The background of the slide is a solid light beige color. On the left side, there are three vertical rectangular panels of Arabic calligraphy. The top panel is dark red with white script. The middle panel is light beige with dark red script. The bottom panel is dark red with white script. On the right side, there is a vertical strip of six small, dark red squares, each containing white calligraphic text.

Natural Language Processing and Categorizing Arabic Articles

by: Abeer Almdani , Sarah Alrashidi

Outlines

1. Introduction and business problem
2. Data
 - EDA
 - Text preprocessing
3. Algorithm
 - Topic Modeling
 - Recommendation-systems
4. Tools
5. Future work



Introduction and Business problem

The goal of this project is to use Natural Language Processing and Unsupervised Learning models to categorize them into sections. After refining the approach. In this project, we will analyze the articles and categorize them into groups such as political



DATA

Data source Kaggle

<https://www.kaggle.com/oussamaseffai/arabic-article-headline-generation>

- 1050 rows
- 3 columns (title,introduction,index)
- Max words number is 1972 words
- Min words number is 116 words
- Max token length is 397 tokens
- Min token length is 19 tokens

DATA

Most Used Words

18.611316	طريق
15.534055	علاج
15.425522	جسم
15.257080	بشر
15.015648	ماء
13.401183	زيت
13.044927	عام
12.822835	اعراض
12.047188	عمل

bigrams n2

(28 , 'زيت زيت')
(27 , 'درج حرار')
(24 , 'احجار كريم')

trigrams n3

(12 , 'زيت جوز هند')
(11 , 'ممارس تمار رياض')
(11 , 'كريا الدم بيضاء')

DATA: Exploratory data analysis (EDA)



Remove NULLS
and Duplicate
values



Merge Title
and
Introduction
columns



Remove
Unnecessary
columns



Removed short
decomments

DATA: Text preprocessing

Remove Non-Arabic Char

Remove extra spaces

Remove Arabic Diacritics

Remove Arabic Numbers

٠،١،٢،٣،٤،٥،٦،٧،٨،٩

Normalize Arabic Letters

ك → گ

Remove Arabic Punctuations

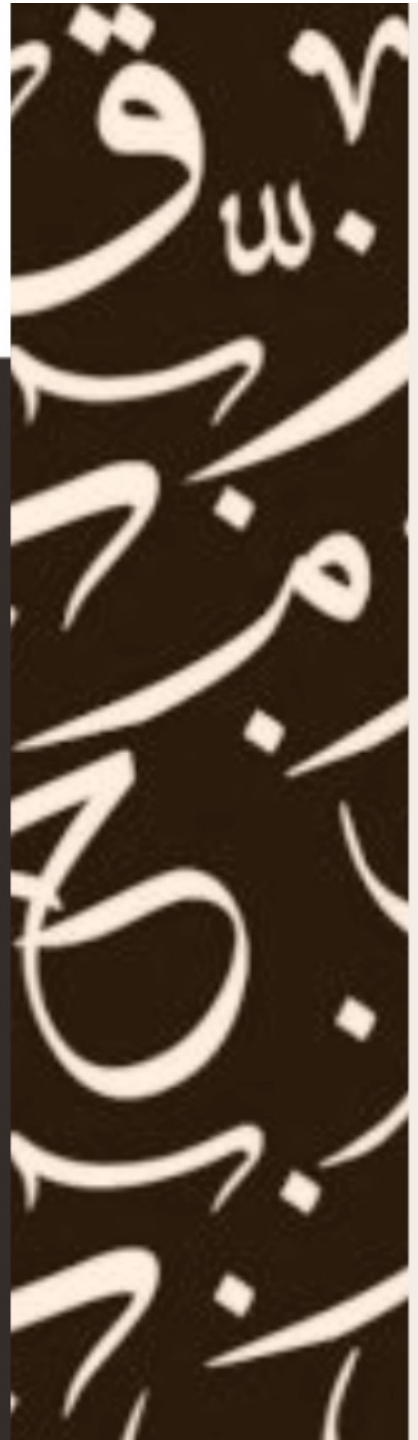
؟ ، . / !

