# NLP Milestone1

Sohaila Ibrahim  
**ID:** 52-21225

Nada Elbehery  
**ID:** 52-8973

March 10, 2025

## Contents

# 1 Overview

In Milestone 1, we focused on conducting a thorough data analysis and review of the dataset. This process involved exploring the structure and characteristics of the data, identifying patterns, trends, and any potential inconsistencies or irregularities. Additionally, we applied preprocessing and cleaning techniques to transform the raw dataset into a structured and well-organized format suitable for machine learning (ML) tasks.

# 2 Dataset Description

Our dataset consists of multiple text files containing transcripts in the Egyptian dialect. These texts are primarily sourced from various episodes of *El Mokhbir El Iqtisadi* (The Economic Informant), a popular Arabic-language YouTube channel that discusses global economic and political issues. The dataset covers a diverse range of topics, including:

- International trade and financial crises.

- Energy policies and resource control.

- Geopolitical conflicts and their economic implications.

- Technological advancements and digital economy dynamics.

- The economic rise and fall of corporations and industries.

- Social and economic trends in different regions.

Some files discuss historical events and economic strategies of major powers such as the United States, China, and Russia, while others focus on inflation, supply chain disruptions, and corporate financial maneuvers.

# 3 Methodology

Natural Language Processing (NLP) for Arabic text involves several preprocessing steps to enhance text representation for further analysis. This report details the sequential flow of preprocessing techniques applied to Arabic text, followed by TF-IDF weighting, Vector Space Model representation, and LDA topic modeling.

# 4 Preprocessing Flow

Farasa is a library designed for Arabic Natural Language Processing (NLP). It provides various tools for tasks such as tokenization, part-of-speech tagging, named entity recognition, stemming, diacritization, and segmentation. The preprocessing pipeline consists of multiple steps that leverage Farasa's modules to prepare Arabic text for analysis.

## 4.1 Segmentation

Segmentation is the process of breaking down words into smaller meaningful units, such as prefixes, stems, and suffixes. The `FarasaSegmenter` module is responsible for segmenting Arabic words into their morphemes.

```
from farasa.segmenter import FarasaSegmenter

def segment(text):
    segmenter = FarasaSegmenter()
    segmented = segmenter.segment(text)
    return segmented
```

## 4.2 Stemming

The `FarasaStemmer` module reduces words to their root forms.

```
from farasa.stemmer import FarasaStemmer

def stem(text):
    stemmer = FarasaStemmer()
    stemmed = stemmer.stem(text)
    return stemmed
```

## 4.3 Part-of-Speech Tagging

The `FarasaPOSTagger` module assigns part-of-speech (POS) tags to words in Arabic text.

```
from farasa.pos import FarasaPOSTagger

def tag(text):
    tagger = FarasaPOSTagger()
    tagged = tagger.tag(text)
    return tagged
```

## 4.4 Named Entity Recognition (NER)

The `FarasaNamedEntityRecognizer` module detects named entities such as persons, organizations, and locations in Arabic text.

```
from farasa.ner import FarasaNamedEntityRecognizer

def recognize(text):
    named_entity_recognizer = FarasaNamedEntityRecognizer()
    named_entity_recognized = named_entity_recognizer.recognize(text)
    return named_entity_recognized
```

## 4.5 Stopword Removal

Eliminating frequently occurring words with little semantic meaning. Predefined Arabic stop words were collected from the library arabicstopwords and extended with additional stop words to handle the Egyptian dialect.

## 4.6 Removing Non-Arabic Characters

Filtering out non-Arabic symbols to clean the text.

