

Arabic News Articles

Alanoud Alosaimi & Mashael Alfahaid



WORKFLOW

Introduction

Tools

Dataset

EDA and Preprocess

Topics modelling



INTRODACTION

TOOLS



Numpy



Jupyter



Pandas



Seaborn



Matplotlib



Scikit-Learn



DATASET



Extract the data from json file



Data Source: GitHub, [Link](#)



Before Cleaning: (31,030 x 6)



After Cleaning: (25,617 x 1)

EDA



Removed
Null

Removed
Duplicated

Removed
unneeded
columns



PREPROCESS

Removed
punctuations and
numbers.

Removed Arabic
tatweel.
(السلام' to 'السلام')

Removed Arabic
tashakeel.
(العربية' to 'العربية')

Removed Arabic
Stop words
(adding more than
200 words).

Removed repeating
char.
(وكان' to 'وكان')

Removed Non
Arabic letters.

Normalize Arabic
text.

(ي' to 'ي')

(آ' to 'أ')

Stemming and
tokenization



EMBEDDING

Count Vectorizer

TF-IDF Vectorizer

TOPICS MODELING

Latent Dirichlet Allocation (LDA)

- Count Vectorizer
- TF-IDF Vectorizer

Latent Semantic Analysis (LSA)

- Count Vectorizer
- TF-IDF Vectorizer

Non-Negative Matrix Factorization (NMF)

- Count Vectorizer
- TF-IDF Vectorizer

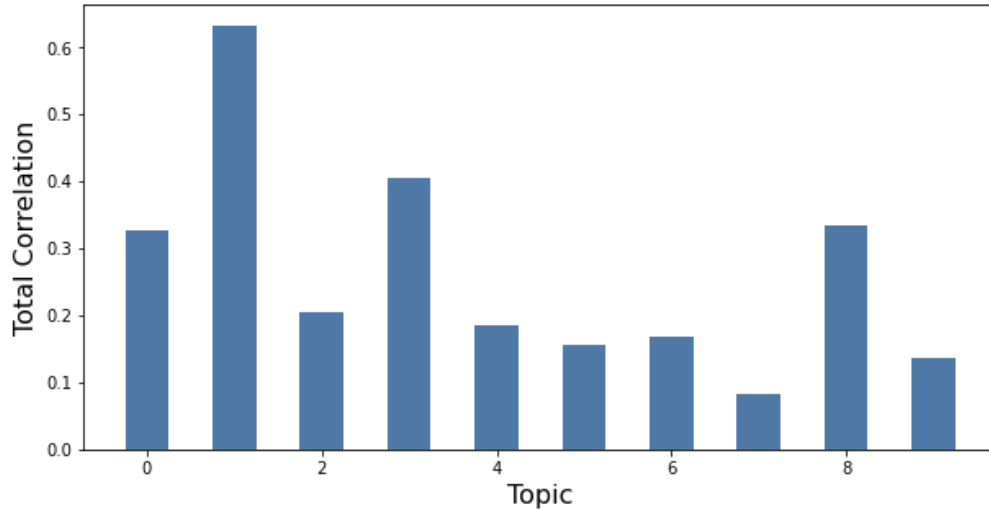
CorEx

- Count Vectorizer
- TF-IDF Vectorizer

CorEx with Anchors

- TF-IDF Vectorizer

CorEx | TF-IDF Vectorizer



- داعش, قتل, تنظيم, هجوم, سوري, ارهاب, تركي, مقاتل, كردستان, مقتل (إرهاب)
- فريق, كره, اعب, ناد, لاعب, مدرب, معسكر, لكر, اتحاد, موسم (رياضة)
- جامعه, طلاب, تعليم, انتخابيه, جامع, بلديه, مدارس, قبول, مرشح, ناخب (تعليم)
- اسعار, سوق, مايه, نفط, اسهم, تراجع, مليار, موشر, دولار, ميه (مالية)
- وزراء, حكومه, عدن, حوث, يمن, مقاومه, لاغات, يمنيه, تحالف, انسانيه (سياسيه)
- مرض, علاج, مستشفى, طبيه, صحيه, طفله, مريض, امراض, فيروس, صحه (صحة)
- حفل, مهرج, عاليه, عيد, فعال, مسرح, تهان, عروض, فطر, زوار (ترفيه)
- امطار, رياح, سطحيه, حراره, طقس, افقيه, هطول, كمساع, لارصاد, غزير (طقس)
- فكر, كتاب, دين, لماد, شيء, مسلم, كثير, كاتب, نفس, لان (ثقافة)
- شهداء, شهيد, طوار, عسير, مسجد, استشهاد, تفجير, تعاز, امن, يتعمد (حرب)

FUTURE WORK



Classification Model



Find similar articles



**Thanks for
your listening**