hub>=0.4.0 in c:\programdata\anaconda3\lib\site-packages (from sentence-transformers>=0.4.1->bertopic) (0.
10.1)
Requirement already satisfied: colorama in c:\programdata\anaconda3\lib\site-packages (from tqdm>=4.41.1->bertopic) (0.4.6)
Requirement already satisfied: numba>=0.49 in c:\users\purus\appdata\roaming\python\python310\site-packages (from umap-learn>=0.5.0->bertopic) (0.57.
1)
Requirement already satisfied: pynndescent>=0.5 in c:\users\purus\appdata\roaming\python\python310\site-packages (from umap-learn>=0.5.0->bertopic)
(0.5.10)
Requirement already satisfied: filelock in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers>=0.4.1->ber
topic) (3.9.0)
Requirement already satisfied: requests in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers>=0.4.1->ber
topic) (2.28.1)
Requirement already satisfied: pyyaml>=5.1 in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers>=0.4.1->
bertopic) (6.0)
Requirement already satisfied: typing-extensions>=3.7.4.3 in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transf
ormers>=0.4.1->bertopic) (4.4.0)
Requirement already satisfied: packaging>=20.9 in c:\programdata\anaconda3\lib\site-packages (from huggingface-hub>=0.4.0->sentence-transformers>=0.
4.1->bertopic) (22.0)
Requirement already satisfied: llvmlite<0.41,>=0.40.0dev0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from numba>=0.49->umap-le
arn>=0.5.0->bertopic) (0.40.1)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil>=2.8.2->pandas>=1.1.5->bertopic) (1.16.0)
Requirement already satisfied: sympy in c:\programdata\anaconda3\lib\site-packages (from torch>=1.6.0->sentence-transformers>=0.4.1->bertopic) (1.11.
1)
Requirement already satisfied: networkx in c:\programdata\anaconda3\lib\site-packages (from torch>=1.6.0->sentence-transformers>=0.4.1->bertopic) (2.
8.4)
Requirement already satisfied: jinja2 in c:\programdata\anaconda3\lib\site-packages (from torch>=1.6.0->sentence-transformers>=0.4.1->bertopic) (3.1.
2)
Requirement already satisfied: regex!=2019.12.17 in c:\programdata\anaconda3\lib\site-packages (from transformers<5.0.0,>=4.6.0->sentence-transformer
s>=0.4.1->bertopic) (2022.7.9)
Requirement already satisfied: sacremoses in c:\programdata\anaconda3\lib\site-packages (from transformers<5.0.0,>=4.6.0->sentence-transformers>=0.4.
1->bertopic) (0.0.53)
Requirement already satisfied: tokenizers!=0.11.3,<0.13,>=0.11.1 in c:\programdata\anaconda3\lib\site-packages (from transformers<5.0.0,>=4.6.0->sent
ence-transformers>=0.4.1->bertopic) (0.11.4)
Requirement already satisfied: click in c:\programdata\anaconda3\lib\site-packages (from nltk->sentence-transformers>=0.4.1->bertopic) (8.0.4)
Requirement already satisfied: pillow!=8.3.*,>=5.3.0 in c:\programdata\anaconda3\lib\site-packages (from torchvision->sentence-transformers>=0.4.1->b
ertopic) (9.4.0)
Requirement already satisfied: MarkupSafe>=2.0 in c:\programdata\anaconda3\lib\site-packages (from jinja2->torch>=1.6.0->sentence-transformers>=0.4.1
->bertopic) (2.1.1)
Requirement already satisfied: charset-normalizer<3,>=2 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentenc
e-transformers>=0.4.1->bertopic) (2.0.4)
Requirement already satisfied: idna<4,>=2.5 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentence-transforme
rs>=0.4.1->bertopic) (3.4)
Requirement already satisfied: urllib3<1.27,>=1.21.1 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentence-t
ransformers>=0.4.1->bertopic) (1.26.14)
Requirement already satisfied: certifi>=2017.4.17 in c:\programdata\anaconda3\lib\site-packages (from requests->huggingface-hub>=0.4.0->sentence-tran
sformers>=0.4.1->bertopic) (2023.5.7)
Requirement already satisfied: mpmath>=0.19 in c:\programdata\anaconda3\lib\site-packages (from sympy->torch>=1.6.0->sentence-transformers>=0.4.1->be
rtopic) (1.2.1)
Requirement already satisfied: pyspellchecker in c:\users\purus\appdata\roaming\python\python310\site-packages (0.7.2)
Requirement already satisfied: numpy in c:\users\purus\appdata\roaming\python\python310\site-packages (1.24.0)
Requirement already satisfied: wordcloud in c:\users\purus\appdata\roaming\python\python310\site-packages (1.9.2)
Requirement already satisfied: numpy>=1.6.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from wordcloud) (1.24.0)
Requirement already satisfied: pillow in c:\programdata\anaconda3\lib\site-packages (from wordcloud) (9.4.0)
Requirement already satisfied: matplotlib in c:\programdata\anaconda3\lib\site-packages (from wordcloud) (3.7.0)
Requirement already satisfied: contourpy>=1.0.1 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.0.5)
Requirement already satisfied: cycler>=0.10 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (0.11.0)
```

Loading [MathJax]/extensions/Safe.js

```
Requirement already satisfied: fonttools>=4.22.0 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (4.25.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.4.4)
Requirement already satisfied: packaging>=20.0 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (22.0)
Requirement already satisfied: pyparsing>=2.3.1 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (3.0.9)
Requirement already satisfied: python-dateutil>=2.7 in c:\programdata\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.8.2)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil>=2.7->matplotlib->wordcloud) (1.16.0)
```

In [9]:
```python
import os
os.environ["TOKENIZERS_PARALLELISM"] = "false"
```

## Supress warnings

In [10]:
```python
import warnings
import numpy as np
!pip install numba

warnings.filterwarnings("ignore")
```

```
Requirement already satisfied: numba in c:\users\purus\appdata\roaming\python\python310\site-packages (0.57.1)
Requirement already satisfied: llvmlite<0.41,>=0.40.0dev0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from numba) (0.40.1)
Requirement already satisfied: numpy<1.25,>=1.21 in c:\users\purus\appdata\roaming\python\python310\site-packages (from numba) (1.24.0)
```

## Imports

In [11]:
```python
from transformers import BertTokenizer
import nltk
import re
import nltk
import pandas as pd
import numpy as np
import warnings
warnings.filterwarnings('ignore')
nltk.download("punkt")
nltk.download('stopwords')
nltk.download('all')
from sklearn.cluster import KMeans
from sklearn.model_selection import KFold
from sklearn.feature_extraction.text import CountVectorizer
```

In [9]:
```python
import os
os.environ["TOKENIZERS_PARALLELISM"] = "false"
```

Loading [MathJax]/extensions/Safe.js

```
[nltk_data] Downloading package punkt to
[nltk_data]     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading collection 'all'
[nltk_data]     |
[nltk_data]     | Downloading package abc to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package abc is already up-to-date!
[nltk_data]     | Downloading package alpino to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package alpino is already up-to-date!
[nltk_data]     | Downloading package averaged_perceptron_tagger to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package averaged_perceptron_tagger is already up-
[nltk_data]     |       to-date!
[nltk_data]     | Downloading package averaged_perceptron_tagger_ru to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package averaged_perceptron_tagger_ru is already
[nltk_data]     |       up-to-date!
[nltk_data]     | Downloading package basque_grammars to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package basque_grammars is already up-to-date!
[nltk_data]     | Downloading package bcp47 to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package bcp47 is already up-to-date!
[nltk_data]     | Downloading package biocreative_ppi to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package biocreative_ppi is already up-to-date!
[nltk_data]     | Downloading package bllip_wsj_no_aux to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package bllip_wsj_no_aux is already up-to-date!
[nltk_data]     | Downloading package book_grammars to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package book_grammars is already up-to-date!
[nltk_data]     | Downloading package brown to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package brown is already up-to-date!
[nltk_data]     | Downloading package brown_tei to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package brown_tei is already up-to-date!
[nltk_data]     | Downloading package cess_cat to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package cess_cat is already up-to-date!
[nltk_data]     | Downloading package cess_esp to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package cess_esp is already up-to-date!
[nltk_data]     | Downloading package chat80 to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package chat80 is already up-to-date!
[nltk_data]     | Downloading package city_database to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package city_database is already up-to-date!
[nltk_data]     | Downloading package cmudict to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package cmudict is already up-to-date!
[nltk_data]     | Downloading package comparative_sentences to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package comparative_sentences is already up-to-
[nltk_data]     |       date!
[nltk_data]     | Downloading package comtrans to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package comtrans is already up-to-date!
[nltk_data]     | Downloading package conll2000 to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package conll2000 is already up-to-date!
[nltk_data]     | Downloading package conll2002 to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package conll2002 is already up-to-date!
[nltk_data]     | Downloading package conll2007 to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package conll2007 is already up-to-date!
[nltk_data]     | Downloading package crubadan to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package crubadan is already up-to-date!
[nltk_data]     | Downloading package dependency_treebank to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package dependency_treebank is already up-to-date!
[nltk_data]     | Downloading package dolch to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package dolch is already up-to-date!
[nltk_data]     | Downloading package europarl_raw to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package europarl_raw is already up-to-date!
[nltk_data]     | Downloading package extended_omw to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package extended_omw is already up-to-date!
[nltk_data]     | Downloading package floresta to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package floresta is already up-to-date!
[nltk_data]     | Downloading package framenet_v15 to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package framenet_v15 is already up-to-date!
[nltk_data]     | Downloading package framenet_v17 to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package framenet_v17 is already up-to-date!
[nltk_data]     | Downloading package gazetteers to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package gazetteers is already up-to-date!
[nltk_data]     | Downloading package genesis to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package genesis is already up-to-date!
[nltk_data]     | Downloading package gutenberg to
[nltk_data]     |     C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]     |   Package gutenberg is already up-to-date!
```

Loading [MathJax]/extensions/Safe.js

```
[nltk_data]    |   Downloading package ieer to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package ieer is already up-to-date!
[nltk_data]    |   Downloading package inaugural to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package inaugural is already up-to-date!
[nltk_data]    |   Downloading package indian to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package indian is already up-to-date!
[nltk_data]    |   Downloading package jeita to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package jeita is already up-to-date!
[nltk_data]    |   Downloading package kimmo to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package kimmo is already up-to-date!
[nltk_data]    |   Downloading package knbc to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package knbc is already up-to-date!
[nltk_data]    |   Downloading package large_grammars to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package large_grammars is already up-to-date!
[nltk_data]    |   Downloading package lin_thesaurus to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package lin_thesaurus is already up-to-date!
[nltk_data]    |   Downloading package mac_morpho to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package mac_morpho is already up-to-date!
[nltk_data]    |   Downloading package machado to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package machado is already up-to-date!
[nltk_data]    |   Downloading package masc_tagged to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package masc_tagged is already up-to-date!
[nltk_data]    |   Downloading package maxent_ne_chunker to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package maxent_ne_chunker is already up-to-date!
[nltk_data]    |   Downloading package maxent_treebank_pos_tagger to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package maxent_treebank_pos_tagger is already up-
[nltk_data]    |         to-date!
[nltk_data]    |   Downloading package moses_sample to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package moses_sample is already up-to-date!
[nltk_data]    |   Downloading package movie_reviews to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package movie_reviews is already up-to-date!
[nltk_data]    |   Downloading package mte_teip5 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package mte_teip5 is already up-to-date!
[nltk_data]    |   Downloading package mwa_ppdb to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package mwa_ppdb is already up-to-date!
[nltk_data]    |   Downloading package names to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package names is already up-to-date!
[nltk_data]    |   Downloading package nombank.1.0 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package nombank.1.0 is already up-to-date!
[nltk_data]    |   Downloading package nonbreaking_prefixes to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package nonbreaking_prefixes is already up-to-date!
[nltk_data]    |   Downloading package nps_chat to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package nps_chat is already up-to-date!
[nltk_data]    |   Downloading package omw to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package omw is already up-to-date!
[nltk_data]    |   Downloading package omw-1.4 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package omw-1.4 is already up-to-date!
[nltk_data]    |   Downloading package opinion_lexicon to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package opinion_lexicon is already up-to-date!
[nltk_data]    |   Downloading package panlex_swadesh to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package panlex_swadesh is already up-to-date!
[nltk_data]    |   Downloading package paradigms to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package paradigms is already up-to-date!
[nltk_data]    |   Downloading package pe08 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package pe08 is already up-to-date!
[nltk_data]    |   Downloading package perluniprops to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package perluniprops is already up-to-date!
[nltk_data]    |   Downloading package pil to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package pil is already up-to-date!
[nltk_data]    |   Downloading package pl196x to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package pl196x is already up-to-date!
[nltk_data]    |   Downloading package porter_test to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package porter_test is already up-to-date!
[nltk_data]    |   Downloading package ppattach to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package ppattach is already up-to-date!
[nltk_data]    |   Downloading package problem_reports to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package problem_reports is already up-to-date!
[nltk_data]    |   Downloading package product_reviews_1 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package product_reviews_1 is already up-to-date!
[nltk_data]    |   Downloading package product_reviews_2 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package product_reviews_2 is already up-to-date!
[nltk_data]    |   Downloading package propbank to
```

Loading [MathJax]/extensions/Safe.js

```
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package propbank is already up-to-date!
[nltk_data]    |   Downloading package pros_cons to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package pros_cons is already up-to-date!
[nltk_data]    |   Downloading package ptb to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package ptb is already up-to-date!
[nltk_data]    |   Downloading package punkt to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package punkt is already up-to-date!
[nltk_data]    |   Downloading package qc to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package qc is already up-to-date!
[nltk_data]    |   Downloading package reuters to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package reuters is already up-to-date!
[nltk_data]    |   Downloading package rslp to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package rslp is already up-to-date!
[nltk_data]    |   Downloading package rte to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package rte is already up-to-date!
[nltk_data]    |   Downloading package sample_grammars to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package sample_grammars is already up-to-date!
[nltk_data]    |   Downloading package semcor to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package semcor is already up-to-date!
[nltk_data]    |   Downloading package senseval to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package senseval is already up-to-date!
[nltk_data]    |   Downloading package sentence_polarity to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package sentence_polarity is already up-to-date!
[nltk_data]    |   Downloading package sentiwordnet to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package sentiwordnet is already up-to-date!
[nltk_data]    |   Downloading package shakespeare to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package shakespeare is already up-to-date!
[nltk_data]    |   Downloading package sinica_treebank to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package sinica_treebank is already up-to-date!
[nltk_data]    |   Downloading package smultron to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package smultron is already up-to-date!
[nltk_data]    |   Downloading package snowball_data to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package snowball_data is already up-to-date!
[nltk_data]    |   Downloading package spanish_grammars to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package spanish_grammars is already up-to-date!
[nltk_data]    |   Downloading package state_union to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package state_union is already up-to-date!
[nltk_data]    |   Downloading package stopwords to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package stopwords is already up-to-date!
[nltk_data]    |   Downloading package subjectivity to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package subjectivity is already up-to-date!
[nltk_data]    |   Downloading package swadesh to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package swadesh is already up-to-date!
[nltk_data]    |   Downloading package switchboard to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package switchboard is already up-to-date!
[nltk_data]    |   Downloading package tagsets to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package tagsets is already up-to-date!
[nltk_data]    |   Downloading package timit to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package timit is already up-to-date!
[nltk_data]    |   Downloading package toolbox to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package toolbox is already up-to-date!
[nltk_data]    |   Downloading package treebank to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package treebank is already up-to-date!
[nltk_data]    |   Downloading package twitter_samples to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package twitter_samples is already up-to-date!
[nltk_data]    |   Downloading package udhr to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package udhr is already up-to-date!
[nltk_data]    |   Downloading package udhr2 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package udhr2 is already up-to-date!
[nltk_data]    |   Downloading package unicode_samples to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package unicode_samples is already up-to-date!
[nltk_data]    |   Downloading package universal_tagset to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package universal_tagset is already up-to-date!
[nltk_data]    |   Downloading package universal_treebanks_v20 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package universal_treebanks_v20 is already up-to-
[nltk_data]    |         date!
[nltk_data]    |   Downloading package vader_lexicon to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package vader_lexicon is already up-to-date!
[nltk_data]    |   Downloading package verbnet to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
[nltk_data]    |     Package verbnet is already up-to-date!
[nltk_data]    |   Downloading package verbnet3 to
[nltk_data]    |       C:\Users\purus\AppData\Roaming\nltk_data...
```

Loading [MathJax]/extensions/Safe.js

