

Sentiment Classification for Arabic Texts

Trainee: Yathrib Alqahtani

Trainer: Prof. Essam Al Daoud

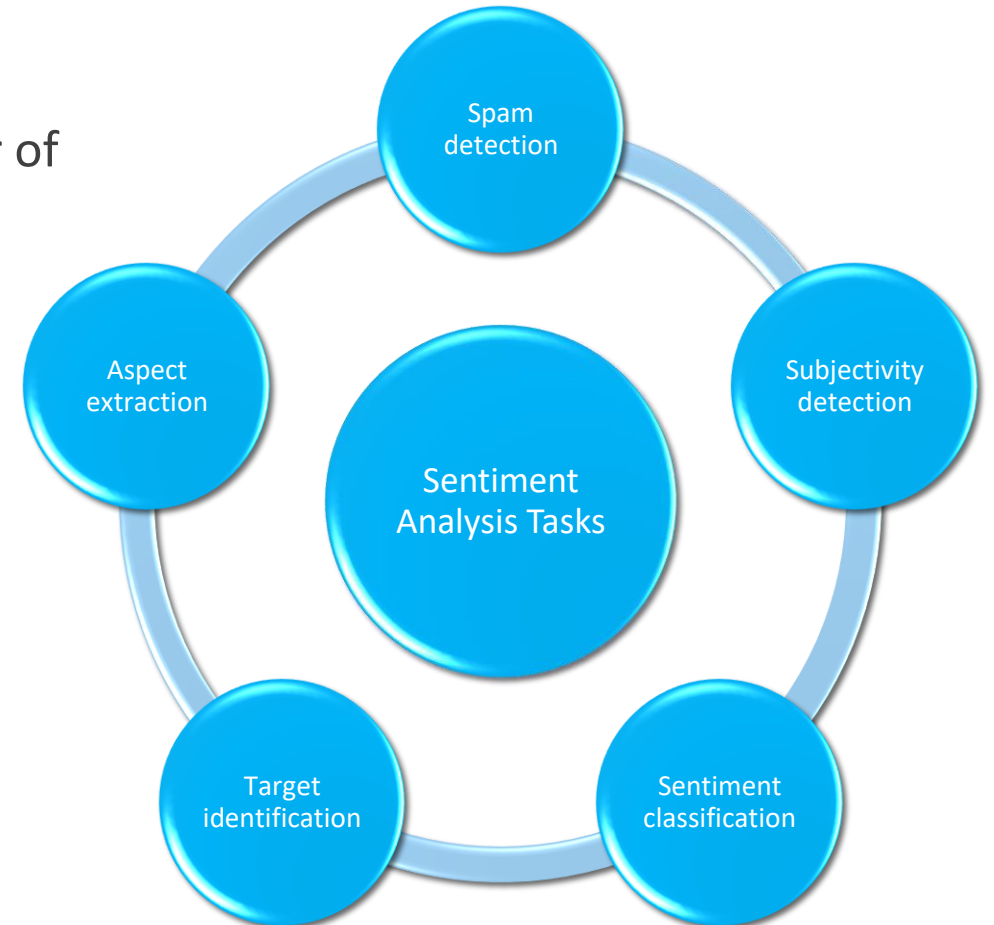


Outline



Introduction

- ❖ The rapid development and spread of Internet technologies has led to a dramatic increase in the number of Internet users viewing and creating content online.
 - A massive amount of opinionated documents.
- ❖ The automatic extraction of opinions has become an active and broad research field in Natural Language Processing (NLP), referred to as **sentiment analysis (SA)**



Sentiment Classification

- ❖ **Sentiment classification** deals with the identification of sentiment as either positive or negative.
- ❖ Applications have reached almost all domains, from finance and healthcare to sports and politics.
- ❖ Recent **supervised** deep learning models have achieved considerable progress.

"I love this movie.
I've seen it many times
and it's still awesome."

