# Text Analytics of Course Reviews on Coursera Platform

Chan Huan Yang, Ramindhran Raja Mohan

*School of Computer Science, University Science of Malaysia, Penang, Malaysia*

## Introduction

We are now entering a new era - the revolution of online learning. From working professionals to recent high school graduates, many of them have found the reasons to take all or some of their courses online in platform such as Coursera, Udemy, and Edx

## Problem Statement

Ratings and reviews are always the major consideration factor by online course seekers before they joining the course. However, it can be time-consuming to read all the information especially the course reviews.
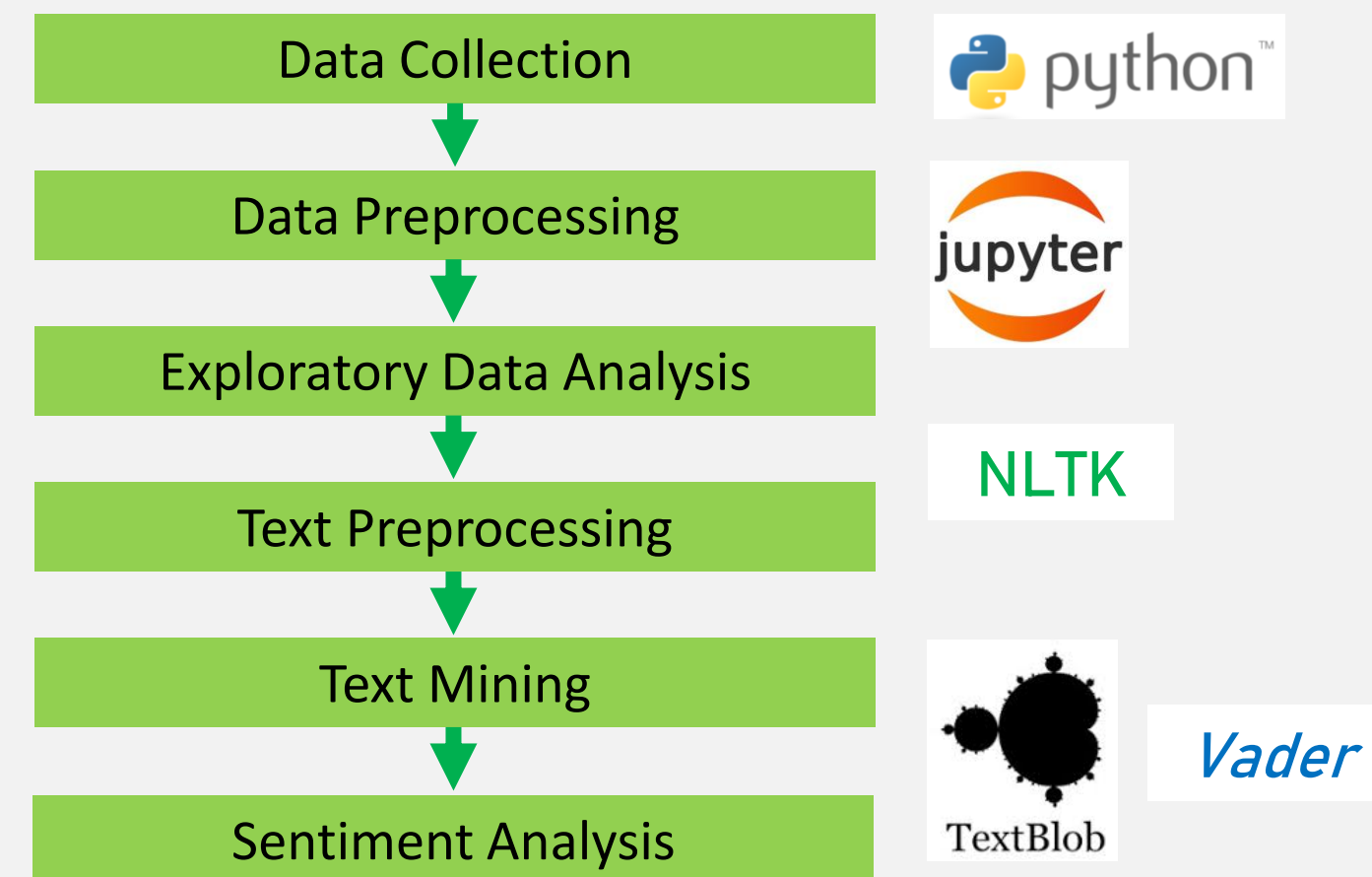
## Research Question

The research questions for this work were:
- How text analytics techniques such as n-gram analysis, word cloud, and sentiment analysis can be applied to improve the online course searching process?
- What insights can be obtained by using text analytics techniques such as n-gram analysis, word cloud, and sentiment analysis?

