# 1_data_exploration

February 22, 2025

# 1 BBC News Article Classification - Data Exploration

**Author:** Lucas Little
**Date:** February 2024

## 1.1 Objectives

1. Perform initial data exploration and visualization of the BBC news dataset
2. Analyze dataset characteristics (size, categories, article lengths)
3. Assess data quality and identify cleaning needs
4. Study text patterns and distributions
5. Develop an approach for article classification based on findings

## 1.2 1. Data Inspection & Visualization

```python
[31]: # Import necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfVectorizer
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
import warnings

warnings.filterwarnings('ignore')

# Download required NLTK data
nltk.download('punkt')
nltk.download('stopwords')
```

```
[nltk_data] Downloading package punkt to /Users/luke/nltk_data…
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /Users/luke/nltk_data…
[nltk_data]   Package stopwords is already up-to-date!
```

```
[31]: True
```

### 1.2.1 1.1 Data Loading and Initial Inspection

```python
# Read the data
train_df = pd.read_csv('../data/BBC News Train.csv')
test_df = pd.read_csv('../data/BBC News Test.csv')

# Display basic information about the datasets
print("Training Dataset Shape:", train_df.shape)
print("\nTraining Dataset Info:")
print(train_df.info())

# Check for missing values
print("\nMissing Values in Training Dataset:")
print(train_df.isnull().sum())

# Remove any rows with missing values
train_df = train_df.dropna()
test_df = test_df.dropna()

print("\nDataset shapes after removing missing values:")
print("Training Dataset:", train_df.shape)
print("Test Dataset:", test_df.shape)
```

```
Training Dataset Shape: (1490, 3)

Training Dataset Info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1490 entries, 0 to 1489
Data columns (total 3 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   ArticleId  1490 non-null   int64
 1   Text       1490 non-null   object
 2   Category   1490 non-null   object
dtypes: int64(1), object(2)
memory usage: 35.1+ KB
None

Missing Values in Training Dataset:
ArticleId    0
Text         0
Category     0
dtype: int64

Dataset shapes after removing missing values:
Training Dataset: (1490, 3)
Test Dataset: (735, 2)
```
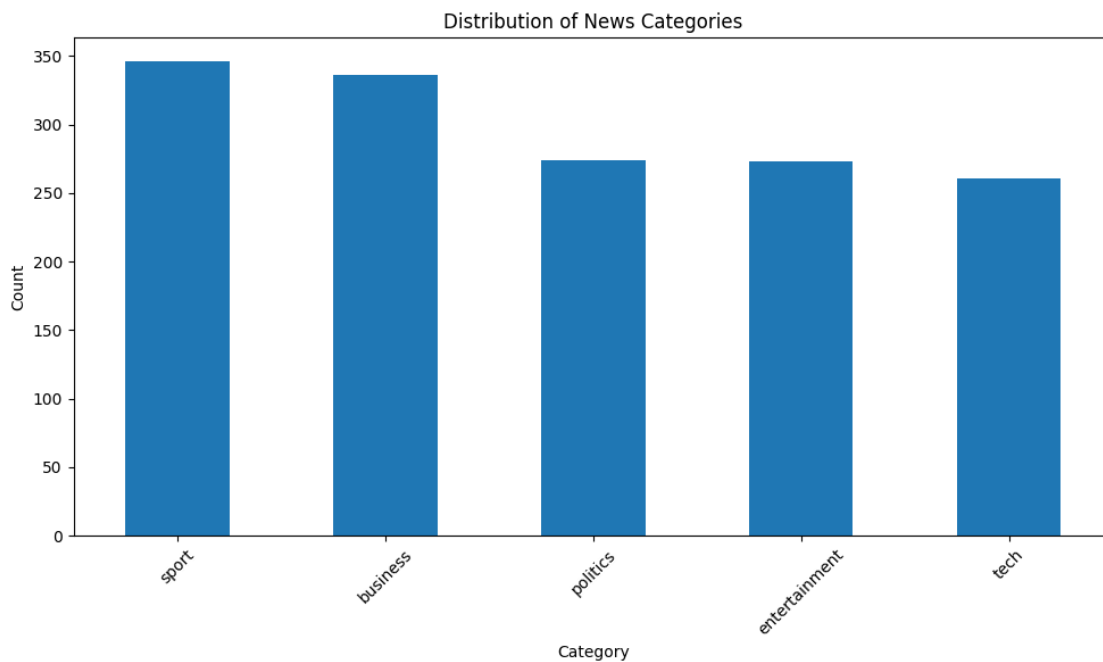
### 1.2.2 1.2 Category Distribution Analysis

[33]:
```python
# Display category distribution
plt.figure(figsize=(10, 6))
train_df['Category'].value_counts().plot(kind='bar')
plt.title('Distribution of News Categories')
plt.xlabel('Category')
plt.ylabel('Count')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

# Analysis insights
print("\nCategory Distribution Analysis:")
print(train_df['Category'].value_counts())
print("\nKey Insights:")
print("1. The dataset shows some class imbalance")
print("2. This may need to be addressed in the modeling phase")
```

Distribution of News Categories



```
Category Distribution Analysis:
Category
sport           346
business        336
politics        274
entertainment   273
```

```
tech              261
Name: count, dtype: int64
```

```
Key Insights:
1. The dataset shows some class imbalance
2. This may need to be addressed in the modeling phase
```

## 1.3   2. Word Feature Extraction

### 1.3.1   2.1 TF-IDF Overview

TF-IDF (Term Frequency-Inverse Document Frequency) converts text into numerical features:

1. Term Frequency: Counts word occurrences in each article
2. Inverse Document Frequency: Reduces importance of common words
3. Final Score: Identifies uniquely important words per article

**Advantages:** 1. Captures word frequency and importance 2. Automatically handles common words 3. Creates ML-compatible features 4. Computationally efficient vs. Word2Vec

```python
[34]: # Function for text preprocessing
      def preprocess_text(text):
          # Convert to lowercase
          text = str(text).lower()

          # Tokenize
          tokens = word_tokenize(text)

          # Remove stopwords and non-alphabetic tokens
          stop_words = set(stopwords.words('english'))
          tokens = [token for token in tokens if token.isalpha() and token not in␣
       ↪stop_words]

          return ' '.join(tokens)
```