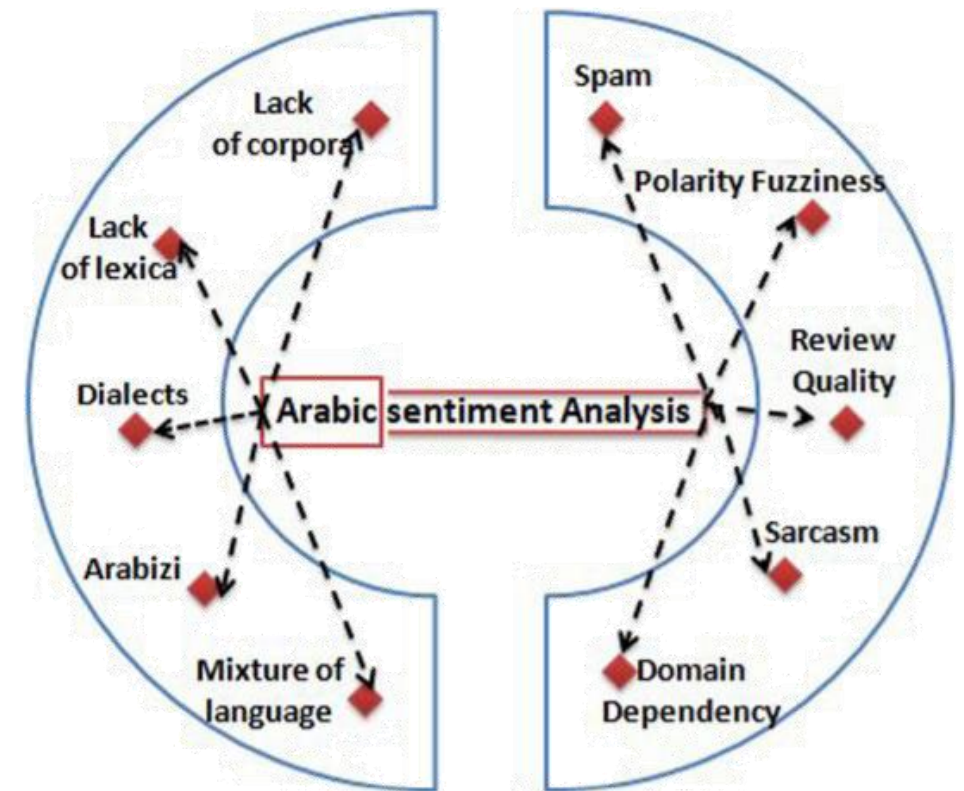


"This movie is bad.
I don't like it it all.
It's terrible."



Sentiment Classification

- ❖ Arabic is the fourth most used Internet language after English, Chinese and Spanish.
- ❖ Around 185 million web users are Arabic speakers.
- ❖ Many Challenges of building models for the Arabic!



Sentiment Classification

❖ Examples of dialects variety in Arabic.

MSA word	Dialectical word	Arabizi	Country
حلو / جميل jameel	حلو	7elew	Lebanon
	حلو	7ilew	Saudi
	حلو	7low/ hlow	Tunisia
جدا jiddan	كثير	ktir	Lebanon
	وايد	wayed	Emirate
	أوي	2awi	Egypt
	برشا	barcha	Tunisia
دراجة Darraja	بسكلات	Besklet	Tunisia
	دراقة	Darraga	Egypt

Sentiment Classification

❖ Arabic vs English!



Dataset

- **Large Multi-Domain Resources for Arabic Sentiment Analysis**

- Modern Standard Arabic (MSA)
- 29444 reviews
- Restaurants, products, hotels, movies.

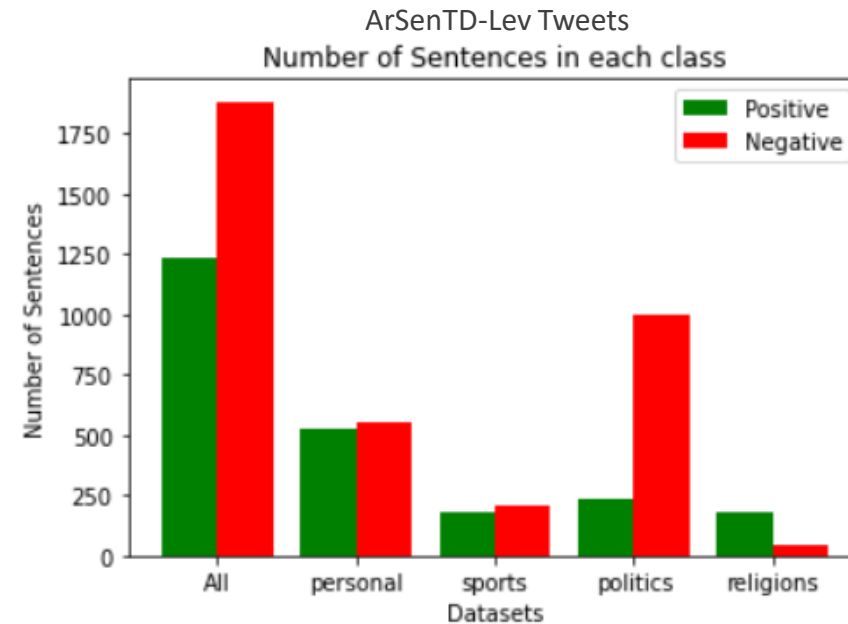
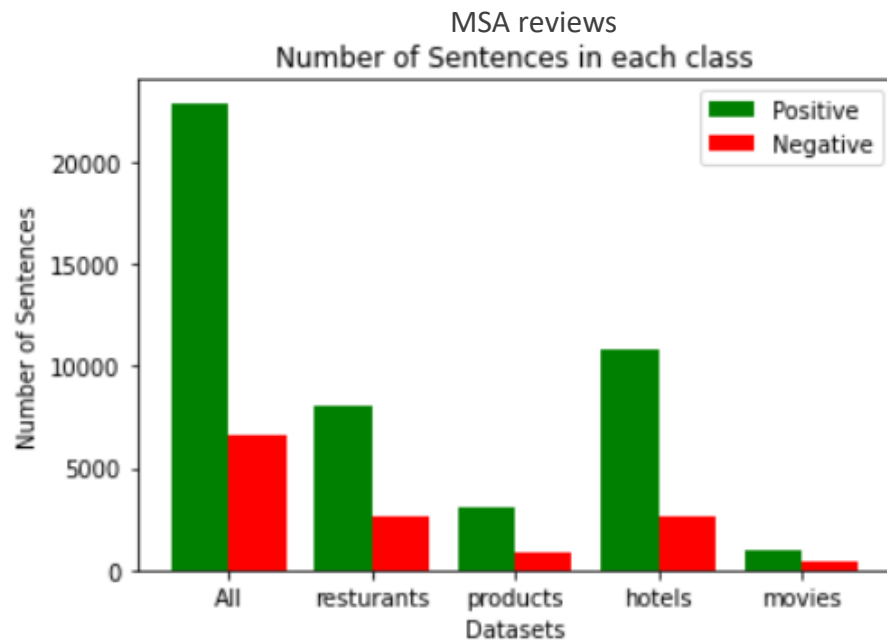
- **ArSenTD-Lev (Arabic Sentiment Twitter Dataset for LEVantine dialect)**