```python
In [12]: from nltk.tokenize import sent_tokenize,word_tokenize
         from nltk.corpus import stopwords
         from string import punctuation
         from string import punctuation
         from nltk.corpus import stopwords, brown
         import re
         from nltk.stem import WordNetLemmatizer
         import matplotlib.pyplot as plt
         %matplotlib inline
         from sklearn.metrics import confusion_matrix
         from sklearn.metrics import roc_curve
         from sklearn.metrics import auc
         import seaborn as sns
         import matplotlib.pyplot as plt
```

## Loading Data

```python
In [13]: from sklearn.datasets import fetch_20newsgroups
```

```python
In [14]: categories = ['alt.atheism', 'soc.religion.christian',
                       'comp.graphics', 'sci.med', 'talk.religion.misc',
                       'sci.space']
```

```python
In [15]: remove = ('headers', 'footers', 'quotes')
```

```python
In [16]: def convert_to_np(dataset):
             return np.asarray(dataset.data), dataset.target
```

```python
In [17]: data_train = fetch_20newsgroups(subset='train', categories=categories,
                                          shuffle=True, random_state=42,
                                          remove=remove)


         data_test = fetch_20newsgroups(subset='test', categories=categories,
                                        shuffle=True, random_state=42,
                                        remove=remove)
         x_validation,y_validation =convert_to_np(data_test)
         x_train,y_train = convert_to_np(data_train)



         print('data loaded')

         data loaded
```

```python
In [18]: def size_mb(docs):
             return sum(len(s.encode('utf-8')) for s in docs) / 1e6
```

```python
In [19]: data_train_size_mb = size_mb(data_train.data)
         data_test_size_mb = size_mb(data_test.data)
```

```python
In [20]: print("%d documents - %0.3fMB (training set)" % (
             len(data_train.data), data_train_size_mb))
         print("%d documents - %0.3fMB (test set)" % (
             len(data_test.data), data_test_size_mb))
         print("%d categories" % len(categories))
         print()

         3227 documents - 4.110MB (training set)
         2147 documents - 3.037MB (test set)
         6 categories
```

## Distribution of data

Loading [MathJax]/extensions/Safe.js

```
In [21]:   # Finding frequency of each category
           targets, frequency = np.unique(data_train.target, return_counts=True)
           targets, frequency
```

```
Out[21]:   (array([0, 1, 2, 3, 4, 5], dtype=int64),
            array([480, 584, 594, 593, 599, 377], dtype=int64))
```

```
In [22]:   targets_str = np.array(data_train.target_names)
           print(list(zip(targets_str, frequency)))
```

```
[('alt.atheism', 480), ('comp.graphics', 584), ('sci.med', 594), ('sci.space', 593), ('soc.religion.christian', 599), ('talk.religion.misc', 377)]
```

```
In [23]:   # Training data class distribution
           fig=plt.figure(figsize=(10, 5), dpi= 80, facecolor='w', edgecolor='k')
           plt.bar(targets_str,frequency)
           plt.xticks(rotation=90)
           plt.title('Class distribution of 20 Newsgroups Training Data')
           plt.xlabel('News Group')
           plt.ylabel('Frequency')
           plt.show()
```



When a dataset has equal samples in all categories and no imbalance issues, it means that each category or class in the dataset has the same number of samples. This balanced distribution ensures that there is no bias towards any particular category during the analysis.

# Preprocessing steps correctly implemented and explained

Below are the steps incorporated for preprocessing Noise removal:digits, characters, and pieces of text that interfere with the process of text analysis Lowercasing:to deal with sparsity issues in the dataset we have done this step Normalization : 1) Stop-word removal 2) We choose Lemmatization over Stemming as it is doing things properly with the use of vocabulary and morphological analysis of words 3) Speller to fix speeling errors

## Preprocessing steps:

### Noise Removal:

```
In this step, you remove unwanted elements from the text that can interfere with text analysis. This typically includes removing
digits, special characters, and other irrelevant pieces of text that don't contribute to the analysis.
```

### Lowercasing:

```
Lowercasing refers to converting all text to lowercase letters. This step helps to address sparsity issues in the dataset by
treating words with different capitalization as the same. For example, "apple" and "Apple" would be considered the same word
after lowercasing.
```

### Normalization:

### Stop-word Removal:

```
Stop words are common words like "a," "the," "and," etc., that don't carry significant meaning in the analysis. Removing stop
words helps reduce noise and focuses on more important words in the text.
```

### Lemmatization:

```
Lemmatization is the process of reducing words to their base or root form. Unlike stemming, which simply chops off word endings,
lemmatization considers vocabulary and the morphological analysis of words to derive their base form. This step helps to maintain
the integrity of words and ensure meaningful analysis.
```

### Spelling Correction:

Loading [MathJax]/extensions/Safe.js

Spelling errors can negatively impact the analysis and interpretation of text. Using a speller, you can correct spelling errors to improve the quality and accuracy of the analysis.

Overall, these preprocessing steps help to clean and prepare the text data for further analysis by removing noise, standardizing text formats, and improving the quality of the text. Each step contributes to enhancing the accuracy and effectiveness of text analysis tasks.

```python
In [24]:    # example didn', "didn't", 'doesn', "doesn't", 'hadn', "hadn't", 'hasn', "hasn't"
            stop_words  = stopwords.words('english')
```

```python
In [25]:    from autocorrect import Speller
            from nltk.tokenize import word_tokenize


            def to_lower(text):

                """
                Converting text to lower case as in, converting "Hello" to  "hello" or "Hi" to "hi".
                """

                # Specll check the words
                spell  = Speller(lang='en')

                texts = spell(text)

                return ' '.join([w.lower() for w in word_tokenize(text)])
```

```python
In [20]:    def clean_text(lower_case):
                # split text phrases into words
                words  = nltk.word_tokenize(lower_case)

                # Create a list of all the punctuations
                punctuations = [ '/', '!', '?', ';', ':', '(',')', '[',']', '-', '_', '%']

                # Remove all the special characters
                punctuations = re.sub(r'\W', ' ', str(lower_case))

                # Initialize the stopwords variable, which is a list of words ('and', 'the', 'i', 'yourself', 'is') that do not hold much values as key words
                stop_words  = stopwords.words('english')

                # Getting rid of all the words that contain numbers in them
                w_num = re.sub('\w*\d\w*', '', lower_case).strip()

                # remove all single characters
                lower_case = re.sub(r'\s+[a-zA-Z]\s+', ' ', lower_case)

                # Substituting multiple spaces with single space
                lower_case = re.sub(r'\s+', ' ', lower_case, flags=re.I)

                # Removing prefixed 'b'
                lower_case = re.sub(r'^b\s+', '', lower_case)

                # Removing non-english characters
                lower_case = re.sub(r'^b\s+', '', lower_case)

                # Return keywords which are not in stop words
                keywords = [word for word in words if not word in stop_words  and word in punctuations and  word in w_num]

                return keywords
```

## Pre-processing of training data

```python
In [21]:    # Training data lemmatization
            # Lemmatize the words
            wordnet_lemmatizer = WordNetLemmatizer()
            for idx, txt in enumerate(data_train['data']):
                lemmatized_word = [wordnet_lemmatizer.lemmatize(word) for word in clean_text(to_lower(txt))]
                clean_data = ' '.join(lemmatized_word)
                data_train['data'][idx]=clean_data
```

## Preprocessing of Test data

```python
In [22]:    # Test data lemmatization
            # Lemmatize the words
            wordnet_lemmatizer = WordNetLemmatizer()
            for idx, txt in enumerate(data_test['data']):
                lemmatized_word = [wordnet_lemmatizer.lemmatize(word) for word in clean_text(to_lower(txt))]
                clean_data = ' '.join(lemmatized_word)
                data_test['data'][idx]=clean_data
```

```python
In [23]:    x_validation,y_validation =convert_to_np(data_test)
            x_train,y_train = convert_to_np(data_train)
```

## Word2vec, GloVe, and OpenAI embedding comparison

I compared the behaviour of two frameworks—BERT and WordVec—I used to build embedding for a small text.

## Word2Vec (CBOW)

```python
In [24]:    import gensim
```
Loading [MathJax]/extensions/Safe.js
```python
            ...els.word2vec import Word2Vec
```

```python
from gensim.test.utils import common_texts
```

```python
In [25]:   text = "After stealing money from the bank vault the bank robber was seen fishing on the Mississippi river bank"
```

```python
In [26]:   # plot word count for news text
           from wordcloud import WordCloud
           wordcloud = WordCloud(background_color='black',
                              max_words=200).generate(text)
           fig = plt.figure(figsize=[10,10])
           plt.title('WordCloud of Sample Sentence')
           plt.axis('off')
           plt.imshow(wordcloud)
           plt.show()
```



WordCloud of Sample Sentence

```python
In [27]:   #Word2Vec
           #training the gensim on the data
           #Using the Cbow architecture for the word2Vec
           #, size = 50 removed from word2Vec
           from gensim.models import Word2Vec
           model_cbow = Word2Vec([text.split(" ")], min_count = 1, workers = 3, vector_size = 50, window = 5, sg = 0)
```

```python
In [28]:   model_cbow.init_sims(replace = True)
           model_cbow.train([["hello", "world"]], total_examples=1, epochs=1)
```

```
Out[28]:   (0, 2)
```

```python
In [29]:   words = list(model_cbow.wv.index_to_key)
           print(words)
```

```
['bank', 'the', 'river', 'Mississippi', 'on', 'fishing', 'seen', 'was', 'robber', 'vault', 'from', 'money', 'stealing', 'After']
```

```python
In [30]:   vector = model_cbow.wv['bank']
           vector
```

```
Out[30]:   array([-0.01259352,  0.00555269,  0.11985431,  0.21158655, -0.21848367,
                  -0.16714124,  0.15168934,  0.2107344 , -0.11778943, -0.08838436,
                   0.17333424, -0.03601422, -0.10654427,  0.15392466, -0.1141429 ,
                  -0.04264994,  0.06755769,  0.02329457, -0.19458175, -0.22190945,
                   0.1717199 ,  0.11907724,  0.15870726,  0.01791622,  0.14915332,
                  -0.07997645, -0.02222663,  0.13547736, -0.17664881, -0.09244104,
                  -0.17641266, -0.02184243,  0.22400671, -0.1718937 , -0.05480954,
                  -0.04550866,  0.189702  , -0.13928957,  0.00106066, -0.11164343,
                  -0.22554341,  0.11759838, -0.20572253, -0.10314386, -0.00082434,
                  -0.00695595, -0.17992744,  0.22580628,  0.11700572,  0.21684425],
                 dtype=float32)
```

```python
In [31]:   !pip install --upgrade pandas
           !pip install --upgrade pyLDAvis
           !pip install gensim
```

Loading [MathJax]/extensions/Safe.js

```
Requirement already satisfied: pandas in c:\users\purus\appdata\roaming\python\python310\site-packages (2.0.3)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\programdata\anaconda3\lib\site-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in c:\programdata\anaconda3\lib\site-packages (from pandas) (2022.7)
Requirement already satisfied: tzdata>=2022.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pandas) (2023.3)
Requirement already satisfied: numpy>=1.21.0 in c:\programdata\anaconda3\lib\site-packages (from pandas) (1.23.5)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)
Requirement already satisfied: pyLDAvis in c:\programdata\anaconda3\lib\site-packages (3.4.0)
Collecting pyLDAvis
  Using cached pyLDAvis-3.4.1-py3-none-any.whl (2.6 MB)
Collecting numpy>=1.24.2 (from pyLDAvis)
  Using cached numpy-1.25.1-cp310-cp310-win_amd64.whl (15.0 MB)
Requirement already satisfied: scipy in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (1.11.1)
Requirement already satisfied: pandas>=2.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (2.0.3)
Requirement already satisfied: joblib>=1.2.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (1.3.1)
Requirement already satisfied: jinja2 in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (3.1.2)
Requirement already satisfied: numexpr in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (2.8.4)
Requirement already satisfied: funcy in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (2.0)
Requirement already satisfied: scikit-learn>=1.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (1.3.0)
Requirement already satisfied: gensim in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (4.3.0)
Requirement already satisfied: setuptools in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (68.0.0)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\programdata\anaconda3\lib\site-packages (from pandas>=2.0.0->pyLDAvis) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in c:\programdata\anaconda3\lib\site-packages (from pandas>=2.0.0->pyLDAvis) (2022.7)
Requirement already satisfied: tzdata>=2022.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pandas>=2.0.0->pyLDAvis) (2023.3)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from scikit-learn>=1.0.0->pyLDAvis) (3.1.0)
Requirement already satisfied: smart-open>=1.8.1 in c:\programdata\anaconda3\lib\site-packages (from gensim->pyLDAvis) (5.2.1)
Requirement already satisfied: FuzzyTM>=0.4.0 in c:\programdata\anaconda3\lib\site-packages (from gensim->pyLDAvis) (2.0.5)
Requirement already satisfied: MarkupSafe>=2.0 in c:\programdata\anaconda3\lib\site-packages (from jinja2->pyLDAvis) (2.1.1)
Requirement already satisfied: pyfume in c:\programdata\anaconda3\lib\site-packages (from FuzzyTM>=0.4.0->gensim->pyLDAvis) (0.2.25)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil>=2.8.2->pandas>=2.0.0->pyLDAvis) (1.16.0)
Requirement already satisfied: simpful in c:\programdata\anaconda3\lib\site-packages (from pyfume->FuzzyTM>=0.4.0->gensim->pyLDAvis) (2.11.0)
Requirement already satisfied: fst-pso in c:\programdata\anaconda3\lib\site-packages (from pyfume->FuzzyTM>=0.4.0->gensim->pyLDAvis) (1.8.1)
Requirement already satisfied: miniful in c:\programdata\anaconda3\lib\site-packages (from fst-pso->pyfume->FuzzyTM>=0.4.0->gensim->pyLDAvis) (0.0.6)
Installing collected packages: numpy, pyLDAvis
  Attempting uninstall: numpy
    Found existing installation: numpy 1.23.5
    Uninstalling numpy-1.23.5:
      Successfully uninstalled numpy-1.23.5
  Rolling back uninstall of numpy
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\entry_points.txt
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\entry_points.txt
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\installer
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\installer
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\license.txt
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\license.txt
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\licenses_bundled.txt
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\licenses_bundled.txt
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\metadata
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\metadata
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\record
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\record
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\top_level.txt
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\top_level.txt
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy-1.23.5.dist-info\wheel
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qly96gw6\wheel
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__config__.py
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__config__.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__init__.cython-30.pxd
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__init__.cython-30.pxd
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__init__.pxd
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__init__.pxd
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__init__.py
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__init__.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__init__.pyi
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__init__.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\__config__.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\__config__.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\__init__.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\__init__.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\_distributor_init.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\_distributor_init.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\_globals.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\_globals.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\_pytesttester.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\_pytesttester.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\_version.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\_version.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\conftest.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\conftest.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\ctypeslib.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\ctypeslib.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\matlib.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\matlib.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\setup.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\setup.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\__pycache__\version.cpython-310.pyc
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\__pycache__\version.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\_distributor_init.py
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\_distributor_init.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\_globals.py
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\_globals.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\_pyinstaller\
   from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-yinstaller
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\_pytesttester.py
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\_pytesttester.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\_pytesttester.pyi
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\_pytesttester.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\_typing\
   from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-yping
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\_version.py
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\_version.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__init__.py
   from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__init__.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\__init__.cpython-310.pyc
 s\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\__init__.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_array_object.cpython-310.pyc
```

Loading [MathJax]/extensions/Safe.js

```
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_array_object.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_constants.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_constants.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_creation_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_creation_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_data_type_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_data_type_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_dtypes.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_dtypes.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_elementwise_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_elementwise_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_manipulation_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_manipulation_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_searching_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_searching_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_set_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_set_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_sorting_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_sorting_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_statistical_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_statistical_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_typing.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_typing.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\_utility_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\_utility_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\linalg.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\linalg.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\__pycache__\setup.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\__pycache__\setup.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_array_object.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_array_object.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_constants.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_constants.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_creation_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_creation_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_data_type_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_data_type_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_dtypes.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_dtypes.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_elementwise_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_elementwise_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_manipulation_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_manipulation_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_searching_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_searching_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_set_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_set_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_sorting_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_sorting_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_statistical_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_statistical_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_typing.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_typing.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\_utility_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\_utility_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\linalg.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\linalg.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\setup.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\setup.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__init__.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__init__.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\__init__.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\__init__.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\test_array_object.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\test_array_object.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\test_creation_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\test_creation_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\test_data_type_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\test_data_type_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\test_elementwise_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\test_elementwise_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\test_set_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\test_set_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\test_sorting_functions.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\test_sorting_functions.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\__pycache__\test_validation.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\__pycache__\test_validation.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\test_array_object.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\test_array_object.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\test_creation_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\test_creation_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\test_data_type_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\test_data_type_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\test_elementwise_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\test_elementwise_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\test_set_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\test_set_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\test_sorting_functions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\test_sorting_functions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\array_api\tests\test_validation.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\array_api\tests\test_validation.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\compat\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-mpat
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\conftest.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\conftest.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\__init__.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\__init__.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\__init__.pyi
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\__init__.pyi
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\__pycache__\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\core\~_pycache__
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_add_newdocs.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_add_newdocs.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_add_newdocs_scalars.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_add_newdocs_scalars.py
    [Loading [MathJax]/extensions/Safe.js] programdata\anaconda3\lib\site-packages\numpy\core\_asarray.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_asarray.py
```

