



Arabic Poem

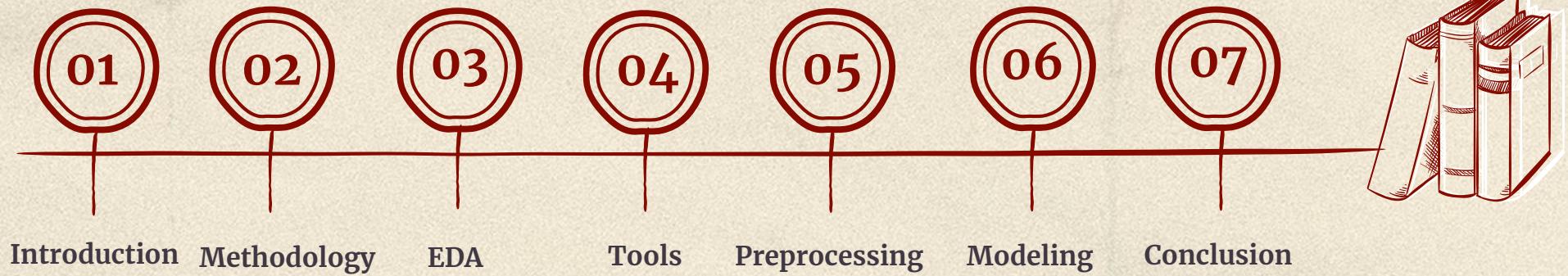
Ohood Soliman
Ajwad Almajnuni
Mashael Asiri





أنا الذي نظر الأعمى إلى أديبي
وأسّمعت كلماتي من به صمم

أبو الطيب المتنبي



Introduction

The purpose of this project is to apply the most important natural language processing techniques and starting with implementing all text preprocessing in "Arabic" then using classification model to classify each Poet name and his poetries and we will use topic modeling to cluster each poetries depend on poetry types or poetry rhymes.

About Data : This dataset contains 1,831,770 observations and 8 columns



Methodology



EDA



Preprocessing



Modeling



EDA

01

Check null values

Duplicates

02

Rename columns

03

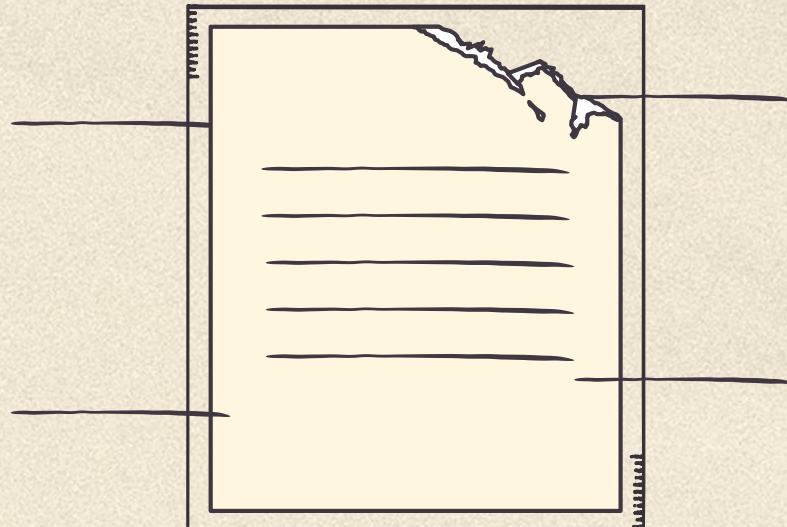
Visualization



Tools

Matplotlib
wordcloud

NLTK
Pyarabic

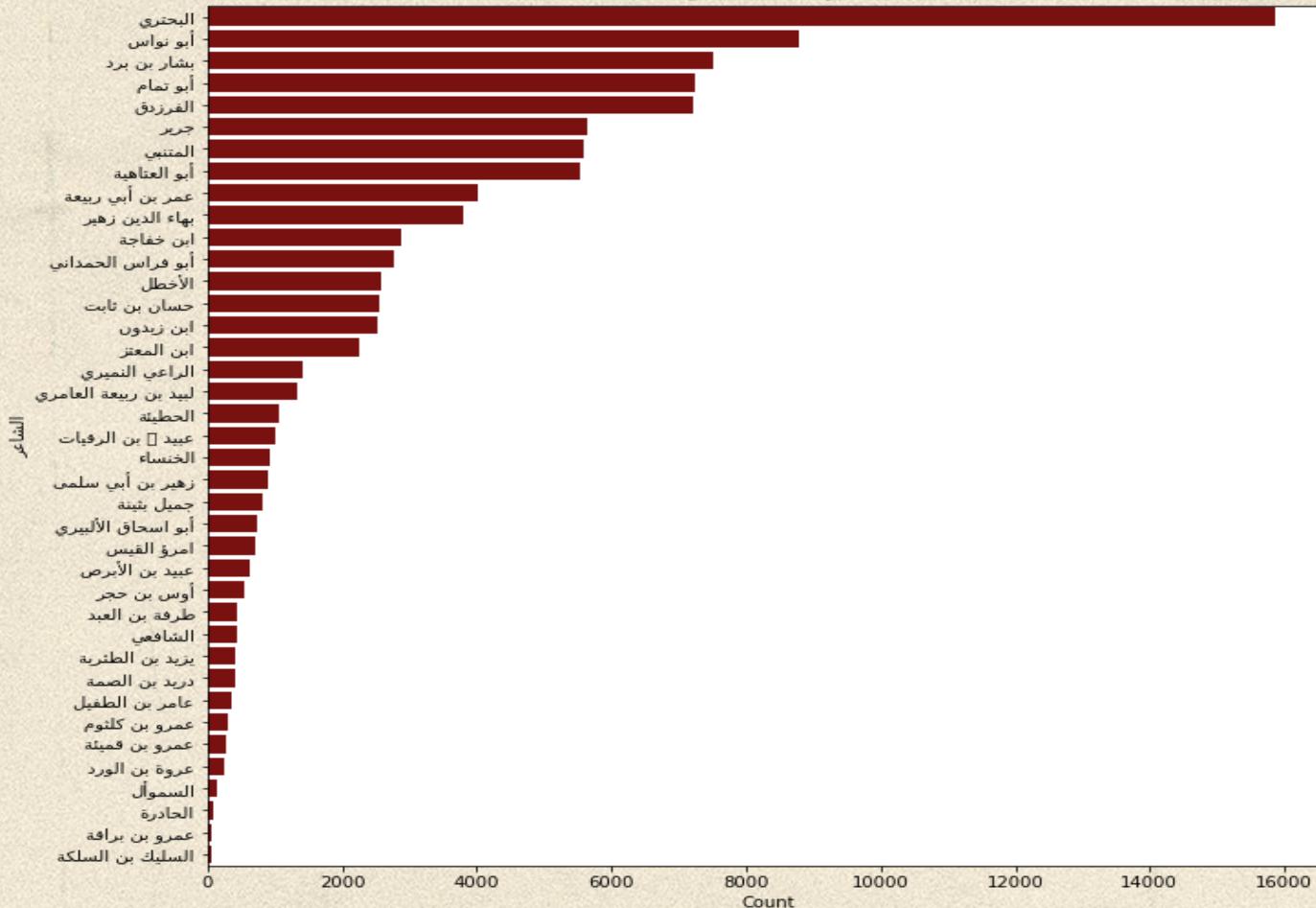


Python
Jupyter
notebook

NumPy
Pandas
SKLearn

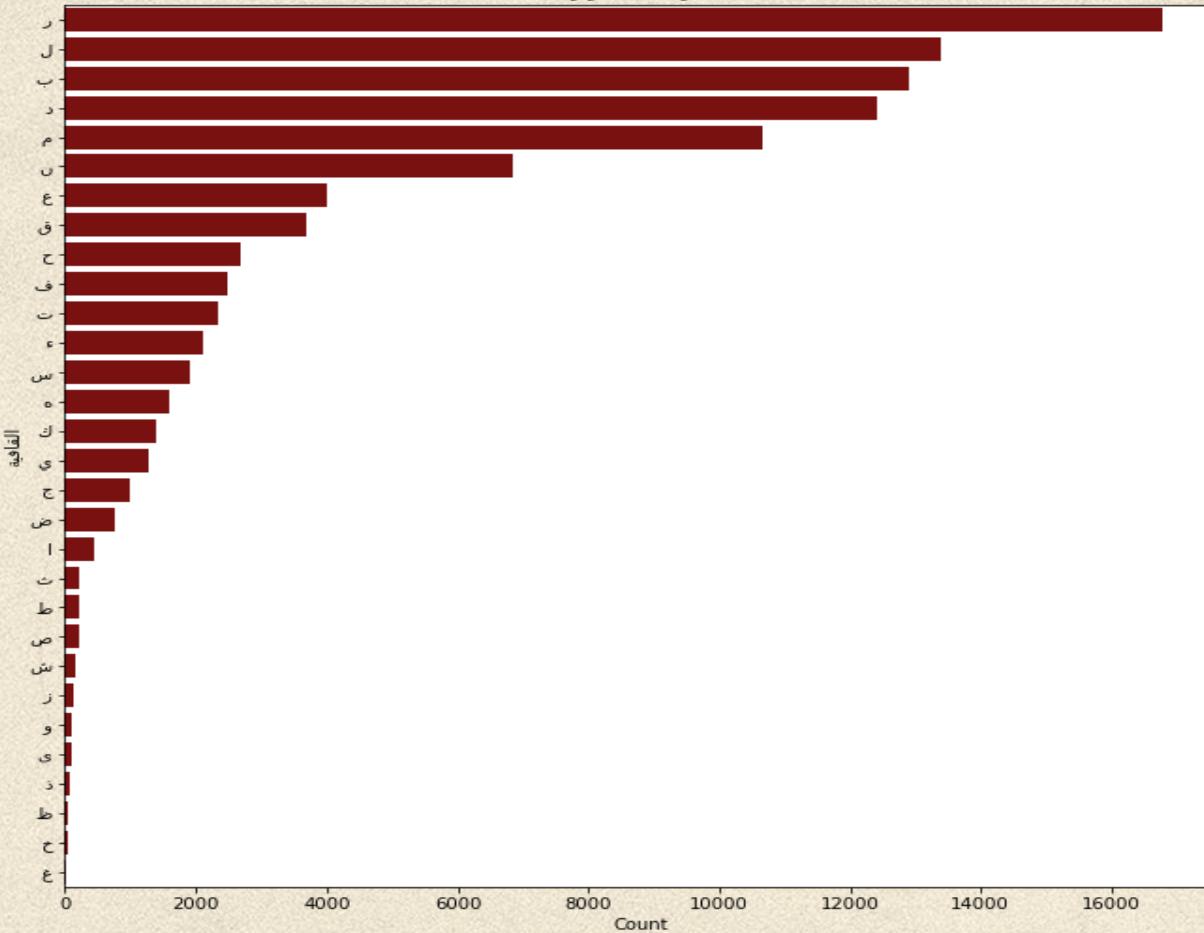
أكثر القصائد كانت للشاعر ...

أكبر القصائد للشاعر



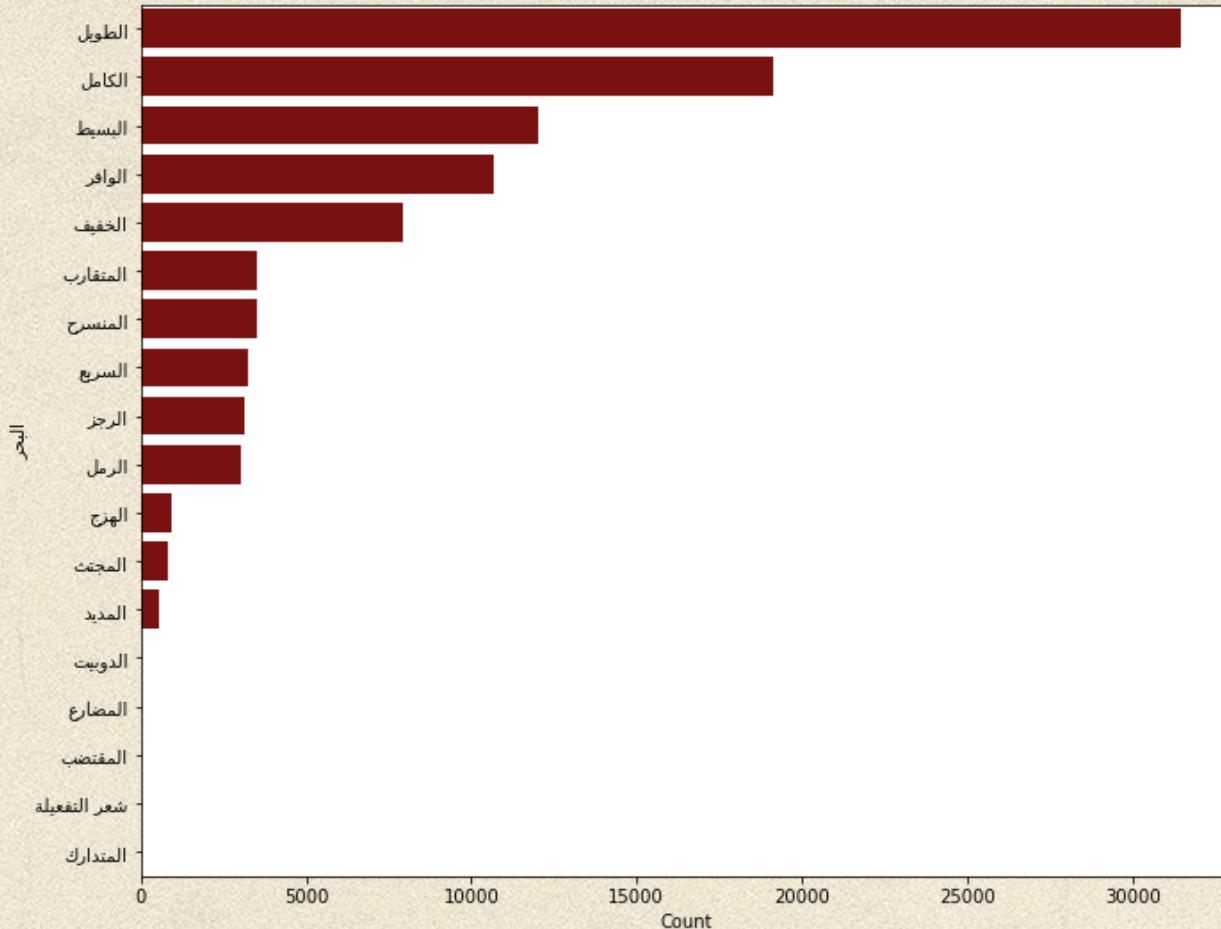
القافية الأكثر تكرارا

أكثر قافية تكرارا



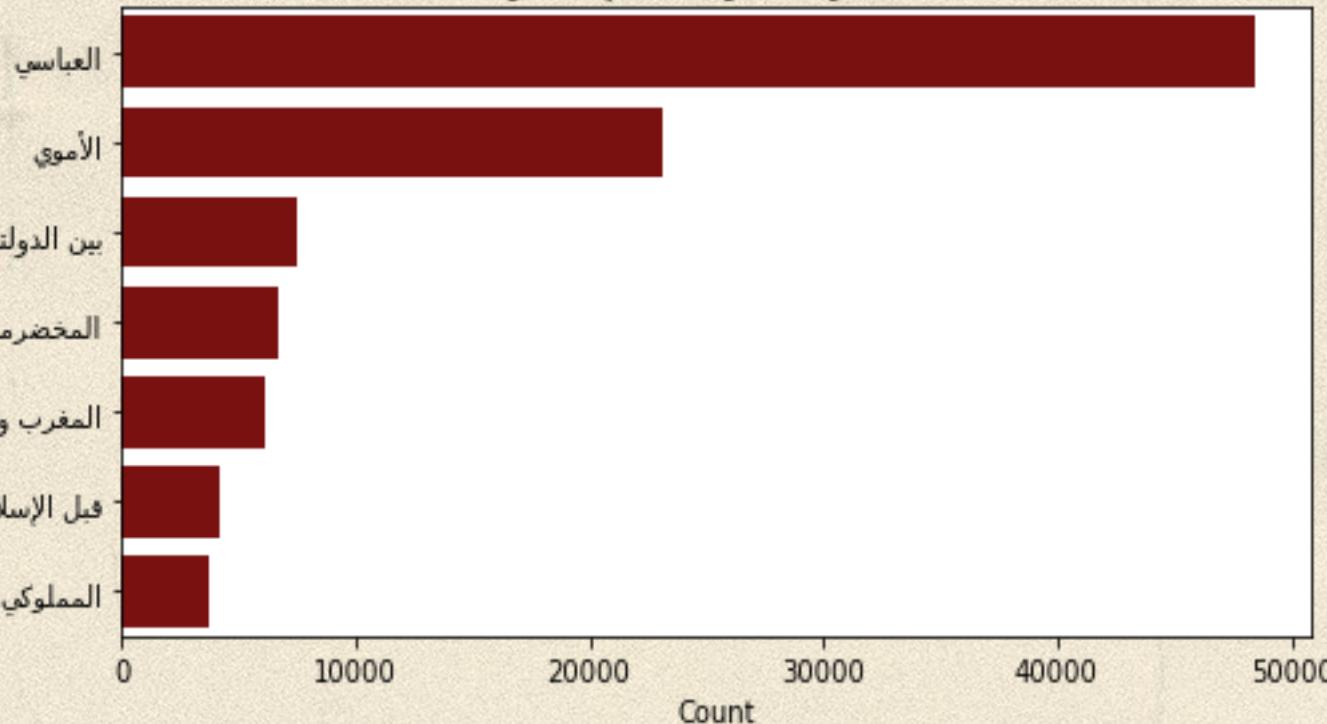
أكثر بحور الشعر استخداما

أكثر بحور الشعر استخداما



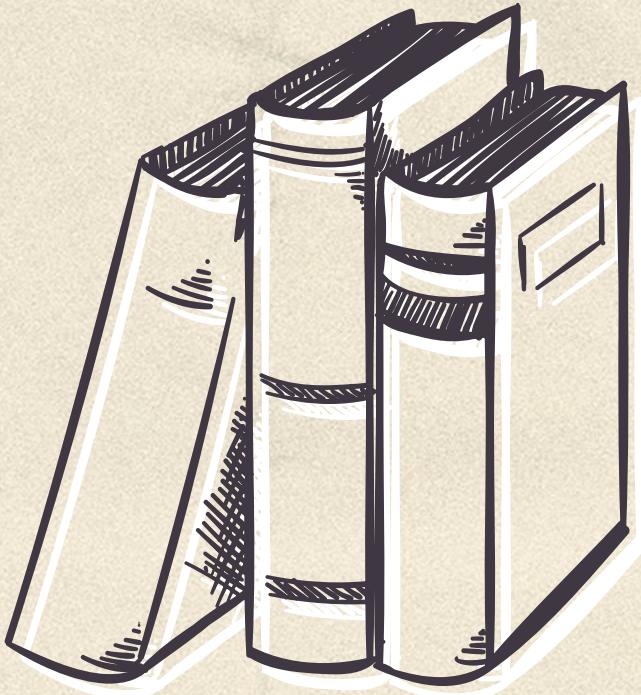
كانت اكثراً الشعر في العصر

....أكثراً الشعر كانت في العصر



NLP Preprocessing

- **Remove**
 - *normalize Arabic*
 - repeated letters remove
 - punctuations
 - special character
 - non-arabic-alphabet



NLP Preprocessing

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



Wordcloud

ML algrathim



- Logistic Regression
- Random Forest
- XGBoost
- Decision Tree

- k -means
- Hierarchical clustering

Logistic Regression

Training

Training Logistic Regression with best
parameters1 **F1 score= 0.7456413443067098**

Test

Test Logistic Regression1 with best
parameters1 **F1 score= 0.6387911437172631**



Decision Tree

Training

Training Decision Tree with best parameters1
F1 score= 0.30552456974887593

Test

Test Decision Tree with best parameters1
F1 score= 0.2918717311288094



Random Forest

Training

Training Random Forest with best parameters1
F1 score= 0.34913956079465364

Test

Test Random Forest with best parameters1
F1 score= 0.3137159705219611



XGBoost

Training

Training XGboost with best parameters1
F1 score= 0.6858045498280895

Test

Test XGboost with best parameters1
F1 score= 0.6529368761859843

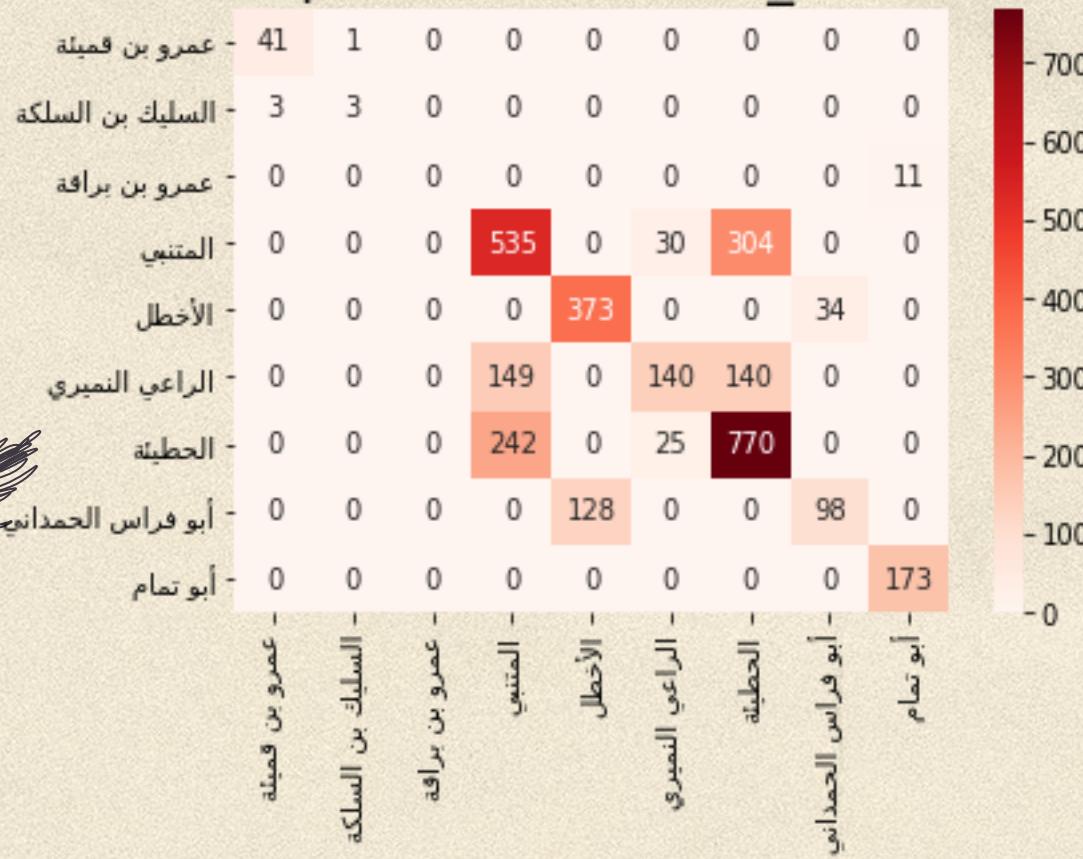


Comparison all classification algorithms scores

Model	Accuracy
Logistic Regression Algorithm	0.64
Decision Tree Algorithm	0.37
Random Forst Algorithm	0.41
Xgboost Algorithm	0.67

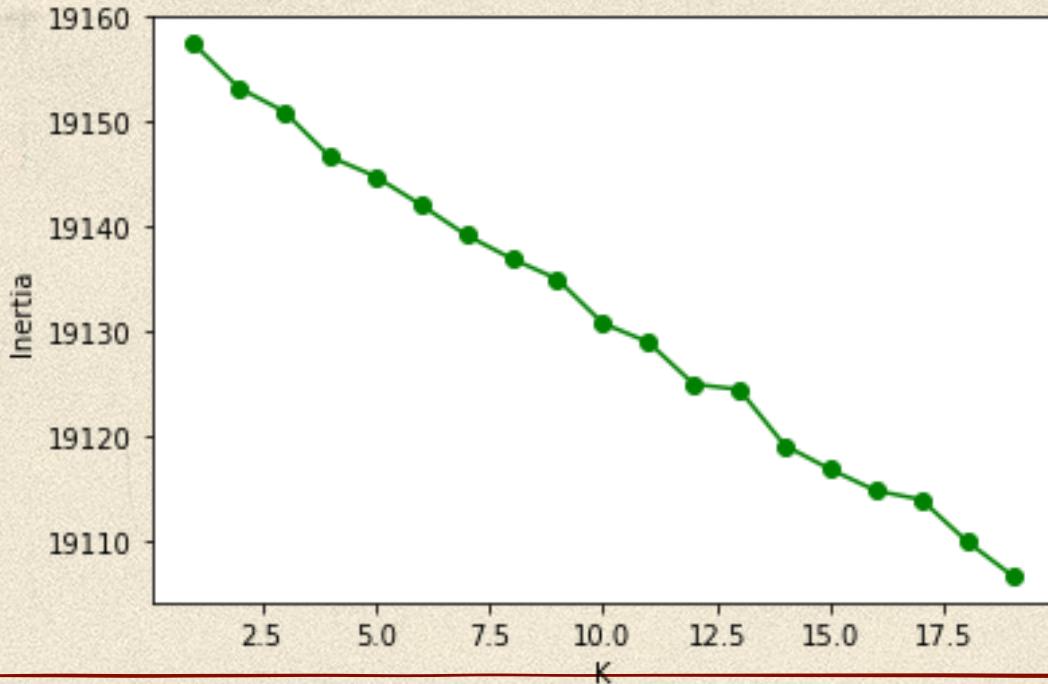


XGboost Mosel to predict each Poet_name with his poems



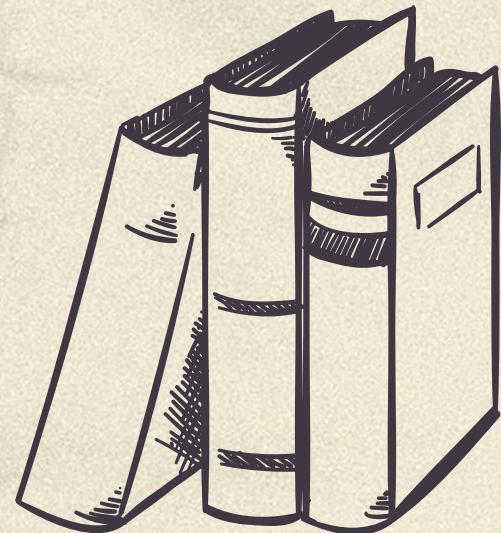
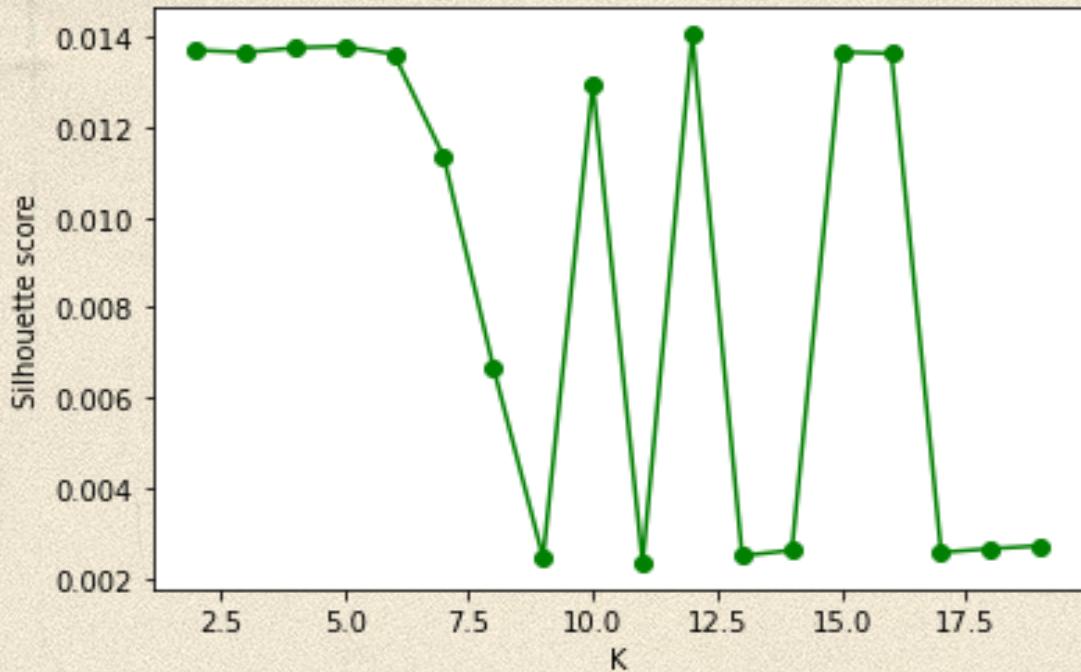
Unsupervised Modeling

- Kmeans



Modeling

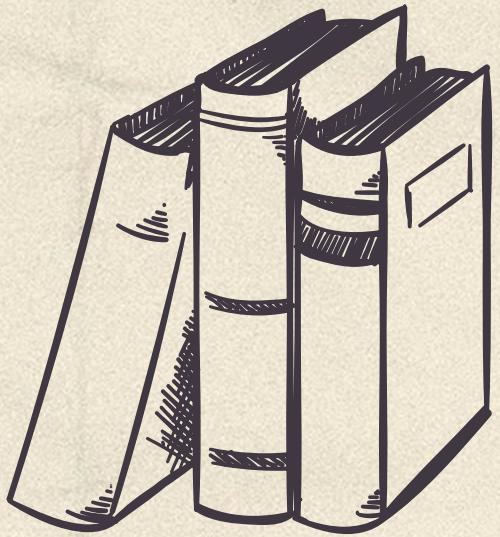
- Kmeans



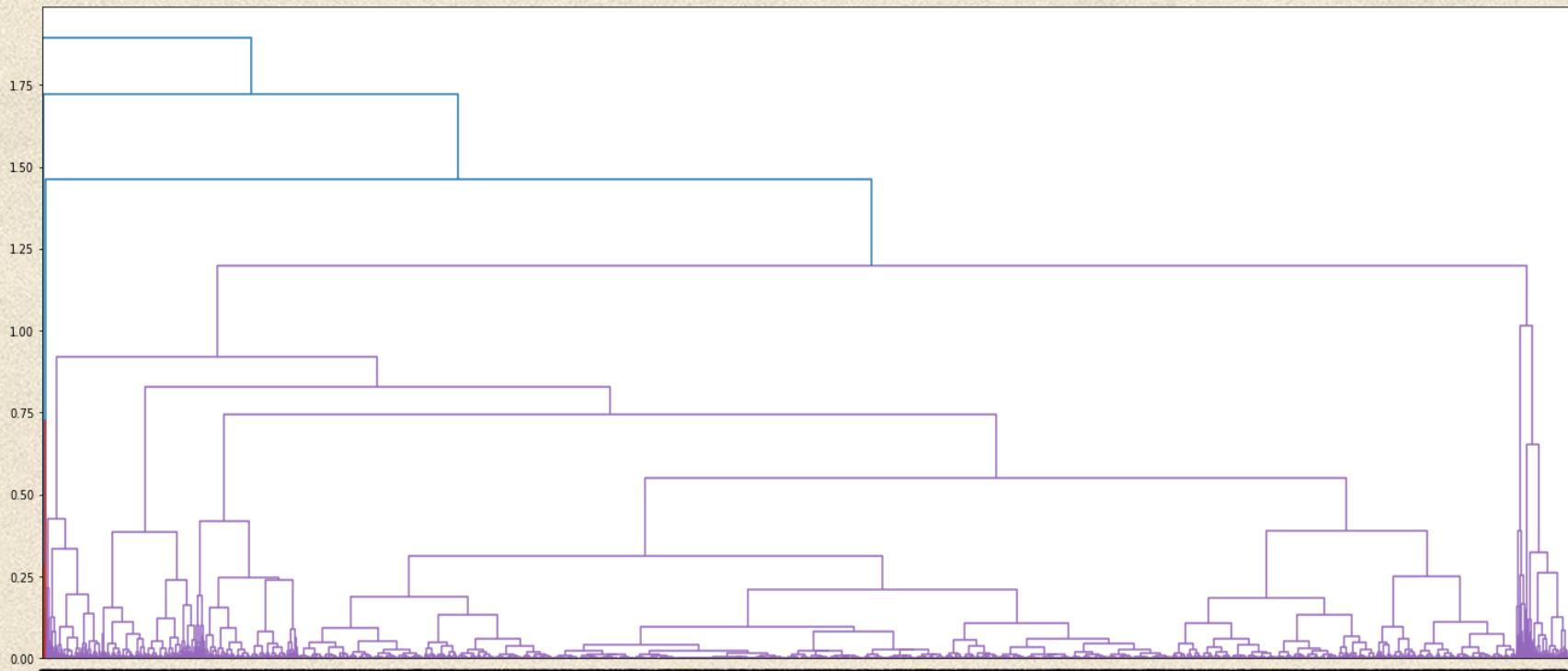
Modeling

- Hierarchical clustering with different linkages

Ward	Complete	average
2.83s	2.60s	2.71s

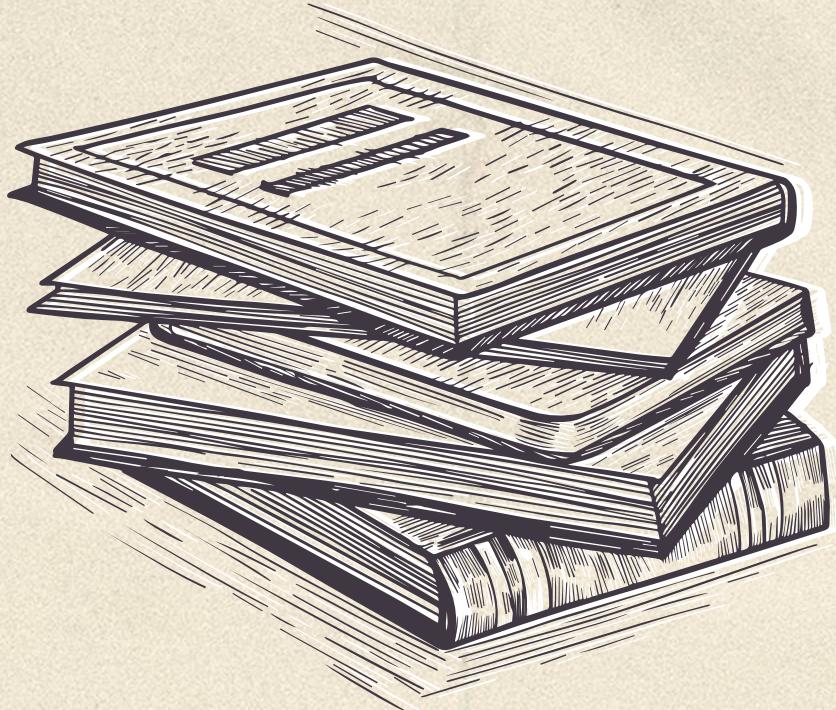


Ward



Conclusion

In attempts to predict the best model , if you enter the poem, the name of the poet will be predicted . We made several models to determine the best model.. . The best result is **XGBoost** .





**THANK
YOU**