- Dialectal Arabic (Levantine)
- 3115 Tweets
- Personal, sports, politics, religion.

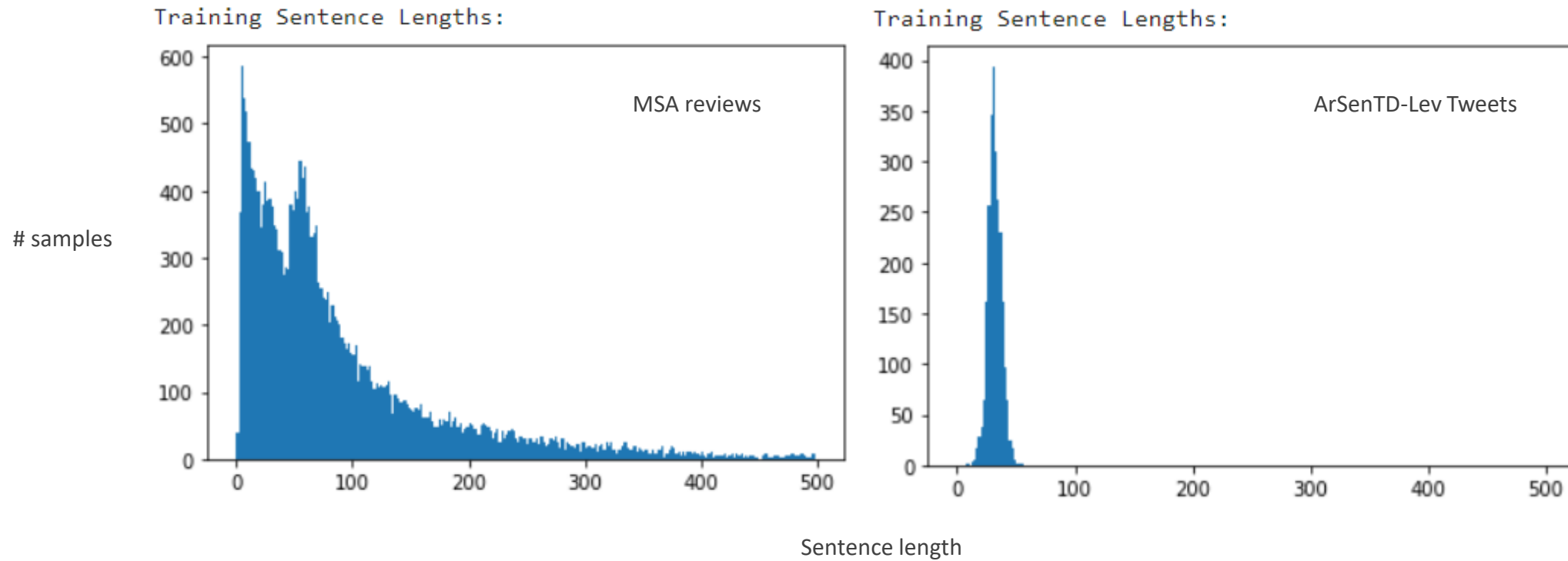
Dataset – Class distributions

❖ Imbalance

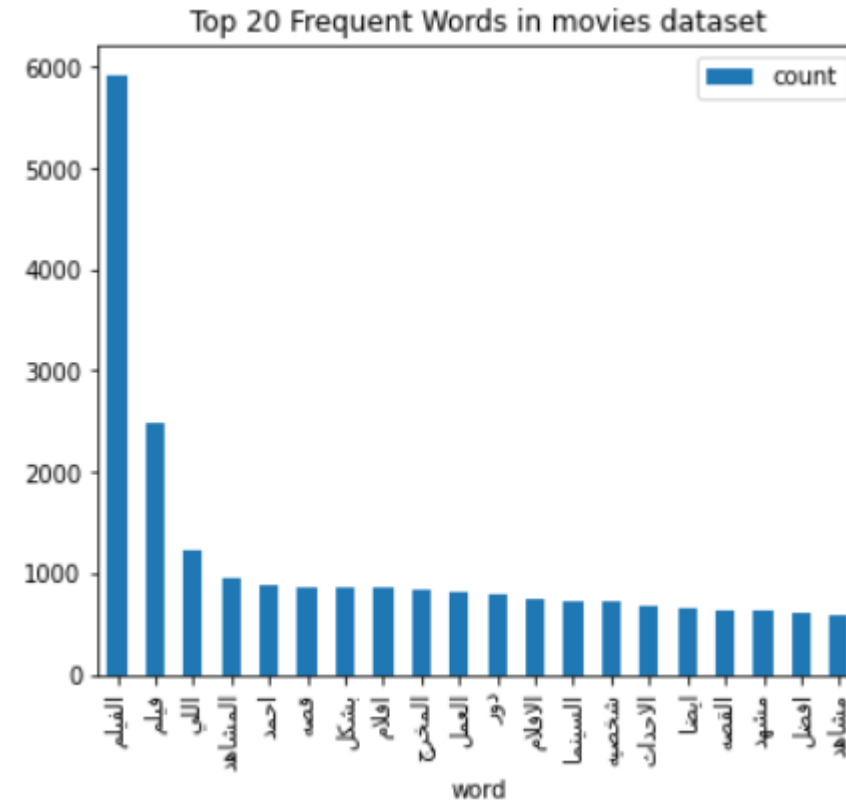