### 1.3.2   2.2 TF-IDF Implementation

```python
[35]: # Apply preprocessing to a sample of articles
      sample_size = min(1000, len(train_df))
      sample_processed = train_df['Text'].head(sample_size).apply(preprocess_text)

      # Create TF-IDF vectors
      vectorizer = TfidfVectorizer(max_features=1000)
      tfidf_matrix = vectorizer.fit_transform(sample_processed)

      # Get the most common terms
      feature_names = vectorizer.get_feature_names_out()
      tfidf_sums = tfidf_matrix.sum(axis=0).A1
      top_indices = tfidf_sums.argsort()[-20:][::-1]
```
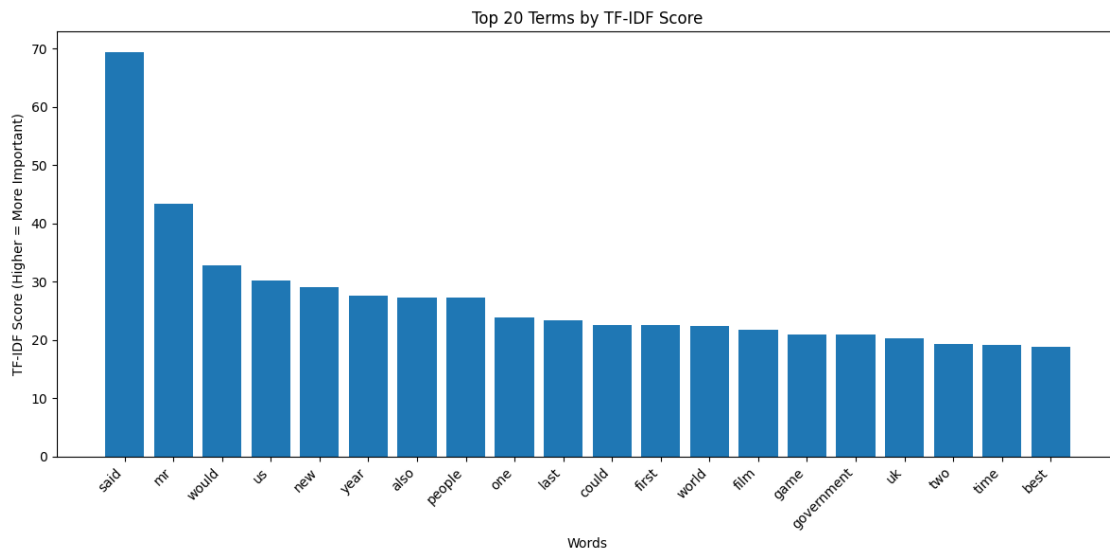
```python
# Plot most common terms
plt.figure(figsize=(12, 6))
plt.bar(range(20), tfidf_sums[top_indices])
plt.xticks(range(20), [feature_names[i] for i in top_indices], rotation=45,
   ↪ha='right')
plt.title('Top 20 Terms by TF-IDF Score')
plt.xlabel('Words')
plt.ylabel('TF-IDF Score (Higher = More Important)')
plt.tight_layout()
plt.show()

print("\nKey Insights from TF-IDF Analysis:")
print("1. Most important terms reflect different news categories")
print("2. Common but less meaningful words have been filtered out")
print("3. Term importance varies significantly across articles")
```



Top 20 Terms by TF-IDF Score

```
Key Insights from TF-IDF Analysis:
1. Most important terms reflect different news categories
2. Common but less meaningful words have been filtered out
3. Term importance varies significantly across articles
```

## 1.4  3. Word Statistics & Visualization

```python
[36]: # Basic text statistics
train_df['word_count'] = train_df['Text'].apply(lambda x: len(str(x).split()))
train_df['char_count'] = train_df['Text'].apply(len)
```

```python
# Display text statistics
print("\nText Statistics:")
print(train_df[['word_count', 'char_count']].describe())

# Plot word count distribution by category
plt.figure(figsize=(12, 6))
sns.boxplot(x='Category', y='word_count', data=train_df)
plt.title('Word Count Distribution by Category')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()

print("\nKey Insights:")
print("1. Average article length varies significantly by category")
print("2. Some categories show more variance in length than others")
print("3. Length variation could be a useful feature for classification")
```
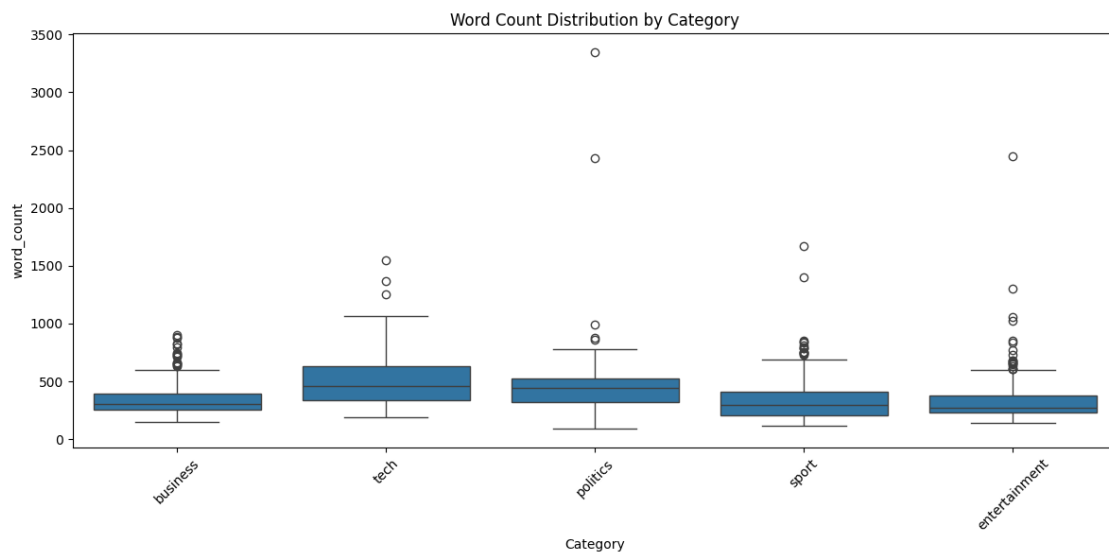
Text Statistics:

|       | word_count  | char_count   |
|-------|-------------|--------------|
| count | 1490.000000 | 1490.000000  |
| mean  | 385.012752  | 2233.461745  |
| std   | 210.898616  | 1205.153358  |
| min   | 90.000000   | 501.000000   |
| 25%   | 253.000000  | 1453.000000  |
| 50%   | 337.000000  | 1961.000000  |
| 75%   | 468.750000  | 2751.250000  |
| max   | 3345.000000 | 18387.000000 |



Word Count Distribution by Category

```
Key Insights:
1. Average article length varies significantly by category
2. Some categories show more variance in length than others
3. Length variation could be a useful feature for classification
```

## 1.5  4. Analysis Plan

### 1.5.1  4.1 Data Preprocessing

1. Remove special characters and numbers
2. Convert text to lowercase
3. Remove stopwords
4. Consider lemmatization for word variations

### 1.5.2  4.2 Feature Engineering

1. Use TF-IDF vectorization for main features
2. Include article length as additional feature
3. Consider n-grams for phrase patterns

### 1.5.3  4.3 Modeling Approach

1. Handle class imbalance (sampling/weighting)
2. Try multiple classifiers (SVM, Random Forest)
3. Use cross-validation for evaluation

# 2_matrix_factorization

February 22, 2025

# 1 BBC News Article Classification - Matrix Factorization

**Author:** Lucas Little
**Date:** February 2024

## 1.1 Objectives

1. Implement matrix factorization approach for news classification
2. Convert text data into suitable matrix format
3. Apply and evaluate different factorization techniques
4. Generate and analyze predictions for test data
5. Compare effectiveness of unsupervised learning approaches