## 4.7 Sample text

For example, given the input text:

إزايك عامل إيه أنا النهاردة كان عندي شغل كتير بس الحمد لله خلصت كل حاجة بعد كده رحت الكافيه شوية عشان أريح دماغي إنت عامل إيه في جديد

- After preprocessing, the text is transformed into:

إزايك عامل أي أنا نهارد ي شغل كتير حمد ل الله خلص كل حاج كد ه رح كافي شوي أريح دماغ إن عامل أي جديد

- After tagging, the text is transformed into:

ال + ADV/بس NOUN-MS/كتير NOUN-MS/شغل PRON/ي+ NOUN-MS/عند V/كان DET+NOUN+NSUFF-FD/ال + نهارد ++ PART/أنا PRON/+أي NOUN-MS/عامل NOUN-MS/إزايك S/S
ال + PRON/رح V/ت+ PRON/ي+ NOUN-MS/ك NOUN-MS/كد PRON/+ت NOUN-MS/خلص NOUN+NSUFF-FS/بعد NOUN-MS/حاج PRON/كل NOUN-MS/الله PREP/ل+ DET+NOUN-MS/حمد
NOUN-MS E/E/جديد PREP/في NOUN-MS/أي NOUN-MS/عامل PRON/عامل V/ت+ إن ADJ-MS/دماغ NOUN-FP/دماغي ADJ-MS/أريح ADJ+NSUFF-FS/عشان شوي+ DET+NOUN-MS/الكافيه

- After entity recognition, the text is transformed into:

O/إيه O/عامل O/إنت O/دماغي O/أريح O/عشان O/شوية O/الكافيه O/رحت O/كده O/بعد O/حاجة O/كل O/خلصت O/لله O/الحمد O/بس O/كتير O/شغل O/عندي O/كان O/النهاردة O/أنا O/إيه O/عامل O/إزايك O/في O/جديد

Another example, input text:

سافر محمد بن سلمان إلى الرياض يوم الإثنين لحضور اجتماع مع ممثلي شركة أرامكو. ناقش الاجتماع مستقبل الطاقة في المملكة العربية السعودية وتأثير الأسعار العالمية. كما التقى محمد بالرئيس التنفيذي لشركة سابك وأكد على أهمية الاستثمارات المستقبلية. في اليوم التالي، زار جامعة الملك سعود والتقى بالطلاب والأساتذة لمناقشة أحدث الأبحاث العلمية.

- After preprocessing, the text is transformed into:

سافر محمد سلمان روضة يوم إثن حضور اجتماع مثل شرك أرامكو ناقش اجتماع مستقبل مملك عربي سعودي تأثير سعر عالمي التقى محمد رئيس تنفيذي شرك سابك أكد أهمي استثمار مستقبلي يوم تالي زار جامع ملك سعود التقى طالب أساتذ مناقش أحدث بحث علمي

- After tagging, the text is transformed into:

NOUN-MS/حضور PREP/ل+ NOUN-MS/محمد V/سافر S/S بين/NOUN-MS سلمان/NOUN-MS إلى/PREP رياض /ال+ ل/DET+NOUN-MS يوم/NOUN-MS إثنين /ل+ DET+NOUN-MS
اجتماع/NOUN-MS مع/NOUN-MS مستقبل/NOUN-MS روي/NOUN-MS شرك /NSUFF ++ NOUN+NSUFF-FS/ أرامكو NOUN-MS /PUNC E/E
ال+ DET+ADJ+NSUFF-FS عربي /PREP+ مملك /ل+ DET+NOUN+NSUFF-FS في/ال+ DET+NOUN-MS طنق ++ NOUN-MS/مستقبل ل/DET+NOUN-MS ناقش/V اجتماع /S/S
سعودي/DET+ADJ+NSUFF-FP CONJ/و تأثير/NOUN-MS أسعار /ل+ DET+NOUN-MP علمي /ال+ DET+ADJ+NSUFF-FS /. PUNC E/E
V/أكد CONJ/و NOUN-MS/سابك NOUN+NSUFF-FS /ل+ PREP ++ شرك/NOUN-MS تنفيذي/DET+NOUN-MS ل/DET+ADJ-MS محمد V/التقى الثقى /PART+PREP كما/S/S
على/NOUN+NSUFF-FS أصمي ++NOUN/استثمار دات /ل+ DET+NOUN+NSUFF-FP مستقبلي /ل+ DET+ADJ+NSUFF-FD /. PUNC E/E
ل+ PREP ب /V التقى CONJ/و NOUN-MS/سعود ل+ DET+ADJ-MS ملك /ل+ DET+NOUN+NSUFF-FS جامع ++ يوم /V/زار PUNC/، دالي /ل+ DET+ADJ-MS في/PREP S/S
ال+ علمي/DET+NOUN-MP ل/+ بحث NOUN+NSUFF-FS أحدث/NOUN-MS مناقش/NOUN+NSUFF-MP /ل+ PREP ++ أساتذ CONJ/و طالب/DET+NOUN-MS
+ال/DET+ADJ+NSUFF-FP /. PUNC E/E

