





```
def preprocess(text):

    text = text.lower()

    tokens = word_tokenize(text)

    stop_words = set(stopwords.words('english')) stemmer = PorterStemmer()

    words = [stemmer.stem(word) for word in tokens if word.isalnum() and word not in
stop_words]

    return words documents = {}

for filename in os.listdir():

    if filename.endswith(".txt"):

        with open(filename, 'r', encoding='utf-8', errors='ignore') as f:

            text = f.read()

            documents[filename] = preprocess(text)

print(f"Total documents loaded: {len(documents)}")

inverted_index = defaultdict(set)

for doc_id, words in documents.items():

    for word in set(words): # avoid duplicates per document

        inverted_index[word].add(doc_id)

vocab_size = len(inverted_index)

print(f"\nVocabulary Size: {vocab_size} words")

print("\nSample inverted index terms:")

for term in list(inverted_index)[:10]:

    print(f"{term}: {sorted(inverted_index[term])}")
```

**OUTPUT:**

Total documents loaded: 11

```
defaultdict(<class 'set'>, {'today': {'HI.txt'}, 'work': {'HI.txt'}, 'warm': {'untitled4.txt', 'HI.txt'},
'hi': {'HI.txt'}, 'professor': {'HI.txt'}, 'aditya': {'HI.txt'}, 'sushuma': {'HI.txt'}, 'assist': {'HI.txt'},
```

[illegible]



```
'sunni': {'untitled4.txt'})
```

Vocabulary Size: 9 words

Sample inverted index terms:

hi: ['HI.txt']

sushuma: ['HI.txt']

today: ['HI.txt']

aditya: ['HI.txt']

work: ['HI.txt']

professor: ['HI.txt']

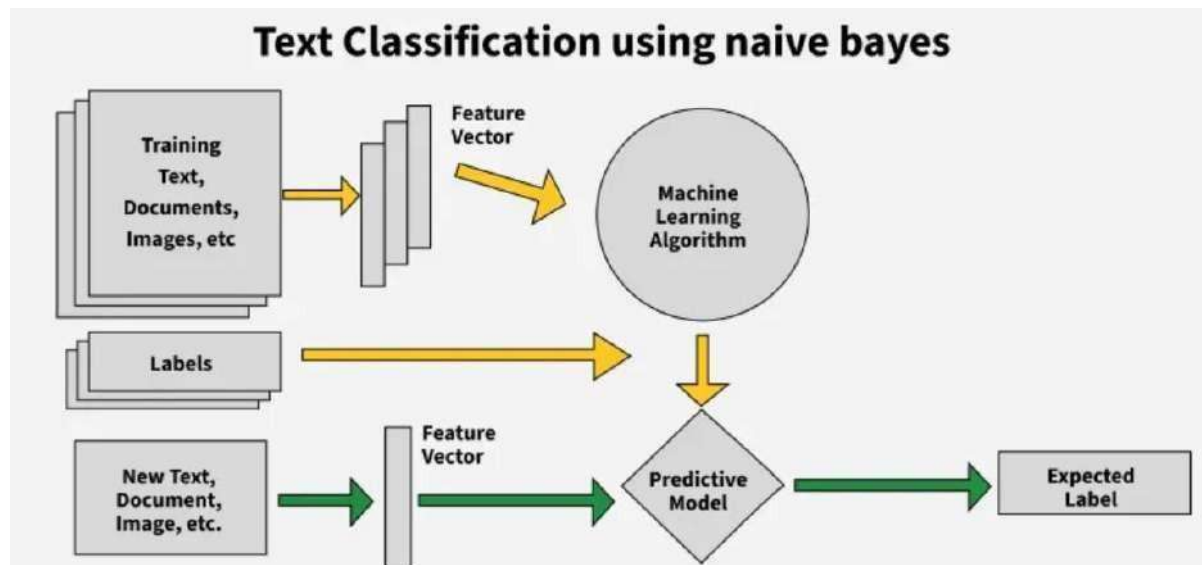
assist: ['HI.txt']

```
warm: ['HI.txt', 'untitled4.txt']
```

sunni: ['untitled4.txt']

