

The background image shows several men in Saudi Arabia, wearing traditional white thobes and red-and-white checkered ghutras. They are holding and reading newspapers. The scene is slightly blurred, focusing on the text overlay. The newspapers have Arabic text and some images on them.

Saudi Newspapers Articles **Topic Extraction**

Rawabi Alharbi

01

Introduction

05

PREPROCESSING

02

Methodology

06

Topic Modeling

03

Tools

07

EDA

04

EDA

08

Conclusion



INTRODUCTION

- Newspapers Articles
- Problem statement



INTRODUCTION

- **Dataset**
 - Downloaded from Github
 - **31,030** Arabic newspaper articles

Methodology



Preprocessing

EDA

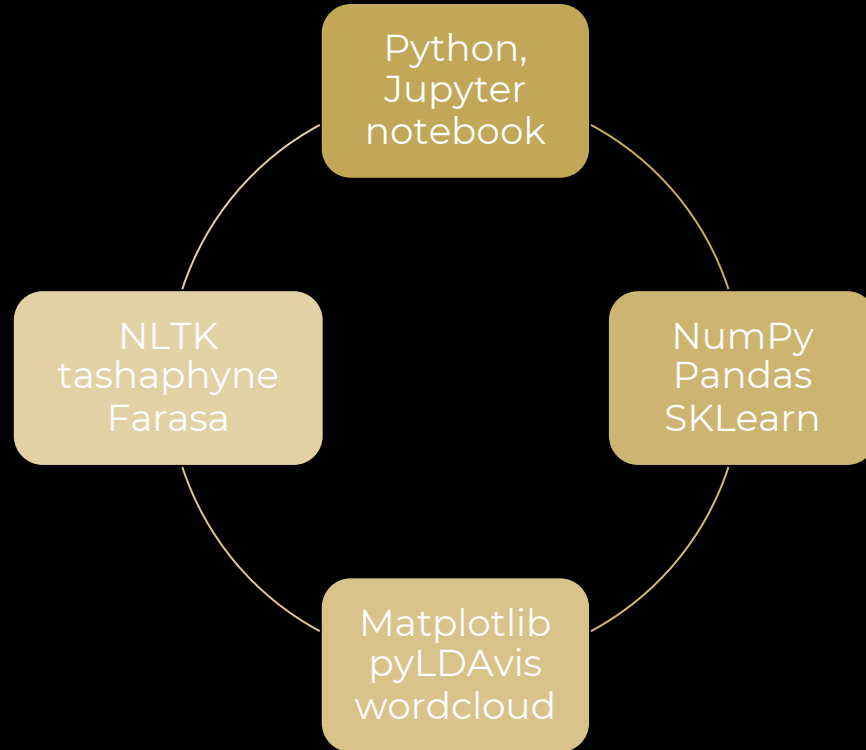


Topic Modeling

EDA



Tools



EDA

Word Cloud of the most frequent words



NLP Preprocessing

- **Remove**

- English letters
- English Numbers
- Special characters
- Arabic Punctuations
- Newline characters