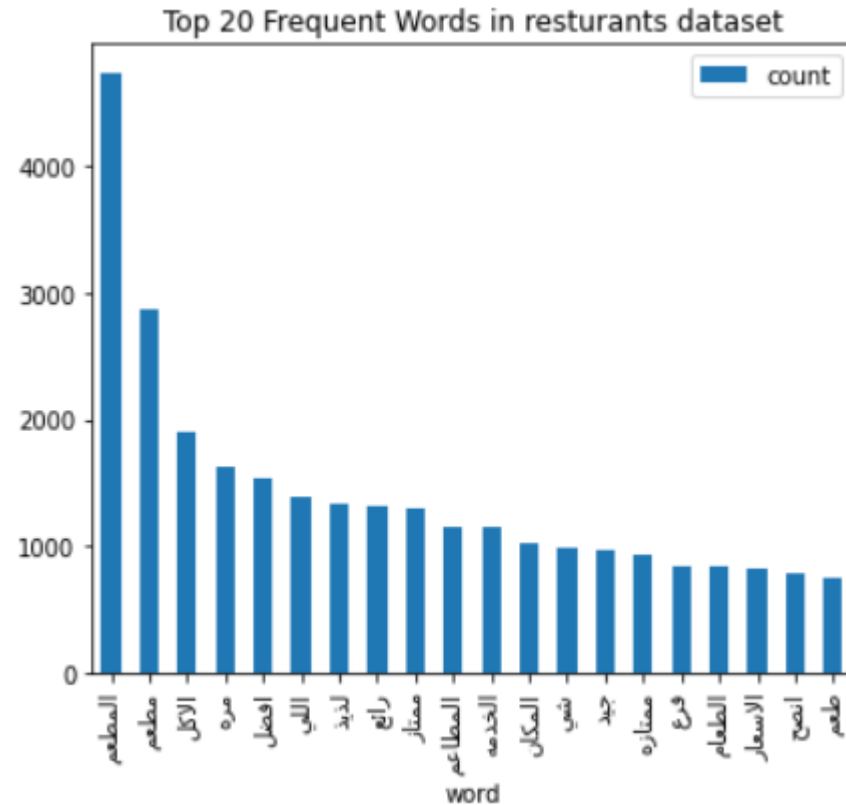
- F1-score is better than accuracy
- Majority differs between reviews and tweets!



Dataset – Sentence lengths



Dataset – Frequent features (MSA)