```
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_asarray.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_asarray.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_dtype.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_dtype.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_dtype_ctypes.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_dtype_ctypes.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_exceptions.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_exceptions.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_internal.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_internal.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_internal.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_internal.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_machar.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_machar.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_methods.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_methods.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_multiarray_tests.cp310-win_amd64.pyd
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_multiarray_tests.cp310-win_amd64.pyd
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_multiarray_umath.cp310-win_amd64.pyd
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_multiarray_umath.cp310-win_amd64.pyd
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_operand_flag_tests.cp310-win_amd64.pyd
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_operand_flag_tests.cp310-win_amd64.pyd
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_rational_tests.cp310-win_amd64.pyd
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_rational_tests.cp310-win_amd64.pyd
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_simd.cp310-win_amd64.pyd
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_simd.cp310-win_amd64.pyd
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_string_helpers.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_string_helpers.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_struct_ufunc_tests.cp310-win_amd64.pyd
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_struct_ufunc_tests.cp310-win_amd64.pyd
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_type_aliases.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_type_aliases.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_type_aliases.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_type_aliases.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_ufunc_config.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_ufunc_config.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_ufunc_config.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_ufunc_config.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\_umath_tests.cp310-win_amd64.pyd
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\_umath_tests.cp310-win_amd64.pyd
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\arrayprint.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\arrayprint.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\arrayprint.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\arrayprint.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\cversions.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\cversions.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\defchararray.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\defchararray.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\defchararray.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\defchararray.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\einsumfunc.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\einsumfunc.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\einsumfunc.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\einsumfunc.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\fromnumeric.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\fromnumeric.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\fromnumeric.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\fromnumeric.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\function_base.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\function_base.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\function_base.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\function_base.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\generate_numpy_api.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\generate_numpy_api.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\getlimits.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\getlimits.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\getlimits.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\getlimits.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\.doxyfile
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\.doxyfile
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\__multiarray_api.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\__multiarray_api.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\__ufunc_api.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\__ufunc_api.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\_neighborhood_iterator_imp.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\_neighborhood_iterator_imp.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\_numpyconfig.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\_numpyconfig.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\arrayobject.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\arrayobject.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\arrayscalars.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\arrayscalars.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\experimental_dtype_api.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\experimental_dtype_api.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\halffloat.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\halffloat.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\libdivide\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\core\include\numpy\~ibdivide
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\ndarrayobject.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\ndarrayobject.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\ndarraytypes.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\ndarraytypes.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\noprefix.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\noprefix.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_1_7_deprecated_api.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_1_7_deprecated_api.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_3kcompat.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_3kcompat.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_common.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_common.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_cpu.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_cpu.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_endian.h
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_endian.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_interrupt.h
  s\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_interrupt.h
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_math.h
```

Loading [MathJax]/extensions/Safe.js

```
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_math.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_no_deprecated_api.h
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_no_deprecated_api.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\npy_os.h
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\npy_os.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\numpyconfig.h
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\numpyconfig.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\old_defines.h
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\old_defines.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\oldnumeric.h
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\oldnumeric.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\random\
    from C:\ProgramData\anaconda3\Lib\site-packages\numpy\core\include\numpy\~andom
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\ufuncobject.h
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\ufuncobject.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\include\numpy\utils.h
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\include\numpy\utils.h
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\lib\
    from C:\ProgramData\anaconda3\Lib\site-packages\numpy\core\~ib
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\memmap.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\memmap.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\memmap.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\memmap.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\multiarray.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\multiarray.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\multiarray.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\multiarray.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\numeric.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\numeric.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\numeric.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\numeric.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\numerictypes.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\numerictypes.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\numerictypes.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\numerictypes.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\overrides.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\overrides.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\records.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\records.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\records.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\records.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\setup.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\setup.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\setup_common.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\setup_common.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\shape_base.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\shape_base.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\shape_base.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\shape_base.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__init__.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__init__.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\__init__.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\__init__.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\_locales.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\_locales.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test__exceptions.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test__exceptions.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_abc.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_abc.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_api.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_api.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_argparse.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_argparse.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_array_coercion.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_array_coercion.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_array_interface.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_array_interface.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_arraymethod.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_arraymethod.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_arrayprint.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_arrayprint.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_casting_unittests.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_casting_unittests.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_conversion_utils.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_conversion_utils.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_cpu_dispatcher.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_cpu_dispatcher.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_cpu_features.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_cpu_features.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_custom_dtypes.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_custom_dtypes.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_cython.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_cython.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_datetime.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_datetime.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_defchararray.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_defchararray.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_deprecations.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_deprecations.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_dlpack.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_dlpack.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_dtype.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_dtype.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_einsum.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_einsum.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_errstate.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_errstate.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_extint128.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_extint128.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_function_base.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_function_base.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_getlimits.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_getlimits.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_half.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_half.cpython-310.pyc
  [Loading [MathJax]/extensions/Safe.js] programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_hashtable.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_hashtable.cpython-310.pyc
```

```
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_indexerrors.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_indexerrors.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_indexing.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_indexing.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_item_selection.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_item_selection.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_limited_api.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_limited_api.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_longdouble.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_longdouble.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_machar.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_machar.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_mem_overlap.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_mem_overlap.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_mem_policy.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_mem_policy.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_memmap.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_memmap.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_multiarray.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_multiarray.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_nditer.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_nditer.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_numeric.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_numeric.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_numerictypes.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_numerictypes.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_overrides.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_overrides.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_print.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_print.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_protocols.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_protocols.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_records.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_records.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_regression.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_regression.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_scalar_ctors.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_scalar_ctors.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_scalar_methods.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_scalar_methods.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_scalarbuffer.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_scalarbuffer.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_scalarinherit.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_scalarinherit.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_scalarmath.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_scalarmath.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_scalarprint.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_scalarprint.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_shape_base.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_shape_base.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_simd.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_simd.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_simd_module.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_simd_module.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_ufunc.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_ufunc.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_umath.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_umath.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_umath_accuracy.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_umath_accuracy.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_umath_complex.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_umath_complex.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\__pycache__\test_unicode.cpython-310.pyc
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\__pycache__\test_unicode.cpython-310.pyc
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\_locales.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\_locales.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\data\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\core\tests\~ata
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\examples\cython\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\core\tests\examples\~ython
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\examples\limited_api\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\core\tests\examples\~imited_api
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test__exceptions.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test__exceptions.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_abc.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_abc.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_api.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_api.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_argparse.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_argparse.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_array_coercion.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_array_coercion.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_array_interface.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_array_interface.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_arraymethod.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_arraymethod.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_arrayprint.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_arrayprint.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_casting_unittests.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_casting_unittests.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_conversion_utils.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_conversion_utils.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_cpu_dispatcher.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_cpu_dispatcher.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_cpu_features.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_cpu_features.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_custom_dtypes.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_custom_dtypes.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_cython.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_cython.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_datetime.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_datetime.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_defchararray.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_defchararray.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_deprecations.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_deprecations.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_dlpack.py
```

Loading [MathJax]/extensions/Safe.js

```
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_dlpack.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_dtype.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_dtype.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_einsum.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_einsum.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_errstate.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_errstate.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_extint128.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_extint128.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_function_base.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_function_base.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_getlimits.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_getlimits.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_half.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_half.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_hashtable.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_hashtable.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_indexerrors.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_indexerrors.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_indexing.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_indexing.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_item_selection.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_item_selection.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_limited_api.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_limited_api.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_longdouble.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_longdouble.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_machar.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_machar.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_mem_overlap.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_mem_overlap.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_mem_policy.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_mem_policy.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_memmap.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_memmap.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_multiarray.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_multiarray.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_nditer.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_nditer.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_numeric.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_numeric.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_numerictypes.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_numerictypes.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_overrides.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_overrides.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_print.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_print.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_protocols.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_protocols.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_records.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_records.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_regression.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_regression.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_scalar_ctors.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_scalar_ctors.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_scalar_methods.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_scalar_methods.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_scalarbuffer.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_scalarbuffer.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_scalarinherit.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_scalarinherit.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_scalarmath.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_scalarmath.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_scalarprint.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_scalarprint.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_shape_base.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_shape_base.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_simd.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_simd.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_simd_module.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_simd_module.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_ufunc.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_ufunc.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_umath.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_umath.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_umath_accuracy.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_umath_accuracy.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_umath_complex.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_umath_complex.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\tests\test_unicode.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\tests\test_unicode.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\umath.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\umath.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\core\umath_tests.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\core\umath_tests.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\ctypeslib.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\ctypeslib.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\ctypeslib.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\ctypeslib.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__config__.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__config__.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__init__.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__init__.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__init__.pyi
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__init__.pyi
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\__config__.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\__config__.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\__init__.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\__init__.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\_shell_utils.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\_shell_utils.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\armccompiler.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\armccompiler.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\ccompiler.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\ccompiler.cpython-310.pyc
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\ccompiler_opt.cpython-310.pyc
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\ccompiler_opt.cpython-310.pyc
```

Loading [MathJax]/extensions/Safe.js

```
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\conv_template.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\conv_template.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\core.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\core.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\cpuinfo.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\cpuinfo.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\exec_command.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\exec_command.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\extension.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\extension.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\from_template.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\from_template.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\intelccompiler.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\intelccompiler.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\lib2def.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\lib2def.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\line_endings.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\line_endings.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\log.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\log.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\mingw32ccompiler.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\mingw32ccompiler.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\misc_util.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\misc_util.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\msvc9compiler.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\msvc9compiler.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\msvccompiler.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\msvccompiler.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\npy_pkg_config.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\npy_pkg_config.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\numpy_distribution.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\numpy_distribution.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\pathccompiler.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\pathccompiler.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\setup.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\setup.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\system_info.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\system_info.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\__pycache__\unixccompiler.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\__pycache__\unixccompiler.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\_shell_utils.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\_shell_utils.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\armccompiler.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\armccompiler.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\ccompiler.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\ccompiler.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\ccompiler_opt.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\ccompiler_opt.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_asimd.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_asimd.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_asimddp.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_asimddp.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_asimdfhm.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_asimdfhm.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_asimdhp.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_asimdhp.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx2.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx2.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512_clx.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512_clx.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512_cnl.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512_cnl.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512_icl.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512_icl.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512_knl.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512_knl.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512_knm.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512_knm.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512_skx.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512_skx.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512cd.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512cd.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_avx512f.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_avx512f.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_f16c.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_f16c.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_fma3.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_fma3.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_fma4.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_fma4.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_neon.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_neon.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_neon_fp16.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_neon_fp16.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_neon_vfpv4.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_neon_vfpv4.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_popcnt.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_popcnt.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_sse.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_sse.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_sse2.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_sse2.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_sse3.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_sse3.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_sse41.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_sse41.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_sse42.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_sse42.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_ssse3.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_ssse3.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_vsx.c
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_vsx.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_vsx2.c
  s\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_vsx2.c
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_vsx3.c
```

Loading [MathJax]/extensions/Safe.js

```
         from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_vsx3.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_vsx4.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_vsx4.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_vx.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_vx.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_vxe.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_vxe.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_vxe2.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_vxe2.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\cpu_xop.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\cpu_xop.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\extra_avx512bw_mask.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\extra_avx512bw_mask.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\extra_avx512dq_mask.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\extra_avx512dq_mask.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\extra_avx512f_reduce.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\extra_avx512f_reduce.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\extra_vsx4_mma.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\extra_vsx4_mma.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\extra_vsx_asm.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\extra_vsx_asm.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\checks\test_flags.c
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\checks\test_flags.c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\command\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\distutils\~-mmand
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\conv_template.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\conv_template.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\core.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\core.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\cpuinfo.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\cpuinfo.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\exec_command.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\exec_command.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\extension.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\extension.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\fcompiler\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\distutils\~-ompiler
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\from_template.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\from_template.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\intelccompiler.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\intelccompiler.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\lib2def.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\lib2def.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\line_endings.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\line_endings.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\log.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\log.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\mingw32ccompiler.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\mingw32ccompiler.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\mingw\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\distutils\~-ngw
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\misc_util.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\misc_util.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\msvc9compiler.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\msvc9compiler.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\msvccompiler.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\msvccompiler.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\npy_pkg_config.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\npy_pkg_config.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\numpy_distribution.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\numpy_distribution.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\pathccompiler.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\pathccompiler.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\setup.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\setup.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\system_info.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\system_info.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\tests\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\distutils\~-sts
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\distutils\unixccompiler.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\distutils\unixccompiler.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\doc\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-c
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\dual.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\dual.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\__init__.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\__init__.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\__init__.pyi
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\__init__.pyi
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\__main__.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\__main__.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\__pycache__\
      from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\~_pycache__
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\__version__.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\__version__.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\auxfuncs.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\auxfuncs.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\capi_maps.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\capi_maps.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\cb_rules.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\cb_rules.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\cfuncs.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\cfuncs.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\common_rules.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\common_rules.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\crackfortran.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\crackfortran.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\diagnose.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\diagnose.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\f2py2e.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\f2py2e.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\f90mod_rules.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\f90mod_rules.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\func2subr.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\func2subr.py
    Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\rules.py
      from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\rules.py
```

Loading [MathJax]/extensions/Safe.js

```
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\setup.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\setup.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\src\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\~rc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\symbolic.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\symbolic.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__init__.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__init__.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\__init__.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\__init__.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_abstract_interface.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_abstract_interface.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_array_from_pyobj.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_array_from_pyobj.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_assumed_shape.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_assumed_shape.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_block_docstring.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_block_docstring.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_callback.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_callback.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_common.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_common.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_compile_function.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_compile_function.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_crackfortran.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_crackfortran.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_f2cmap.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_f2cmap.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_f2py2e.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_f2py2e.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_kind.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_kind.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_mixed.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_mixed.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_module_doc.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_module_doc.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_parameter.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_parameter.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_quoted_character.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_quoted_character.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_regression.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_regression.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_return_character.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_return_character.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_return_complex.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_return_complex.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_return_integer.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_return_integer.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_return_logical.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_return_logical.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_return_real.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_return_real.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_semicolon_split.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_semicolon_split.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_size.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_size.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_string.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_string.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\test_symbolic.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\test_symbolic.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\__pycache__\util.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\__pycache__\util.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\abstract_interface\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~bstract_interface
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\array_from_pyobj\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~rray_from_pyobj
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\assumed_shape\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~ssumed_shape
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\block_docstring\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~lock_docstring
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\callback\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~allback
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\cli\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~li
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\common\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~ommon
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\accesstype.f90
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\accesstype.f90
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\foo_deps.f90
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\foo_deps.f90
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\gh15035.f
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\gh15035.f
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\gh17859.f
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\gh17859.f
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\gh2848.f90
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\gh2848.f90
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\operators.f90
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\operators.f90
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\privatemod.f90
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\privatemod.f90
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\crackfortran\publicmod.f90
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\crackfortran\publicmod.f90
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\f2cmap\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~2cmap
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\kind\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~ind
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\mixed\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~ixed
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\module_data\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~odule_data
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\negative_bounds\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~egative_bounds
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\parameter\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~arameter
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\quoted_character\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~uoted_character
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\regression\
```

Loading [MathJax]/extensions/Safe.js