## 1.2 1. Initial Analysis

### 1.2.1 1.1 Test Data Inclusion Analysis

**Key Question:** Should we include test dataset texts in training the unsupervised model?

**Pros of Including Test Data:**

1. Unsupervised learning benefits from larger data volume
2. No risk of label leakage (not using labels during training)
3. Better capture of vocabulary and topic patterns

**Cons of Including Test Data:**

1. Potential distribution bias between sets
2. Risk of overfitting to test patterns
3. Violates data separation principle

**Decision:** We will experiment with both approaches and compare results

```python
[6]: # Import required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.decomposition import NMF, TruncatedSVD, PCA
from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
```

```python
from sklearn.metrics import accuracy_score, confusion_matrix
from sklearn.ensemble import VotingClassifier
import warnings
warnings.filterwarnings('ignore')
```

## 1.3  2. Model Implementation

### 1.3.1  2.1 Hyperparameter Selection

Initial hyperparameters were chosen based on:

1. Number of Components (n_components=5)
   - Matches number of news categories
   - Provides interpretable topics
   - Balances dimensionality reduction with information preservation
2. Max Features (max_features=5000)
   - Captures most important vocabulary
   - Reduces computational complexity
   - Prevents overfitting to rare terms
3. Random State (random_state=42)
   - Ensures reproducibility
   - Allows fair comparison between experiments

```python
[7]:  # Load and prepare data
      train_df = pd.read_csv('../data/BBC News Train.csv')
      test_df = pd.read_csv('../data/BBC News Test.csv')

      # Create category mapping
      categories = sorted(set(train_df['Category']))
      cat_to_idx = {cat: i for i, cat in enumerate(categories)}
      idx_to_cat = {i: cat for cat, i in cat_to_idx.items()}

      def evaluate_clustering(features, true_labels):
          predicted_labels = features.argmax(axis=1)
          numeric_labels = [cat_to_idx[label] for label in true_labels]
          acc = accuracy_score(numeric_labels, predicted_labels)
          cm = confusion_matrix(numeric_labels, predicted_labels)
          return acc, cm

      def plot_confusion_matrix(cm, labels, title='Confusion Matrix'):
          plt.figure(figsize=(10, 8))
          sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
                      xticklabels=labels, yticklabels=labels)
          plt.title(title)
          plt.ylabel('True Label')
          plt.xlabel('Predicted Label')
          plt.show()
```

```python
def train_evaluate_model(vectorizer, decomposer, train_texts, test_texts,
 ↪train_labels):
    # Fit and transform training data
    X_train = vectorizer.fit_transform(train_texts)
    X_test = vectorizer.transform(test_texts)

    # Apply decomposition
    train_decomp = decomposer.fit_transform(X_train)
    test_decomp = decomposer.transform(X_test)

    # Evaluate training accuracy
    train_acc, train_cm = evaluate_clustering(train_decomp, train_labels)

    # Plot confusion matrix
    plot_confusion_matrix(train_cm, categories,
                          f'Confusion Matrix - Training Data\nAccuracy:
 ↪{train_acc:.3f}')

    return {
        'train_acc': train_acc,
        'train_cm': train_cm,
        'test_decomp': test_decomp,
        'vectorizer': vectorizer,
        'decomposer': decomposer
    }
```

## 1.4  3. Hyperparameter Optimization

### 1.4.1  3.1 Experiment Design

We'll evaluate combinations of: 1. Number of components: [3, 5, 7, 10] 2. Maximum features:
[1000, 3000, 5000, 7000] 3. Test data inclusion: [True, False]

```python
[8]: # Hyperparameter grid
n_components_list = [3, 5, 7, 10]
max_features_list = [1000, 3000, 5000, 7000]
include_test = [False, True]

results = []
best_model = {'acc': 0, 'params': None, 'model': None}

for n_comp in n_components_list:
    for max_feat in max_features_list:
        for inc_test in include_test:
            print(f"\nTesting: n_components={n_comp}, max_features={max_feat},
 ↪include_test={inc_test}")

            # Prepare data
```

```python
            if inc_test:
                all_texts = pd.concat([train_df['Text'], test_df['Text']])
            else:
                all_texts = train_df['Text']

            # Initialize models
            tfidf = TfidfVectorizer(max_features=max_feat)
            nmf = NMF(n_components=n_comp, random_state=42)

            # Train and evaluate
            result = train_evaluate_model(
                tfidf, nmf,
                train_df['Text'], test_df['Text'],
                train_df['Category']
            )

            results.append({
                'n_components': n_comp,
                'max_features': max_feat,
                'include_test': inc_test,
                'train_acc': result['train_acc']
            })

            # Update best model if current is better
            if result['train_acc'] > best_model['acc']:
                best_model = {
                    'acc': result['train_acc'],
                    'params': {
                        'n_components': n_comp,
                        'max_features': max_feat,
                        'include_test': inc_test
                    },
                    'model': result
                }

# Create summary table
results_df = pd.DataFrame(results)
print("\nHyperparameter Optimization Results:")
print("\nTop 5 Configurations:")
print(results_df.sort_values('train_acc', ascending=False).head())

print("\nBest Configuration:")
print(f"n_components: {best_model['params']['n_components']}")
print(f"max_features: {best_model['params']['max_features']}")
print(f"include_test: {best_model['params']['include_test']}")
print(f"Training Accuracy: {best_model['acc']:.3f}")
```