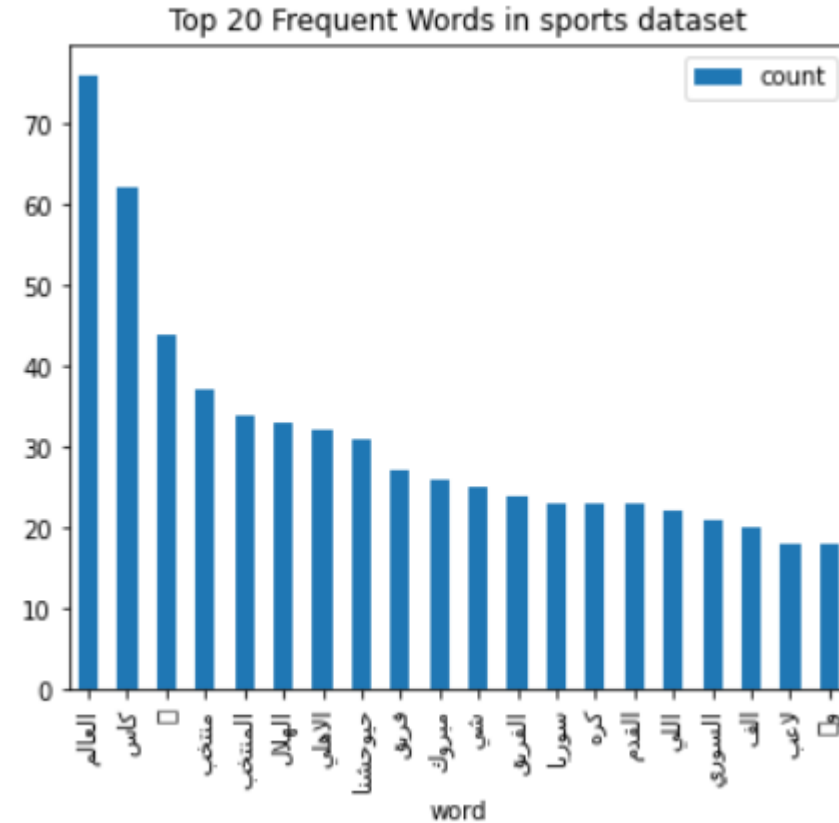
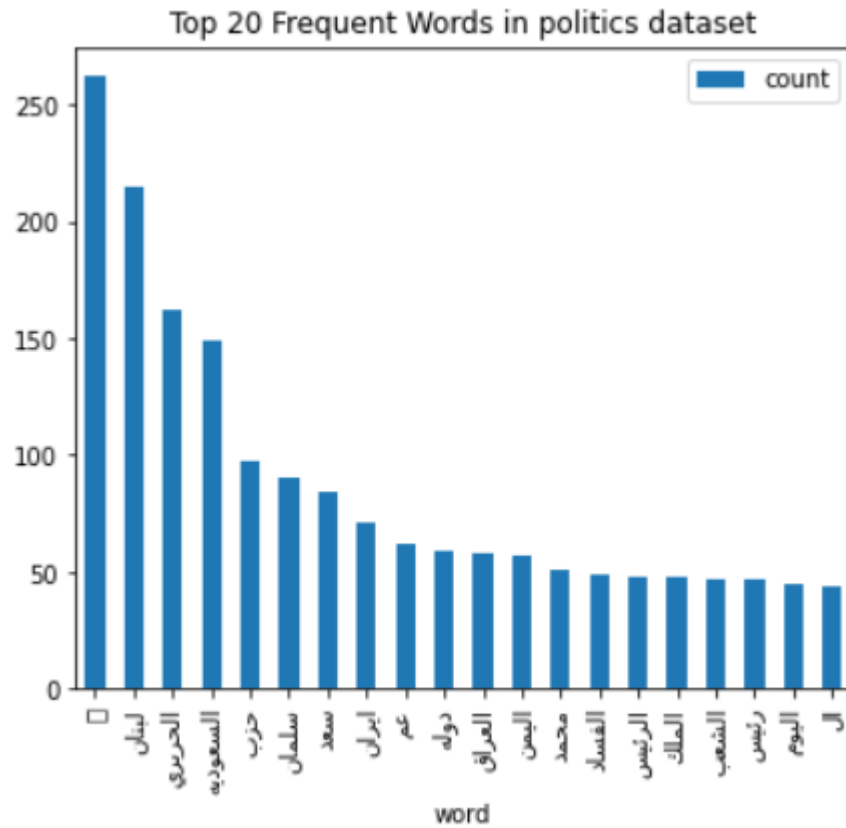


Dataset - Frequent features (MSA)

WordCloud of top 500 words in hotels dataset



Dataset - Frequent features (ArSenTD-Lev)



Dataset - Frequent features (ArSenTD-Lev)

WordCloud of top 2000 words in religions dataset



Cleaning and pre-processing

- Remove unwanted parts
 - Tashkel, URLs, punctuations, emojis, numbers and English letters, twitter hashtags and mentions.
- Remove Arabic stop words.
- Normalization.
- Stemming was tested but didn't show improvements.

