

HW2 Report

I. Introduction

This report presents the results of mining frequent itemsets from a dataset of tweets related to flu shots collected in the year 2014. Utilizing the Apriori algorithm, the study aims to uncover common patterns and associations within public discussions about flu vaccinations on social media.

II. Methodology

Dataset Description

The dataset comprises tweets containing various keywords associated with flu shots. Each row in the dataset represents a tweet, with the "text_keywords" attribute capturing the essence of each tweet through a set of keywords.

Apriori Algorithm Implementation

The Apriori algorithm was implemented to identify frequent itemsets among the keywords in the dataset. This method iteratively finds groups of items (keywords) that appear together in tweets equally or more frequently than a specified threshold, min_support.

III. Minimum Support Threshold

A minimum support threshold of 1000 was chosen for this analysis. This threshold was selected to ensure that only the most significant and prevalent itemsets were considered, reducing noise and focusing on patterns that represent widespread topics of discussion. Support counts lower than this, would include transactions that hold no clear significance, while a minimum support higher than this would prune too many itemsets.

IV. Algorithmic Optimizations

To enhance the efficiency of the Apriori algorithm, several optimizations were employed:

- Transaction Reduction: Transactions not containing any candidate itemsets for subsequent iterations were excluded from further analysis.
- Efficient Counting: A hashing technique was utilized to count occurrences of itemsets efficiently, minimizing computational overhead.

- Pruning: Itemsets unlikely to meet the minimum support in future iterations were pruned early in the process.

V. Experiences and Lessons Learned

Challenges encountered included handling large datasets and ensuring the algorithm's efficiency. Overcoming these obstacles required careful optimization and a deep understanding of the algorithm's workings. This project underscored the importance of preprocessing and optimizing data analysis algorithms for handling big data effectively.

VI. Analysis of Results

The analysis revealed several key itemsets, including "flu", "shot", "flu shot", "got shot", "get shot", "flu shot today", and combinations thereof. Notably, the itemset "flu shot" with a support count of 23150 indicates a high level of discussion around flu vaccinations. Combinations with "get" and "got" highlight actions or intentions related to receiving flu shots, reflecting public sentiment towards obtaining a flu shot.

VII. Knowledge Gained

The frequent itemsets identified through this analysis offer insights into public perceptions and behaviors regarding flu shots. The prominence of itemsets related to receiving the flu shot suggests a proactive stance toward flu prevention among the Twitter population.

VIII. Conclusion

This project demonstrated the use of the Apriori algorithm in extracting meaningful patterns from social media data. By analyzing discussions around flu shots on Twitter, we gained valuable insights into public attitudes toward flu vaccinations. The experiences and lessons learned from this assignment underscore the potential of data mining techniques in understanding and leveraging social media discourse for public health purposes.