```
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~egression
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\return_character\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~eturn_character
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\return_complex\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~eturn_complex
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\return_integer\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~eturn_integer
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\return_logical\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~eturn_logical
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\return_real\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~eturn_real
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\size\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\f2py\tests\src\~ize
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\string\char.f90
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\string\char.f90
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\string\fixed_string.f90
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\string\fixed_string.f90
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\src\string\string.f
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\src\string\string.f
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_abstract_interface.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_abstract_interface.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_array_from_pyobj.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_array_from_pyobj.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_assumed_shape.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_assumed_shape.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_block_docstring.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_block_docstring.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_callback.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_callback.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_common.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_common.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_compile_function.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_compile_function.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_crackfortran.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_crackfortran.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_f2cmap.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_f2cmap.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_f2py2e.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_f2py2e.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_kind.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_kind.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_mixed.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_mixed.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_module_doc.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_module_doc.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_parameter.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_parameter.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_quoted_character.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_quoted_character.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_regression.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_regression.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_return_character.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_return_character.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_return_complex.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_return_complex.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_return_integer.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_return_integer.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_return_logical.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_return_logical.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_return_real.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_return_real.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_semicolon_split.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_semicolon_split.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_size.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_size.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_string.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_string.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\test_symbolic.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\test_symbolic.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\tests\util.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\tests\util.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\f2py\use_rules.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\f2py\use_rules.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\fft\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-t
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\lib\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-b
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\license.txt
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\license.txt
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\linalg\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-nalg
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\ma\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~ma
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\matlib.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\matlib.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\matrixlib\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-trixlib
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\__init__.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\__init__.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\__init__.pyi
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\__init__.pyi
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\__pycache__\
           from C:\ProgramData\anaconda3\Lib\site-packages\numpy\polynomial\~_pycache__
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\_polybase.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\_polybase.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\_polybase.pyi
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\_polybase.pyi
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\chebyshev.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\chebyshev.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\chebyshev.pyi
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\chebyshev.pyi
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\hermite.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\hermite.py
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\hermite.pyi
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\hermite.pyi
         Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\hermite_e.py
           from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\hermite_e.py
```

Loading [MathJax]/extensions/Safe.js

```
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\hermite_e.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\hermite_e.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\laguerre.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\laguerre.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\laguerre.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\laguerre.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\legendre.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\legendre.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\legendre.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\legendre.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\polynomial.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\polynomial.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\polynomial.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\polynomial.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\polyutils.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\polyutils.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\polyutils.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\polyutils.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\setup.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\setup.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__init__.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__init__.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\__init__.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\__init__.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_chebyshev.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_chebyshev.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_classes.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_classes.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_hermite.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_hermite.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_hermite_e.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_hermite_e.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_laguerre.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_laguerre.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_legendre.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_legendre.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_polynomial.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_polynomial.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_polyutils.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_polyutils.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\__pycache__\test_printing.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\__pycache__\test_printing.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_chebyshev.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_chebyshev.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_classes.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_classes.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_hermite.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_hermite.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_hermite_e.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_hermite_e.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_laguerre.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_laguerre.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_legendre.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_legendre.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_polynomial.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_polynomial.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_polyutils.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_polyutils.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\polynomial\tests\test_printing.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\polynomial\tests\test_printing.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\py.typed
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\py.typed
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\random\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-ndom
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\setup.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\setup.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\__init__.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\testing\__init__.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\__init__.pyi
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\testing\__init__.pyi
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\__pycache__\__init__.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\testing\__pycache__\__init__.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\__pycache__\print_coercion_tables.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\testing\__pycache__\print_coercion_tables.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\__pycache__\setup.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\testing\__pycache__\setup.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\_private\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\testing\~private
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\print_coercion_tables.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\testing\print_coercion_tables.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\setup.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\testing\setup.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\testing\tests\
  from C:\ProgramData\anaconda3\Lib\site-packages\numpy\testing\~ests
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__init__.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__init__.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\__init__.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\__init__.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test__all__.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test__all__.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test_ctypeslib.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test_ctypeslib.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test_matlib.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test_matlib.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test_numpy_version.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test_numpy_version.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test_public_api.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test_public_api.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test_reloading.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test_reloading.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test_scripts.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test_scripts.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\__pycache__\test_warnings.cpython-310.pyc
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\__pycache__\test_warnings.cpython-310.pyc
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test__all__.py
  from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test__all__.py
Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test_ctypeslib.py
```

Loading [MathJax]/extensions/Safe.js

```
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test_ctypeslib.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test_matlib.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test_matlib.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test_numpy_version.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test_numpy_version.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test_public_api.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test_public_api.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test_reloading.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test_reloading.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test_scripts.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test_scripts.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\tests\test_warnings.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\tests\test_warnings.py
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\typing\
    from C:\ProgramData\anaconda3\Lib\site-packages\numpy\~-ping
  Moving to c:\programdata\anaconda3\lib\site-packages\numpy\version.py
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-gqyeen0f\version.py
  Moving to c:\programdata\anaconda3\scripts\f2py.exe
    from C:\Users\purus\AppData\Local\Temp\pip-uninstall-qlvrf_kt\f2py.exe
```

ERROR: Could not install packages due to an OSError: [WinError 5] Access is denied: 'C:\\ProgramData\\anaconda3\\Lib\\site-packages\\numpy\\.libs\\li
bopenblas64__v0.3.23-gcc_10_3_0.dll'
Consider using the `--user` option or check the permissions.

```
Requirement already satisfied: gensim in c:\programdata\anaconda3\lib\site-packages (4.3.0)
Requirement already satisfied: numpy>=1.18.5 in c:\programdata\anaconda3\lib\site-packages (from gensim) (1.23.5)
Requirement already satisfied: scipy>=1.7.0 in c:\programdata\anaconda3\lib\site-packages (from gensim) (1.11.1)
Requirement already satisfied: smart-open>=1.8.1 in c:\programdata\anaconda3\lib\site-packages (from gensim) (5.2.1)
Requirement already satisfied: FuzzyTM>=0.4.0 in c:\programdata\anaconda3\lib\site-packages (from gensim) (2.0.5)
Requirement already satisfied: pandas in c:\users\purus\appdata\roaming\python\python310\site-packages (from FuzzyTM>=0.4.0->gensim) (2.0.3)
Requirement already satisfied: pyfume in c:\programdata\anaconda3\lib\site-packages (from FuzzyTM>=0.4.0->gensim) (0.2.25)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\programdata\anaconda3\lib\site-packages (from pandas->FuzzyTM>=0.4.0->gensim) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in c:\programdata\anaconda3\lib\site-packages (from pandas->FuzzyTM>=0.4.0->gensim) (2022.7)
Requirement already satisfied: tzdata>=2022.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pandas->FuzzyTM>=0.4.0->gensim)
(2023.3)
Requirement already satisfied: simpful in c:\programdata\anaconda3\lib\site-packages (from pyfume->FuzzyTM>=0.4.0->gensim) (2.11.0)
Requirement already satisfied: fst-pso in c:\programdata\anaconda3\lib\site-packages (from pyfume->FuzzyTM>=0.4.0->gensim) (1.8.1)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil>=2.8.2->pandas->FuzzyTM>=0.4.0->gensim)
(1.16.0)
Requirement already satisfied: miniful in c:\programdata\anaconda3\lib\site-packages (from fst-pso->pyfume->FuzzyTM>=0.4.0->gensim) (0.0.6)
```

In [35]:
```
!pip install pyLDAvis
import pyLDAvis
```

```
Requirement already satisfied: pyLDAvis in c:\programdata\anaconda3\lib\site-packages (3.4.0)
Requirement already satisfied: numpy>=1.22.0 in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (1.23.5)
Requirement already satisfied: scipy in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (1.11.1)
Requirement already satisfied: pandas>=1.3.4 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (2.0.3)
Requirement already satisfied: joblib>=1.2.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (1.3.1)
Requirement already satisfied: jinja2 in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (3.1.2)
Requirement already satisfied: numexpr in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (2.8.4)
Requirement already satisfied: funcy in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (2.0)
Requirement already satisfied: scikit-learn>=1.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pyLDAvis) (1.3.0)
Requirement already satisfied: gensim in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (4.3.0)
Requirement already satisfied: setuptools in c:\programdata\anaconda3\lib\site-packages (from pyLDAvis) (68.0.0)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\programdata\anaconda3\lib\site-packages (from pandas>=1.3.4->pyLDAvis) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in c:\programdata\anaconda3\lib\site-packages (from pandas>=1.3.4->pyLDAvis) (2022.7)
Requirement already satisfied: tzdata>=2022.1 in c:\users\purus\appdata\roaming\python\python310\site-packages (from pandas>=1.3.4->pyLDAvis) (2023.
3)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from scikit-learn>=1.0.0->pyLDA
vis) (3.1.0)
Requirement already satisfied: smart-open>=1.8.1 in c:\programdata\anaconda3\lib\site-packages (from gensim->pyLDAvis) (5.2.1)
Requirement already satisfied: FuzzyTM>=0.4.0 in c:\programdata\anaconda3\lib\site-packages (from gensim->pyLDAvis) (2.0.5)
Requirement already satisfied: MarkupSafe>=2.0 in c:\programdata\anaconda3\lib\site-packages (from jinja2->pyLDAvis) (2.1.1)
Requirement already satisfied: pyfume in c:\programdata\anaconda3\lib\site-packages (from FuzzyTM>=0.4.0->gensim->pyLDAvis) (0.2.25)
Requirement already satisfied: six>=1.5 in c:\programdata\anaconda3\lib\site-packages (from python-dateutil>=2.8.2->pandas>=1.3.4->pyLDAvis) (1.16.0)
Requirement already satisfied: simpful in c:\programdata\anaconda3\lib\site-packages (from pyfume->FuzzyTM>=0.4.0->gensim->pyLDAvis) (2.11.0)
Requirement already satisfied: fst-pso in c:\programdata\anaconda3\lib\site-packages (from pyfume->FuzzyTM>=0.4.0->gensim->pyLDAvis) (1.8.1)
Requirement already satisfied: miniful in c:\programdata\anaconda3\lib\site-packages (from fst-pso->pyfume->FuzzyTM>=0.4.0->gensim->pyLDAvis) (0.0.6)
```

In [36]:
```python
import gensim
import pyLDAvis.gensim as gensimvis
from gensim import corpora

# list of documents called 'words'
# Each document in 'words' is a string containing the text

# Preprocess the data to tokenize and remove stopwords
preprocessed_data = []
for doc in words:
    tokens = gensim.utils.simple_preprocess(doc)
    preprocessed_data.append(tokens)

# Create a dictionary from the preprocessed data
dictionary = corpora.Dictionary(preprocessed_data)

# Create a document-term matrix
doc_term_matrix = [dictionary.doc2bow(doc) for doc in preprocessed_data]

# Build the LDA model
lda_model = gensim.models.LdaModel(doc_term_matrix, num_topics=10, id2word=dictionary)

# Visualize the results
vis = gensimvis.prepare(lda_model, doc_term_matrix, dictionary)
pyLDAvis.display(vis)
```

Loading [MathJax]/extensions/Safe.js

| Selected Topic: 0 | Previous Topic | Next Topic | Clear Topic |

Slide to adjust relevance metric:(2)

λ = 1

0.0  0.2  0.4  0.6

## Intertopic Distance Map (via multidimensional scaling)

## Top-14 Most Salient Terms(1)



PC2

PC1

Marginal topic distribution

2%
5%
10%

1

mississippi
on
was
money
from
river
stealing
fishing
seen
bank
the
after
robber
vault

Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see C
2. relevance(term w | topic t) = λ * p(w | t) + (1 - λ) * p(w | t)/p(w); see Sievert & Shirley (

In [37]:
```python
print('the array representation of the word \'bank\'\n:',model_cbow.wv['bank'], '\n the array representation of the word \'robber\'\n:', model_cbow.w
```

```
the array representation of the word 'bank'
: [-0.01259352  0.00555269  0.11985431  0.21158655 -0.21848367 -0.16714124
  0.15168934  0.2107344  -0.11778943 -0.08838436  0.17333424 -0.03601422
 -0.10654427  0.15392466 -0.1141429  -0.04264994  0.06755769  0.02329457
 -0.19458175 -0.22190945  0.1717199   0.11907724  0.15870726  0.01791622
  0.14915332 -0.07997645 -0.02222663  0.13547736 -0.17664881 -0.09244104
 -0.17641266 -0.02184243  0.22400671 -0.1718937  -0.05480954 -0.04550866
  0.189702   -0.13928957  0.00106066 -0.11164343 -0.22554341  0.11759838
 -0.20572253 -0.10314386 -0.00082434 -0.00695595 -0.17992744  0.22580628
  0.11700572  0.21684425]
 the array representation of the word 'robber'
: [-0.18355328  0.03191495 -0.18450207 -0.05770317  0.09562866  0.14996925
  0.03079692  0.05406347 -0.10566583  0.18576439 -0.16216265  0.11947794
 -0.21133634  0.05234125 -0.12794407 -0.10921183 -0.07994997  0.14537773
  0.14907673 -0.12788668  0.019877   -0.21842225  0.20076823  0.2379963
 -0.07050339  0.02056307  0.01919113  0.14083774 -0.22125545  0.01502903
  0.17659806  0.05736633  0.0289223  -0.23967153  0.21808493 -0.16104512
 -0.07691811  0.08984208 -0.01986157  0.0362968   0.04580666 -0.17556296
 -0.25003898  0.23244233  0.159337  -0.17774613  0.08750858  0.00531226
  0.12221606 -0.18304008]
```

In [38]:
```python
print(model_cbow.wv.most_similar('bank', 'vault'))
```

```
[('robber', 0.15167683362960815), ('river', 0.12220965325832367), ('was', 0.09942746162414551), ('stealing', 0.08457799255847931), ('money', 0.08098877966403961), ('Mississippi', 0.08072850853204727), ('from', 0.07423313707113266), ('the', 0.030568802729249), ('After', 0.01814216934144497), ('on', -0.04082303121685982)]
```

In [39]:
```python
from sklearn.manifold import TSNE
import numpy as np
import matplotlib.pyplot as plt

def plot_tsne(model, num):
    labels = model.wv.index_to_key[:num]
    tokens = model.wv[model.wv.index_to_key[:num]]

    tsne = TSNE(perplexity=4, n_components=2, init='pca', n_iter=250, random_state=42)
    data = tsne.fit_transform(tokens)

    x = data[:, 0]
    y = data[:, 1]

    plt.figure(figsize=(10, 10))
    plt.scatter(x, y)

    for i, label in enumerate(labels):
        plt.annotate(label, xy=(x[i], y[i]), xytext=(5, 2), textcoords='offset points', ha='right', va='bottom')

    plt.show()

plot_tsne(model_cbow, 300)
```

Loading [MathJax]/extensions/Safe.js

## Word2Vec Skip Gram

```python
In [40]: model_skipgram = Word2Vec([text.split(" ")], min_count = 1, vector_size = 50, workers = 3, window = 5, sg = 1)
```

```python
In [41]: model_skipgram.init_sims(replace = True)
         model_skipgram.train([["hello", "world"]], total_examples=1, epochs=1)
```

C:\Users\purus\AppData\Local\Temp\ipykernel_16696\30155734.py:1: DeprecationWarning: Call to deprecated `init_sims` (Gensim 4.0.0 implemented interna
l optimizations that make calls to init_sims() unnecessary. init_sims() is now obsoleted and will be completely removed in future versions. See http
s://github.com/RaRe-Technologies/gensim/wiki/Migrating-from-Gensim-3.x-to-4).
  model_skipgram.init_sims(replace = True)
C:\ProgramData\anaconda3\lib\site-packages\gensim\models\word2vec.py:913: DeprecationWarning: Call to deprecated `init_sims` (Use fill_norms() instea
d. See https://github.com/RaRe-Technologies/gensim/wiki/Migrating-from-Gensim-3.x-to-4).
  self.wv.init_sims(replace=replace)

```
Out[41]: (0, 2)
```

```python
In [42]: words = list(model_skipgram.wv.index_to_key)
         print(words)
```

['bank', 'the', 'river', 'Mississippi', 'on', 'fishing', 'seen', 'was', 'robber', 'vault', 'from', 'money', 'stealing', 'After']

```python
In [43]: vector = model_skipgram.wv['bank']
         vector
```

```
Out[43]: array([-0.01259352,  0.00555269,  0.11985431,  0.21158655, -0.21848367,
               -0.16714124,  0.15168934,  0.2107344 , -0.11778943, -0.08838436,
                0.17333424, -0.03601422, -0.10654427,  0.15392466, -0.1141429 ,
               -0.04264994,  0.06755769,  0.02329457, -0.19458175, -0.22190945,
                0.1717199 ,  0.11907724,  0.15870726,  0.01791622,  0.14915332,
               -0.07997645, -0.02222663,  0.13547736, -0.17664881, -0.09244104,
               -0.17641266, -0.02184243,  0.22400671, -0.1718937 , -0.05480954,
               -0.04550866,  0.189702  , -0.13928957,  0.00106066, -0.11164343,
               -0.22554341,  0.11759838, -0.20572253, -0.10314386, -0.00082434,
               -0.00695595, -0.17992744,  0.22580628,  0.11700572,  0.21684425],
              dtype=float32)
```

