

# News Processing Script Documentation

## Purpose:

The purpose of this script is to fetch news articles from various RSS feeds, classify them into categories based on their titles and descriptions, store them in a SQLite database, and asynchronously process them using Celery.

## Components:

### 1. Feed Parsing and Data Extraction:

- The script utilizes the `feedparser` library to parse RSS feeds and extract relevant information from each article entry, including the title, description, link, published date, and source.

### 2. Category Classification:

- Articles are classified into categories based on keywords present in their titles or descriptions. The `classify\_category` function implements this logic by checking for specific keywords associated with different categories.

### 3. Database Storage:

- A SQLite database is used to store the fetched articles. The `store\_articles` function inserts articles into the database, ensuring uniqueness based on the combination of title and source. Duplicate articles are logged, but not re-inserted.

### 4. Task Queue and News Processing:

- Celery is employed to process articles asynchronously. The `process\_articles` task retrieves articles from the database that have not been classified yet and updates their categories using the `classify\_category` function.

### 5. Logging:

- The script logs informational messages to a file (*news\_processing.log*) using the Python `logging` module. This includes messages related to duplicate articles found during insertion.

## **6. Main Function:**

- The `main` function orchestrates the entire process by fetching feeds, storing articles in the database, and initiating asynchronous processing using Celery.

## **Design Choices:**

### **1. Asynchronous Processing with Celery:**

- Celery is chosen for asynchronous task execution to handle potentially time-consuming tasks like article processing without blocking the main execution flow.

### **2. SQLite Database:**

- SQLite is selected for its simplicity and portability, making it suitable for this lightweight application. It provides a self-contained database solution without requiring a separate database server.

### **3. Keyword-based Classification:**

- Categories are assigned to articles based on predefined keywords present in their titles or descriptions. This approach offers a straightforward method for classification, although it may lack nuance compared to more sophisticated natural language processing techniques.

### **4. Logging for Duplicate Articles:**

- Duplicate articles are detected during insertion into the database, and corresponding messages are logged for record-keeping purposes. This helps maintain data integrity and provides visibility into potential issues with the data sources.

### **5. External Libraries:**

- The script utilizes external libraries such as ***feedparser***, ***nltk***, and ***Celery*** to leverage existing functionality and streamline development. These libraries offer robust solutions for tasks like RSS feed parsing, text tokenization, and asynchronous task management.

#### ### Usage Instructions:

1. Ensure that all required dependencies (***feedparser***, ***nltk***, ***Celery***, ***redis***) are installed.
2. Start a Redis server locally, as Celery is configured to use Redis as its message broker.
3. Run the script (***news\_processor.py***) using a Python interpreter.
4. Monitor the progress and logs to track the fetching, storing, and processing of news articles.

#### Future Enhancements:

##### 1. ***Improved Category Classification:***

- Enhance the category classification logic using machine learning or natural language processing techniques for more accurate and nuanced classification.

##### 2. ***Enhanced Error Handling:***

- Implement comprehensive error handling to gracefully handle exceptions and edge cases, ensuring robustness and reliability.

##### 3. ***Performance Optimization:***

- Optimize the code for improved performance, especially for large-scale data processing, by implementing caching mechanisms or parallel processing techniques.

##### 4. ***User Interface:***

- Develop a user interface or command-line interface to provide users with more control over the fetching, storing, and processing of news articles.

This documentation provides an overview of the implemented logic, design choices, and potential areas for future enhancements in the news processing script. If you have any further questions or need additional clarification, feel free to ask!