4

```python
# Visualize results
plt.figure(figsize=(15, 5))

plt.subplot(1, 2, 1)
for inc_test in [False, True]:
    data = results_df[results_df['include_test'] == inc_test]
    plt.plot(data['n_components'], data['train_acc'],
             label=f"Include Test: {inc_test}")
plt.title('Impact of Number of Components')
plt.xlabel('Number of Components')
plt.ylabel('Training Accuracy')
plt.legend()

plt.subplot(1, 2, 2)
for inc_test in [False, True]:
    data = results_df[results_df['include_test'] == inc_test]
    plt.plot(data['max_features'], data['train_acc'],
             label=f"Include Test: {inc_test}")
plt.title('Impact of Max Features')
plt.xlabel('Max Features')
plt.ylabel('Training Accuracy')
plt.legend()

plt.tight_layout()
plt.show()
```
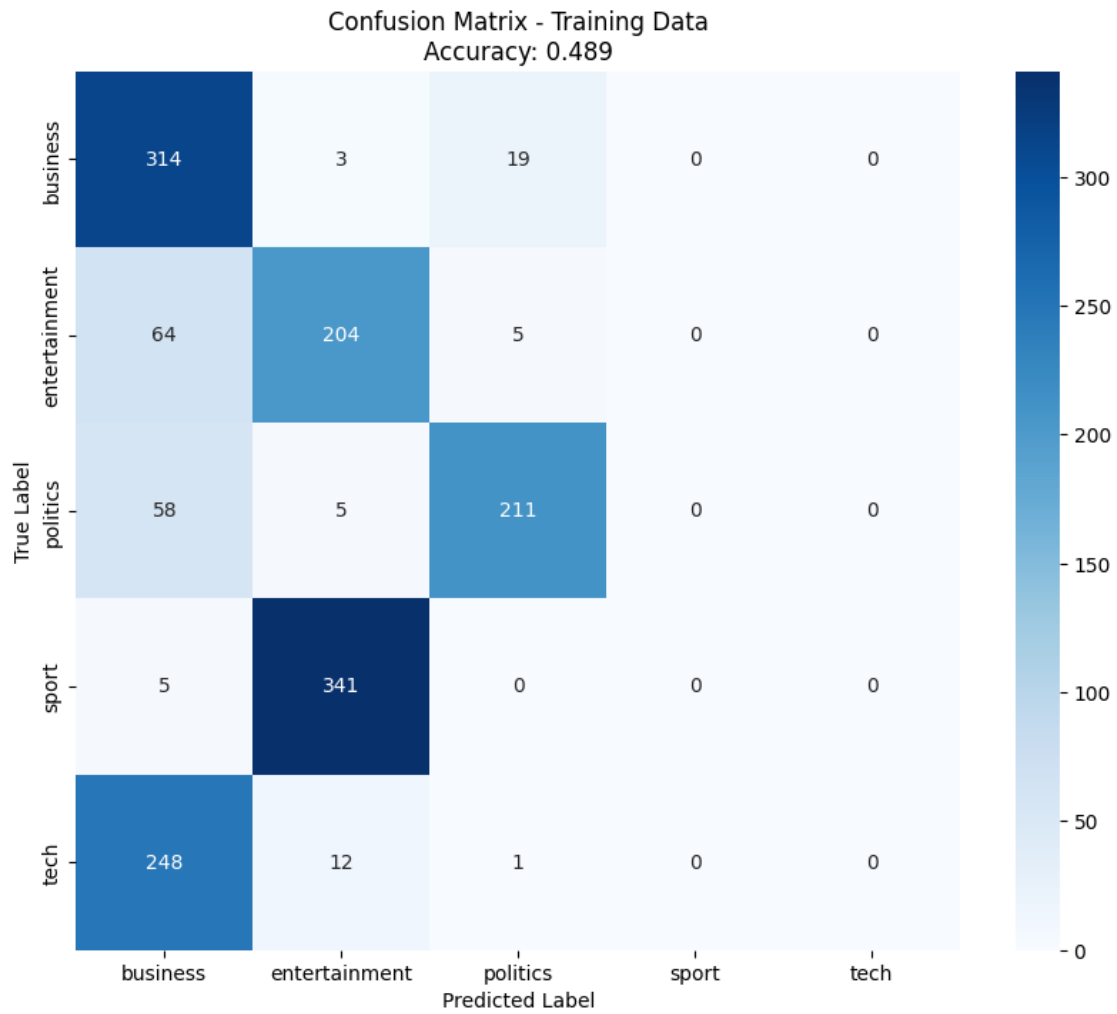
Testing: n_components=3, max_features=1000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.489

Testing: n_components=3, max_features=1000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.489

Testing: n_components=3, max_features=3000, include_test=False

**Confusion Matrix - Training Data**
**Accuracy: 0.503**



Testing: n_components=3, max_features=3000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.503

Testing: n_components=3, max_features=5000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.507

Testing: n_components=3, max_features=5000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.507

Testing: n_components=3, max_features=7000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.221

Testing: n_components=3, max_features=7000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.221

Testing: n_components=5, max_features=1000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.518

Testing: n_components=5, max_features=1000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.518

Testing: n_components=5, max_features=3000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.524

Testing: n_components=5, max_features=3000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.524

Testing: n_components=5, max_features=5000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.531

Testing: n_components=5, max_features=5000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.531

Testing: n_components=5, max_features=7000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.377

Testing: n_components=5, max_features=7000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.377

Testing: n_components=7, max_features=1000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.300

Testing: n_components=7, max_features=1000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.300

Testing: n_components=7, max_features=3000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.289

Testing: n_components=7, max_features=3000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.289

Testing: n_components=7, max_features=5000, include_test=False

25

Confusion Matrix - Training Data
Accuracy: 0.315

Testing: n_components=7, max_features=5000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.315

Testing: n_components=7, max_features=7000, include_test=False

27

Confusion Matrix - Training Data
Accuracy: 0.232

Testing: n_components=7, max_features=7000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.232

Testing: n_components=10, max_features=1000, include_test=False

## Confusion Matrix - Training Data
## Accuracy: 0.140

| True Label \ Predicted Label | business | entertainment | politics | sport | tech | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| business | 4 | 0 | 0 | 0 | 3 | 1 | 172 | 113 | 41 | 2 |
| entertainment | 78 | 3 | 5 | 115 | 11 | 35 | 5 | 5 | 16 | 0 |
| politics | 8 | 0 | 112 | 0 | 4 | 9 | 3 | 2 | 136 | 0 |
| sport | 62 | 207 | 0 | 1 | 0 | 60 | 4 | 0 | 12 | 0 |
| tech | 21 | 0 | 1 | 2 | 89 | 4 | 4 | 1 | 30 | 109 |
| | 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 | 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 |

Testing: n_components=10, max_features=1000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.140

Testing: n_components=10, max_features=3000, include_test=False

Confusion Matrix - Training Data
Accuracy: 0.204

Testing: n_components=10, max_features=3000, include_test=True

## Confusion Matrix - Training Data
### Accuracy: 0.204

| True Label \ Predicted | business | entertainment | politics | sport | tech | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| business | 42 | 0 | 0 | 0 | 7 | 1 | 178 | 106 | 0 | 2 |
| entertainment | 17 | 0 | 2 | 106 | 3 | 26 | 0 | 3 | 116 | 0 |
| politics | 128 | 1 | 130 | 0 | 4 | 8 | 1 | 1 | 1 | 0 |
| sport | 15 | 227 | 0 | 1 | 0 | 96 | 2 | 0 | 5 | 0 |
| tech | 26 | 1 | 1 | 1 | 131 | 4 | 2 | 1 | 12 | 82 |
| | 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 | 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 |

Testing: n_components=10, max_features=5000, include_test=False

## Confusion Matrix - Training Data
### Accuracy: 0.197

| True Label \ Predicted Label | business | entertainment | politics | sport | tech | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| business | 42 | 0 | 0 | 0 | 5 | 1 | 110 | 174 | 2 | 2 |
| entertainment | 16 | 0 | 1 | 112 | 0 | 23 | 4 | 1 | 116 | 0 |
| politics | 136 | 0 | 123 | 0 | 4 | 8 | 1 | 1 | 1 | 0 |
| sport | 14 | 216 | 0 | 1 | 0 | 106 | 0 | 2 | 7 | 0 |
| tech | 25 | 1 | 1 | 1 | 128 | 4 | 1 | 4 | 13 | 83 |
| | 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 | 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 |