```python
In [44]: print(model_skipgram.wv.most_similar('river', 'bank'))
```

[('on', 0.21477800607681274), ('robber', 0.1479063332080841), ('seen', 0.11235106736421585), ('from', 0.1012321263551712), ('money', -0.0089773815125
22697), ('the', -0.023684391751885414), ('fishing', -0.030451636761426926), ('stealing', -0.05601589381694794), ('After', -0.09814958274364471), ('va
ult', -0.10375411063432693)]

```python
In [45]: print('the array representation of the word \'river\'\n:',model_cbow.wv['river'], '\n the array representation of the word \'bank\'\n:', model_cbow.wv
```

Loading [MathJax]/extensions/Safe.js

```
 the array representation of the word 'river'
: [-0.20642233  0.08778626  0.12428615  0.13750665  0.17881607 -0.1477021
  0.02647698  0.14481895 -0.06801289 -0.14784212 -0.00982393 -0.20041768
 -0.13410783  0.17013788  0.08028586  0.17303869  0.16285078  0.1803445
 -0.09074181 -0.013454    0.05623839 -0.10822076  0.20089145 -0.23608108
  0.1619981   0.0697938  -0.11813033  0.1053268  -0.04165894  0.16072272
  0.23863597 -0.10447081 -0.01435281 -0.13639784  0.09221862  0.06673351
  0.16502593  0.14610766  0.22842576  0.2220777   0.18914114 -0.16738306
 -0.21926259 -0.00851949 -0.07423428  0.18905132  0.14221562 -0.03701518
  0.03618421  0.04286749]
 the array representation of the word 'bank'
: [-0.01259352  0.00555269  0.11985431  0.21158655 -0.21848367 -0.16714124
  0.15168934  0.2107344  -0.11778943 -0.08838436  0.17333424 -0.03601422
 -0.10654427  0.15392466 -0.1141429  -0.04264994  0.06755769  0.02329457
 -0.19458175 -0.22190945  0.1717199   0.11907724  0.15870726  0.01791622
  0.14915332 -0.07997645 -0.02222663  0.13547736 -0.17664881 -0.09244104
 -0.17641266 -0.02184243  0.22400671 -0.1718937  -0.05480954 -0.04550866
  0.189702   -0.13928957  0.00106066 -0.11164343 -0.22554341  0.11759838
 -0.20572253 -0.10314386 -0.00082434 -0.00695595 -0.17992744  0.22580628
  0.11700572  0.21684425]
```

In [46]: `plot_tsne(model_skipgram,100)`



## Continuous Bag of Words (CBOW) and Skip-gram models in the context of generating embeddings and topic modeling:

| Continuous Bag of Words | skip-gram |
|---|---|
| The CBOW architecture, as the name suggests, predicts the target word based on the context words surrounding it. It takes a window of context words as input and tries to maximize the probability of predicting the target word. Here are some characteristics of the CBOW model | The Skip-gram architecture aims to predict the context words given a target word. It takes a target word as input and tries to maximize the probability of predicting the surrounding context words. Here are some characteristics of the Skip-gram model |
| Eg: In this sentence "After stealing money from the bank vault the bank robber was seen fishing on the Mississippi river bank" river, bank, Mississippi appears closer with less distance to each other | Eg: In this sentence "After stealing money from the bank vault the bank robber was seen fishing on the Mississippi river bank" robber, fishing, stealing appears closer with less distance to each other |

## BERT Vector Embeddings (With Pretrained)

In [47]:
```python
from pytorch_pretrained_bert import BertTokenizer
import sys
import numpy as np
import random as rn
import torch
from pytorch_pretrained_bert import BertModel
from torch import nn
from pytorch_pretrained_bert import BertTokenizer
from torch.utils.data import TensorDataset, DataLoader, RandomSampler, SequentialSampler
from torch.optim import Adam
from torch.nn.utils import clip_grad_norm_
        splay import clear_output
```

Loading [MathJax]/extensions/Safe.js

```python
import pandas as pd
import matplotlib.pyplot as plt
```

In [48]:
```python
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from_pretrained('bert-base-uncased', do_lower_case=True)
```

In [49]:
```python
tokenize_ = tokenizer.tokenize(text)
print("Text after tokenization: ")
print(tokenize_)
max_len = 25

textLst = tokenize_[:max_len-2]
input_sequence = ["[CLS]"] + textLst + ["[SEP]"]
pad_len = max_len - len(input_sequence)

print("After adding [CLS] and [SEP]: ")
print(input_sequence)
tokens = tokenizer.convert_tokens_to_ids(input_sequence)
print("After converting Tokens to Id: ")
print(tokens)
tokens += [0] * pad_len
print("tokens: ")
print(tokens)
pad_masks = [1] * len(input_sequence) + [0] * pad_len
print("Pad Masking: ")
print(pad_masks)
segment_ids = [0] * max_len
print("Segment Ids: ")
print(segment_ids)
```

```
Text after tokenization:
['after', 'stealing', 'money', 'from', 'the', 'bank', 'vault', 'the', 'bank', 'robber', 'was', 'seen', 'fishing', 'on', 'the', 'mississippi', 'rive
r', 'bank']
After adding [CLS] and [SEP]:
['[CLS]', 'after', 'stealing', 'money', 'from', 'the', 'bank', 'vault', 'the', 'bank', 'robber', 'was', 'seen', 'fishing', 'on', 'the', 'mississipp
i', 'river', 'bank', '[SEP]']
After converting Tokens to Id:
[101, 2044, 11065, 2769, 2013, 1996, 2924, 11632, 1996, 2924, 27307, 2001, 2464, 5645, 2006, 1996, 5900, 2314, 2924, 102]
tokens:
[101, 2044, 11065, 2769, 2013, 1996, 2924, 11632, 1996, 2924, 27307, 2001, 2464, 5645, 2006, 1996, 5900, 2314, 2924, 102, 0, 0, 0, 0, 0]
Pad Masking:
[1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0]
Segment Ids:
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
```

In [50]:
```python
marked_text = "[CLS] " + text + " [SEP]"

# Tokenize our sentence with the BERT tokenizer.
tokenized_text = tokenizer.tokenize(marked_text)
segments_ids = [1] * len(tokenized_text)

# Map the token strings to their vocabulary indeces.
indexed_tokens = tokenizer.convert_tokens_to_ids(tokenized_text)

# Print out the tokens.
print (tokenized_text)
```

```
['[CLS]', 'after', 'stealing', 'money', 'from', 'the', 'bank', 'vault', 'the', 'bank', 'robber', 'was', 'seen', 'fishing', 'on', 'the', 'mississipp
i', 'river', 'bank', '[SEP]']
```

In [51]:
```python
# Convert inputs to PyTorch tensors
tokens_tensor = torch.tensor([indexed_tokens])
segments_tensors = torch.tensor([segments_ids])

# Load pre-trained model (weights)
model = BertModel.from_pretrained('bert-base-multilingual-cased')

# Put the model in "evaluation" mode, meaning feed-forward operation.
model.eval()
```

Loading [MathJax]/extensions/Safe.js

```
Out[51]: BertModel(
           (embeddings): BertEmbeddings(
             (word_embeddings): Embedding(119547, 768, padding_idx=0)
             (position_embeddings): Embedding(512, 768)
             (token_type_embeddings): Embedding(2, 768)
             (LayerNorm): BertLayerNorm()
             (dropout): Dropout(p=0.1, inplace=False)
           )
           (encoder): BertEncoder(
             (layer): ModuleList(
               (0-11): 12 x BertLayer(
                 (attention): BertAttention(
                   (self): BertSelfAttention(
                     (query): Linear(in_features=768, out_features=768, bias=True)
                     (key): Linear(in_features=768, out_features=768, bias=True)
                     (value): Linear(in_features=768, out_features=768, bias=True)
                     (dropout): Dropout(p=0.1, inplace=False)
                   )
                   (output): BertSelfOutput(
                     (dense): Linear(in_features=768, out_features=768, bias=True)
                     (LayerNorm): BertLayerNorm()
                     (dropout): Dropout(p=0.1, inplace=False)
                   )
                 )
                 (intermediate): BertIntermediate(
                   (dense): Linear(in_features=768, out_features=3072, bias=True)
                 )
                 (output): BertOutput(
                   (dense): Linear(in_features=3072, out_features=768, bias=True)
                   (LayerNorm): BertLayerNorm()
                   (dropout): Dropout(p=0.1, inplace=False)
                 )
               )
             )
           )
           (pooler): BertPooler(
             (dense): Linear(in_features=768, out_features=768, bias=True)
             (activation): Tanh()
           )
         )
```

```python
In [52]: # Predict hidden states features for each layer
         with torch.no_grad():
             encoded_layers, _ = model(tokens_tensor, segments_tensors)

         # Concatenate the tensors for all layers. We use `stack` here to
         # create a new dimension in the tensor.
         token_embeddings = torch.stack(encoded_layers, dim=0)

         # Remove dimension 1, the "batches".
         token_embeddings = torch.squeeze(token_embeddings, dim=1)

         # Swap dimensions 0 and 1.
         token_embeddings = token_embeddings.permute(1,0,2)

         token_embeddings.size()
```

```
Out[52]: torch.Size([20, 12, 768])
```

```python
In [53]: # Stores the token vectors, with shape [23 x 768]
         token_vecs_sum = []

         # `token_embeddings` is a [23 x 12 x 768] tensor.

         # For each token in the sentence...
         for token in token_embeddings:

             # `token` is a [12 x 768] tensor

             # Sum the vectors from the last four layers.
             sum_vec = torch.sum(token[-4:], dim=0)

             # Use `sum_vec` to represent `token`.
             token_vecs_sum.append(sum_vec)

         print ('Shape is: %d x %d' % (len(token_vecs_sum), len(token_vecs_sum[0])))
```

```
Shape is: 20 x 768
```

```python
In [54]: # Stores the token vectors, with shape [23 x 768]
         token_vecs_sum = []

         # `token_embeddings` is a [23 x 12 x 768] tensor.

         # For each token in the sentence...
         for token in token_embeddings:

             # `token` is a [12 x 768] tensor

             # Sum the vectors from the last four layers.
             sum_vec = torch.sum(token[-4:], dim=0)

             # Use `sum_vec` to represent `token`.
             token_vecs_sum.append(sum_vec)

         print ('Shape is: %d x %d' % (len(token_vecs_sum), len(token_vecs_sum[0])))
```

```
Shape is: 20 x 768
```

```python
In [55]: # `encoded_layers` has shape [12 x 1 x 23 x 768]

         # `token_vecs` is a tensor with shape [23 x 768]
         token_vecs = encoded_layers[11][0]

         # Calculate the average of all 23 token vectors.
         ing = torch.mean(token_vecs, dim=0)
```

Loading [MathJax]/extensions/Safe.js

In [56]: sentence_embedding[0]

Out[56]: tensor(0.1697)

# Comparison of word2vec(Continuous Bag of Words,skip-gram), BERT

We have choose these two techniques and BERT is the most efficient among all

- word2vec without pre trained model (Building Corpus on Your Own)
- BERT (pre trained model) . Reason for choosing BERT is BERT's sensitivity to single-word cues in context, we draw on data from semantic priming observed in humans. BERT's ability to capture contextual meaning, bidirectional training, pretraining, fine-tuning, and task-agnostic nature make it a more powerful and flexible model compared to GloVe. However, it's worth noting that GloVe still has its merits and can be useful in certain scenarios where contextual information may not be crucial or where computational resources are limited.
- Advantage of using pre trained model is save time, resources, and money compared to building and training your own model. And they are often as effective and more efficient than custom models

| WordVec | Glove | BERT |
|---|---|---|
| Performs well on syntactic and semantic analogies and generate efficient and simple word embeddings.Predictive Model | Generates static word embeddings | Captures contextual word representations and bidirectional training for deeper word understanding |
| Pretrained models (corpus) available | Pretrained models (corpus) available | Pretraining and fine-tuning framework |
| Requires substantial computational resources | Requires substantial computational resources | Requires more computational resources. |
| Limited understanding of word relationships and dependencies. | Ignores contextual information. | BERT models look at the surrounding words to understand the context. |

# Discussion on the embeddings that provide better semantic understanding (BERT has better Understanding)

| Aspect | BERT | Word2Vec |
|---|---|---|
| Model Type | Contextualized word embeddings | Static word embeddings |
| Semantic | Captures contextualized meaning, considering surrounding context | Captures distributional patterns and linear semantic relations |
| Understanding | Considers entire surrounding context | Focuses on local context within a fixed window |
| Context | Performs well on various NLP tasks, including semantic tasks | Focuses on local context within a fixed window |
| Task Performance | Provides fine-grained understanding of word meaning | Effective for tasks involving semantic associations and analogies |
| Granularity | Provides fine-grained understanding of word meaning | Captures broader semantic associations |
| Pre-training Time | Requires significant pre-training time | Faster pre-training compared to BERT |
| Fine-tuning | Allows fine-tuning on specific downstream tasks | Limited fine-tuning capabilities |

# Part 2: Text Classification Model

- Correct implementation of chosen classification model
- Comprehensive discussion on the choice of model, including its advantages and disadvantages

# The following steps were performed:

### Classification Algorithm:

The Multinomial Naive Bayes (MultinomialNB) algorithm was chosen as the classification model for the task.

### The decision to choose the Multinomial Naive Bayes (MultinomialNB) algorithm as the classification model for the task offers several benefits:

Simplicity: MultinomialNB is a straightforward and easy-to-understand algorithm. It is based on the Bayes' theorem and assumes independence between features, making it computationally efficient.

Efficiency: MultinomialNB performs well on large datasets with high-dimensional feature spaces. It is particularly suitable for text classification tasks where the number of features (e.g., word counts or TF-IDF values) can be significant.

Handling of Discrete Features: MultinomialNB is designed to handle discrete features, such as word frequencies or presence/absence indicators in text classification. It calculates the probability distribution over the classes given the feature values.

Robustness to Irrelevant Features: MultinomialNB can handle irrelevant features gracefully. It can still provide reasonably accurate predictions even when there are irrelevant or redundant features in the dataset.

Interpretability: MultinomialNB provides interpretable results by estimating the probability of each class given the input features. This can be useful in understanding the decision-making process and explaining the classification outcomes.

Low Training Time: MultinomialNB has a fast training time, especially compared to more complex algorithms such as deep learning models. This makes it suitable for situations where quick model training is desired.

Availability and Community Support: MultinomialNB is a widely used classification algorithm implemented in various machine learning libraries. It has extensive documentation and is well-supported by the machine learning community, ensuring access to resources, tutorials, and assistance when needed.

Overall, the Multinomial Naive Bayes algorithm offers simplicity, efficiency, robustness, and interpretability, making it a suitable choice for classification tasks, especially in text classification scenarios.

## Hyperparameter Tuning:

The best set of hyperparameters for the MultinomialNB model was identified using RandomizedSearchCV. This technique helps to automatically search for the optimal combination of hyperparameters by randomly sampling from the hyperparameter space.

## Evaluation Metrics:

To assess the performance of the classification model, several evaluation metrics were used. These included the Classification Report, which provides precision, recall, F1-score, and support for each class; the Confusion Matrix, which shows the counts of true positive, true negative, false positive, and false negative predictions; and the ROC Curve, which illustrates the trade-off between true positive rate and false positive rate at different classification thresholds.

By employing these steps, the classification algorithm was trained and evaluated, and the performance was assessed using various metrics.