## Purpose of The Study

Our objective is to propose a text analytics pipeline that includes text cleaning, text lemmatization, sentiment analysis, text mining, and visualization that can help course seekers to gain a quick insight into the courses as well as enables them to make a quick comparison between multiple courses

## Research Method



## Data Collection

The data used in this work is from Kaggle.

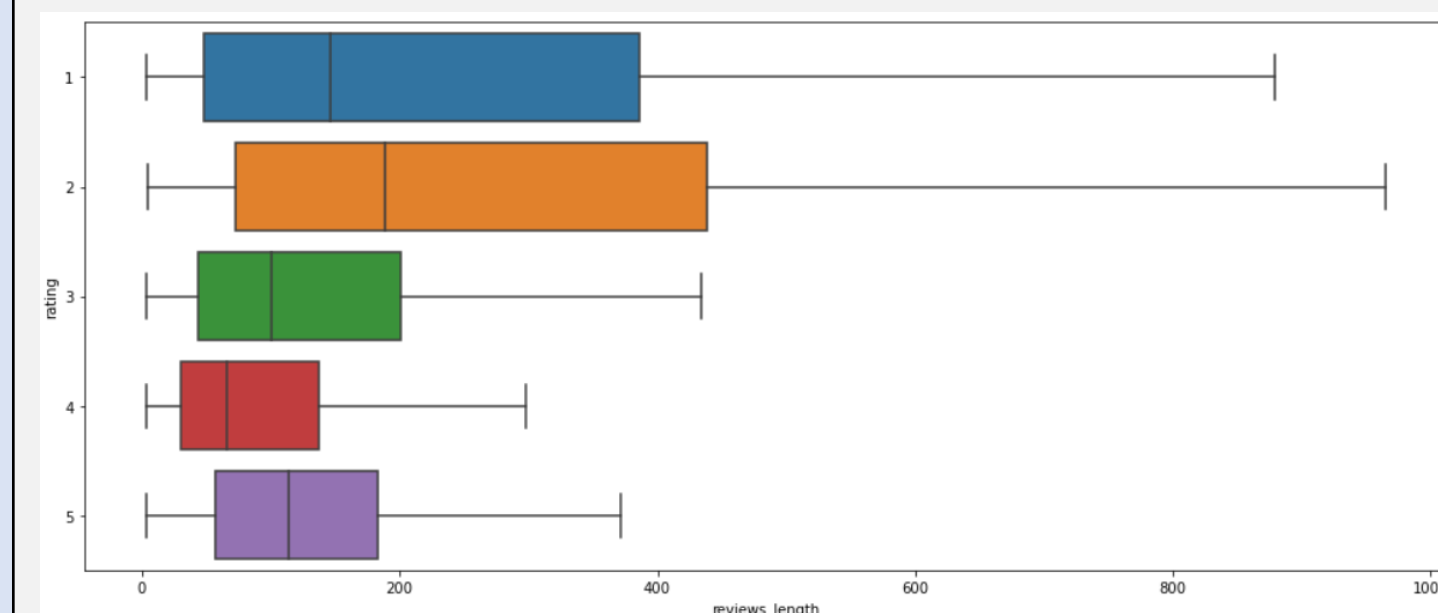| Course Name | URL | No of Review |
|---|---|---|
| Programming for Everybody (Getting Started with Python) | https://www.coursera.org/learn/python | 45218 |
| Python Data Structures | https://www.coursera.org/learn/python-data | 33543 |
| Introduction to Data Science in Python | https://www.coursera.org/learn/python-data-analysis | 14289 |

