

Chapter 47


Improvisation of Cleaning Process on Tweets for Opinion Mining

Arpita Grover

 <https://orcid.org/0000-0001-5273-686X>

Kurukshetra University, India

Pardeep Kumar

 <https://orcid.org/0000-0003-3755-1837>

Kurukshetra University, Kurukshetra, India

Kanwal Garg

Kurukshetra University, Kurukshetra, India

ABSTRACT

In the current scenario, high accessibility to computational facilities encourage generation of a large volume of electronic data. Expansion of the data has persuaded researchers towards critical analyzation so as to extract the maximum possible patterns for wiser decisiveness. Such analysis requires curtailing of text to a better structured format by pre-processing. This scrutiny focuses on implementing pre-processing in two major steps for textual data generated by dint of Twitter API. A NoSQL, document-based database named as MongoDB is used for accumulating raw data. Thereafter, cleaning followed by data transformation is executed on accumulated tweets related to Narendra Modi, Honorable Prime Minister of India.

This chapter published as an Open Access Chapter distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>) which permits unrestricted use, distribution, and production in any medium, provided the author of the original work and original publication source are properly credited.

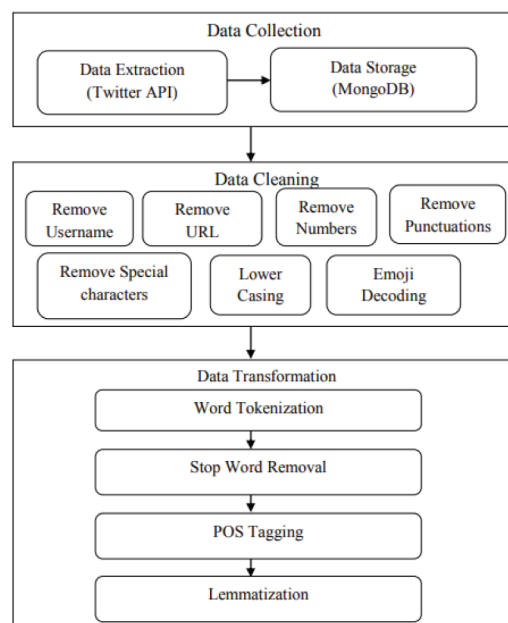
1. INTRODUCTION

Social media brings people together so that they can generate ideas or share their experiences with each other. The information generated through such sites can be utilized in many ways to discover fruitful patterns. But, accumulation of data via such sources create a huge unstructured textual data with numerous unwanted formats. Henceforth, the first step of text mining involves pre-processing of gathered reviews.

The journey of transforming dataset into a form, an algorithm may digest, takes a complicated road. The task embraces four differentiable phases: Cleaning, Annotation, Normalization and Analysis. The step of cleaning comprehends extrication of worthless text, tackling with capitalization and other similar details. Stop words, Punctuations marks, URLs, numbers are some of the instances which can be discarded at this phase. Annotation is a step of applying some scheme over text. In context to natural language processing, this includes part-of-speech tagging. Normalization demonstrates reduction of linguistic. In other words, it is a process that maps terms to a scheme. Basically, standardization of text through lemmatization and stemming are the part of normalization. Finally, text undergoes manipulation, generalization and statistical probing to interpret features.

For this study, pre-processing is accomplished in three major steps, as signified in Figure 1, keeping process of sentiment analysis in consideration. Foremost step included collection of tweets from Twitter by means of Twitter API. Captured data is then stored in a NoSQL database known to be MongoDB. Thereafter, collected tweets underwent cleaning (Zainol et al., 2018) process. Cleaning phase incorporated removal of user name, URLs, numbers, punctuations, special characters along in addition to lower casing and emoji decoding. The first two phases of data collection and cleaning were demonstrated in previous research. Also, it was shown that application of cleaning process still left data with anomalies and that is why the endmost stage of data transformation is introduced in this research. Data transfor-

Figure 1. Preprocessing steps



8 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/improvisation-of-cleaning-process-on-tweets-for-opinion-mining/308525

Related Content

Resilient Supply Chains to Improve the Integrity of Accounting Data in Financial Institutions Worldwide Using Blockchain Technology

Yu Yang and Zecheng Yin (2023). *International Journal of Data Warehousing and Mining* (pp. 1-20).

www.irma-international.org/article/resilient-supply-chains-to-improve-the-integrity-of-accounting-data-in-financial-institutions-worldwide-using-blockchain-technology/320648

Data Mining-Driven Detection of Banking URL Phishing Using Hybrid CNN and Machine Learning Models

Karthik Vanna, Mosiur Rahaman, Akshat Gaurav, Varsha Arya, Razaz Waheeb Attar, Brij B. Gupta and Kwok Tai Chui (2026). *International Journal of Data Warehousing and Mining* (pp. 1-17).

www.irma-international.org/article/data-mining-driven-detection-of-banking-url-phishing-using-hybrid-cnn-and-machine-learning-models/411870

Large-Scale System for Social Media Data Warehousing: The Case of Twitter-Related Drug Abuse Events Integration

Jenhani Ferdaous and Mohamed Salah Gouider (2022). *International Journal of Data Warehousing and Mining* (pp. 1-18).

www.irma-international.org/article/large-scale-system-for-social-media-data-warehousing/290890

Evolution of Spatial Data Templates for Object Classification

Neil Dunstan and Michael de Raadt (2002). *Data Mining: A Heuristic Approach* (pp. 143-156).

www.irma-international.org/chapter/evolution-spatial-data-templates-object/7587

Dynamic View Management System for Query Prediction to View Materialization

Negin Daneshpour and Ahmad Abdollahzadeh Barfouroush (2013). *Developments in Data Extraction, Management, and Analysis* (pp. 132-161).

www.irma-international.org/chapter/dynamic-view-management-system-query/70796