```python
In [57]: from sklearn.feature_extraction.text import TfidfVectorizer
         from sklearn.naive_bayes import MultinomialNB
         from sklearn.pipeline import Pipeline
         from imblearn.under_sampling import RandomUnderSampler

         pipeline = Pipeline(
             [
                 ("vect", TfidfVectorizer()),
                 ("clf", MultinomialNB()),
             ]
         )
         pipeline
```

Out[57]:
```
   ▸      Pipeline

   ▸ TfidfVectorizer

     ▸ MultinomialNB
```

# Hyperparameters

```python
In [58]: import numpy as np

         parameter_grid = {
             "vect__max_df": (0.2, 0.4, 0.6, 0.8, 1.0),
             "vect__min_df": (1, 3, 5, 10),
             "vect__ngram_range": ((1, 1), (1, 2)),  # unigrams or bigrams
             "vect__norm": ("l1", "l2"),
             "clf__alpha": np.logspace(-6, 6, 13),
         }
```

```python
In [59]: from pprint import pprint

         from sklearn.model_selection import RandomizedSearchCV

         random_search = RandomizedSearchCV(
             estimator=pipeline,
             param_distributions=parameter_grid,
             n_iter=40,
             random_state=0,
             n_jobs=2,
             verbose=1
         )

         print("Performing grid search...")
         print("Hyperparameters to be evaluated:")
         pprint(parameter_grid)
```

```
Performing grid search...
Hyperparameters to be evaluated:
{'clf__alpha': array([1.e-06, 1.e-05, 1.e-04, 1.e-03, 1.e-02, 1.e-01, 1.e+00, 1.e+01,
       1.e+02, 1.e+03, 1.e+04, 1.e+05, 1.e+06]),
 'vect__max_df': (0.2, 0.4, 0.6, 0.8, 1.0),
 'vect__min_df': (1, 3, 5, 10),
 'vect__ngram_range': ((1, 1), (1, 2)),
 'vect__norm': ('l1', 'l2')}
```

```python
In [60]: from time import time

         t0 = time()
         X_train = random_search.fit(data_train.data, data_train.target)
         print(f"Done in {time() - t0:.3f}s")
```

```
Fitting 5 folds for each of 40 candidates, totalling 200 fits
Done in 65.155s
```

# Best Estimator

```python
In [61]: print("Best parameters combination found:")
         best_parameters = random_search.best_estimator_.get_params()
         for param_name in sorted(parameter_grid.keys()):
             print(f"{param_name}: {best_parameters[param_name]}")
```

```
Best parameters combination found:
clf__alpha: 0.001
vect__max_df: 0.8
vect__min_df: 3
vect__ngram_range: (1, 2)
vect__norm: l2
```

```python
In [62]: random_search.best_estimator_
```

Loading [MathJax]/extensions/Safe.js

Out[62]: 
```
▸   Pipeline
    ▸ TfidfVectorizer
        ▸ MultinomialNB
```

In [63]:
```python
preds = random_search.best_estimator_.predict(data_test.data)

y_train = data_train.target
y_test = data_test.target
```

In [64]:
```python
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import accuracy_score, cohen_kappa_score, f1_score, classification_report
```

In [65]:
```python
print(classification_report(data_test.target, preds))
```

```
              precision    recall  f1-score   support

           0       0.66      0.50      0.57       319
           1       0.88      0.85      0.86       389
           2       0.86      0.77      0.81       396
           3       0.80      0.83      0.81       394
           4       0.55      0.85      0.67       398
           5       0.50      0.31      0.38       251

    accuracy                           0.72      2147
   macro avg       0.71      0.68      0.68      2147
weighted avg       0.72      0.72      0.71      2147
```

In [66]:
```python
from sklearn import metrics



classification_report = metrics.classification_report(data_test.target,
                                                       preds,
                                                       target_names=data_test.target_names)
print(classification_report)

# Access precision and recall scores
scores = metrics.precision_recall_fscore_support(data_test.target, preds)
precision = scores[0]
recall = scores[1]


print("Precision:", precision)
print("Recall:", recall)
```

```
                        precision    recall  f1-score   support

           alt.atheism       0.66      0.50      0.57       319
         comp.graphics       0.88      0.85      0.86       389
               sci.med       0.86      0.77      0.81       396
             sci.space       0.80      0.83      0.81       394
soc.religion.christian       0.55      0.85      0.67       398
    talk.religion.misc       0.50      0.31      0.38       251

              accuracy                           0.72      2147
             macro avg       0.71      0.68      0.68      2147
          weighted avg       0.72      0.72      0.71      2147

Precision: [0.6557377  0.87798408 0.85875706 0.8009828  0.55482815 0.5       ]
Recall: [0.5015674  0.85089974 0.76767677 0.82741117 0.85175879 0.30677291]
```

In [67]:
```python
pred_ls = preds
test_ls = data_test.target

conf_arr = confusion_matrix(test_ls, pred_ls)

plt.figure(figsize=(8, 6), dpi=80, facecolor='w', edgecolor='k')

CLASSES = categories
ax = sns.heatmap(conf_arr, cmap="Blues_r", annot=True, fmt='d', xticklabels=CLASSES, yticklabels=CLASSES)

plt.title('newsgroup 20 ')
plt.xlabel('Prediction')
plt.ylabel('Actual')
plt.show(ax)
```

Loading [MathJax]/extensions/Safe.js

## newsgroup 20



```
In [68]: x_validation,y_validation =convert_to_np(data_test)
         x_train,y_train = convert_to_np(data_train)
```

```
In [69]: #ROC curve

         from sklearn.metrics import roc_curve
         from sklearn.preprocessing import label_binarize
         from sklearn.linear_model import LogisticRegression
         from sklearn.model_selection import train_test_split
         from sklearn.metrics import roc_curve, auc

         X_train = x_train
         X_test=x_validation
         y_test=y_validation
         # Split the dataset into training and testing sets

         # Convert the text data into TF-IDF features
         vectorizer = TfidfVectorizer()
         X_train = vectorizer.fit_transform(X_train)
         X_test = vectorizer.transform(X_test)

         #Train a MultinomialNB classifier
         classifier =  MultinomialNB()
         classifier.fit(X_train, y_train)

         # Predict probabilities for the test set
         y_scores =  classifier.predict_proba(X_test)

         # Binarize the true labels for ROC curve calculation
         y_test_bin = label_binarize(y_test, classes=np.unique(y_test))

         # Calculate the ROC curve and AUC for each class
         fpr = dict()
         tpr = dict()
         roc_auc = dict()
         n_classes = len(np.unique(y_test))
         for i in range(n_classes):
             fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], y_scores[:, i])
             roc_auc[i] = auc(fpr[i], tpr[i])

         # Plot the ROC curve for each class
         plt.figure(figsize=(10, 8))
         colors = ['blue', 'red', 'green', 'orange', 'purple']  # Add more colors if needed
         for i, color in zip(range(n_classes), colors):
             plt.plot(fpr[i], tpr[i], color=color, lw=2, label=categories[i]+' (AUC = {1:.2f})'.format(i, roc_auc[i]))

         plt.plot([0, 1], [0, 1], color='black', lw=2, linestyle='--')
         plt.xlim([0.0, 1.0])
         plt.ylim([0.0, 1.05])
         plt.xlabel('False Positive Rate')
         plt.ylabel('True Positive Rate')
         plt.title('Receiver Operating Characteristic (ROC)')
         plt.legend(loc='lower right')
         plt.show()
```

Loading [MathJax]/extensions/Safe.js

Receiver Operating Characteristic (ROC)

Legend:
- alt.atheism (AUC = 0.88)
- soc.religion.christian (AUC = 0.98)
- comp.graphics (AUC = 0.97)
- sci.med (AUC = 0.97)
- talk.religion.misc (AUC = 0.91)

In [70]:
```python
import plotly.express as px
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import TfidfVectorizer

data = fetch_20newsgroups(subset='train', categories=categories, shuffle=True, random_state=42)


vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(data.data)
y = data.target


from sklearn.decomposition import PCA

pca = PCA(n_components=2)
X_pca = pca.fit_transform(X.toarray())


fig = px.scatter(x=X_pca[:, 0], y=X_pca[:, 1], color=y, labels={'color': 'Category'})
fig.update_layout(title='Scatter Plot of 20 Newsgroups Dataset (PCA)',
                  xaxis_title='Principal Component 1',
                  yaxis_title='Principal Component 2')
fig.show()
```

C:\ProgramData\anaconda3\lib\site-packages\plotly\express\imshow_utils.py:24: DeprecationWarning: `np.bool8` is a deprecated alias for `np.bool_`.
(Deprecated NumPy 1.24)
  np.bool8: (False, True),

C:\ProgramData\anaconda3\lib\site-packages\plotly\io\_renderers.py:395: DeprecationWarning:

distutils Version classes are deprecated. Use packaging.version instead.

C:\ProgramData\anaconda3\lib\site-packages\plotly\io\_renderers.py:395: DeprecationWarning:

distutils Version classes are deprecated. Use packaging.version instead.

Loading [MathJax]/extensions/Safe.js

Scatter Plot of 20 Newsgroups Dataset (PCA)



## The choice of the model for generating word embeddings and performing multi-class prediction involves considering several factors, including their advantages and disadvantages.

### Word Embeddings:

Word embeddings can be generated using various techniques, such as TF-IDF. TF-IDF assigns weights to words based on their frequency in a document and across the entire corpus. The advantages of word embeddings include:Capturing semantic relationships: Word embeddings can capture the semantic meaning and relationships between words, allowing for better understanding and representation of textual data.

Dimensionality reduction: Word embeddings reduce the dimensionality of the data, making it more manageable for machine learning algorithms. Transferable knowledge: Pre-trained word embeddings can be used as a starting point for new text analysis tasks, leveraging knowledge from large pre-existing corpora.

### Naive Bayes Classifier:

Naive Bayes is a popular choice for multi-class prediction tasks. It is based on Bayes' theorem and assumes independence between features. Some advantages of Naive Bayes Classifier include:

Simple and fast: Naive Bayes is computationally efficient and scales well to large datasets, making it suitable for real-time or high-volume applications. Effective with high-dimensional data: Naive Bayes performs well even with a high number of features, such as word embeddings, and can handle sparse data efficiently. Interpretable results: Naive Bayes provides interpretable results by estimating the conditional probabilities of each class given the input features. However, it is important to note that there are also limitations and considerations to keep in mind:

Word embeddings may not capture all nuances of language: While word embeddings are powerful, they may not fully capture all nuances of language and context, leading to some loss of information. Naive Bayes assumes feature independence: The independence assumption in Naive Bayes may not hold true in some cases, which can impact its performance.

Model selection should consider the specific task and data: The choice of model should be based on the specific task requirements, data characteristics, and available resources. Overall, the use of word embeddings and Naive Bayes Classifier can provide effective results in multi-class prediction tasks, allowing for the exploration of semantic relationships and probabilistic predictions across multiple classes. However, it is crucial to carefully evaluate the suitability of these models based on the specific context and objectives of the analysis.

## Part 3: Text Clustering Model

- Correct implementation of chosen clustering model
- Comprehensive discussion on the choice of model, including its pros and cons

**We selected the KNeighborsClassifier as the clustering algorithm for our analysis. Additionally, we incorporated the K-fold Cross Validation technique to evaluate the performance of our model.**

In [71]:
```python
from sklearn.neighbors import KNeighborsClassifier
import plotly.express as px
```

In [72]:
```python
from __future__ import print_function

from sklearn import __version__ as sklearn_version
print('Sklearn version:', sklearn_version)
from sklearn import model_selection
from imblearn.under_sampling import RandomUnderSampler
```

Sklearn version: 1.3.0

In [73]:
```python
from __future__ import print_function

from comet_ml import Experiment

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.pipeline import Pipeline
```

Loading [MathJax]/extensions/Safe.js

```python
from sklearn.datasets import fetch_20newsgroups
from sklearn.linear_model import SGDClassifier
from sklearn.metrics import accuracy_score
from sklearn.model_selection import KFold

import numpy as np

kf = KFold(n_splits=10)
curr_fold = 0
acc_dict = {}

for train_idx, test_idx in kf.split(x_train):
    text_clf = Pipeline([('vect', CountVectorizer()),  # Counts occurrences of each word
                         ('tfidf', TfidfTransformer()),  # Normalize the counts based on document length
                         ('clf', KNeighborsClassifier(n_neighbors=3, weights='uniform')),
                         ])


    text_clf.fit(x_train[train_idx].tolist(), y_train[train_idx])

    # Predict unseen test data based on fitted classifer
    predicted = text_clf.predict(x_train[test_idx])

    # Compute accuracy
    acc = accuracy_score(y_train[test_idx].tolist(), predicted)
    acc_dict[curr_fold]=acc
    print("accuracy_fold_%s" % curr_fold, acc)

    curr_fold += 1
```

C:\ProgramData\anaconda3\lib\site-packages\requests_toolbelt\_compat.py:56: DeprecationWarning:

'urllib3.contrib.pyopenssl' module is deprecated and will be removed in a future release of urllib3 2.x. Read more in this issue: https://github.com/urllib3/urllib3/issues/2680

C:\ProgramData\anaconda3\lib\site-packages\comet_ml\monkey_patching.py:19: DeprecationWarning:

the imp module is deprecated in favour of importlib and slated for removal in Python 3.12; see the module's documentation for alternative uses

```
accuracy_fold_0 0.25386996904024767
accuracy_fold_1 0.21362229102167182
accuracy_fold_2 0.21981424148606812
accuracy_fold_3 0.19814241486068113
accuracy_fold_4 0.19504643962848298
accuracy_fold_5 0.20743034055727555
accuracy_fold_6 0.22910216718266255
accuracy_fold_7 0.2236024844720497
accuracy_fold_8 0.27639751552795033
accuracy_fold_9 0.2701863354037267
```

In [74]:
```python
from sklearn.decomposition import LatentDirichletAllocation as LDA

pipe_knn = Pipeline([('vect', TfidfVectorizer(stop_words='english')),
                     ('lda', LDA(n_components=6, max_iter=25,
                                 learning_method='online',
                                 learning_offset=200.,
                                 random_state=0)),
                     ('clf',  KNeighborsClassifier()),
                     ])
```

In [75]:
```python
from pprint import pprint

from sklearn.model_selection import RandomizedSearchCV

param_range=[1,2]
knn_param_grid=[{
                'clf__n_neighbors': param_range,
                'clf__weights': ['uniform', 'distance'],
                'clf__metric': ['euclidean', 'manhattan']
                }]

random_search = RandomizedSearchCV(
    pipe_knn,
    param_distributions=knn_param_grid,
    n_iter=40,
    random_state=0,
    n_jobs=2,
    verbose=1,
)

print("Performing grid search...")
print("Hyperparameters to be evaluated:")
pprint(knn_param_grid)
```

```
Performing grid search...
Hyperparameters to be evaluated:
[{'clf__metric': ['euclidean', 'manhattan'],
  'clf__n_neighbors': [1, 2],
  'clf__weights': ['uniform', 'distance']}]
```

In [76]:
```python
from time import time

t0 = time()
X_train = random_search.fit(data_test.data, data_test.target)
print(f"Done in {time() - t0:.3f}s")
```

```
Fitting 5 folds for each of 8 candidates, totalling 40 fits
Done in 98.245s
```

In [77]:
```python
# Predicting our test data
predicted = random_search.best_estimator_.predict(data_test.data)
print('We got an accuracy of',np.mean(predicted == data_test.target)*100, '% over the test data.')
```

Loading [MathJax]/extensions/Safe.js  acy of 97.48486259897533 % over the test data.

Accuracy of the model is 97% for the test data...

BERTopic follows a series of steps to process documents and perform clustering:

==Step 1==: Preprocessing BERTopic preprocesses the input documents by performing text cleaning, tokenization, and encoding. This step ensures that the text is in a suitable format for further analysis.

==Step 2==: BERT Embeddings BERTopic utilizes a pre-trained BERT model to obtain contextualized embeddings for each document. These embeddings capture the semantic meaning of the text by considering the surrounding context.

==Step 3==: Dimensionality Reduction To handle high-dimensional embeddings, BERTopic applies the UMAP (Uniform Manifold Approximation and Projection) algorithm. UMAP reduces the dimensionality of the embeddings while preserving the local structure of the data. This step helps in visualizing and interpreting the document representations.

==Step 4==: Clustering BERTopic employs the HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise) algorithm to cluster the reduced embeddings. HDBSCAN is capable of identifying clusters of varying densities, allowing for more flexible and accurate clustering.

==Step 5==: Topic Assignment BERTopic assigns a topic label to each document based on the cluster it belongs to. This assignment is determined by the proximity of the document's embedding to the centroid of each topic cluster.

Probability Calculation The probabilities of topic assignments are calculated based on the distances between the document embeddings and the centroid of each topic cluster. These probabilities are then normalized to ensure they sum up to 1 across all topics, providing a measure of confidence in the topic assignments.

By following these steps, BERTopic is able to effectively process documents, generate meaningful embeddings, and perform clustering to identify topics within the data.

