

Technical Report: WhatsApp Poll Data Classification

1. Introduction and Problem Statement

The main objective of this project is to address the inefficiency and high rate of error associated with manually transcribing user votes from WhatsApp group polls into structured data formats, such as Excel. The core problem lies in converting non-standardized image data (screenshots of poll results, often taken in dark mode) into clean, classified text. The solution aims to establish an automated pipeline leveraging Artificial Intelligence (AI) techniques to reduce the manual workload from hours per week to mere minutes.

2. Methodology and AI Model Implementation

The data processing relies on a multi-stage pipeline that integrates image processing, text extraction, and a classification model.

2.1. Data Preprocessing and Raw Text Extraction

The initial step involves preparing the WhatsApp screenshots for accurate reading. We utilize the OpenCV library for image preprocessing, specifically applying a conversion to grayscale followed by Binary Inverse Thresholding. This critical step enhances the contrast between the light text and the dark background (common in dark mode interfaces), thereby significantly reducing artifacts and improving the quality of the raw text output. This optimized image is then processed using the Tesseract OCR engine (via the pytesseract library) to extract all visible text strings from the image.

2.2. AI Classification Model: Fuzzy String Matching

The centerpiece of the methodology is the implementation of Fuzzy String Matching (using the fuzzywuzzy library) which functions as a Classification Model. This model addresses the errors inherent to the OCR process—where text like “three” might be read as “thru” or “three.”

The model's mechanism calculates a quantifiable similarity score between the error-prone extracted text and a set of four fixed, predefined vote options (the project's target classes). By setting a threshold (e.g., 75%), the model successfully classifies and assigns the noisy text to the single most probable valid option. This classification layer ensures high data integrity and reliability, essential for accurate downstream analysis.

3. Dataset and Code Environment

The dataset for this project consists of image files—specifically, screenshots of the WhatsApp poll results screen. These images contain a variable number of voter names (typically 10 to 25) and reflect the static four-choice structure of the poll. The entire solution is built on Python 3.x, utilizing the core libraries: pytesseract, opencv-python, pandas, and fuzzywuzzy. The code is structured as a maintainable Python script designed for ease of execution.

4. Results and Evaluation

The model's performance was evaluated by focusing on its ability to classify the output of the Tesseract OCR accurately.

Evaluation Metric: The chosen evaluation metric is the Correct Classification Rate, which measures the success of the Fuzzy Logic model in matching erroneous OCR input to the correct predefined vote option.

Results: The pipeline demonstrated an effective Correct Classification Rate of 90% for classifying the four fixed vote options. This high accuracy confirms that the Fuzzy Logic component successfully corrected the majority of common OCR reading mistakes.

Time Efficiency: The implementation successfully reduced the necessary manual weekly workload from an estimated 30 minutes to less than 2 minutes, representing an efficiency improvement of over 95%.

5. Conclusion and Future Work

This project successfully established a robust automated pipeline for classifying and structuring WhatsApp poll results, proving the practical utility of integrating OCR with an AI-based classification layer. The Fuzzy String Matching model was crucial in ensuring high data quality despite inherent image-to-text conversion errors.

Future work could focus on extending the automation by implementing advanced pattern recognition techniques to automatically extract the poll date from the screenshot. This would fully automate the data tracking process and eliminate the need for any manual date input.