Cleaning and pre-processing

الذاكرة اعتقد كريستيانو افضل لاعب العالم كاكا وميسي ثالثاً تم حدث المكس . FinQeمن الذاكرة... 3@
اكثر ما يزعجنا بعد مستوانا خارج ارضنا . هو تمثيل كونتي بعد فوزه علينا @Yousef_MUFC
الف الف مبروك للمنتخب السوري كسب احترام الجميع لعب بروح والعزيمة شكرا شئ قدمته مدار التصنيفات النهائية برغم الظروف #استراليا سوريا
منرجع للمنتخب اصنقائي بالتاتوي اللي كنت فكرهم رمادين او حتى شبيحة طلحو بيكرها المنتخب و شمتائين .. لك هالكم واحد وعامل جموع مئين جي ؟
والله الاكئين غالين لانو تقريبا اول نادين حبيتهم روح سعودي جيت تي شيرت الهلال @yara_CR7
{ اليابان* . _جرحت الهلال_ . ب { نيتسيمورا } . _حرمت الهلال_ من . { ادواردو*
https://t.co/Qjuyf9vztz ...كاس_العالم_بيوحشنا_٢. القوة القوة لا بارك الله بضعف هذا جمهور خط النار لا تورط نفسك معهم . اتمني ي_@sh0uqa2
استسيو خرافي بكلشي اليوم . ما عنده قتالية الدون , بس كثير بيحاول يقدم
https://t.co/jrTPJJjs4h ! لو كان (نصف) مزيج بين الظاهرة وزيدان .. لاحتفلنا بالهف قبل ان يستلم الكرة , ولكن الغياء يلعب دور -
https://t.co/ougMBGzZiN ...تقام عصر اليوم على ملعب #نادي_صنمك مباراة مهمة جداً جداً بالنسبة للنادي وإدارته وجمهوره الكريم . الفوز ولاتشي غير ا
منتخب فرنسا حالياً منتخب جبار بكل خاتة تحسلك نجم , اتوقع يسون شئ بالمونديال القادم
منتخب لبنان يسجل هدفه الأول في الدقيقة الواحدة والعشرين على منتخب كوريا الشمالية ببطولة كأس آسيا على ملعب بيروت البلدي
https://t.co/FFV7ew7jhM هذا منتخب الاقليات الذي يحتر ان الموت فقط هو ما يستحقه الاكثرية كيف كنت تؤيد هذا الفريق يا سني سوري شكرا استراليا
حتى لو خسرنآ. المهم إنكن رفعتوا راسنا قدام العالم كلو و الكل عرف مين المنتخب السوري و مين نسور قاسيون. #نسور_قاسيون شكراً لكن على هالفرة
https://t.co/xKokrpvihn ...في حزة العصر جاني رساله رسالت مؤلمه يا الله تعيني. شئ يخص الهلال وبقي احواله . إصابة إد
https://t.co/ZmVM14s0qA شوف ابن اختي شوف oki_q8e# . LEITOT . هدف لا نراه كل يوم
https://t.co/Ke41GlVB10 ...الخليوي يستاهل فوق ماهو شابل الفريق هالموسم ظهره تبرع للجمهور ب١٠٠ تذكره على حسابه .. واتحدى اذا في لاعب اجنب
والله اتني بقهر انو كل هل الانتقادات ع بنزيمآ اتني واخر اتني زيدان يحكي انو هو المهاجم الافضل ذ . @iAbdullah_24
https://t.co/d05NSSo1DZ ...للعلی عنوان. ملء كل جنا & الشقيق التأهل لكأس العالم في #روسيا-2018 . عشت في الاوطانAM نبارك من القلب❤ لمنتخب #المغرب
ديري_مديري #زيال_مديري شينآ فشينآ المرينقي يودع الدوري الاسباني فارق النقاط مع المتصدر 10 نقاط .. اجساد تتحرك في الميدان دون هوية وروح#
https://t.co/U3BUxVgKKX ...اه يعرف فيار. احضر هادا الفيديو من واحد من داخل الاتحاد الإسباني بيتكلم بشكل واضح عن المهر التحكيمي لصال @OBalah1998
https://t.co/Hggj8BVjr6 ...الحمد لله... ماقصرتوا بالجوم الزعيم.. مباراة قوية اضاع فيها الزعيم فرص كثيرة وسهلة ولكنها تدل على رغبة قوية في الفوز
الف الف مبروك للمنتخب السوري كسب احترام الجميع لعب بروح والعزيمة شكرا على كل شئ قدمته على مدار التصنيفات النهائية برغم الظروف #استراليا سوريا
. كاس_العالم_حيوحشنا_٥. الجمهور لن يتوقف أبداً و سيكون ضد سياسة تركي الشيك #الأهلي خط احمر و لن نرضى بذلك ابداً لن نتوقف#
تالتا أنتم في كل موسم تجملون الدوري العراقي أطول دوري في العالم وهو يتكون من 20فريق. فكيف اذا كان الدوري 30فريق. تحياتي لكم يا فاشلين @IRAQFA.
https://t.co/d05NSSo1DZ ...للعلی عنوان. ملء كل جنا & الشقيق التأهل لكأس العالم في #روسيا-2018 . عشت في الاوطانAM نبارك من القلب❤ لمنتخب #المغرب
الله أنصر نيمار و باوليڤيڤيو والشباب المحترمة على تشيلي، اللهم تعادل بيرو و كولومبيا، اللهم ثبت خطي ميسي و ماسكيرانو والشباب الوردة
والله اتني بقهر انو كل هل الانتقادات ع بنزيمآ وكل اتني . واخر اتني زيدان يحكي انو هو المهاجم الافضل ذ . @iAbdullah_24
تلمو يا #لبنان شو ناقصنا لتصير زي مصر .. الاعب اللبناني مفكر حالي ميسي كثير علينا اسيا يا حرام ... لك نحن وين و المنتخبات الباقية وين
https://t.co/xKokrpvihn ...في حزة العصر جاني رساله رسالت مؤلمه يا الله تعيني. شئ يخص الهلال وبقي احواله . إصابة إد
https://t. ...❤ . كاس_العالم_حيوحشنا_٥.❤❤ : الحين تركي يقول ان المبلغ اودع تاريخ ٢٨-١٠.❤❤❤ : وقبل دخول فترة الحرة.❤❤❤ : صح#
ماعتقد نادم نسعد اء نسه عنا، نحتاج اكثر من ذلك نحتاج مودة فر سافحه . قدمنا: @BarcaArabNet