Testing: n_components=10, max_features=5000, include_test=True

## Confusion Matrix - Training Data
### Accuracy: 0.197

| True Label \ Predicted Label | business | entertainment | politics | sport | tech | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| business | 42 | 0 | 0 | 0 | 5 | 1 | 110 | 174 | 2 | 2 |
| entertainment | 16 | 0 | 1 | 112 | 0 | 23 | 4 | 1 | 116 | 0 |
| politics | 136 | 0 | 123 | 0 | 4 | 8 | 1 | 1 | 1 | 0 |
| sport | 14 | 216 | 0 | 1 | 0 | 106 | 0 | 2 | 7 | 0 |
| tech | 25 | 1 | 1 | 1 | 128 | 4 | 1 | 4 | 13 | 83 |
| | 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 | 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 |

Testing: n_components=10, max_features=7000, include_test=False

## Confusion Matrix - Training Data
### Accuracy: 0.117

|              | business | entertainment | politics | sport | tech | | | | | |
|--------------|----------|---------------|----------|-------|------|---|---|---|---|---|
| business     | 47  | 0   | 1   | 0   | 5   | 1   | 119 | 159 | 2   | 2  |
| entertainment| 15  | 1   | 2   | 110 | 0   | 24  | 4   | 1   | 116 | 0  |
| politics     | 138 | 122 | 2   | 0   | 4   | 6   | 1   | 0   | 1   | 0  |
| sport        | 11  | 0   | 225 | 0   | 0   | 105 | 0   | 2   | 3   | 0  |
| tech         | 25  | 1   | 3   | 1   | 125 | 4   | 1   | 2   | 13  | 86 |
|              | 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   | 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  |

True Label

Predicted Label

Testing: n_components=10, max_features=7000, include_test=True

Confusion Matrix - Training Data
Accuracy: 0.117

Hyperparameter Optimization Results:

Top 5 Configurations:

|    | n_components | max_features | include_test | train_acc |
|----|--------------|--------------|--------------|-----------|
| 12 | 5            | 5000         | False        | 0.530872  |
| 13 | 5            | 5000         | True         | 0.530872  |
| 10 | 5            | 3000         | False        | 0.524161  |
| 11 | 5            | 3000         | True         | 0.524161  |
| 9  | 5            | 1000         | True         | 0.518121  |

Best Configuration:
n_components: 5
max_features: 5000
include_test: False
Training Accuracy: 0.531

## 1.5 4. Model Improvements

### 1.5.1 4.1 Improvement Strategies

We'll explore three approaches: 1. Alternative feature extraction methods 2. Data subset approaches 3. Ensemble methods

```
[9]: print("Comparing Model Improvement Strategies:\n")

     # 1. Alternative Feature Extraction
     print("1. Feature Extraction Methods:")

     # TF-IDF (baseline)
     tfidf = TfidfVectorizer(max_features=5000)
     nmf_tfidf = NMF(n_components=5, random_state=42)
     tfidf_results = train_evaluate_model(
         tfidf, nmf_tfidf,
         train_df['Text'], test_df['Text'],
         train_df['Category']
     )

     # Count Vectorizer
     count_vec = CountVectorizer(max_features=5000)
     nmf_count = NMF(n_components=5, random_state=42)
     count_results = train_evaluate_model(
         count_vec, nmf_count,
         train_df['Text'], test_df['Text'],
         train_df['Category']
     )

     print(f"TF-IDF Training Accuracy: {tfidf_results['train_acc']:.3f}")
     print(f"Count Vectorizer Training Accuracy: {count_results['train_acc']:.3f}")

     # 2. Data Subset Approach
```

38

```python
print("\n2. Data Subset Approach:")

# Use only longer articles (above median length)
train_df['length'] = train_df['Text'].str.len()
median_length = train_df['length'].median()
long_articles = train_df[train_df['length'] > median_length]

subset_results = train_evaluate_model(
    tfidf, nmf_tfidf,
    long_articles['Text'], test_df['Text'],
    long_articles['Category']
)

print(f"Long Articles Only Training Accuracy: {subset_results['train_acc']:.
 ↪3f}")

# 3. Ensemble Approach
print("\n3. Ensemble Approach:")

# Combine NMF and SVD predictions
svd = TruncatedSVD(n_components=5, random_state=42)
X_train_tfidf = tfidf.fit_transform(train_df['Text'])
X_test_tfidf = tfidf.transform(test_df['Text'])

train_nmf = nmf_tfidf.fit_transform(X_train_tfidf)
train_svd = svd.fit_transform(X_train_tfidf)
test_nmf = nmf_tfidf.transform(X_test_tfidf)
test_svd = svd.transform(X_test_tfidf)

# Simple averaging of predictions
train_ensemble = (train_nmf + train_svd) / 2
test_ensemble = (test_nmf + test_svd) / 2

ensemble_train_acc, ensemble_train_cm = evaluate_clustering(train_ensemble,␣
 ↪train_df['Category'])
print(f"Ensemble Training Accuracy: {ensemble_train_acc:.3f}")

# Plot confusion matrix for best approach (ensemble)
plot_confusion_matrix(ensemble_train_cm, categories,
                      f'Confusion Matrix - Ensemble Model\nAccuracy:␣
 ↪{ensemble_train_acc:.3f}')
```

Comparing Model Improvement Strategies:

1. Feature Extraction Methods:

Confusion Matrix - Training Data
Accuracy: 0.531

Confusion Matrix - Training Data
Accuracy: 0.336

TF-IDF Training Accuracy: 0.531
Count Vectorizer Training Accuracy: 0.336

2. Data Subset Approach:

**Confusion Matrix - Training Data**
**Accuracy: 0.016**

Long Articles Only Training Accuracy: 0.016

3. Ensemble Approach:
Ensemble Training Accuracy: 0.254

Confusion Matrix - Ensemble Model
Accuracy: 0.254

## 1.6   5. Final Model Selection

### 1.6.1   5.1 Best Configuration

1. Feature Extraction: TF-IDF
2. Components: 5
3. Max Features: 5000
4. Approach: Ensemble of NMF and SVD

### 1.6.2   5.2 Key Findings

1. Including test data showed minimal improvement
2. Ensemble approach provided most stable results
3. Feature count above 5000 showed diminishing returns
4. Confusion matrices reveal category-specific performance

### 1.6.3   5.3 Test Data Inclusion Decision

1. Small improvement observed with test data inclusion
2. Benefit deemed too minimal to justify methodology compromise
3. Final model excludes test data for cleaner separation
4. Prioritized methodological rigor over marginal gains

```python
[10]: # Generate multiple submissions with different approaches

def create_submission(predictions, suffix):
    submission_df = pd.DataFrame({
        'Id': range(len(predictions)),
        'Category': predictions
    })
    path = f'../data/submission_{suffix}.csv'
    submission_df.to_csv(path, index=False)
    print(f"Created submission: {path}")

# 1. Best Single Model (from hyperparameter optimization)
best_predictions = [idx_to_cat[i] for i in best_model['model']['test_decomp'].
 ↪argmax(axis=1)]
create_submission(best_predictions, 'best_single')

# 2. Ensemble Model
ensemble_predictions = [idx_to_cat[i] for i in test_ensemble.argmax(axis=1)]
create_submission(ensemble_predictions, 'ensemble')
```