Remove Arabic Stop Word
(350 word)

Remove Repeating Char

وو <--- و

يي <--- ي



Algorithm: Topic Modelling

Embedding

- CountVectorizer
- TF-IDF

Modelling

- LatentDirichletAllocation(LDA)
- LatentSemanticAnalysis(LSA)
- Non-NegativeMatrixFactorization(NMF)
- CorEx
- Clustering and (Elbow curve and TSNE)



NMF | Count Vectorizer

Topic 0

عام , عمل , عالم , شخص , حيا , تم , تاريخ , فتر , بدء , عمر , يوم , خاص

Topic 1

مدين , دول , تقع , جنوب , تعتبر , مساح , بحر , شرق , عرب , جزير , سكان , غرب

Topic 2

جسم , غذاء , تناول , انس , الصح , طبيع , انواع , حرار , يعتبر , لذل , مواد , لان

Topic 3

اصاب , حال , علاج , اعراض , مرض , مصاب , شخص , عدوي , فيروس , جهاز , الدم , تظهر

Topic 4

طريق , ماء , يتم , استخدام , بشر , شعر , لمد , مكون , زيت , وصف , غسل , حصول

Results

We can notice that topics are:

Topic 0: تاريخ

Topic 1: جغرافيا

Topic 2: الغذاء والطعام

Topic 3: الطب

Topic 4: الصحة والجمال



Recommendation system

- Document Similarity with Count vectorizer

مقالك ينتمي إلى المقالات الجغرافية

***** صيد الأسماك في مصر *****

مصر من البلدان العربية الغنية بمصادرها السمكية سواء التي تعيش في المياه العذبة أو بالمياه المالحة ، وقد تم إنشاء الكثير من المشاريع السمكية بسبب توفر المياه على طول مجرى نهر النيل و قناة السويس ، حيث تتم تربية أنواع محددة من الأسماك بقصد تصديرها إلى الخارج .

Recommendation system

1

***** موقع طروادة الأثري و الجغرافي *****

أجريت الحفريات الأولى للمدينة سنة 1870م ، و تعتبر بقاياها من أهم و أبرز مظاهر الاتصال بين حضارتي الأناضول ، و البحر الأبيض المتوسط ، و تقع طروادة على تلة تسمى هيسارليك (بالإنجليزية : Hisarlık) ، و تطل هذه التلة على السهل الممتد على طول ساحل بحر إيجه التركي ، و على بعد 8 . 4 كم من المدخل الجنوبي إلى مضيق الدردنيل ، و قد كشفت الحفريات عن كونه منطقة كانت قد سكنت منذ 8000 عام ؛ و بسبب موقعها ، فقد عملت كجسر ثقافي يربط بين منطقة ترواس ، و البلقان ، و الأناضول ، و بحر إيجه ، و البحر الأسود ؛ و ذلك بسبب الهجرة و التنقل . يعتبر موقع طروادة الأثري ذا أهمية بالغة ؛ لفهم تطور الحضارة الأوروبية ؛ حيث إنه شاهد على تتابع الحضارات التي احتلت المنطقة لأكثر من 4000 سنة ، كما أن العديد من الآثار اليونانية ، و الرومانية ، تعكس خصائص المستوطنات التي أقامت فيها ، إضافة إلى أن موقعها يعد ذا أهمية ثقافية كبيرة ؛ و ذلك بسبب تأثيره في الأعمال الأدبية ، مثل : ملحمة إلياذة للمؤلف هوميروس (بالإنجليزية : Homer's Illiad) ، و الإنيادة لفرجيل (بالإنجليزية : Virgil's Aeneid) . و تعتبر طروادة نقطة استراتيجية للمدخل الجنوبي للمضيق ، الذي يربط البحر الأسود مع بحر إيجه ، عبر بحر مرمرة ، كما أنها مصدر لطريق بري يمتد شمالا على طول ساحل الأناضول الغربي ، و عبر أضيق نقطة من الدردنيل إلى الشاطئ الأوروبي ، و قد ظل الموقع الدقيق لهذه المدينة غير محدد حتى العصر الحديث ؛ حيث بدأت الحفريات لاكتشاف المدينة منذ عام 1870 ، و حتى عام 2005 م ، من خلال العديد من الحملات كان آخرها بقيادة عالم الآثار في جامعة توبنغن ، مانفريد كورفمان (بالإنجليزية : Manfred Korfmann) .

