**Introduction to NLTK**

The Natural Language Toolkit (NLTK) is a Python library used for processing and analyzing natural language textual data.

**Advantages:**

Extensive Natural Language Processing (NLP) Capabilities: NLTK offers a rich set of natural language processing functionalities, including tokenization, part-of-speech tagging, named entity recognition, syntactic analysis, sentiment analysis, text classification, and more. This makes it a powerful tool for handling textual data.

Widespread Use in Education and Academia: NLTK is widely employed in educational and academic settings and has an abundance of tutorials and documentation available for learning and reference.

Open Source and Free: NLTK is open-source and freely available, enabling its free usage in various projects.

Community Support: NLTK has an active community where users can get support and contribute to its development.

**Disadvantages:**

Performance: NLTK might lack efficiency when handling large-scale textual data as certain functionalities could be relatively slow.

Steep Learning Curve: For beginners, NLTK might have a relatively steep learning curve due to its extensive functionalities and options, requiring some time to master.

Not Suitable for All NLP Tasks: Despite offering a wide range of NLP tools, NLTK might not be suitable for all types of NLP tasks. Some specific domain problems might require other specialized tools or libraries.

**NLTK excels in handling various natural language processing tasks, including but not limited to:**

Text analysis and mining

Text classification

Semantic analysis

Text generation

Information retrieval

Language models and machine translation

**Here is a classic example of using NLTK for text classification, focusing on sentiment analysis:**

Issue: Perform sentiment analysis on a set of movie reviews to determine whether the comments are positive or negative.

First, ensure that NLTK library is installed and necessary datasets and resources are downloaded. Then, you can follow these steps for sentiment analysis:

```
import nltk

nltk.download('movie_reviews')

nltk.download('punkt')


from nltk.corpus import movie_reviews

from nltk.tokenize import word_tokenize

from nltk.classify import NaiveBayesClassifier

from nltk.classify.util import accuracy


# Prepare the dataset

documents = [(list(word_tokenize(review)), category)

        for category in movie_reviews.categories()

        for review in movie_reviews.fileids(category)]


# Build feature extraction function

def extract_features(document):

    words = set(document)

    features = {}

    for word in word_features:

        features[f'contains({word})'] = (word in words)

    return features


# Build feature set
```

```python
all_words = nltk.FreqDist(w.lower() for w in movie_reviews.words())
word_features = list(all_words.keys())[:2000]
featuresets = [(extract_features(d), c) for (d, c) in documents]

# Split into training and testing sets
train_set = featuresets[:1500]
test_set = featuresets[1500:]

# Train the Naive Bayes classifier
classifier = NaiveBayesClassifier.train(train_set)

# Evaluate the classifier's performance
print("Accuracy:", accuracy(classifier, test_set))

# Perform sentiment analysis
text_to_analyze = "This movie was really great! I loved it."
words_to_analyze = word_tokenize(text_to_analyze)
features_to_analyze = extract_features(words_to_analyze)
sentiment = classifier.classify(features_to_analyze)
print("Sentiment:", sentiment)
```

The above code demonstrates how to use NLTK to build a Naive Bayes classifier for sentiment analysis. First, it prepares a dataset of movie reviews, then constructs a feature extraction function and feature set. Next, it divides the dataset into training and testing sets and proceeds to train the classifier. Finally, the classifier can be used for sentiment analysis on new text.