- After entity recognition, the text is transformed into:

B-ORG ,/B-ORG. أر امكو O/شركة O/مثلي O/مع O/اجتماع O/لحضور O/الإثنين O/يوم B-LOC/الرياض O/إلى I-PERS/سلمان I-PERS/بين B-PERS/محمد O/سافر
O ,/العلمية O/الأسعار O/وتأثير I-LOC/السعودية I-LOC/العربية B-LOC/المملكة O/في O/الطاقة O/مستقبل O/الاجتماع O/ناقش
O ,/المستقبلية O/الاستثمار O/أهمية O/على O/وأكد O/سابك O/شركة O/التنفيذي O/بالرئيس B-PERS/محمد O/التقى O/كما
O ,/O/العلمية O/الأبحاث O/أحدث O/لمناقشة O/وا الأساتذة O/بالطلاب O/والتقى I-ORG/سعود I-ORG/الملك B-ORG/جامعة O/زار O ,/O/التالي O/اليوم O/في

# 5 TF-IDF Weighting

TF-IDF (Term Frequency-Inverse Document Frequency) measures the importance of words in a corpus.

# 6 Vector Space Model

The Vector Space Model represents documents numerically, facilitating similarity measurements. For example, given a query vector:

```
Query vector: [0.19348638238589283]
Document _____.txt vector: [0]
```

This allows us to identify the most relevant documents.

# 7 LDA Topic Modeling

Latent Dirichlet Allocation (LDA) is an unsupervised machine learning algorithm used for topic modeling, which aims to discover hidden topics within a collection of documents. It assumes that each document is a mixture of multiple topics, and each topic is characterized by a distribution of words. The function

**perform_lda(dic)** handles the creation of the model. It takes as an input a dictionary of word counts per document, structured as:

$$\texttt{dic} = \{\text{file\_name}_i : \{\text{word}_j : \text{count}_{ij}\}\} \tag{1}$$

where each document (file) contains words mapped to their frequency.

## 7.1 Processing Steps

1. Extracts unique words from each document:

$$\texttt{processed\_docs} = [\text{list of words from each document}] \tag{2}$$

2. Creates a Gensim dictionary mapping words to unique IDs.

3. Converts documents into a Bag-of-Words (BoW) representation:

$$\texttt{corpus} = [(\text{word\_id}, \text{frequency}) \text{ tuples}] \tag{3}$$

4. Trains an LDA model with:

   - $\texttt{num\_topics} = 6$ (number of topics)
   - $\texttt{passes} = 10$ (iterations for better topic assignment)
   - $\texttt{random\_state} = 46$ (for reproducibility)

## 7.2 Enhancing Topic Modeling with Named Entity Recognition (NER)

Named Entity Recognition (NER) can improve topic modeling by refining the textual representation of documents before applying Latent Dirichlet Allocation (LDA). We explore three different approaches to integrating NER:

1. **Use Only Named Entities:** In this approach, topic modeling is performed solely on named entities extracted from the documents.

2. **Boost Named Entities:** Here, named entities are given higher weights in the document representation, making them more influential in topic generation.

3. **Entities with Top 20 Words:** This approach combines named entities with the 20 most frequent words per document to balance contextual and specific information.

# 8 Conclusion

This NLP pipeline successfully processes Arabic text, extracts key information using TF-IDF and NER, and represents documents numerically for further analysis using the Vector Space Model and LDA topic modeling.