2

***** مميزات جزيرة المرجان *****

تحتوي هذه الجزيرة على مساكن مميزة تعرف بباب البحر ، و هي في معظمها عبارة عن وحدات فاخرة للسكن ، إضافة إلى فنادق عالمية كفندق ريكسوس التركي ، و الذي يضم ستمئة و خمسين غرفة ، و هو يعتبر من المنتجعات الضخمة ذات النجوم الخمسة ، و الذي بني على مبانٍ ثلاثة ذات شكل هرمي ، و فيها قاعات خاصة بالمؤتمرات ، إضافة إلى المسرح ، و يوجد أيضا فندق هيلتون دبل تريه ، و كراون بلازا بفرعه الوحيد في دولة الإمارات على أرض هذه الجزيرة . يأتي في مقدمة هذه الفنادق منتجع ذا مديرة ، الذي تم منحه لأول مرة في العالم من قبل فريق ذا مديرة ، و يأتي على مساحة تقدر بأربعمئة ألف متر مربع ، و يتمتع بإطلالة ساحرة ، و كذلك الأمر بالنسبة لفندق سبا جزيرة المرجان ، حيث يوجد فيه حوالي ثلاثمئة غرفة و جناح ، و قد صمم بطراز أندلسي ، و يتبع له كورنيش يبلغ طوله حوالي كيلومترا . تحوي أيضا جزيرة المرجان على عدد من المطاعم العالمية ، و تتوزع على شاطئ الجزيرة و كورنيشها المجهز بعدة مسارات ، بعضها مخصص لهواة رياضة المشي و حتى هواة المشي ، و بعضها الآخر لهواة ركوب الدراجات الهوائية ، إضافة إلى العديد من المراكز الترفيهية في الجزيرة . الجدير بالذكر أنه حسب تقرير لهيئة تنمية السياحة في دولة الإمارات السنوي ، قد بلغت عائدات الدولة من إمارة رأس الخيمة في عام ألفين و أربعة عشر للميلاد ما يزيد عن مليار درهم إماراتي ، و لعل لجزيرة المرجان اليد الطولى في الدعم الاقتصادي و الاستثماري في البلاد .

3

***** أكبر مدينة في العراق *****

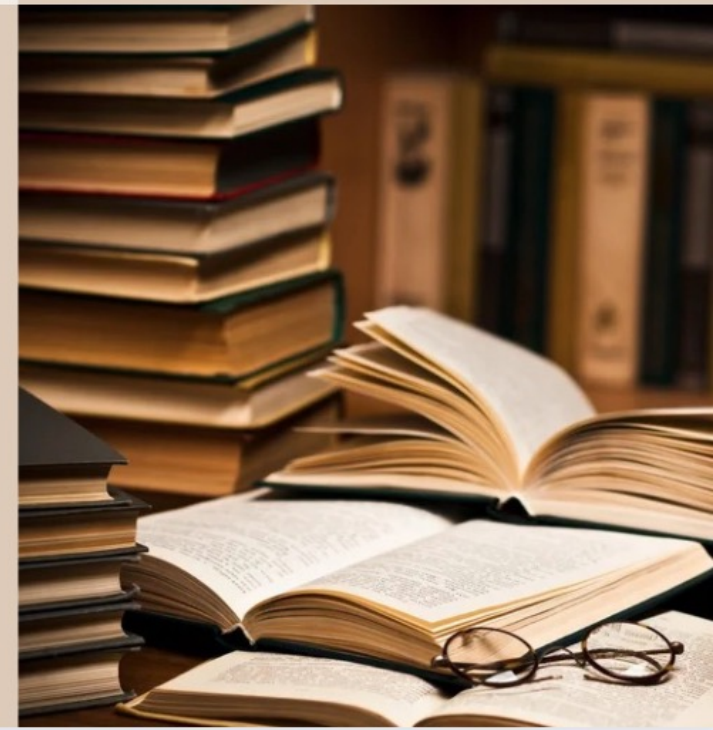
تعد بغداد عاصمة دولة العراق و أكبر مدينة فيها و ثاني أكبر مدينة على مستوى الوطن العربي و غرب قارة اسيا ، حيث يبلغ عدد سكانها حوالي 7,180,889 نسمة ، تقع بغداد على طول نهر دجلة ، و بسبب موقعها المميز ازدهرت و تطورت المدينة بشكل سريع ، كما ساعدها موقعها من السيطرة على طرق التجارة الواقعة على طول نهر دجلة ، و يقسم نهر دجلة المدينة الى قسمين هما : الرصافة الواقعة في الجهة الشرقية ، و الكرخ الواقعة في الجزء الغربي ، و تعتبر مدينة بغداد مركزا للثقافة العربية ، و تضم المدينة مجموعة من المعالم التاريخية و الأثرية السياحية أهمها : نصب الشهيد ، حديقة الحيوانات ، المتحف الوطني .

