

الشعر العربي

NLP project



Presented by:
Mada Abudahish
Mozah AlKhaldi

Workflow

introduction

01

Tools

02

EDA

03

Text
prepressing

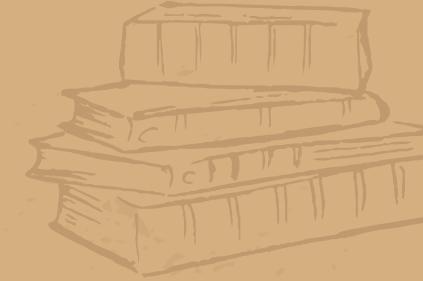
04

Topic
Modeling

05

Classification

06





Introduction

Arabic poems are the oldest and the most prominent form of Arabic literature today.

Goal

In this project, we will develop an **Arabic poems** topic modelling to predict the correct phase (**العصر**) of poems that can correctly **classify** the topic.

Dataset

The dataset that will be used in
this project has been downloaded
from (الديوان)

- 1,831,770 **rows**
- 8 **columns**

Tools

- Software Platform :Jupyter Notebook
- Programming Language: Python
- Python Libraries:
 - Statistics libraries:
 - Sklearn
 - Nltk
 - Data manipulation libraries:
 - Pandas
 - Numpy
 - byArabic package
 - Camel tool
 - Visualization libraries
 - Matplotlib
 - Seaborn
 - wordcloud

Exploratory data analysis (EDA)



**Remove
Unnecessary
columns**



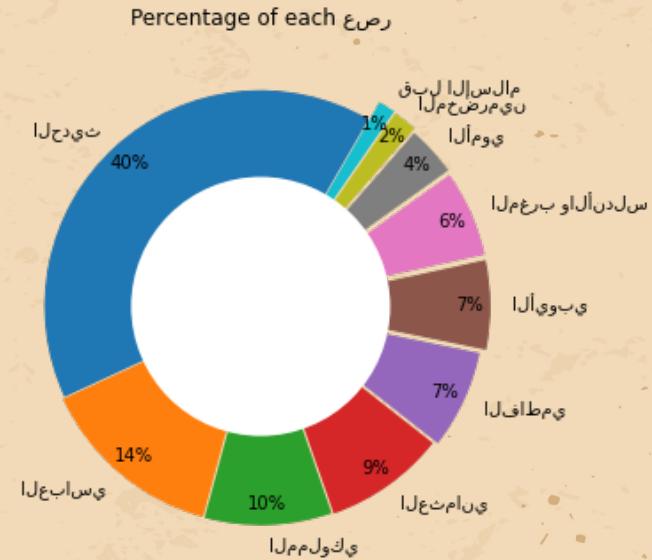
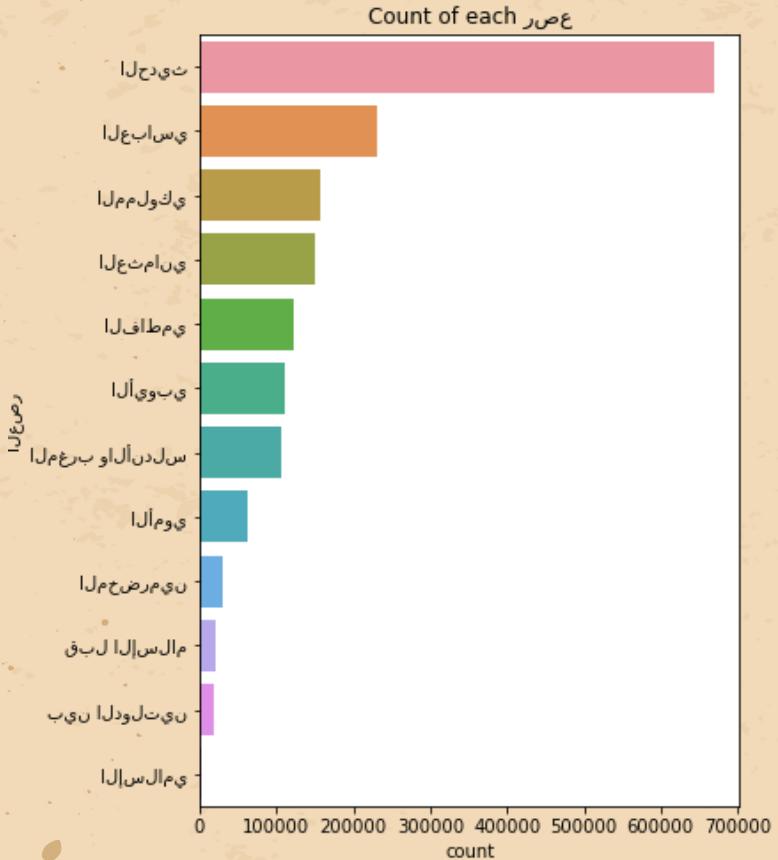
**Remove NULLS
and Duplicate
values**



Split Data by
الديوان



Distribution of العصر



Word Clouds



Text preressing

01

Remove harakat and tashkeel , tatweel.