الذاكرة اعتقد كريستيانو افضل لاعب العالم كاكا وميسي ثالثاً حدث المكس
يزعجنا مستوانا خارج ارضنا تمثيل كونتي فوزه
الف الف مبروك للمنتخب السوري كسب احترام الجميع لعب بروح والعزيمة شكرا شئ قدمته مدار التصنيفات النهائية برغم الظروف استراليا سوريا
منرجع للمنتخب اصنقائي بالتاتوي اللي كنت فكرهم رمادين شبيحه طلحو بيكرها المنتخب شمتائين هالكم وعامل جموع مئين جي
والله الاكئين غالين لانو تقريبا اول نادين حبيتهم روح سعودي جيت شيرت الهلال
اليابان جرحت الهلال نيتسيمورا حرمت الهلال ادواردو
...كاس العالم بيوحشنا القوة القوة بارك الله بضعف جمهور خط النار تورط اتمني ي
استسيو خرافي بكلشي اليوم قتاليه الدون كثير بيحاول يقدم
نصف مزيج الظاهره وزيدان لاحتفلنا بالهف يستلم الكرة الغياء يلعب دور
...تقام عصر اليوم ملعب نادي صنمك مباراة مهمه بالنسبه للنادي وادارته وجمهوره الكريم الفوز ولاتشي ا
منتخب فرنسا حالياً منتخب جبار خاتة تحسلك نجم اتوقع يسون شئ بالمونديال القادم
منتخب لبنان يسجل هدفه الاول الدقيقة الواحدة والعشرين منتخب كوريا الشمالية ببطولة كأس اسيا ملعب بيروت البلدي
منتخب الاقليات يعتبر الموت يستحقه الاكثرية كنت تؤيد الفريق سني سوري شكرا استراليا
خسرنآ المهم انكن رفعتوا راسنا قدام العالم كلو الكل عرف مين المنتخب السوري مين نسور قاسيون نسور قاسيون شكرا لكن هالفرة
...حزه العصر جاني رساله رسالت مؤلمه الله تعيني شئ يخص الهلال وبقي احواله اصابه اد
هدف نراه يوم شوف اختي شوف
...الخليوي يستاهل ماهو شابل الفريق هالموسم ظهره تبرع للجمهور ب١٠٠ تذكره حسابه واتحدى لاعب اجنب
والله اتني بقهر انو الانتقادات بنزيمآ اتني واخر اتني زيدان يحكي انو المهاجم الافضل
...نبارك القلب لمنتخب المغرب الشقيق التأهل لكاس العالم روسيا عشت الاوطان للطي عنوان ملء جنا
ديري مديري مديري شينآ فشينآ المرينقي يودع الدوري الاسباني فارق النقاط المتصدر نقاط اجساد تتحرك الميدان هوية وروح
...يعرف فيار احضر هادا الفيديو داخل الاتحاد الإسباني بيتكلم بشكل واضح المهر التحكيمي لصال
...الحمد لله ماقصرتوا بالجوم الزعيم مباراة قوية اضاع الزعيم فرص كثيرة وسهلة تدل على رغبة قوية الفوز
الف الف مبروك للمنتخب السوري كسب احترام الجميع لعب بروح والعزيمة شكرا شئ قدمته مدار التصنيفات النهائية برغم الظروف استراليا سوريا
كاس العالم حيوحشنا الجمهور يتوقف ابدآ سيكون ضد سياسة تركي الشيك الاهلي خط احمر نرضي ابدآ نتوقف
تالتا موسم تجملون الدوري العراقي أطول دوري العالم يتكون فريق الدوري فريق تحياتي فاشلين
نسور قاسيون بلا يا بنذا غول ثاني سومه هالفريق بعزيمتو وشجاعته يستاهل يكون بكأس العالم حارب الظروف قلوبنا معكو وانتالله حتربحو
...نبارك القلب لمنتخب المغرب الشقيق التأهل لكاس العالم روسيا عشت الاوطان للطي عنوان ملء جنا
الله أنصر نيمار باوليڤيڤيو والشباب المحترمة تشيلي، اللهم تعادل بيرو كولومبيا، اللهم ثبت خطي ميسي ماسكيرانو والشباب الوردة
والله اتني بقهر انو الانتقادات بنزيمآ اتني واخر اتني زيدان يحكي انو المهاجم الافضل
تلمو لبنان شو ناقصنا لتصير زي مصر الاعب اللبناني مفكر حالي ميسي كثير اسيا حرام وين المنتخبات الباقية وين
...حزه العصر جاني رساله رسالت مؤلمه الله تعيني شئ يخص الهلال وبقي احواله اصابه اد
كاس العالم حيوحشنا الحين تركي يقول المبلغ اودع تاريخ ٢٨ صح دخول فترة الحرة صح

Modelling

Linear SVM

TF-IDF with unigrams and bigrams

- Default parameters

AraBERT

AraBERT word embeddings

- Version v02

Bi-LSTM

Arabic word2vec

- MSA reviews → a version pretrained on Wikipedia
- ArSenTD-Lev tweets → a version pretrained on Tweets

Bi-LSTM

Keres tokenizing and padding (indexing)

- **Tools:** Google Collab for cloud processing / Pretrained language model (Hugging face) / Pretrained words embeddings (GitHub)