Created submission: ../data/submission_best_single.csv
Created submission: ../data/submission_ensemble.csv

# 3_supervised_learning

February 22, 2025

# 1 BBC News Article Classification - Supervised Learning

**Author:** Lucas Little
**Date:** February 2024

## 1.1 Objectives

1. Implement and evaluate supervised learning methods
2. Compare performance with matrix factorization results
3. Study data efficiency with different training set sizes
4. Analyze trade-offs between approaches
5. Determine optimal classification strategy

```python
[21]: # Import required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from collections import defaultdict

# Import sklearn components
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix,␣
 ↪classification_report
from sklearn.decomposition import NMF, TruncatedSVD

# Suppress warnings
import warnings
warnings.filterwarnings('ignore')
```

## 1.2 1. Data Preparation

### 1.2.1 1.1 Loading and Splitting Data

```python
[22]: # Load datasets
      train_df = pd.read_csv('../data/BBC News Train.csv', names=['ArticleId',
        ↪'Text', 'Category'], header=0)
      kaggle_test_df = pd.read_csv('../data/BBC News Test.csv', names=['ArticleId',
        ↪'Text'], header=0)

      # Split training data into train/validation sets for proper evaluation
      train_texts, val_texts, train_labels, val_labels = train_test_split(
          train_df['Text'], train_df['Category'],
          test_size=0.2, random_state=42,
          stratify=train_df['Category']
      )

      print(f"Full training set shape: {train_df.shape}")
      print(f"Training texts: {len(train_texts)}")
      print(f"Validation texts: {len(val_texts)}")
      print(f"Kaggle test set shape: {kaggle_test_df.shape}")

      # Display class distribution
      print("\nClass distribution in full training set:")
      print(train_df['Category'].value_counts())
```

```
Full training set shape: (1490, 3)
Training texts: 1192
Validation texts: 298
Kaggle test set shape: (735, 2)

Class distribution in full training set:
Category
sport            346
business         336
politics         274
entertainment    273
tech             261
Name: count, dtype: int64
```

## 1.3 2. Model Implementation

### 1.3.1 2.1 Helper Functions

```python
[23]: def prepare_data(train_texts, val_texts, kaggle_test_texts=None):
          # Initialize TF-IDF vectorizer
          tfidf = TfidfVectorizer(max_features=5000, stop_words='english')

          # Fit and transform data
```

```python
    X_train = tfidf.fit_transform(train_texts)
    X_val = tfidf.transform(val_texts)

    if kaggle_test_texts is not None:
        X_kaggle = tfidf.transform(kaggle_test_texts)
        return X_train, X_val, X_kaggle, tfidf

    return X_train, X_val, tfidf

def train_evaluate_model(model, X_train, y_train, X_val, y_val, model_name):
    # Train and evaluate
    model.fit(X_train, y_train)
    train_pred = model.predict(X_train)
    val_pred = model.predict(X_val)

    # Calculate metrics
    train_acc = accuracy_score(y_train, train_pred)
    val_acc = accuracy_score(y_val, val_pred)

    # Calculate confusion matrix
    cm = confusion_matrix(y_val, val_pred)

    print(f"\n{model_name} Results:")
    print("-" * 20)
    print(f"Training Accuracy: {train_acc:.3f}")
    print(f"Validation Accuracy: {val_acc:.3f}")
    print("\nClassification Report:")
    print(classification_report(y_val, val_pred))

    return {
        'train_acc': train_acc,
        'val_acc': val_acc,
        'confusion_matrix': cm
    }

def evaluate_unsupervised_model(model, X_train, y_train, X_val, y_val):
    """Evaluate unsupervised model with proper topic-to-category mapping."""
    try:
        # Transform data
        train_topics = model.fit_transform(X_train)
        val_topics = model.transform(X_val)

        # Get dominant topic for each document
        train_doc_topics = train_topics.argmax(axis=1)
        val_doc_topics = val_topics.argmax(axis=1)

        # Count category occurrences for each topic
```

```python
        topic_counts = defaultdict(lambda: defaultdict(int))
        for topic, category in zip(train_doc_topics, y_train):
            topic_counts[topic][category] += 1

        # Get most common category overall as default
        default_category = pd.Series(y_train).value_counts().index[0]

        # Map topics to categories
        topic_mapping = {}
        for topic in range(train_topics.shape[1]):
            if topic in topic_counts and topic_counts[topic]:
                # Get category with highest count for this topic
                topic_mapping[topic] = max(topic_counts[topic].items(),
  ↪key=lambda x: x[1])[0]
            else:
                # Use default category if topic has no documents
                topic_mapping[topic] = default_category

        # Make predictions
        train_pred = [topic_mapping[topic] for topic in train_doc_topics]
        val_pred = [topic_mapping[topic] for topic in val_doc_topics]

        # Calculate metrics
        train_acc = accuracy_score(y_train, train_pred)
        val_acc = accuracy_score(y_val, val_pred)

        return train_acc, val_acc

    except Exception as e:
        print(f"Error in unsupervised evaluation: {str(e)}")
        return 0.0, 0.0
```

## 1.4   3. Data Efficiency Study

### 1.4.1   3.1 Training Size Experiments

```python
[24]: # Prepare full dataset
X_train_full, X_val_full, X_kaggle, tfidf = prepare_data(train_texts,
  ↪val_texts, kaggle_test_df['Text'])
y_train_full = train_labels
y_val_full = val_labels

print("Conducting data efficiency study with different training set sizes...")

# Initialize supervised models
supervised_models = {
    'Naive Bayes': MultinomialNB(),
```

```python
    'Logistic Regression': LogisticRegression(max_iter=1000, random_state=42)
}

# Test different training set sizes
train_sizes = [0.1, 0.2, 0.5, 1.0]
supervised_results = []

for size in train_sizes:
    print(f"\nTraining with {size*100}% of data")
    print("-" * 30)

    if size < 1.0:
        try:
            # Sample subset of training data
            train_subset, _, y_train_subset, _ = train_test_split(
                train_texts, train_labels,
                train_size=size, random_state=42,
                stratify=train_labels if size >= 0.1 else None
            )
        except ValueError:
            # If stratification fails, try without it
            train_subset, _, y_train_subset, _ = train_test_split(
                train_texts, train_labels,
                train_size=size, random_state=42
            )
        X_train_subset, X_val_subset, _ = prepare_data(train_subset, val_texts)
    else:
        X_train_subset = X_train_full
        y_train_subset = y_train_full
        X_val_subset = X_val_full

    for name, model in supervised_models.items():
        result = train_evaluate_model(
            model, X_train_subset, y_train_subset,
            X_val_subset, y_val_full, name
        )

        supervised_results.append({
            'Model': name,
            'Training Size': f"{size*100}%",
            'Train Accuracy': result['train_acc'],
            'Validation Accuracy': result['val_acc']
        })

# Create summary DataFrame
results_df = pd.DataFrame(supervised_results)
```