02

Remove extra spaces and punctuations

03

Tokenization
Stemming.

04

Replace أ أ أ with أ And و و و with و
And repeated letter with single letter like .[ا][ا][ا], [س][س][س], [و][و][و]

05

Remove Stop Word.



Topic Modelling



Modelling

- Latent Dirichlet Allocation (**LDA**)
- Latent Semantic Analysis (**LSA**)
- Non-Negative Matrix Factorization (**NMF**)
- **CorEx**
- **Clustering**



Embedding

- **CountVectorizer**
- **TF-IDF**



NMF / Count Vectorizer

Topic 0

يوم، الدهر، الناس، الهوى، خير، الزمان، قلبي، الدنيا، الورى، الارض

Topic 1

الهوى، قلبي، الحب، القلب، حب، قلب، الفواد، النفس، الحياة، فوادي

Topic 2

الزمان، الايام، الناس، الدهر، يوم، الملك، الورى، العلي، الهدى، الدنيا

Topic 3

الارض، الاسى، دمي، الموت، الحزن، دم، المدى، الدماء، الدموع، الاولى

Topic 4

الحق، الناس، الاسلام، الله، الحياة، خير، اليوم، العباد، محمد، الحر

الدهر والزمان Topic 0:

غزل Topic 1:

فخر Topic 2:

رثاء Topic 3:

اسلامي Topic 4:

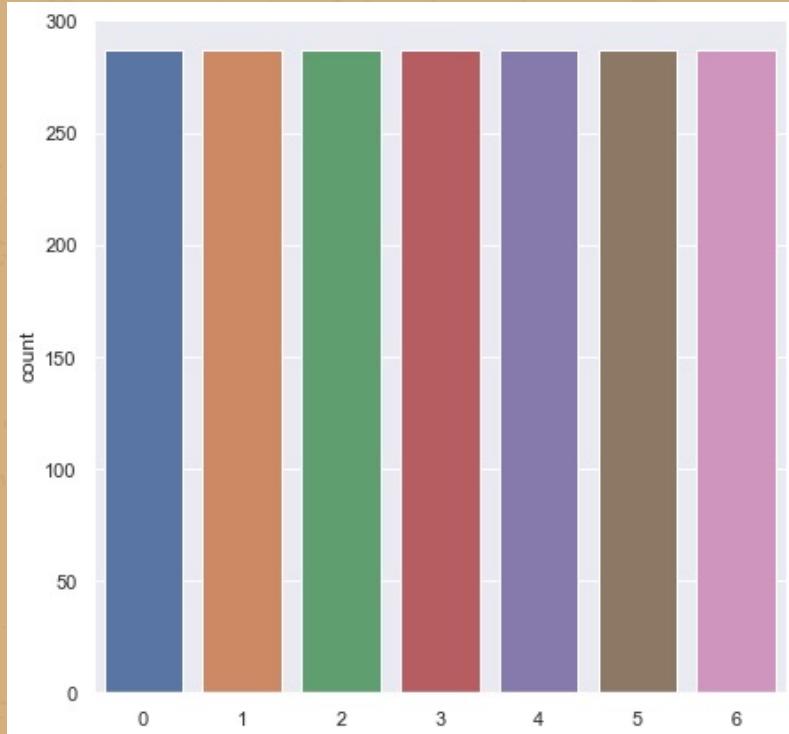
Classification

- XGBClassifier
- Gaussain Naive Bayes
- Support vector machines (SVMs)
- Random Forest
- Multinomial Naive Bayes
- Decision Tree
- Logistic Regression

Handling Data Imbalanced

SMOTE

The total observations before 733
The total observations after: 2009



Best Model

XGBClassifier



F1

0.44

Accurecy

0.25

Future Work:

1.

Feed more features to the classifiers

2.

Recommendation system

Thanks!

