# Natural Language Processing using Python Programming

## Notebook 08.1: Introduction to Text Classification and the ML Pipeline

`Python` `3.8+`  `NLTK` `Latest`  `SpaCy` `Latest`  `Scikit-learn` `Latest`  `License` `MIT`

---

**Part of the comprehensive learning series:** Natural Language Processing using Python Programming

**Learning Objectives:**

- Master the fundamentals of supervised text classification for multi-class problems
- Understand the standard ML pipeline workflow for text classification tasks
- Implement scikit-learn Pipeline objects for robust, production-ready workflows
- Learn best practices for managing vectorization and classification steps
- Build foundation for advanced classification algorithms and evaluation methods

---

- **Text Classification** is a supervised learning task where we assign predefined categories (labels) to text documents.

- This is the foundation for sentiment analysis, spam detection, topic labeling, and intent recognition.

- This notebook introduces the standard classification pipeline and, critically, the **Scikit-learn** `Pipeline` **object**, which simplifies and formalizes the workflow.

## 1. Setting up: Multi-Class Dataset

- We will use the **20 Newsgroups dataset**, a classic multi-class problem where documents are classified into 20 different topics (e.g., 'comp.graphics', 'rec.sport.baseball').

In [1]:
```python
# Import necessary libraries
# This code snippet demonstrates how to load and explore the 20 Newsgroups da
# a popular dataset for text classification tasks.
# It fetches a subset of the dataset for faster loading and prints out some b
import pandas as pd
from sklearn.datasets import fetch_20newsgroups
from sklearn.model_selection import train_test_split

# Fetching a subset of the 20 Newsgroups data for faster loading and demonstr
```

```python
categories = ['alt.atheism', 'soc.religion.christian', 'comp.graphics', 'rec.

newsgroups_train = fetch_20newsgroups(
    subset='train',                    # Using the training subset
    categories=categories,             # Categories to include
    shuffle=True,                      # Shuffling the data
    random_state=42                    # Reproducibility
)

# Splitting data into features and labels
# X contains the text data, and y contains the corresponding labels
X = newsgroups_train.data
y = newsgroups_train.target

print("20 Newsgroups Subset Loaded.")
print(f"Total documents: {len(X)}")
print(f"Number of classes: {len(newsgroups_train.target_names)}")
print(f"Example Class Names: {newsgroups_train.target_names}")
```

```
20 Newsgroups Subset Loaded.
Total documents: 2260
Number of classes: 4
Example Class Names: ['alt.atheism', 'comp.graphics', 'rec.sport.baseball', 's
oc.religion.christian']
```

## 2. The Text Classification Pipeline (Conceptual)

- The overall workflow remains consistent:

    1. **Data Split:** Separate $X$ (features/text) and $y$ (labels) into training and testing sets.

    2. **Feature Extraction:** Convert text to numerical vectors (e.g., TF-IDF).

    3. **Model Training:** Fit a classifier (e.g., Naive Bayes) to the vectors.

    4. **Prediction:** Use the model on the test vectors.

    5. **Evaluation:** Calculate performance metrics.

### Problem with Manual Steps:

- Manually managing the vectorizer ( `fit_transform` on train, `transform` on test) and ensuring consistency between steps is error-prone, especially during hyperparameter tuning.

---

## 3. Formalizing the Workflow with Scikit-learn `Pipeline`

- The **`Pipeline` object** chains multiple estimators into one.