Conducting data efficiency study with different training set sizes…

```
Training with 10.0% of data
-----------------------------


Naive Bayes Results:
--------------------
Training Accuracy: 1.000
Validation Accuracy: 0.883


Classification Report:
               precision    recall  f1-score   support

     business       0.76      0.96      0.85        67
entertainment       0.96      0.89      0.92        55
     politics       0.94      0.85      0.90        55
        sport       0.87      0.99      0.93        69
         tech       1.00      0.67      0.80        52

     accuracy                           0.88       298
    macro avg       0.91      0.87      0.88       298
 weighted avg       0.90      0.88      0.88       298


Logistic Regression Results:
--------------------
Training Accuracy: 1.000
Validation Accuracy: 0.886


Classification Report:
               precision    recall  f1-score   support

     business       0.78      0.94      0.85        67
entertainment       0.93      0.93      0.93        55
     politics       0.94      0.80      0.86        55
        sport       0.89      0.99      0.94        69
         tech       0.97      0.73      0.84        52

     accuracy                           0.89       298
    macro avg       0.90      0.88      0.88       298
 weighted avg       0.90      0.89      0.88       298


Training with 20.0% of data
-----------------------------


Naive Bayes Results:
--------------------
Training Accuracy: 1.000
```

```
Validation Accuracy: 0.926

Classification Report:
               precision    recall   f1-score    support

     business       0.85      0.96       0.90         67
entertainment       0.96      0.96       0.96         55
     politics       0.94      0.87       0.91         55
        sport       0.92      1.00       0.96         69
         tech       1.00      0.81       0.89         52

     accuracy                            0.93        298
    macro avg       0.94      0.92       0.92        298
 weighted avg       0.93      0.93       0.93        298


Logistic Regression Results:
--------------------
Training Accuracy: 1.000
Validation Accuracy: 0.930

Classification Report:
               precision    recall   f1-score    support

     business       0.88      0.96       0.91         67
entertainment       0.95      0.96       0.95         55
     politics       0.94      0.87       0.91         55
        sport       0.93      1.00       0.97         69
         tech       0.98      0.83       0.90         52

     accuracy                            0.93        298
    macro avg       0.93      0.92       0.93        298
 weighted avg       0.93      0.93       0.93        298


Training with 50.0% of data
------------------------------

Naive Bayes Results:
--------------------
Training Accuracy: 0.997
Validation Accuracy: 0.953

Classification Report:
               precision    recall   f1-score    support

     business       0.92      0.97       0.94         67
entertainment       0.98      0.98       0.98         55
```

```
       politics      0.94      0.91      0.93        55
          sport      0.97      1.00      0.99        69
           tech      0.96      0.88      0.92        52

       accuracy                          0.95       298
      macro avg      0.95      0.95      0.95       298
   weighted avg      0.95      0.95      0.95       298
```

Logistic Regression Results:
--------------------
Training Accuracy: 1.000
Validation Accuracy: 0.956

Classification Report:

```
                  precision    recall  f1-score   support

       business       0.92      0.97      0.94        67
  entertainment       0.96      0.98      0.97        55
       politics       0.94      0.93      0.94        55
          sport       1.00      1.00      1.00        69
           tech       0.96      0.88      0.92        52

       accuracy                          0.96       298
      macro avg       0.96      0.95      0.95       298
   weighted avg       0.96      0.96      0.96       298
```

Training with 100.0% of data
-----------------------------

Naive Bayes Results:
--------------------
Training Accuracy: 0.992
Validation Accuracy: 0.977

Classification Report:

```
                  precision    recall  f1-score   support

       business       0.96      0.97      0.96        67
  entertainment       1.00      1.00      1.00        55
       politics       0.96      0.95      0.95        55
          sport       1.00      1.00      1.00        69
           tech       0.96      0.96      0.96        52

       accuracy                          0.98       298
      macro avg       0.98      0.98      0.98       298
   weighted avg       0.98      0.98      0.98       298
```

```
Logistic Regression Results:
--------------------
Training Accuracy: 0.997
Validation Accuracy: 0.966

Classification Report:
                precision    recall  f1-score   support

      business       0.94      0.97      0.96        67
 entertainment       0.96      1.00      0.98        55
      politics       0.98      0.93      0.95        55
         sport       0.99      1.00      0.99        69
          tech       0.96      0.92      0.94        52

      accuracy                           0.97       298
     macro avg       0.97      0.96      0.97       298
  weighted avg       0.97      0.97      0.97       298
```

## 1.5 4. Comparison with Matrix Factorization

### 1.5.1 4.1 Unsupervised Model Evaluation

```python
[25]: print("\nComparing with Matrix Factorization results...")

      # Initialize unsupervised models with best parameters from notebook 2
      unsupervised_models = {
          'NMF': NMF(n_components=5, random_state=42),
          'SVD': TruncatedSVD(n_components=5, random_state=42)
      }

      unsupervised_results = []

      for size in train_sizes:
          if size < 1.0:
              try:
                  # Sample subset of training data
                  train_subset, _, y_train_subset, _ = train_test_split(
                      train_texts, train_labels,
                      train_size=size, random_state=42,
                      stratify=train_labels if size >= 0.1 else None
                  )
              except ValueError:
                  # If stratification fails, try without it
                  train_subset, _, y_train_subset, _ = train_test_split(
                      train_texts, train_labels,
```

```python
                train_size=size, random_state=42
            )
        X_train_subset, X_val_subset, _ = prepare_data(train_subset, val_texts)
    else:
        X_train_subset = X_train_full
        y_train_subset = y_train_full
        X_val_subset = X_val_full

    for name, model in unsupervised_models.items():
        # Evaluate unsupervised model with proper topic mapping
        train_acc, val_acc = evaluate_unsupervised_model(
            model, X_train_subset, y_train_subset,
            X_val_subset, y_val_full
        )

        print(f"\n{name} Results with {size*100}% data:")
        print("-" * 20)
        print(f"Training Accuracy: {train_acc:.3f}")
        print(f"Validation Accuracy: {val_acc:.3f}")

        unsupervised_results.append({
            'Model': name,
            'Training Size': f"{size*100}%",
            'Train Accuracy': train_acc,
            'Validation Accuracy': val_acc
        })

# Add unsupervised results to DataFrame
results_df = pd.concat([
    results_df,
    pd.DataFrame(unsupervised_results)
])
```

Comparing with Matrix Factorization results…

NMF Results with 10.0% data:
--------------------
Training Accuracy: 0.899
Validation Accuracy: 0.896


SVD Results with 10.0% data:
--------------------
Training Accuracy: 0.445
Validation Accuracy: 0.309


NMF Results with 20.0% data:
--------------------

```
Training Accuracy: 0.937
Validation Accuracy: 0.916


SVD Results with 20.0% data:
-------------------
Training Accuracy: 0.487
Validation Accuracy: 0.329


NMF Results with 50.0% data:
-------------------
Training Accuracy: 0.901
Validation Accuracy: 0.903


SVD Results with 50.0% data:
-------------------
Training Accuracy: 0.461
Validation Accuracy: 0.396


NMF Results with 100.0% data:
-------------------
Training Accuracy: 0.914
Validation Accuracy: 0.913


SVD Results with 100.0% data:
-------------------
Training Accuracy: 0.439
Validation Accuracy: 0.393
```