```
In [1]: pip install --upgrade numba
```

```
Requirement already satisfied: numba in c:\users\purus\appdata\roaming\python\python310\site-packages (0.57.1)
Requirement already satisfied: llvmlite<0.41,>=0.40.0dev0 in c:\users\purus\appdata\roaming\python\python310\site-packages (from numba) (0.40.1)
Requirement already satisfied: numpy<1.25,>=1.21 in c:\users\purus\appdata\roaming\python\python310\site-packages (from numba) (1.24.0)
Note: you may need to restart the kernel to use updated packages.
```

```python
In [26]: from bertopic import BERTopic

topic_model = BERTopic(language="english", calculate_probabilities=True, verbose=True)

topics, probs = topic_model.fit_transform(data_train.data)
```

```
Batches:   0%|          | 0/101 [00:00<?, ?it/s]
2023-07-10 03:51:21,606 - BERTopic - Transformed documents to Embeddings
2023-07-10 03:51:37,991 - BERTopic - Reduced dimensionality
2023-07-10 03:51:38,439 - BERTopic - Clustered reduced embeddings
```

```python
In [27]: freq = topic_model.get_topic_info()
freq
```

| | Topic | Count | Name | Representation | Representative_Docs |
|---|---|---|---|---|---|
| 0 | -1 | 1000 | -1_the_of_to_and | [the, of, to, and, in, is, that, for, you, it] | [[text deleted]\n\nI wish that you had followe... |
| 1 | 0 | 493 | 0_the_of_and_in | [the, of, and, in, to, it, is, that, my, with] | [As promised, below is a personal critique of ... |
| 2 | 1 | 370 | 1_image_jpeg_for_and | [image, jpeg, for, and, to, file, you, the, it... | [Archive-name: typing-injury-faq/software\nVer... |
| 3 | 2 | 85 | 2_hello___ | [hello, , , , , , , , , ] | [\n, Hello,, Hello,] |
| 4 | 3 | 70 | 3_you_that_your_jim | [you, that, your, jim, context, to, out, it, o... | [New in this version: challenge #5, plus an a... |
| 5 | 4 | 69 | 4_god_atheists_is_that | [god, atheists, is, that, atheism, of, to, the... | [Archive-name: atheism/logic\nAlt-atheism-arch... |
| 6 | 5 | 59 | 5_the_judas_of_bible | [the, judas, of, bible, and, greek, in, book, ... | [I produced an error last week about CHORION:\... |
| 7 | 6 | 51 | 6_islam_islamic_of_the | [islam, islamic, of, the, muslims, quran, muha... | [I apologize for the long delay in getting a r... |
| 8 | 7 | 50 | 7_moon_lunar_prize_the | [moon, lunar, prize, the, to, would, billion, ... | [\nWishful thinking mostly. It's more likely t... |
| 9 | 8 | 50 | 8_koresh_they_the_fbi | [koresh, they, the, fbi, was, that, to, he, no... | [\nTrue. At first, the news media seemed entr... |
| 10 | 9 | 44 | 9_hell_eternal_that_he | [hell, eternal, that, he, god, to, heaven, the... | [\n\nAnd yet, Jayne, as we read the Gospels a... |
| 11 | 10 | 44 | 10_science_of_the_is | [science, of, the, is, that, one, to, as, in, ... | [Avoiding mistakes is certainly highly desirab... |
| 12 | 11 | 43 | 11_polygon_edge_xxxx_algorithm | [polygon, edge, xxxx, algorithm, points, line,... | [About a year ago I started work on a problem ... |
| 13 | 12 | 40 | 12_tobacco_health_smokeless_coli | [tobacco, health, smokeless, coli, o157h7, amo... | [The following is a survey we are conducting f... |
| 14 | 13 | 35 | 13_probe_spacecraft_satellite_orbit | [probe, spacecraft, satellite, orbit, the, ven... | [Forwarded from Neal Ausman, Galileo Mission D... |
| 15 | 14 | 32 | 14_homosexuality_gay_to_people | [homosexuality, gay, to, people, that, homosex... | [Tony-\n\nI read your post, it was nothing new... |
| 16 | 15 | 29 | 15_data_ftp_available_for | [data, ftp, available, for, images, image, and... | [------------------------------------\n\t+ ..... |
| 17 | 16 | 27 | 16_station_redesign_ssf_space | [station, redesign, ssf, space, the, billion, ... | [In the April edition of "One Small Step for a... |
| 18 | 17 | 27 | 17_ra_satan_god_the | [ra, satan, god, the, of, lucifer, do, that, a... | [ \n \nDefine perfect then.\n \n \n \nTake you... |
| 19 | 18 | 27 | 18_penalty_system_punishment_cruel | [penalty, system, punishment, cruel, is, murde... | [My turn to jump in! :)\n\n\nI think you mean ... |
| 20 | 19 | 26 | 19_god_you_we_to | [god, you, we, to, and, that, our, christ, the... | [\n\n\nAre you your own master? Do you have a... |
| 21 | 20 | 26 | 20_space_astronaut_and_nasa | [space, astronaut, and, nasa, candidates, aero... | [Sorry for asking a question that's not entire... |
| 22 | 21 | 26 | 21_tomb_he_the_jesus | [tomb, he, the, jesus, resurrection, was, that... | [[much of the excellent post deleted for space... |
| 23 | 22 | 26 | 22_truth_that_absolute_bible | [truth, that, absolute, bible, is, are, you, t... | [Dean Velasco quoted a letter from James M Sto... |
| 24 | 23 | 25 | 23_moral_animals_morality_is | [moral, animals, morality, is, you, do, omnisc... | [#In <1qvabj$g1j@horus.ap.mchp.sni.de> frank@D... |
| 25 | 24 | 24 | 24_war_in_the_bosnia | [war, in, the, bosnia, serbs, that, religion, ... | [\n\nThe Bible does tell us that governments a... |
| 26 | 25 | 23 | 25_den_sphere_points_radius | [den, sphere, points, radius, plane, ellipse, ... | [\n\n\n\n\n\n\n\n\nGood I had a bad feeling ... |
| 27 | 26 | 23 | 26_advertising_space_billboard_inflatable | [advertising, space, billboard, inflatable, wo... | [From the article "What's New" Apr-16-93 in sc... |
| 28 | 27 | 22 | 27_marriage_married_ceremony_couple | [marriage, married, ceremony, couple, to, comm... | [\n\tI originally wrote to the person who aske... |
| 29 | 28 | 22 | 28_pope_church_the_schism | [pope, church, the, schism, catholic, of, litu... | [Here is some material by Michael Davies on th... |
| 30 | 29 | 22 | 29_objective_morality_morals_compromise | [objective, morality, morals, compromise, mora... | [\nI'll take a wild guess and say Freedom is o... |
| 31 | 30 | 20 | 30_jews_antisemitism_jewish_casual | [jews, antisemitism, jewish, casual, was, they... | [\nI think the problem here is that I pretty m... |
| 32 | 31 | 19 | 31_de_van_het_een | [de, van, het, een, en, te, utrecht, op, orbit... | [Hiya \n\nI'm a VERY amuture astronomer in Ade... |
| 33 | 32 | 19 | 32_why_is_exist_existence | [why, is, exist, existence, that, not, to, it,... | [This kind of argument cries for a comment...\... |
| 34 | 33 | 18 | 33_order_oto_reuss_amorc | [order, oto, reuss, amorc, rosicrucian, ordo, ... | [930420\n\nDo what thou wilt shall be the whol... |
| 35 | 34 | 18 | 34_group_split_newsgroup_graphics | [group, split, newsgroup, graphics, aspects, g... | [Concerning the proposed newsgroup split, I pe... |
| 36 | 35 | 17 | 35_mary_her_she_maria | [mary, her, she, maria, was, the, sin, he, his... | [Biblical basis for the Immaculate Conception:... |
| 37 | 36 | 16 | 36_42_question_answer_alice | [42, question, answer, alice, discovered, ever... | [: Well,\n: \n: 42 is 101010 binary, and who w... |
| 38 | 37 | 15 | 37_lds_mormons_mormon_the | [lds, mormons, mormon, the, to, church, and, s... | [:\n (lots of stuff about the Nicene Creed del... |
| 39 | 38 | 15 | 38_software_process_level_shuttle | [software, process, level, shuttle, maturity, ... | [\n My understanding is that the 'expected e... |
| 40 | 39 | 14 | 39_he_ungodly_we_him | [he, ungodly, we, him, whatever, his, will, in... | [The parable of the Prodigal Son is not about ... |
| 41 | 40 | 14 | 40_photography_kirlian_pictures_object | [photography, kirlian, pictures, object, field... | [I think that's the correct spelling..\n\tI am... |
| 42 | 41 | 14 | 41_easter_resurrection_pagan_goddess | [easter, resurrection, pagan, goddess, celebra... | [for SRC\n\nIn most languages, the Feast of th... |
| 43 | 42 | 13 | 42_sabbath_ceremonial_day_law | [sabbath, ceremonial, day, law, saturday, wors... | [[In response to some of the discussions on th... |
| 44 | 43 | 13 | 43_revelation_scripture_prophecy_god | [revelation, scripture, prophecy, god, the, he... | [\n\nThis is one of the differences between OT... |
| 45 | 44 | 12 | 44_law_jesus_the_not | [law, jesus, the, not, is, that, paul, god, to... | [\nOK, here's at least one Christian's answer:... |
| 46 | 45 | 12 | 45_centaur_proton_stage_tanks | [centaur, proton, stage, tanks, payload, mile,... | [Reading from a Amoco Performance Products dat... |
| 47 | 46 | 12 | 46_mining_miners_right_basically | [mining, miners, right, basically, space, go, ... | [===\nI aint talking the large or even the "mi... |
| 48 | 47 | 12 | 47_cancer_center_centers_research | [cancer, center, centers, research, medical, u... | [\nThat's ridiculous!\n\n\nThey aren't designe... |
| 49 | 48 | 11 | 48_cview_temp_file_files | [cview, temp, file, files, disk, floppy, direc... | [: >over where it places its temp files: it ju... |
| 50 | 49 | 11 | 49_constant_mass_km_velocity | [constant, mass, km, velocity, radius, orbit, ... | [I have the "osculating elements at perigee" o... |
| 51 | 50 | 11 | 50_church_churches_there_that | [church, churches, there, that, people, to, is... | [Here are some notes about what the church is ... |
| 52 | 51 | 11 | 51_motto_things_state_worse | [motto, things, state, worse, think, it, anthe... | [\n\nIn this era of AIDS, isn't someone's fuck... |

```
In [28]: freq['Name'].head()
```

```
0         -1_the_of_to_and
1         0_the_of_and_in
2     1_image_jpeg_for_and
3                2_hello___
4     3_you_that_your_jim
Name: Name, dtype: object
```

```
In [29]: freq['Representation'].head()
```

Loading [MathJax]/extensions/Safe.js

```
Out[29]:  0          [the, of, to, and, in, is, that, for, you, it]
          1          [the, of, and, in, to, it, is, that, my, with]
          2          [image, jpeg, for, and, to, file, you, the, it...
          3                                 [hello, , , , , , , , , ]
          4       [you, that, your, jim, context, to, out, it, o...
          Name: Representation, dtype: object
```

```
In [30]:  freq['Representative_Docs'].head()
```

```
Out[30]:  0       [[text deleted]\n\nI wish that you had followe...
          1       [As promised, below is a personal critique of ...
          2       [Archive-name: typing-injury-faq/software\nVer...
          3                                 [\n, Hello,, Hello,]
          4       [New in this version:  challenge #5, plus an a...
          Name: Representative_Docs, dtype: object
```

```
In [31]:  topic_model.get_topic(0)   # Selecting the most frequent topic
```

```
Out[31]:  [('the', 0.014611941381660297),
           ('of', 0.014491629430234026),
           ('and', 0.014130160488989085),
           ('in', 0.01395276868505395),
           ('to', 0.013822016719459363),
           ('it', 0.013376969465908337),
           ('is', 0.013252788336954275),
           ('that', 0.011239821494151961),
           ('my', 0.010714789994410856),
           ('with', 0.010551194063474238)]
```

```
In [32]:  topic_model.visualize_topics()
```



**Intertopic Distance Map**

```
In [33]:  topic_model.visualize_distribution(probs[1], min_probability=0.001)
```

## Topic Probability Distribution

Topic 51: motto_things_state_wors...
Topic 49: constant_mass_km_veloci...
Topic 47: cancer_center_centers_r...
Topic 45: centaur_proton_stage_ta...
Topic 43: revelation_scripture_pr...
Topic 41: easter_resurrection_pag...
Topic 39: he_ungodly_we_him_whate...
Topic 37: lds_mormons_mormon_the_...
Topic 35: mary_her_she_maria_was
Topic 33: order_oto_reuss_amorc_r...
Topic 31: de_van_het_een_en
Topic 29: objective_morality_mora...
Topic 27: marriage_married_ceremo...
Topic 25: den_sphere_points_radiu...
Topic 23: moral_animals_morality_...
Topic 21: tomb_he_the_jesus_resur...
Topic 19: god_you_we_to_and
Topic 17: ra_satan_god_the_of
Topic 15: data_ftp_available_for_...
Topic 13: probe_spacecraft_satell...
Topic 11: polygon_edge_xxxx_algor...
Topic 9: hell_eternal_that_he_god
Topic 7: moon_lunar_prize_the_to
Topic 5: the_judas_of_bible_and
Topic 3: you_that_your_jim_contex...
Topic 0: the_of_and_in_to

Probability (x-axis: 0, 0.05, 0.1, 0.15)

```
In [34]: topic_model.visualize_hierarchy(top_n_topics=50)
```

## Hierarchical Clustering



33_order_oto_reuss
39_he_ungodly_we
35_mary_her_she
17_ra_satan_god
22_truth_that_absolute
4_god_atheists_is
32_why_is_exist
19_god_you_we
9_hell_eternal_that
43_revelation_scripture_pro...
6_islam_islamic_of
5_the_judas_of
14_homosexuality_gay_to
0_the_of_and
10_science_of_the
3_you_that_your
24_war_in_the
8_koresh_they_the
37_lds_mormons_mormon
21_tomb_he_the
28_pope_church_the
27_marriage_married_ceremony
46_mining_miners_right
30_jews_antisemitism_jewish
42_sabbath_ceremonial_day
44_law_jesus_the
41_easter_resurrection_pagan
18_penalty_system_punishment
23_moral_animals_morality
29_objective_morality_morals
31_de_van_het
49_constant_mass_km
45_centaur_proton_stage
26_advertising_space_billbo...
7_moon_lunar_prize
16_station_redesign_ssf
13_probe_spacecraft_satellite
20_space_astronaut_and
47_cancer_center_centers
12_tobacco_health_smokeless
48_cview_temp_file
15_data_ftp_available
1_image_jpeg_for
11_polygon_edge_xxxx
25_den_sphere_points
38_software_process_level
34_group_split_newsgroup
36_42_question_answer
40_photography_kirlian_pict...
2_hello__

(x-axis: 0, 0.2, 0.4, 0.6, 0.8, 1, 1.2, 1.4, 1.6)

```
In [35]: topic_model.visualize_barchart(top_n_topics=8)
```

Loading [MathJax]/extensions/Safe.js

## Topic Word Scores

**Topic 0**
- the
- of
- and
- in
- to

(x-axis: 0, 0.005, 0.01, 0.015)

**Topic 1**
- image
- jpeg
- for
- and
- to

(x-axis: 0, 0.005, 0.01, 0.015, 0.02)

**Topic 2**
- hello

(x-axis: 0, 2, 4)

**Topic 3**
- you
- that
- your
- jim
- context

(x-axis: 0, 0.01, 0.02)

**Topic 4**
- god
- atheists
- is
- that
- atheism

(x-axis: 0, 0.005, 0.01, 0.015, 0.02)

**Topic 5**
- the
- judas
- of
- bible
- and

(x-axis: 0, 0.005, 0.01, 0.015, 0.02)

**Topic 6**
- islam
- islamic
- of
- the
- muslims

(x-axis: 0, 0.01, 0.02, 0.03)

**Topic 7**
- moon
- lunar
- prize
- the
- to

(x-axis: 0, 0.01, 0.02, 0.03)

```
In [36]:  topic_model.visualize_heatmap(n_clusters=20, width=1000, height=1000)
```

## Similarity Matrix



```
In [37]:  topic_model.visualize_term_rank()
```

Loading [MathJax]/extensions/Safe.js

## Term score decline per Topic

In [38]: topic_model.reduce_topics(data_train.data, nr_topics=60)

2023-07-10 03:51:41,688 - BERTopic - Reduced number of topics from 53 to 53

Out[38]: <bertopic._bertopic.BERTopic at 0x1e745ec1270>

In [39]: topic_model.visualize_topics()

## Intertopic Distance Map



In [40]: topic_model.visualize_heatmap(n_clusters=3, width=1000, height=1000)

Loading [MathJax]/extensions/Safe.js

**Similarity Matrix**

```
topic_model.visualize_hierarchy(top_n_topics=50)
```

Loading [MathJax]/extensions/Safe.js

## Hierarchical Clustering

```python
from __future__ import print_function
import pyLDAvis
import pyLDAvis.sklearn
pyLDAvis.enable_notebook()
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from sklearn.decomposition import LatentDirichletAllocation
```

**BERTopic is a topic modeling technique that automatically summarizes a large corpus of texts.**

## Imbalance in the data set

When introducing imbalance in a dataset, it means creating a scenario where there is an unequal distribution of samples across different categories or classes. This imbalance can occur when certain categories have significantly fewer instances compared to others.

In the context of clustering, introducing imbalance can impact the clustering results. Clustering algorithms aim to group similar data points together based on their features or characteristics. However, when there is an imbalance in the dataset, the clustering algorithm may be biased towards the majority class or category. This means that the clusters formed may be dominated by the majority class, while the minority classes may be underrepresented or overlooked.

The impact of imbalance on clustering can lead to several outcomes. First, the minority classes may be merged with the majority class in the clustering process, making it difficult to distinguish between different minority groups. Second, the clusters formed may not accurately represent the underlying structure of the data, as the algorithm may prioritize capturing the patterns and characteristics of the majority class.

To address the impact of imbalance on clustering, various techniques can be employed. These include oversampling the minority class to balance the dataset, applying clustering algorithms that are specifically designed to handle imbalanced data, or using evaluation metrics that consider the class imbalance, such as Adjusted Rand Index (ARI) or Fowlkes-Mallows Index (FMI). By addressing the imbalance, it is possible to obtain more accurate and meaningful clusters that reflect the true distribution of the data.

