# Arabic Dialect Sentimental Analysis

**Prepared by** :     Hazem Mohamed Hosny
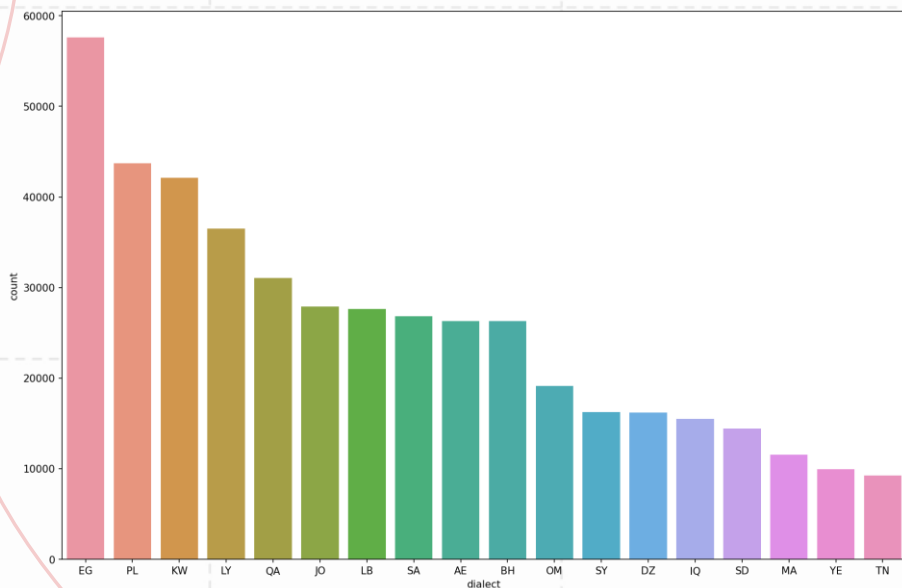
# CONTENTS

# 01  Introduction

As the official language of 22 countries spread across the Middle-East North Africa (MENA) region, <u>Arabic is the 4th most used language on the Internet</u>. <u>Statistics from 2018</u> show 164 million internet users in the Middle East and 121 million internet users in North Africa.

As a language, Arabic has complex morphology and various dialects. The complexity increases significantly when considering the informal nature of social-media text and the distinction between Modern Standard Arabic (MSA) and Dialectical Arabic (DA). Further complicating matters, there are numerous dialects, for example the Egyptian dialect is different from the Levantine dialect which is used in Palestine, Jordan, Syria, and Lebanon. Both of these dialects are also distinct from the Gulf dialect used in Kuwait, Bahrain, Qatar, and the United Arab Emirates.

# 02 DataSet

Our Dataset contains 458197 rows where labels are imbalanced. "EG" is the most major label, Our as "TN" is least label in occurrence. With a distribution as shown in the figure.



| | id | dialect | text |
|---|---|---|---|
| 0 | 1175358310087892992 | IQ | @Nw8ieJUwaCAAreT يغير .. ينتفض .. لكن بالنهاية . |
| 1 | 1175416117793349632 | IQ | @7zNqXP0yrODdRjK ح .. يعني هذا محسوب على البشر... |
| 2 | 1175450108898565888 | IQ | @KanaanRema خليجي كلامه من مبين |
| 3 | 1175471073770573824 | IQ | @HAIDER76128900 الحلوه وروحك مرورك يسلملي💐 |
| 4 | 1175496913145217024 | IQ | @hmo2406 محمد اخ الغيبه هل وين 🌸💐 |
| ... | ... | ... | ... |
| 458192 | 1019484980282580992 | BH | @Al_mhbaa_7 باسطانا اللي منك مبسوطين😜 |
| 458193 | 1021083283709407232 | BH | @Zzainabali @P_ameerah يختي ابش ماينده والله |
| 458194 | 1017477537889431552 | BH | @Al_mhbaa_7 مس احنا مننا تهربي حنا لك عملنا شو... |
| 458195 | 1022430374696239232 | BH | @haneenalmwla وبالعافيه فيها يبارك الله 😊😊😊 |
| 458196 | 1022409931029458944 | BH | @jolnar121 سحليه لك بتطلع ي ضيفي السحله😅😅 |

