

# Detailed Documentation for Log Categorization Script using LLM

## Introduction

This script processes log files, categorizes log entries using a pre-trained BERT model, and saves the structured output to a CSV file. The script is designed to handle multiple log formats and includes error handling to manage unexpected formats or missing data.

## Libraries Used

1. `argparse`: A library for parsing command-line arguments. It allows the script to accept input file paths and model paths as arguments.
2. `pandas`: A powerful data manipulation library used for handling data in a DataFrame format.
3. `torch`: The PyTorch library, which provides the tools for building and training machine learning models.
4. `transformers`: A library by Hugging Face that provides pre-trained transformer models like BERT for natural language processing tasks.
5. `sys`: A library that provides access to some variables used or maintained by the interpreter and to functions that interact with the interpreter.
6. `re`: A library for regular expressions, used for parsing log entries from text.

## Function Descriptions

### 1. `load_model(model_path)`

**Purpose:** This function loads a pre-trained BERT model and its tokenizer from the specified path.

**Parameters:**

- `model_path` (str): The path to the directory containing the pre-trained model and tokenizer.

**Returns:**

- `model`: The loaded BERT model for sequence classification.
- `tokenizer`: The tokenizer corresponding to the BERT model.

**2. parse\_log\_file(file\_path)**

**Purpose:** This function reads and parses the log file, extracting log entries in different formats and structuring them into a DataFrame.

**Parameters:**

- `file_path` (str): The path to the log file.

**Returns:**

- `df` (DataFrame): A pandas DataFrame containing parsed log entries with columns `timestamp`, `message_id`, and `readable_message`.

**Process:**

- Reads the log file content line by line.
- Uses regular expressions to match and extract log entries in different formats.
- Extracts timestamps, message IDs, and readable messages.
- Stores the extracted information in a list of dictionaries, which is then converted to a DataFrame.

**Error Handling:**

- Raises a `ValueError` if no valid log entries are found.

**3. process\_csv(file\_path, tokenizer)**

**Purpose:** This function processes the log entries by tokenizing the readable\_message using the provided tokenizer.

**Parameters:**

- file\_path (str): The path to the log file.
- tokenizer: The tokenizer for encoding log messages.

**Returns:**

- df (DataFrame): A DataFrame with an additional column tokens containing tokenized log messages.

**Process:**

- Calls parse\_log\_file to get the initial DataFrame.
- Applies the tokenizer to the readable\_message column to create a tokens column.

**Error Handling:**

- Catches and raises exceptions related to missing columns or other issues during processing.

#### **4. infer(model, tokenizer, df)**

**Purpose:** This function performs inference using the BERT model to categorize log entries.

**Parameters:**

- model: The pre-trained BERT model for sequence classification.
- tokenizer: The tokenizer for encoding log messages.
- df (DataFrame): The DataFrame containing tokenized log messages.

**Returns:**

- df (DataFrame): The DataFrame with an additional column category containing predicted categories.

**Process:**

- Converts the tokenized messages into a format suitable for the model.
- Uses the model to predict categories for each log message.
- Adds the predictions as a new column category in the DataFrame.

## **5. human\_readable\_output(df)**

**Purpose:** This function converts numerical categories into human-readable labels (Error or Non-Error).

**Parameters:**

- df (DataFrame): The DataFrame containing predicted categories.

**Returns:**

- df (DataFrame): The DataFrame with human-readable categories.

**Process:**

- Maps numerical predictions (0 or 1) to Non-Error and Error.

## **6. main(csv\_file=None, model\_path=None)**

**Purpose:** This is the main function that orchestrates the entire process from loading the model, processing the log file, performing inference, and saving the output to a CSV file.

**Parameters:**

- csv\_file (str, optional): The path to the log file. Default is None.
- model\_path (str, optional): The path to the pre-trained model. Default is None.

**Process:**

- If csv\_file and model\_path are provided, it processes the log file, categorizes the logs, and saves the output.
- If not provided, it parses command-line arguments to get the file paths.
- Includes error handling to catch and print exceptions.

**Command-line Arguments:**

- `csv_file`: Path to the log file.
- `model_path`: Path to the pre-trained model.

## Execution

### Interactive Mode:

- The script checks if it's running in an interactive environment (like Jupyter) and uses default paths if so.

### Command-line Mode:

- The script can be executed from the command line, accepting file paths as arguments.
- Command line to run this script if you are running in the python interpreter...

`“python log_categorizer.py /path/to/log_file.txt /path/to/model_folder ”`