[illegible]





**PROGRAM:**

```
from sklearn.datasets import fetch_20newsgroups

from sklearn.model_selection import train_test_split

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.naive_bayes import MultinomialNB

from sklearn.metrics import accuracy_score, classification_report

categories = ['sci.space', 'rec.sport.hockey', 'comp.graphics', 'alt.atheism']

newsgroups = fetch_20newsgroups(subset='all', categories=categories, shuffle=True, random_state=42)

print(f"Total documents: {len(newsgroups.data)}")

print(f"Target classes: {newsgroups.target_names}")

X_train, X_test, y_train, y_test = train_test_split(newsgroups.data, newsgroups.target, test_size=0.2,
                                                    random_state=42)

vectorizer = TfidfVectorizer(stop_words='english')

X_train_tfidf = vectorizer.fit_transform(X_train)

X_test_tfidf = vectorizer.transform(X_test)

nb = MultinomialNB()
```











```
websites = input("Enter comma-separated websites to limit crawling (e.g., bbc.com,cnn.com):")
websites = websites.strip().split(',')

SERP_API_KEY = '8d6bc2b3eef2e66c277a5a34be29b70d490834e929934539b15ae91c71dd569c'

search_url = 'https://serpapi.com/search.json'

def search_news(topic, websites):
    all_results = []

    for site in websites:
        params = {
            "engine": "google",
            "q": f'{topic} site:{site.strip()}',
            "api_key": SERP_API_KEY
        }

        response = requests.get(search_url, params=params)
        data = response.json()

        if "organic_results" in data:
            for result in data["organic_results"]:
                title = result.get("title")
                link = result.get("link")
                snippet = result.get("snippet", "")
                all_results.append((title, link, snippet))

    return all_results

def display_results(results):
    for idx, (title, link, snippet) in enumerate(results, start=1):
        print(f"\nNews {idx}:")
        print(f"Title : {title}")
        print(f"Link : {link}")
        print(f"Snippet : {snippet}")
```

[illegible]



```
print(f"URL : {link}")

print(f"Summary : {snippet}")

results = search_news(topic, websites)

if results:

    display_results(results)

else:

    print("No results found.")
```

**OUTPUT:**

Enter the news topic to search for: AI in healthcare

Enter comma-separated websites to limit crawling (e.g., `bbc.com,cnn.com`): `bbc.com,cnn.com`

News 1:

Title : AI in healthcare: what are the risks for the NHS?

URL : <https://www.bbc.com/news/articles/c6233x9k4dlo>

Summary : Generative AI will be transformative for NHS patient outcomes, a senior government advisor says.

News 2:

Title : How AI can spot diseases that doctors aren't looking for

URL : <https://www.bbc.com/news/articles/c9q7zgy1xlpo>

Summary : AI can take a second look at medical scans and flag up potential problems that doctors might not see.

News 3:

Title : How AI Has Transformed Healthcare's Future

URL : <https://www.bbc.com/storyworks/hpe-greenlake/how-ai-has-transformed-healthcares-future>

Summary : AI can link seemingly unrelated information to reveal new research pathways that yield better results. For example, AI models have identified potential ...

[illegible]



News 4:

Title : Hospitals will use AI to speed up patient care

URL : <https://www.bbc.com/news/articles/cye0yywdegdo>

Summary : Hospitals across the region are to use artificial intelligence (AI) technology to reduce unnecessary admissions and lengthy stays, ...

News 5:

Title : Can AI help modernise Ireland's healthcare system?

URL : <https://www.bbc.com/news/articles/cly7yxm3py5o>

Summary : Ireland is investing billions of euros to revamp its healthcare service - will AI help?

News 6:

Title : How artificial intelligence is matching drugs to patients

URL : <https://www.bbc.com/news/business-65260592>

Summary : Health-tech firms around the world are increasingly using AI to help tailor drugs for patients.

[illegible]