NLP Preprocessing

- **Stemming**

- ISRIStemmer ✗
- ArabicLightStemmer ✗
- FarasaStemmer



WAITING



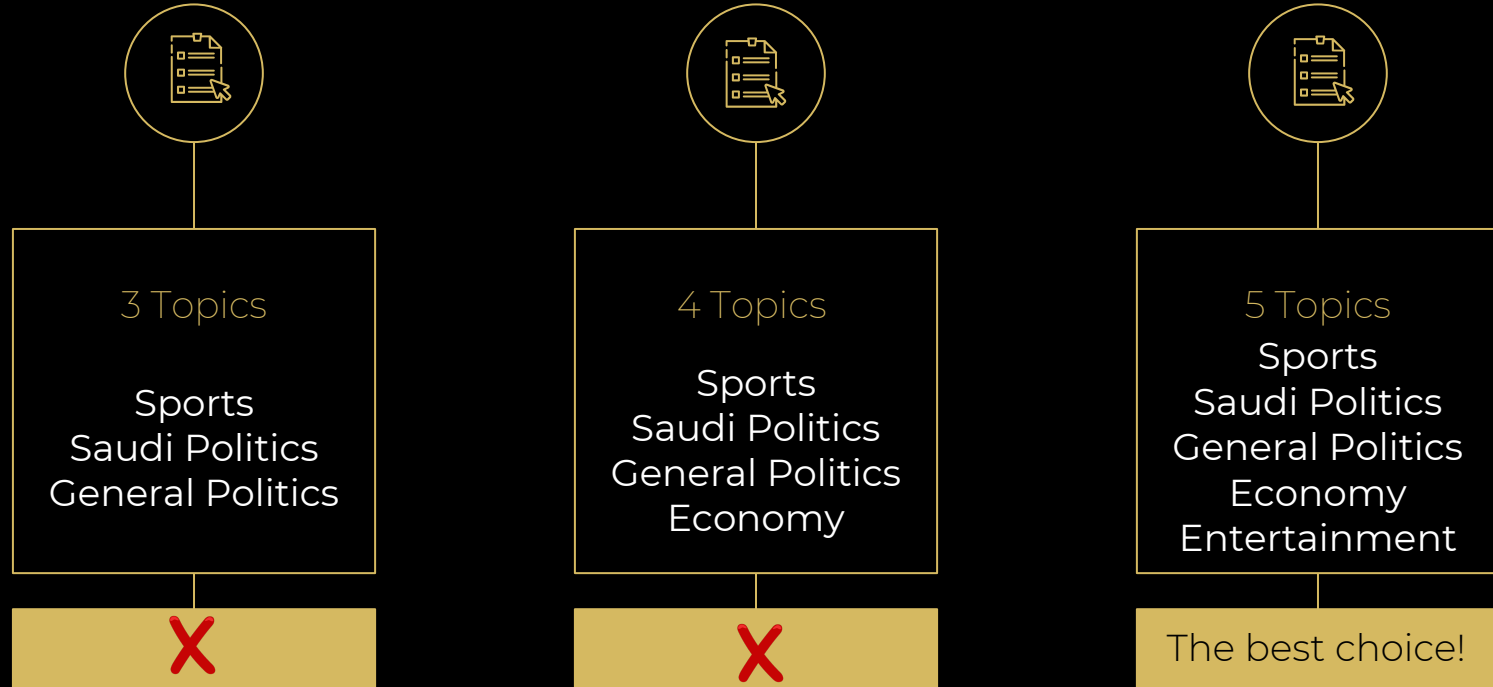
NLP Preprocessing

- **TF-IDF Vectorizer**
 - Remove Arabic stop words

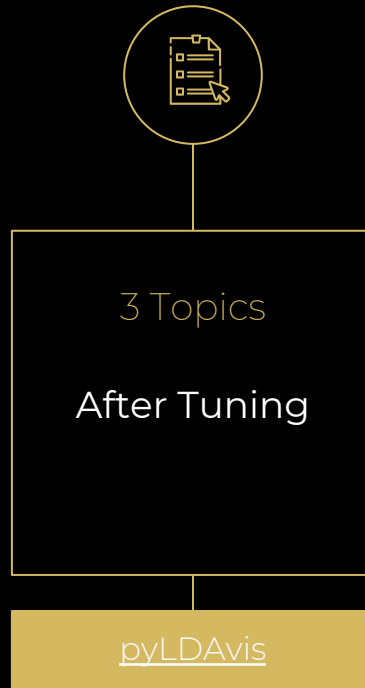


Topic Modeling - NMF

Tuning



Topic Modeling - LDA



EDA – After Topic Modeling

Word Cloud of the most frequent words in Topic 1



الرياضة

EDA – After Topic Modeling

Word Cloud of the most frequent words in Topic 2



السياسة السعودية

EDA – After Topic Modeling

Word Cloud of the most frequent words in Topic 3



السياسة العامة

EDA – After Topic Modeling

Word Cloud of the most frequent words in Topic 4



الاقتصاد

EDA – After Topic Modeling

Word Cloud of the most frequent words in Topic 5



السياحة/الترفيه



Result

5 topics were extracted from the articles

- Sports
- Saudi Politics
- General Politics
- Economy
- Entertainment



Conclusion

The model works with good performance to extract the most meaningful topics from the Saudi newspapers.

Future work

- Apply a good stemming & lemmatization
- Apply clustering

Thank
You!