- Crucially, it ensures that the data transformation (e.g., vectorization) is **fitted only on the training data** and automatically **applied to all subsequent data** (test data, cross-validation data).

```
In [2]:  # Import necessary libraries for building the pipeline
         from sklearn.pipeline import Pipeline                           # Pipeline cl
         from sklearn.feature_extraction.text import TfidfVectorizer     # TF-IDF Vect
         from sklearn.naive_bayes import MultinomialNB                   # Naive Bayes
         from sklearn.metrics import classification_report              # For evaluat

         # 1. Split the data
         X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, ran

         # 2. Define the Pipeline steps
         text_clf = Pipeline([
             ('tfidf', TfidfVectorizer()),  # Step 1: Feature Extraction
             ('clf', MultinomialNB()),      # Step 2: Classifier
         ])

         # 3. Train the Pipeline (The Pipeline automatically handles fit_transform ->
         text_clf.fit(X_train, y_train)

         # 4. Predict (The Pipeline automatically handles transform -> predict)
         predicted = text_clf.predict(X_test)

         print("Pipeline Trained Successfully.\n")
```

```
Pipeline Trained Successfully.
```

```
In [4]:  print("Example Predictions:")
         for i, (doc, pred_label) in enumerate(zip(X_test[:2], predicted[:2])):
             true_label = newsgroups_train.target_names[y_test[i]]
             predicted_class = newsgroups_train.target_names[pred_label]
             print(f" - Actual: {true_label:<20} | Predicted: {predicted_class}")
```

```
Example Predictions:
 - Actual: alt.atheism          | Predicted: alt.atheism
 - Actual: rec.sport.baseball   | Predicted: rec.sport.baseball
```

```
In [5]:  print("\n--- Classification Report (Preview) ---")
         print(classification_report(y_test, predicted, target_names=newsgroups_train.
```

```
--- Classification Report (Preview) ---
                        precision    recall  f1-score   support

          alt.atheism       1.00      0.74      0.85       122
        comp.graphics       0.99      0.95      0.97       141
    rec.sport.baseball       0.98      0.97      0.98       155
soc.religion.christian       0.78      0.99      0.87       147

             accuracy                           0.92       565
            macro avg       0.94      0.91      0.92       565
         weighted avg       0.93      0.92      0.92       565
```

**Advantage of the `Pipeline`**

- The `Pipeline` is an immutable, single object representing the entire process. This is vital for:

    - **Consistency:** Eliminates the risk of transforming test data incorrectly.

    - **Cross-Validation:** Simplifies tuning by allowing grid search over *both* vectorizer and classifier parameters simultaneously.

    - **Deployment:** The entire process (Vectorizer + Classifier) can be saved/loaded as one file ( `.pkl` file, Chapter 10.3), ready for production.

## 4. Summary and Next Steps

- We established the need for supervised text classification and, most importantly, introduced the **Scikit-learn `Pipeline`** as the best practice for managing the text ML workflow.

- We successfully trained a multi-class classifier.

- In the next notebook (**8.2**), we will use this `Pipeline` structure to compare different classification algorithms: **Logistic Regression, Naive Bayes, and Support Vector Machines**.

### Key Takeaways

- **Text Classification Fundamentals:** We mastered the supervised learning approach to text classification, understanding how to assign predefined categories to text documents.

- **Pipeline Architecture Mastery:** We learned the critical importance of scikit-learn's Pipeline object for managing complex ML workflows with consistency and reliability.

- **Multi-Class Implementation:** We successfully implemented a 4-class text classifier using the 20 Newsgroups dataset, demonstrating scalability beyond binary classification.

- **Production Best Practices:** We established the foundation for robust, deployment-ready text classification systems using standardized pipeline structures.

---

### *Next Notebook Preview*

- With pipeline fundamentals mastered, we're ready to explore **algorithm comparison** and performance optimization.

- The next notebook will compare multiple classification algorithms (Logistic Regression, Naive Bayes, SVM) within the pipeline framework for comprehensive performance analysis.

---

## About This Project

This notebook is part of the **Natural Language Processing using Python Programming for Beginners** repository - a comprehensive, beginner-friendly guide for mastering NLP using Python, NLTK, and SpaCy.

**Repository:** `NLP`

## Author

**Prakash Ukhalkar**

GitHub  prakash-ukhalkar

---

Built with ❤️ for the Python community