

The provided Jupyter Notebook outlines a thorough and structured approach to analyzing the changes between initial and final versions of a dataset containing questions and answers. The analysis begins with the setup and data loading, where essential libraries such as `pandas`, `numpy`, `matplotlib`, `nltk`, `spacy`, `textblob`, `gensim`, and `sklearn` are imported. These libraries are crucial for data manipulation, visualization, and natural language processing (NLP) tasks. The initial and final datasets are loaded from CSV files, setting a strong foundation for the subsequent analysis.

In the data cleaning section, irrelevant columns like `Original ID` and `identifier` are dropped, streamlining the dataset. Missing values in `section_heading` and `control_heading` are forward-filled to maintain data continuity, while rows with missing `answer` values are removed to ensure data integrity. A check for duplicates confirms that the datasets contain unique entries. This thorough cleaning process ensures a clean and consistent dataset, ready for in-depth analysis.

The exploratory data analysis (EDA) phase combines questions and answers into single text strings for both initial and final versions. This holistic view of the text allows for a comprehensive word frequency analysis, identifying the most common words and their frequencies. This step provides an initial understanding of the dataset's themes and vocabulary, setting the stage for deeper analysis.

Text preprocessing involves tokenization, lemmatization, and the removal of stop words. These steps standardize the text, making it easier to analyze by focusing on meaningful words and removing common but uninformative ones. This standardization is essential for improving the accuracy of subsequent analyses.

The text analysis section reveals several insights. A comparison of word frequencies shows slight changes, indicating minor content adjustments. For instance, the most common words in the dataset might have slightly shifted in their frequency of occurrence, suggesting subtle changes in the focus or terminology used. An analysis of answer lengths reveals differences between the initial and final versions, suggesting content refinement. For example, the average length of answers might have increased or decreased, indicating a change in the level of detail provided. The cosine similarity score of 0.9994 indicates minimal changes in the overall text content, implying that the core content remains largely unchanged. Sentiment analysis, performed using `TextBlob`, shows variations in sentiment scores, revealing subtle shifts in text tone or perspective. These findings suggest that while the overall content remains similar, there are nuanced changes in word usage, answer length, and sentiment.

Advanced NLP techniques further confirm these subtle changes. Text complexity analysis, including metrics like average word length, reveals efforts to modify readability. For instance, a decrease in average word length might suggest an attempt to make the text more accessible. Named Entity Recognition (NER) using `spaCy` extracts named entities, which remain largely consistent between versions, highlighting stable core content. High Jaccard and semantic

similarity scores reaffirm minor changes in text content, indicating that the vocabulary and semantics have not significantly diverged. These advanced techniques provide a deeper understanding of the text, confirming subtle refinements while maintaining core content consistency.

Feature engineering extracts various features, such as the number of words, sentences, entities, and noun chunks, providing a detailed comparison. Visual differences in these features between initial and final versions highlight specific areas of content modification. For example, an increase in the number of noun chunks might indicate a more descriptive text, while changes in the number of entities could reflect a shift in focus. This granular view of textual changes reveals specific areas of refinement, offering deeper insights into the dataset's evolution.

Dimensionality reduction and visualization techniques effectively illustrate the high degree of textual similarity and stable thematic content. Topic modeling using Latent Dirichlet Allocation (LDA) identifies consistent topics across versions, suggesting stable thematic content. For instance, the primary topics in both versions might revolve around similar themes, such as technology or education. Truncated Singular Value Decomposition (SVD) visualizes document similarities in 2D space, offering an intuitive understanding of textual changes. Semantic similarity using sentence transformers is visualized through heatmaps and boxplots, confirming the high similarity between initial and final texts. These visualizations provide a clear, intuitive representation of the text's consistency and minor changes.

In conclusion, the notebook presents a detailed and methodical approach to comparing textual data across two versions. The findings reveal high consistency in text, with minor changes in word usage, answer length, and sentiment. Changes in text complexity metrics suggest efforts to improve readability, while consistent topics across versions indicate stable thematic content. These insights offer a framework for assessing textual changes, useful in content review and quality assurance processes. The comprehensive analysis and critical insights foster a deeper understanding of the nuances in textual data, offering a rich basis for further exploration and discussion. The notebook's methodologies and findings provide valuable contributions to the field of NLP and content analysis, demonstrating the application of various techniques for text analysis and highlighting the importance of textual refinement and quality assurance in content development.