## Data Preprocessing
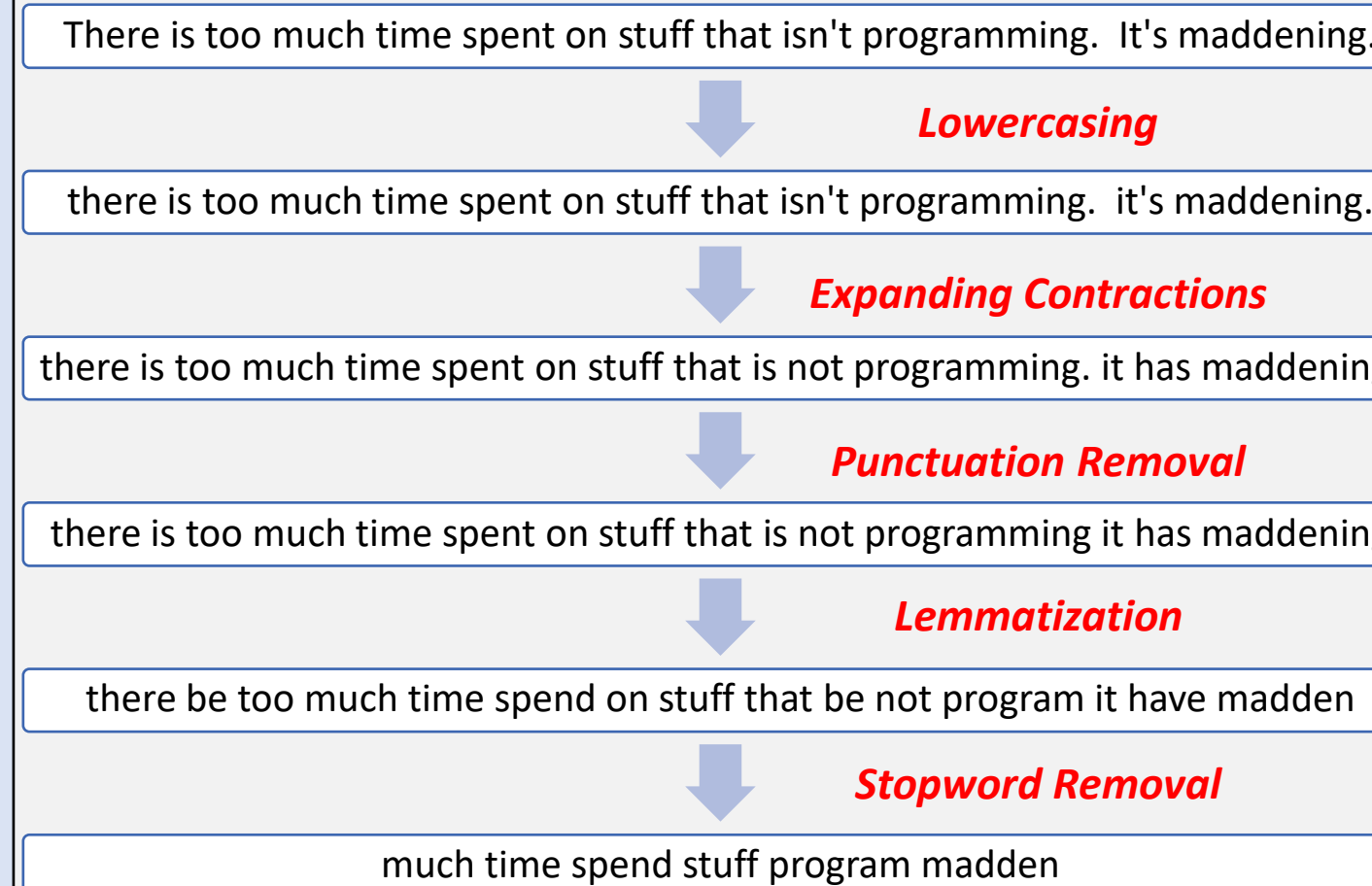
Data preprocessing steps include:
- removed the duplicate reviews
- removed the reviews with string's length less than three
- selected the English labeled reviews only using package "langid"