```python
newsgroups_train = fetch_20newsgroups(subset="train", categories=categories)
newsgroups_test = fetch_20newsgroups(subset="test", categories=categories)

X_train = data_train.data
X_test = data_test.data

y_train = data_train.target
y_test = data_test.target
```

```python
from sklearn.naive_bayes import MultinomialNB
```

Loading [MathJax]/extensions/Safe.js

```
In [46]:  from imblearn.pipeline import make_pipeline as make_pipeline_imb
          from imblearn.under_sampling import RandomUnderSampler

          model = make_pipeline_imb(TfidfVectorizer(), RandomUnderSampler(), MultinomialNB())

          model.fit(X_train, y_train)
          y_pred = model.predict(X_test)
```

```
In [47]:  from imblearn.metrics import classification_report_imbalanced

          # Assuming you have defined y_test and y_pred
          report = classification_report_imbalanced(y_test, y_pred)
          print(report)
```

```
                   pre       rec       spe        f1       geo       iba       sup

            0      0.47      0.28      0.95      0.35      0.51      0.25       319
            1      0.90      0.77      0.98      0.83      0.87      0.74       389
            2      0.95      0.58      0.99      0.72      0.76      0.55       396
            3      0.93      0.57      0.99      0.70      0.75      0.54       394
            4      0.35      0.94      0.60      0.51      0.75      0.58       398
            5      0.62      0.18      0.99      0.27      0.42      0.16       251

    avg / total    0.72      0.59      0.91      0.59      0.70      0.50      2147
```

## Dealing with abbreviations and misspelled words

It is an important aspect of text preprocessing. It involves identifying and resolving abbreviated forms of words and correcting misspelled words to improve the accuracy and consistency of the text analysis process.

Abbreviations can pose a challenge in text analysis because they can vary widely and may not be recognizable by standard language models or algorithms. Therefore, it is crucial to develop strategies for handling abbreviations in order to ensure accurate analysis results. One common approach is to create a dictionary or lookup table that maps abbreviations to their expanded forms. This allows the algorithm to replace the abbreviated form with the corresponding full word, enabling better understanding and interpretation of the text.

Misspelled words are another issue that needs to be addressed in text preprocessing. Misspellings can occur due to typographical errors, lack of knowledge or attention, or variations in spelling conventions. Correcting misspelled words involves employing techniques such as spell checking, which compares the input word against a dictionary of correctly spelled words and suggests possible corrections.

In addition to spell checking, advanced methods like fuzzy matching or phonetic algorithms can be utilized to handle misspellings. Fuzzy matching algorithms consider the similarity between words based on their character patterns and suggest alternative corrections. Phonetics algorithms focus on the pronunciation of words and can assist in identifying potential correct spellings based on their phonetic similarity.

Overall, dealing with abbreviations and misspelled words requires a combination of techniques such as creating abbreviation dictionaries, employing spell checking algorithms, and utilizing fuzzy matching or phonetic algorithms to ensure accurate and consistent text analysis results.

The dataframe contains instances of misspelled words and abbreviations that will be corrected using a spellchecker.

```
In [48]:  import pandas as pd
          from spellchecker import SpellChecker
```

```
In [49]:  df = pd.DataFrame(['swtch', 'cola', 'FBI', 'smsng', 'BCA', 'MIB'], columns=['misspelled'])
          abbreviations = {
              'FBI': 'Federal Bureau of Investigation',
              'BCA': 'Bank Central Asia',
              'MIB': 'Men In Black',
              'cola': 'Coca Cola'
          }

          spell = SpellChecker()
          df['fixed'] = df['misspelled'].apply(spell.correction).replace(abbreviations)
          print("Misspelled and short form been corrected and abbreviated")
          print(df)
```

```
Misspelled and short form been corrected and abbreviated
  misspelled                            fixed
0      swtch                           switch
1       cola                        Coca Cola
2        FBI  Federal Bureau of Investigation
3      smsng                             sing
4        BCA                              bra
5        MIB                              mix
```

## Word Sense Disambiguation

In natural language processing, the phenomenon of words having multiple meanings in different contexts is known as "polysemy." This polysemy can introduce ambiguity and hinder accurate text analysis. To address this challenge, we employ a technique called "word sense disambiguation."

Word sense disambiguation involves determining the intended meaning or sense of a word within a particular context. By disambiguating words based on their intended senses, we can treat them as distinct entities during analysis, enabling more accurate and meaningful interpretation of the text.

The goal of word sense disambiguation is to assign the most appropriate sense or meaning to each occurrence of a word, considering the surrounding context and the specific task or application at hand. This process may involve leveraging linguistic resources such as dictionaries, semantic networks, or corpus-based methods that analyze word usage patterns in large text collections.

By disambiguating words and resolving their multiple meanings, we can enhance the precision and relevance of text analysis tasks such as information retrieval, machine translation, sentiment analysis, and many others. Word sense disambiguation is a fundamental step in understanding and accurately interpreting natural language text in various contexts.

```
In [50]:  import nltk
          nltk.download('wordnet')
          Loading [MathJax]/extensions/Safe.js  omw-1.4')
```

```python
from nltk.wsd import lesk
from nltk import word_tokenize
```

We took 'Bank' as example to diffentiate sentence1: "Keep your savings in the bank" sentence2: "It's so risky to drive over the banks of the road"

In [51]:
```python
sentence1 = "Keep your savings in the bank"
sentence2 = "It's so risky to drive over the banks of the road"
```

Sentence 1

In [52]:
```python
def get_synset(sentence, word):
    return lesk(word_tokenize(sentence), word)
get_synset(sentence1,'bank')
```

Out[52]:
```
Synset('savings_bank.n.02')
```

Here, savings_bank.n.02 refers to a container for keeping money safely at home

In [53]:
```python
get_synset(sentence2,'bank')
```

Out[53]:
```
Synset('bank.v.07')
```

By utilizing the Lesk algorithm, we successfully determined the specific meaning of a word within its given context. For example, in this case, the term "bank.v.07" specifically denotes a slope or incline in the curve of a road. This algorithm enables us to disambiguate words and accurately discern their intended sense based on the surrounding context.

# Dealing with Slang

To address the challenge of SMS slangs, which involve abbreviations and informal language commonly used in text messages, we adopt a method of transforming these slangs into their corresponding full words. Manually replacing each abbreviation is impractical due to the large number of slangs present. Instead, we utilize a combination of pre-built dictionaries and online resources to convert the input text into a more elaborate and comprehensible format.

By leveraging existing dictionaries and online word resources, we can automatically map SMS slangs to their expanded forms. For example, abbreviations like "bt" can be replaced with "but," "hv" can be transformed into "have," and "nt" can be converted to "not." This approach allows us to effectively handle a wide range of slangs without the need for manual intervention.

The use of dictionaries and online word databases enables us to process text inputs efficiently and accurately. These resources contain a vast collection of words and their corresponding meanings, allowing us to convert SMS slangs into their appropriate interpretations. By making our input text more elaborate and closer to standard language, we enhance the accuracy and clarity of subsequent text analysis tasks.

## Create Dictionary words

In [54]:
```python
CONTRACTION_MAP = {
"ain't": "is not",
"aren't": "are not",
"can't": "cannot",
"can't've": "cannot have",
"'cause": "because",
"could've": "could have",
"couldn't": "could not",
"couldn't've": "could not have",
"didn't": "did not",
"doesn't": "does not",
"don't": "do not",
"hadn't": "had not",
"hadn't've": "had not have",
"hasn't": "has not",
"haven't": "have not",
"he'd": "he would",
"he'd've": "he would have",
"he'll": "he will",
"he'll've": "he he will have",
"he's": "he is",
"how'd": "how did",
"how'd'y": "how do you",
"how'll": "how will",
"how's": "how is",
"I'd": "I would",
"I'd've": "I would have",
"I'll": "I will",
"I'll've": "I will have",
"I'm": "I am",
"I've": "I have",
"i'd": "i would",
"i'd've": "i would have",
"i'll": "i will",
"i'll've": "i will have",
"i'm": "i am",
"i've": "i have",
"isn't": "is not",
"it'd": "it would",
"it'd've": "it would have",
"it'll": "it will",
"it'll've": "it will have",
"it's": "it is",
"let's": "let us",
"ma'am": "madam",
"mayn't": "may not",
"might've": "might have",
"mightn't": "might not",
```

```
    "mightn't've": "might not have",
    "must've": "must have",
    "mustn't": "must not",
    "mustn't've": "must not have",
    "needn't": "need not",
    "needn't've": "need not have",
    "o'clock": "of the clock",
    "oughtn't": "ought not",
    "oughtn't've": "ought not have",
    "shan't": "shall not",
    "sha'n't": "shall not",
    "shan't've": "shall not have",
    "she'd": "she would",
    "she'd've": "she would have",
    "she'll": "she will",
    "she'll've": "she will have",
    "she's": "she is",
    "should've": "should have",
    "shouldn't": "should not",
    "shouldn't've": "should not have",
    "so've": "so have",
    "so's": "so as",
    "that'd": "that would",
    "that'd've": "that would have",
    "that's": "that is",
    "there'd": "there would",
    "there'd've": "there would have",
    "there's": "there is",
    "they'd": "they would",
    "they'd've": "they would have",
    "they'll": "they will",
    "they'll've": "they will have",
    "they're": "they are",
    "they've": "they have",
    "to've": "to have",
    "wasn't": "was not",
    "we'd": "we would",
    "we'd've": "we would have",
    "we'll": "we will",
    "we'll've": "we will have",
    "we're": "we are",
    "we've": "we have",
    "weren't": "were not",
    "what'll": "what will",
    "what'll've": "what will have",
    "what're": "what are",
    "what's": "what is",
    "what've": "what have",
    "when's": "when is",
    "when've": "when have",
    "where'd": "where did",
    "where's": "where is",
    "where've": "where have",
    "who'll": "who will",
    "who'll've": "who will have",
    "who's": "who is",
    "who've": "who have",
    "why's": "why is",
    "why've": "why have",
    "will've": "will have",
    "won't": "will not",
    "won't've": "will not have",
    "would've": "would have",
    "wouldn't": "would not",
    "wouldn't've": "would not have",
    "y'all": "you all",
    "y'all'd": "you all would",
    "y'all'd've": "you all would have",
    "y'all're": "you all are",
    "y'all've": "you all have",
    "you'd": "you would",
    "you'd've": "you would have",
    "you'll": "you will",
    "you'll've": "you will have",
    "you're": "you are",
    "you've": "you have",


}
```

To extract keywords from a web page, we can create a request to retrieve the HTML content of the page. In this case, we will be extracting slangs starting with different alphabets from the website "https://www.noslang.com/dictionary/". Each alphabet has a specific URL pattern associated with it.

To retrieve slangs starting with 'A', we can use the URL "https://www.noslang.com/dictionary/a". Similarly, for slangs starting with 'B', we can use the URL "https://www.noslang.com/dictionary/b". The pattern continues for other alphabets as well.

Once we have retrieved the HTML content from the web page, we can extract the desired keywords, in this case, the slangs, from the HTML using techniques like web scraping or using appropriate libraries such as BeautifulSoup.

To store the extracted slangs in a file in JSON format, we can create a JSON object or dictionary and populate it with the extracted slangs as key-value pairs. Each key represents an alphabet, and the corresponding value is a list of slangs starting with that alphabet. We can then write this JSON object to a file using the appropriate functions provided by the programming language or JSON libraries.

By following this approach, we can create a request for each web page corresponding to different alphabets, retrieve the HTML content, extract the slangs, and finally store them in a file in JSON format.

```
In [55]:  from bs4 import BeautifulSoup
          import urllib3
          import json
          http=urllib3.PoolManager()
          Abbr_dict={}
                            t the Slangs from https://www.noslang.com/dictionary/
          def getAbbr(alpha):
```
Loading [MathJax]/extensions/Safe.js

```python
    global Abbr_dict
    r=http.request('GET','https://www.noslang.com/dictionary/'+alpha)
    soup=BeautifulSoup(r.data,'html.parser')
    # print(soup.findAll('div'))
    for i in soup.findAll('div', class_="dictonary-word"):
        abbr=i.find('abbr')['title']
        #print(abbr)
        Abbr_dict[str(i.text[:2])]=abbr
        #Abbr_dict[str(i.text)]=abbr[:2]
        #print(Abbr_dict)


linkDict=[]
#Generating a-z
for one in range(97,123):
    linkDict.append(chr(one))
#Creating Links for https://www.noslang.com/dictionary/a...https://www.noslang.com/dictionary/b....etc
for i in linkDict:
    getAbbr(i)
# finally writing into a json file
with open("ShortendText.json","w") as file:
    jsonDict = json.dump(Abbr_dict,file)
    print("File Created and Content added")
```

```
File Created and Content added
```

Based on the provided text and available information, the system will search the dictionary and file for the corresponding slang terms and retrieve the appropriate results.

In [56]:
```python
import json

with open('ShortendText.json','r') as file:
    Abbr_dict=json.loads(file.read())

#print(Abbr_dict)
line = "can't've you, i'm always stuck and Ay . DG ."
print("Input: ",line)
splitLine=line.split()
#print("SplitLine1", splitLine)

for i in line.split():
    if i in CONTRACTION_MAP:
        #print(i,CONTRACTION_MAP[i])
        splitLine[splitLine.index(i)]=CONTRACTION_MAP[i]
        #print("SplitLine2", splitLine)
    if i in Abbr_dict :
        #print(i,Abbr_dict[i])
        splitLine[splitLine.index(i)]=Abbr_dict[i]
        #print("SplitLine3", splitLine)

result = ' '.join(splitLine)
print("Result:",result)
```

```
Input:  can't've you, i'm always stuck and Ay . DG .
Result: cannot have you, i am always stuck and Are you bored because I am . Don't get me wrong .
```

When working with the 20 Newsgroups dataset, there are several potential improvements and considerations to address various challenges:

Handling Unbalanced Data: The 20 Newsgroups dataset is relatively balanced, with an equal number of samples in each category. However, if you encounter an unbalanced dataset in topic modeling, techniques such as oversampling or undersampling can be employed to balance the representation of different topics.

Dealing with Slang, Abbreviations, or Typos: Preprocessing steps can be implemented to handle slang, abbreviations, and typos. This may involve creating a dictionary or file of commonly used slang terms or abbreviations and mapping them to their full forms. Additionally, spell-checking or correction algorithms can be applied to fix typographical errors in the text.

Addressing Context and Word Disambiguation Challenges: Topic modeling algorithms, including BERTopic, leverage contextual information to extract topics. However, the challenge of word disambiguation still exists. Techniques such as word sense disambiguation or context-aware word embeddings can be used to address this challenge by considering the surrounding context of ambiguous words.

Incorporating Contextual Word Embeddings: BERTopic utilizes BERT-based embeddings to capture contextual information within the text. This helps in understanding the meaning of words in different contexts. Using pre-trained contextual word embeddings like BERT, ELMO, or GPT can improve the accuracy and relevance of the extracted topics.

Fine-tuning BERT for Topic Modeling: While BERTopic uses pre-trained BERT models, fine-tuning BERT specifically for topic modeling on domain-specific datasets can further enhance its performance. Fine-tuning allows the model to adapt to the specific characteristics of the dataset and extract more meaningful topics.

Domain-Specific Topic Labeling: To improve the interpretability of topics, manual labeling or post-processing techniques can be employed. This involves assigning meaningful labels to the extracted topics based on the knowledge of the specific domain or context of the dataset.

By considering these improvements and addressing the challenges mentioned, the topic modeling process on the 20 Newsgroups dataset can yield more accurate and relevant results.

Thank You !!!

In [ ]: