



KSA Ministry of Justice

NLP project and cosine similarity

Problem statement

The data contains judicial information, including final commercial judgments issued by commercial courts, this is where our project will help find the most similar judgments.

contents

1

collect data

Web scrapping
and tools

2

Preprocessing

- Data cleaning
- Text preprocessing

3

Convert to vector

TF –IDF
SVD

4

Distances matrix

Cosine similarity

5

Tools

- Python and Jupiter Notebook
- NumPy and Pandas for data manipulation
- Matplotlib and Seaborn for plotting visualization
- SkLearn

1

Collect data

Web scrapping

الصفحة الرئيسية > التصنيف > مجموعة الأحكام القضائية

● بيانات الحكم

المحكمة: المحكمة العامة

المدينة: بريدة

تاريخها: ٢٦/١٠/١٤٤٣

رقم القضية - القرار: ٣٣٣

محكمة الاستئناف: المحكمة التجارية

المدينة: منطقة الرياض

تاريخه: ٢٤/١٠/١٤٤٣

رقم القرار: ٥٥٢٨

● التصنيف

● عنوان الحكم

● نص الحكم

● الاستئناف



691



طباعة



مشاركة



-



=



+

حجم الخط

Cleaning data

1

Rename columns

2

Drop duplicate values

3

Fill null values

2

Preprocessing

NLP on Arabic data

Preprocessing

1

Remove punctuations

2

remove \n

3

remove numbers

Preprocessing

1

Remove stop words

2

Correct words

3

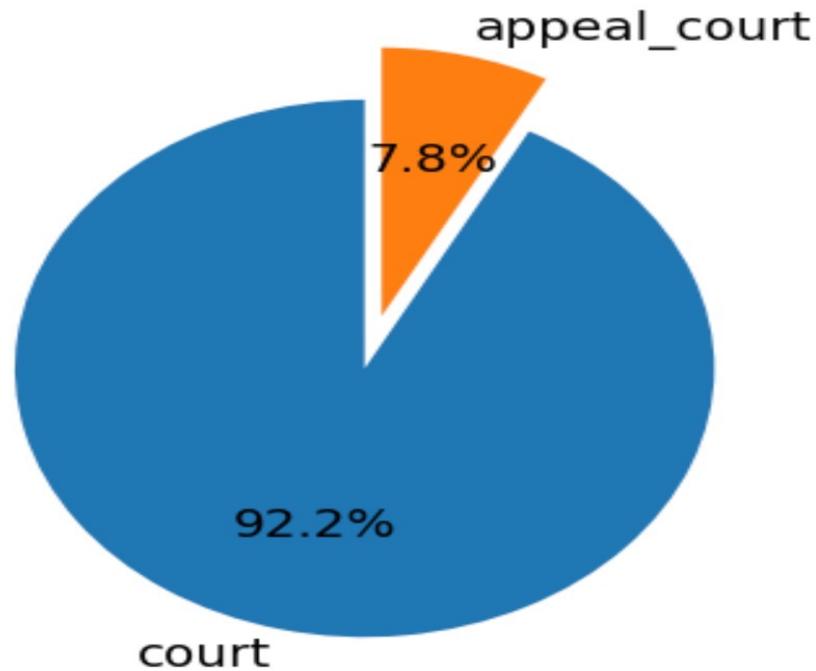
remove English words

3

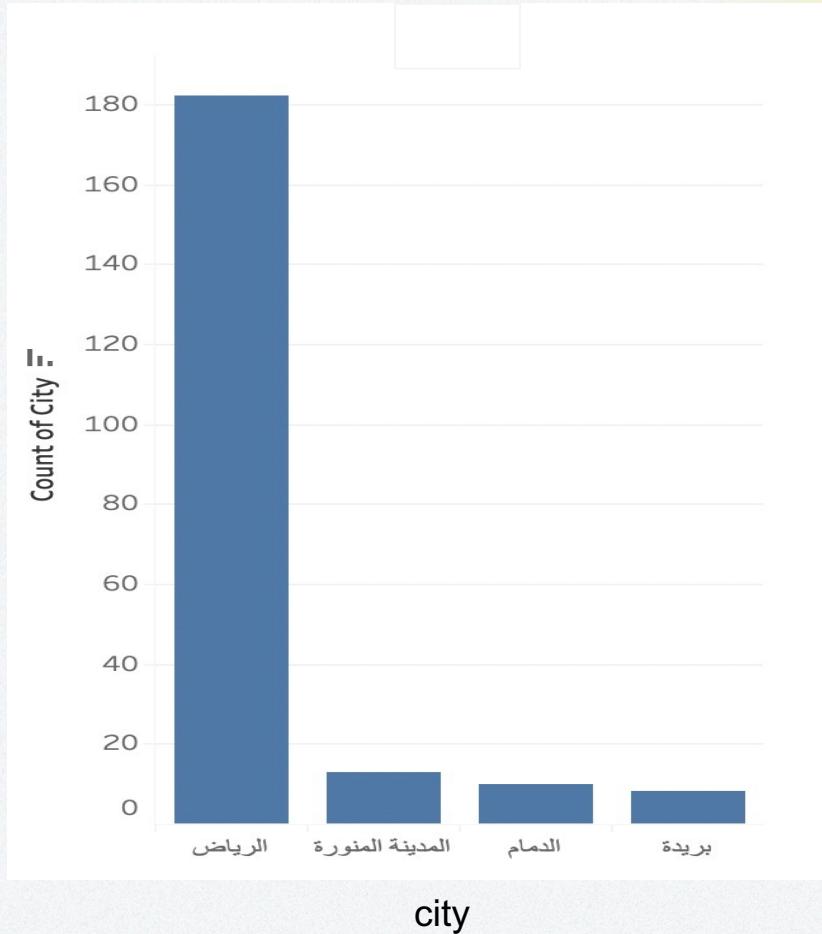
EDA

Shape data : (2043, 8)

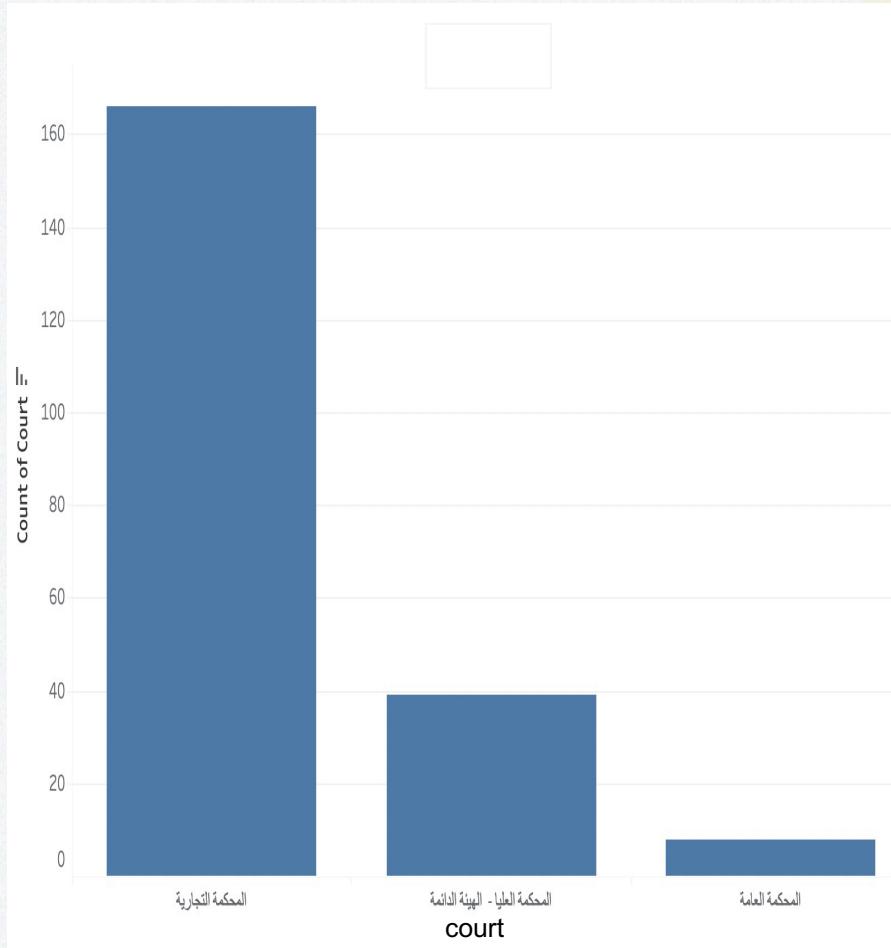
Appeal court percentage according to Court percentage