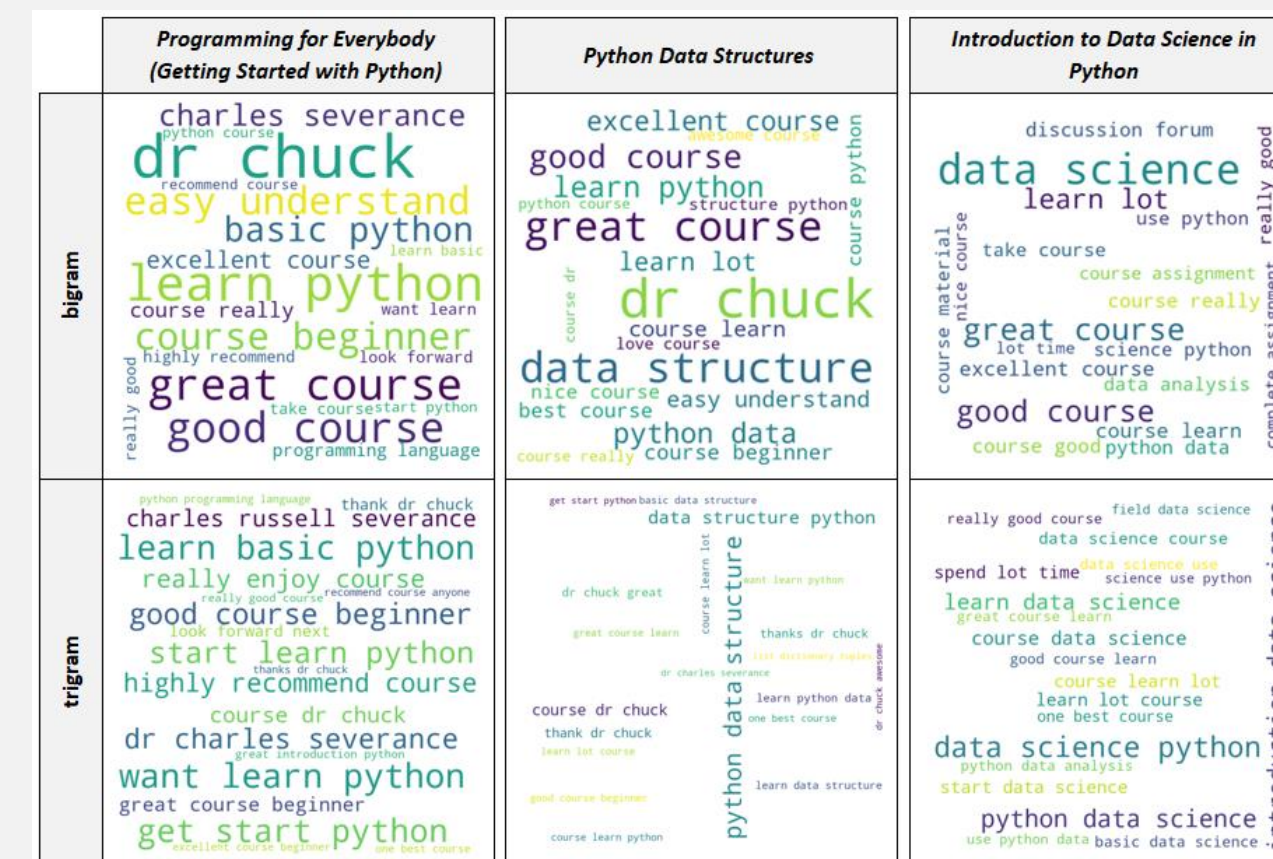## Exploratory Data Analysis

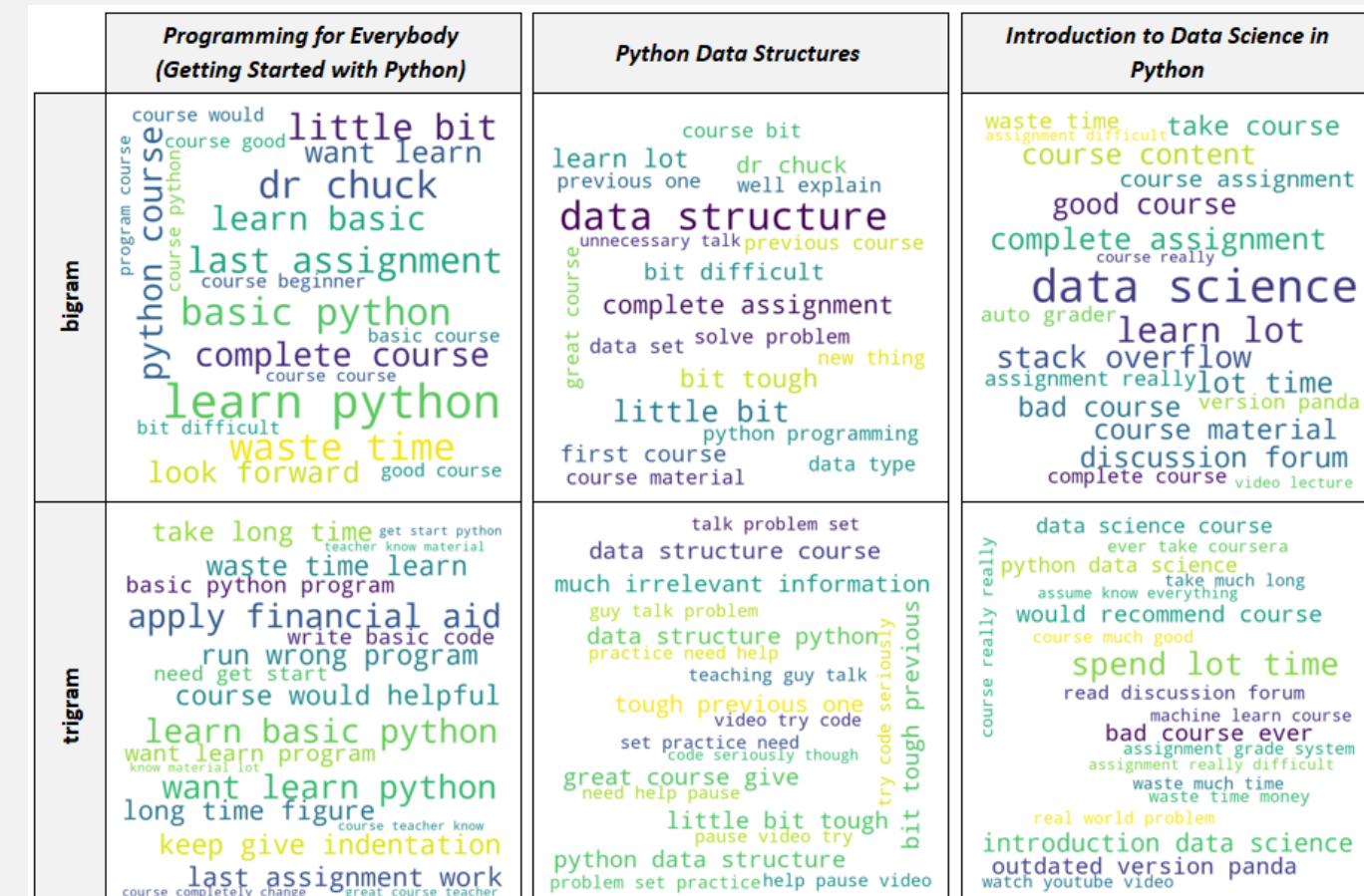### Reviews length of each rating



## Text Preprocessing

There is too much time spent on stuff that isn't programming. It's maddening.

⬇ *Lowercasing*

there is too much time spent on stuff that isn't programming.  it's maddening.

⬇ *Expanding Contractions*

there is too much time spent on stuff that is not programming. it has maddening.

⬇ *Punctuation Removal*

there is too much time spent on stuff that is not programming it has maddening

⬇ *Lemmatization*

there be too much time spend on stuff that be not program it have madden

⬇ *Stopword Removal*

much time spend stuff program madden

## Text Mining

### Word cloud of overall reviews



## Sentiment Analysis

| name | Average of Polarity | | |
|---|---|---|---|
| | Textblob | Vader | Overall |
| Programming for Everybody (Getting Started with Python) | 0.40 | 0.68 | 0.54 |
| Python Data Structures | 0.49 | 0.62 | 0.55 |
| Introduction to Data Science in Python | 0.34 | 0.45 | 0.40 |

### Word cloud of negative reviews



## Discussion and Finding

### Evaluation of accuracy and usability of n-gram



### Evaluation of sentiment polarity score

*Sentiment Polarity Score vs Actual Rating*



### Confusion matrix and classification report

| | Predicted negative | Predicted positive |
|---|---|---|
| Observed negative | 326 | 385 |
| Observed positive | 772 | 28324 |

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| negative | 30% | 46% | 36% | 711 |
| positive | 99% | 97% | 98% | 29096 |
| | | | | |
| accuracy | | 96% | | 29807 |
| macro average | 64% | 72% | 67% | 29807 |
| weightage avg | 97% | 96% | 97% | 29807 |

### Example review from the dataset

★ By Aayush D • Dec 7, 2018

Too easy of a course. completed in a day without much effort... And didn't really get as much out of it as I thought I would.

⬇

*"easy course complete day without much effort really get much think would"*

| Textblob | = 0.17 (positive) | |
| Vader | = 0.44 (positive) | ✗ |
| Actual | = 1 star (negative) | |

## Conclusion

The n-gram analysis and word cloud are sufficient enough to provide an accurate and informative glance into the course. However, it falls short on sentiment analysis especially in detecting the negative reviews.