Results

Dataset	Domain	Size	SVM		AraBERT		Bi-LSTM-W2V		Bi-LSTM tokenizing/padding	
			F1	Acc	F1	Acc	F1	Acc	F1	Acc
MSA	Restaurants	10705	0.92	0.87	0.94	0.91	0.90	0.84	0.90	0.85
	Products	3964	0.91	0.86	0.94	0.90	0.91	0.85	0.89	0.82
	Hotels	13422	0.97	0.95	0.99	0.98	0.92	0.87	0.95	0.93
	Movies	1353	0.89	0.83	0.90	0.85	0.84	0.71	0.83	0.72
	All	29444	0.94	0.90	0.96	0.94	0.59	0.49	0.93	0.88
ArSenTD -Lev	personal	1078	0.79	0.79	0.87	0.90	0.82	0.87	0.75	0.81
	Sports	384	0.82	0.83						
	Politics	1237	0.67	0.90						
	Religions	227	0.92	0.87						
	All	3115	0.76	0.83						

- AraBERT gave the best results in MSA and dialects and in all domains.

Results

Dataset	Domain	Size	SVM		AraBERT		Bi-LSTM-W2V		Bi-LSTM tokenizing/padding	
			F1	Acc	F1	Acc	F1	Acc	F1	Acc
MSA	Restaurants	10705	0.92	0.87	0.94	0.91	0.90	0.84	0.90	0.85
	Products	3964	0.91	0.86	0.94	0.90	0.91	0.85	0.89	0.82
	Hotels	13422	0.97	0.95	0.99	0.98	0.92	0.87	0.95	0.93
	Movies	1353	0.89	0.83	0.90	0.85	0.84	0.71	0.83	0.72
	All	29444	0.94	0.90	0.96	0.94	0.59	0.49	0.93	0.88
ArSenTD -Lev	personal	1078	0.79	0.79	0.87	0.90	0.82	0.87	0.75	0.81
	Sports	384	0.82	0.83						
	Politics	1237	0.67	0.90						
	Religions	227	0.92	0.87						
	All	3115	0.76	0.83						

- Bi-LSTM-w2v didn't outperform simple SVM's in MSA reviews ,but it increased f1 score in ArSenTD-Lev tweets from 0.76 to 0.82

Results

Dataset	Domain	Size	SVM		AraBERT		Bi-LSTM-W2V		Bi-LSTM tokenizing/padding	
			F1	Acc	F1	Acc	F1	Acc	F1	Acc
MSA	Restaurants	10705	0.92	0.87	0.94	0.91	0.90	0.84	0.90	0.85
	Products	3964	0.91	0.86	0.94	0.90	0.91	0.85	0.89	0.82
	Hotels	13422	0.97	0.95	0.99	0.98	0.92	0.87	0.95	0.93
	Movies	1353	0.89	0.83	0.90	0.85	0.84	0.71	0.83	0.72
	All	29444	0.94	0.90	0.96	0.94	0.59	0.49	0.93	0.88
ArSenTD -Lev	personal	1078	0.79	0.79	0.87	0.90	0.82	0.87	0.75	0.81
	Sports	384	0.82	0.83						
	Politics	1237	0.67	0.90						
	Religions	227	0.92	0.87						
	All	3115	0.76	0.83						

- Bi-LSTM-tokenizing/padding tends to beat or perform better than Bi-LSTM-w2v only for large datasets.

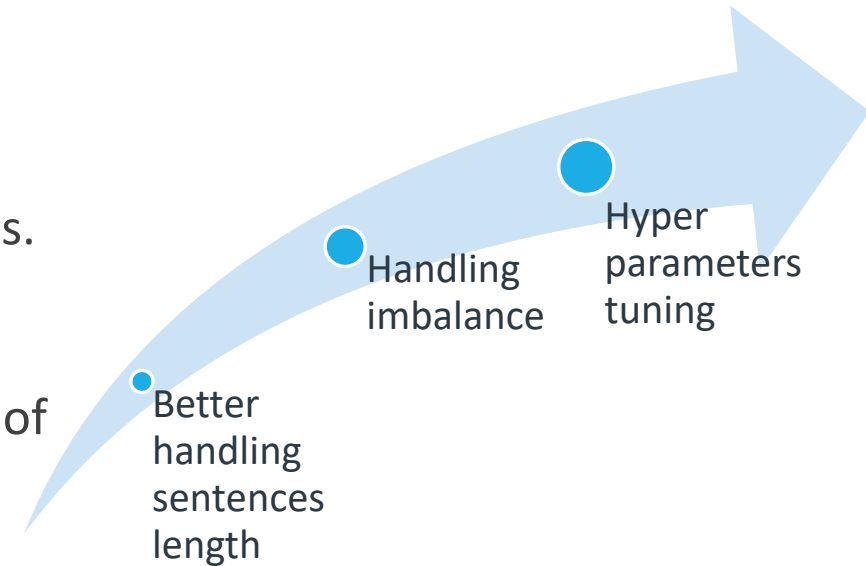
Conclusion

❖ This project focus on binary sentiment classification on Arabic texts.

- MSA and Dialects
- Reviews and tweets
- Different domains

❖ There is substantial room for improvement for the reporting results.

❖ Extra efforts to be done to overcome challenges caused by variety of dialects in Arabic and limited availability of its resources.



Thanks for listening