Most city where judgments issued



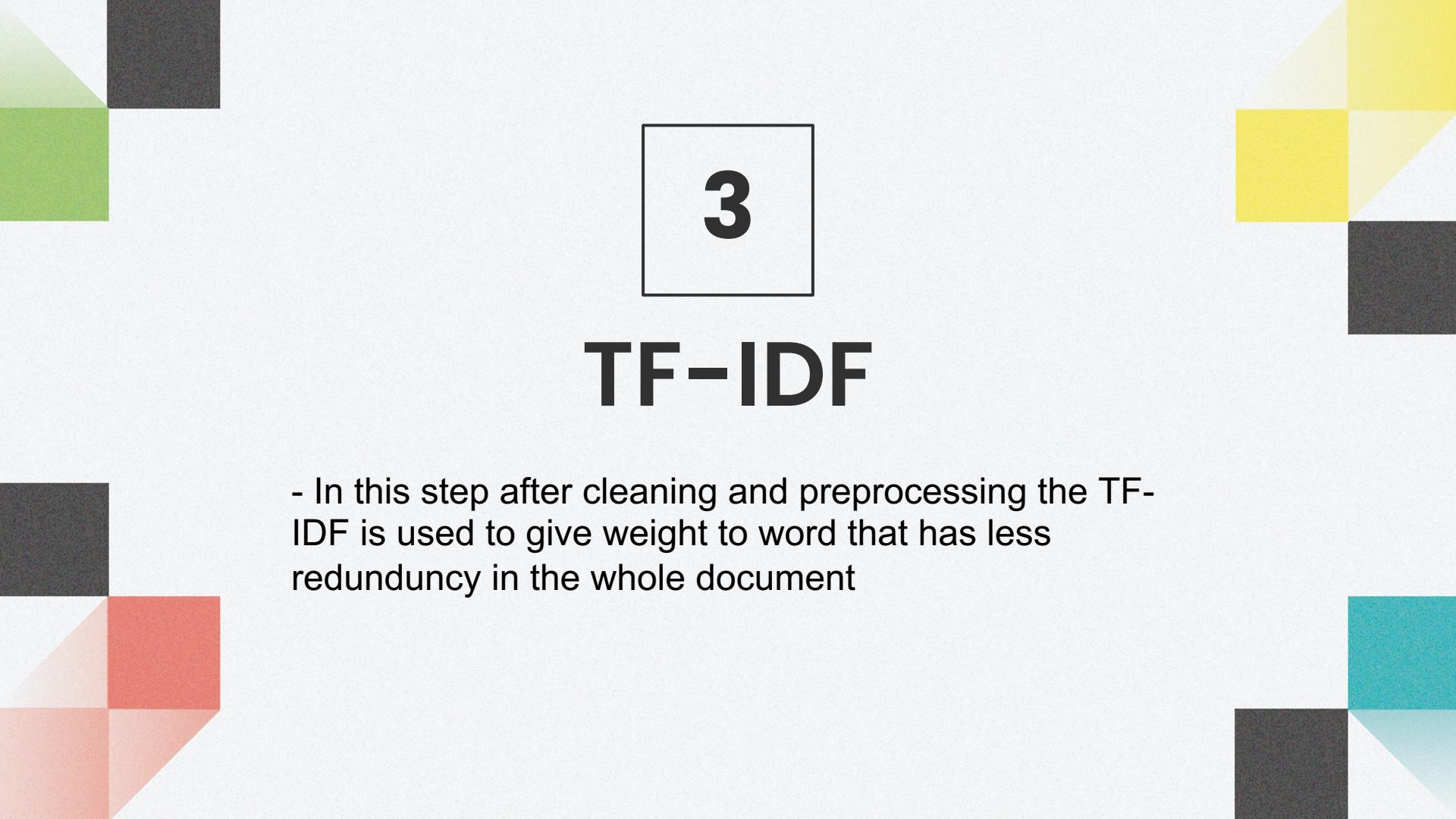
What is the most redundant court



Word cloud

Word cloud judgment_text





3

TF-IDF

- In this step after cleaning and preprocessing the TF-IDF is used to give weight to word that has less redundancy in the whole document

4

Cosine similarity

Output distance Cosine matrix

```
cosine_similarity(dt_tfidf.loc[25:25,:],dt_tfidf).argsort()[0][-6:]
```

```
array([ 235,  982,  735,    87, 1108,   25])
```

```
cosine_similarity(dt_tfidf.loc[25:25,:],dt_tfidf.loc[1108:1108,:])
```

```
array([[0.70741606]])
```

Exsample:

```
df train.judgment text.iloc[25]
```

```
df['train-judgment-text'].iloc[1108]
```

Thanks!

Do you have any questions?