458197 rows × 3 columns

# 03   Proposed Methodology

**1. PreProcessing "text" which includes:**

- Remove URLs
- Remove RT and cc
- Remove @username
- Remove new line in both Windows or Apple
- Clean hashtags
- Clean emojis where translating native emojis (such as: ':)' , ':(', …etc to words), and remove emojis belongs to emojis library or other utf-8 format.
- Remove punctuations
- For Arabic normalizations: remove stopwords, normalize Arabic letters (such as: transform ى to ي, transform ة to ه, etc…), Using pyarabic Package to remove Tashkel, and finally remove longation.

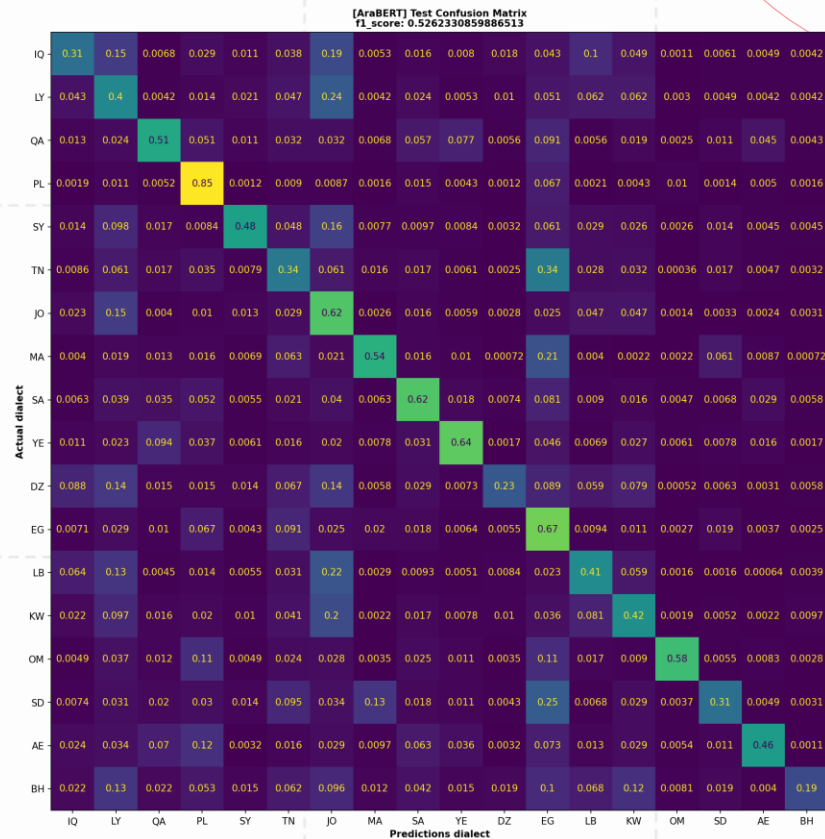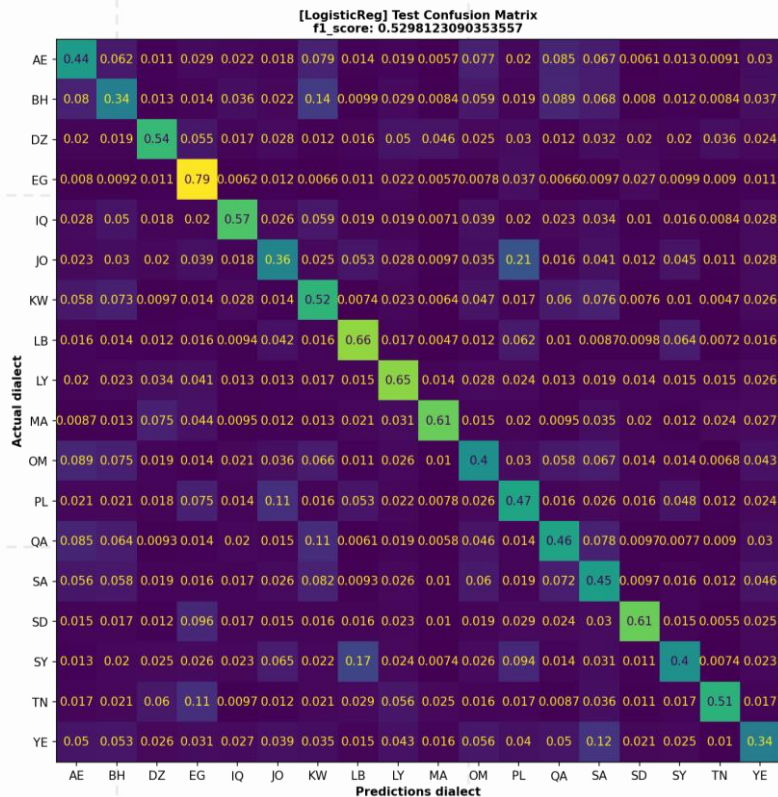| | id | dialect | text | word_count | char_count | avg_char_per_word | stopwords | emoji_count |
|---|---|---|---|---|---|---|---|---|
| 0 | 1175358310087892992 | IQ | بالنهايه ينتفض يعير | 3 | 19 | 5.666667 | 0 | 0 |
| 1 | 1175416117793349632 | IQ | يعني محسوب البشر حيونه ووحشيه وتطلبون العرب يحترمكم وي،من بدينكم ولاينعتكم بالارهاب | 12 | 83 | 6.000000 | 0 | 0 |
| 2 | 1175450108898565888 | IQ | مبين كلامه خليجي | 3 | 16 | 4.666667 | 0 | 0 |
| 3 | 1175471073770573824 | IQ | يسلملي مرورك وروحك الحلوه | 4 | 25 | 5.500000 | 0 | 0 |
| 4 | 1175496913145217024 | IQ | وين الغيبه اخ محمد | 4 | 18 | 3.750000 | 0 | 0 |