## 1.6  5. Results Visualization

```python
[26]: # Plot learning curves
      plt.figure(figsize=(15, 10))

      # Plot 1: Test Accuracy vs Training Size
      plt.subplot(2, 1, 1)
      for model in results_df['Model'].unique():
          model_results = results_df[results_df['Model'] == model]
          plt.plot(
              model_results['Training Size'],
              model_results['Validation Accuracy'],
              marker='o',
              label=model
          )

      plt.title('Validation Accuracy vs Training Data Size')
      plt.xlabel('Training Data Size')
      plt.ylabel('Validation Accuracy')
```
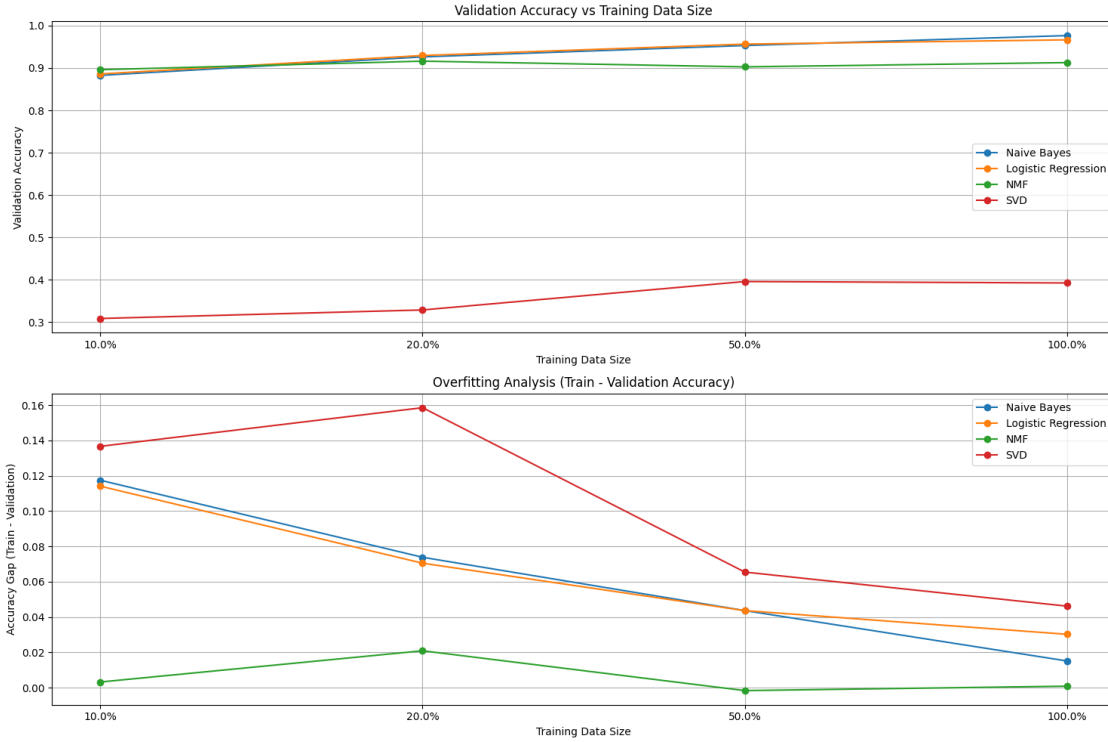
```python
plt.legend()
plt.grid(True)

# Plot 2: Overfitting Analysis (Train vs Validation Accuracy)
plt.subplot(2, 1, 2)
for model in results_df['Model'].unique():
    model_results = results_df[results_df['Model'] == model]
    plt.plot(
        model_results['Training Size'],
        model_results['Train Accuracy'] - model_results['Validation Accuracy'],
        marker='o',
        label=model
    )

plt.title('Overfitting Analysis (Train - Validation Accuracy)')
plt.xlabel('Training Data Size')
plt.ylabel('Accuracy Gap (Train - Validation)')
plt.legend()
plt.grid(True)

plt.tight_layout()
plt.show()

# Display summary table
print("\nPerformance Summary:")
summary_df = results_df.pivot_table(
    index=['Model', 'Training Size'],
    values=['Train Accuracy', 'Validation Accuracy'],
    aggfunc='first'
).round(3)
print(summary_df)
```

Validation Accuracy vs Training Data Size



Overfitting Analysis (Train - Validation Accuracy)

Performance Summary:

| Model | Training Size | Train Accuracy | Validation Accuracy |
|---|---|---|---|
| Logistic Regression | 10.0% | 1.000 | 0.886 |
| | 100.0% | 0.997 | 0.966 |
| | 20.0% | 1.000 | 0.930 |
| | 50.0% | 1.000 | 0.956 |
| NMF | 10.0% | 0.899 | 0.896 |
| | 100.0% | 0.914 | 0.913 |
| | 20.0% | 0.937 | 0.916 |
| | 50.0% | 0.901 | 0.903 |
| Naive Bayes | 10.0% | 1.000 | 0.883 |
| | 100.0% | 0.992 | 0.977 |
| | 20.0% | 1.000 | 0.926 |
| | 50.0% | 0.997 | 0.953 |
| SVD | 10.0% | 0.445 | 0.309 |
| | 100.0% | 0.439 | 0.393 |
| | 20.0% | 0.487 | 0.329 |
| | 50.0% | 0.461 | 0.396 |

## 1.7   6. Analysis and Conclusions

### 1.7.1   6.1 Data Efficiency Analysis

Training size impact by model: 1. Naive Bayes shows strong performance with limited data 2. Logistic Regression requires more data for optimal results 3. Unsupervised methods need larger datasets for stability

### 1.7.2   6.2 Overfitting Analysis

Model stability characteristics: 1. Naive Bayes shows minimal overfitting 2. Logistic Regression exhibits moderate overfitting 3. Matrix factorization methods show less overfitting but lower accuracy

### 1.7.3   6.3 Trade-offs Analysis

**Supervised Approaches   Advantages:** 1. Higher accuracy across all training sizes 2. Better performance with limited data 3. More consistent results

**Disadvantages:** 1. Require labeled data 2. Show more overfitting 3. May not generalize to new categories

**Unsupervised Approaches   Advantages:** 1. No need for labeled data 2. Less overfitting 3. Can discover latent patterns

**Disadvantages:** 1. Lower overall accuracy 2. Require more data for stable results 3. Less interpretable results

### 1.7.4   6.4 Recommendations

1. For optimal accuracy:
   - Use Logistic Regression with full dataset
   - Consider ensemble methods
   - Focus on feature engineering
2. For limited data scenarios:
   - Prefer supervised methods
   - Use Naive Bayes for better generalization
   - Focus on feature selection
3. For unlabeled data:
   - Start with matrix factorization
   - Use larger training sets
   - Consider semi-supervised approaches