FUTURE WORK:

More
Recommendation
systems

Remove More
Arabic Stop Word

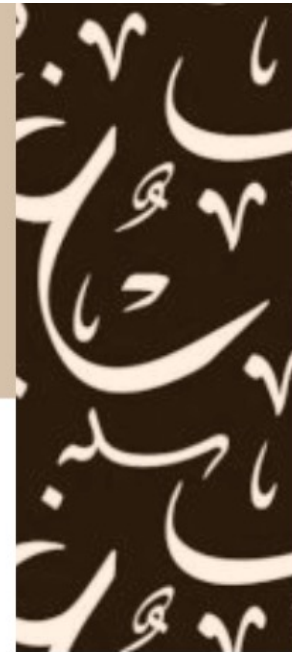
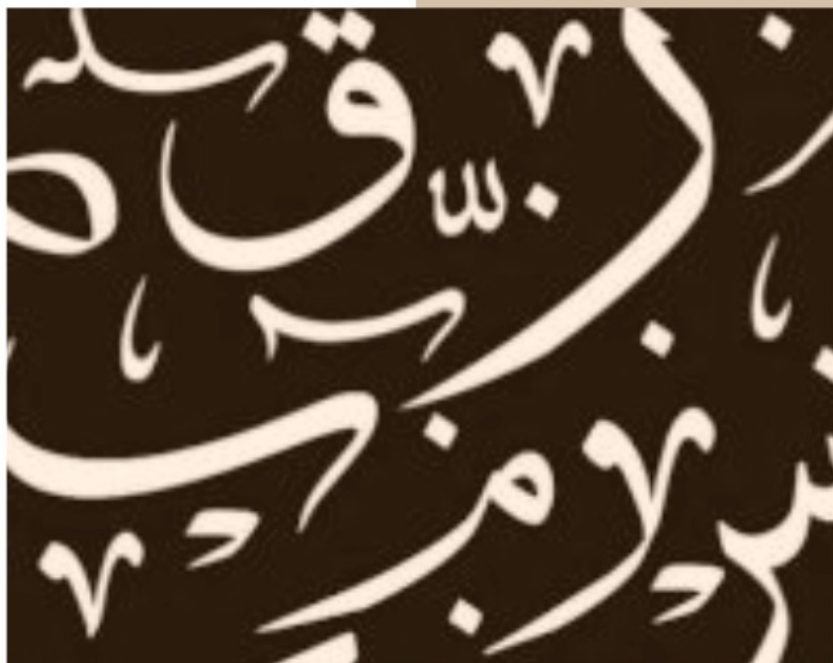
Optimizing
Arabic Python
Libraries



Tools:

- **Software Platform:**
 - Jupyter Notebook
- **Programming Language:**
 - Python
- **Python Libraries:**
 - Sklearn
 - Nltk
 - pyarabic.araby
- **Data manipulation libraries**
 - Pandas
 - Numpy
- **Visualization**
 - Matplotlib
 - Seaborn
 - Canva





THANKS!