# 03 Proposed Methodology

**2. For ML Models:**

- TF-IDF Vectorizer with unigram word analyzer extraction.
- A Stratified Train test split is considered to be more fit for unbalanced data for more generalized validation score.
- A Lasso Logistic Regression were tested along the TF-IDF Vectorizer on the cleaned data and with class weight to overcome unbalance data
- Complement Naive Bayes classifier were tested too. It is particularly suited for imbalanced data sets.
- Stochastic gradient descent Classifier (SGDClassifier) with Hinge loss as parameter that performs as Linear SVM were tested with class weights setted.

**3. For Deep Learning Model:**

- An AraBERT version 2 model were tested from "aubmindlab/bert-base-arabertv02" without FarasaSegmenter over the original data with "ArabertPreprocessor()" object, offered Transformers Package over Pytorch.

# 04  Metrics and Results



[LogisticReg] Test Confusion Matrix
f1_score: 0.5298123090353557

[AraBERT] Test Confusion Matrix
f1_score: 0.5262330859886513

# 05 Deployment

**In this section, 5 main libraries were used :**
1. scikit-learn for tfidf vectorizer and Logistic regression Model
2. Pyarabic for striping tashkeel
3. Nltk for Arabic stop words
4. Transformers from Hugging Face to add our AraBERT model to the pipeline to predict dialect of text.
5. Flask / Fast API it's a two technologies that used to implement our webserver app

FastAPI

HUGGING FACE

## default

| GET | / Home |

| POST | /Predict_ml Arabic Tweet Dialect Prediction (Logistic Regression Model) |

| POST | /Predict Arabic Tweet Dialect Prediction (AraBERT best label) |

| POST | /Predict_all_labels Arabic Tweet Dialect Prediction with (AraBERT all labels) |

# THANKS!

Do you have any questions?