

## import libraries

```
In [4]: import pandas as pd
import matplotlib.pyplot as plt
import os
import nltk
from nltk.tokenize import word_tokenize
from nltk.util import ngrams
import matplotlib.pyplot as plt
from collections import Counter
import string
```

```
In [5]: nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\LENOVO\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

```
Out[5]: True
```

## read data and Concatenate all csv files

```
In [6]: # Assuming all CSV files are in the same directory
data_directory = "C:/Users/LENOVO/Downloads/archive/stories/"
csv_files = [file for file in os.listdir(data_directory) if file.endswith('.csv')]

# Initialize an empty list to store DataFrames from each CSV
dfs = []

for file in csv_files:
    file_path = os.path.join(data_directory, file)
    df = pd.read_csv(file_path)
    dfs.append(df)

# Concatenate all DataFrames into a single DataFrame
combined_data = pd.concat(dfs, ignore_index=True)
```

```
In [7]: combined_data.head()
```

Out[7]:

	Unnamed: 0	id	title	date	author	story	topic
0	0	f06aa998054e11eba66e646e69d991ea	بيت الشعر" يسائل وزير الثقافة" عن كوابيس سوداء	الجمعة 02 أكتوبر 2020 - 23:19	هسبريس من الرباط	وجه "بيت الشعر في المغرب" ...إلى وزير الثقافة وال	art-et-culture
1	1	f1cf1b9c054e11ebb718646e69d991ea	مهرجان "سينما المؤلف" يستحضر روح ثريا جبران	الجمعة 02 أكتوبر 2020 - 07:26	هسبريس من الرباط	في ظلّ استمرار حالة الطوارئ...الصحية المرتبطة بج	art-et-culture
2	2	f2d282a4054e11eb800f646e69d991ea	فيلم "بدون عنف" لهشام...العسري ..كعب الحذاء ووا	الجمعة 02 أكتوبر 2020 - 04:00	عفيفة الحسينات	تشير مشاهدة فيلم قصير ضمن...الثلاثية الأخيرة للم	art-et-culture
3	3	f3f46cac054e11eba403646e69d991ea	تنين ووهان" .. مريم أيت أحمد" ...توقع أولى "روا	الجمعة 02 أكتوبر 2020 - 02:00	حاوزها: وائل بورشاشن	من قلب أيام "الحجر"، رأت التّور ... الفصول	art-et-culture
4	4	f50f0476054e11eba31b646e69d991ea	مسكر يتخلّى عن دعم "الوزارة" "بسبب" الجمهور	الخميس 01 أكتوبر 2020 - 19:40	هسبريس من الرباط	أعلن الفنان المغربيّ سعيد...مسكر تخليه عن مبلغ ا	art-et-culture

In [8]:

combined\_data.columns

Out[8]:

Index(['Unnamed: 0', 'id', 'title', 'date', 'author', 'story', 'topic'], dtype='object')

In [9]:

print(combined\_data['story'][0])

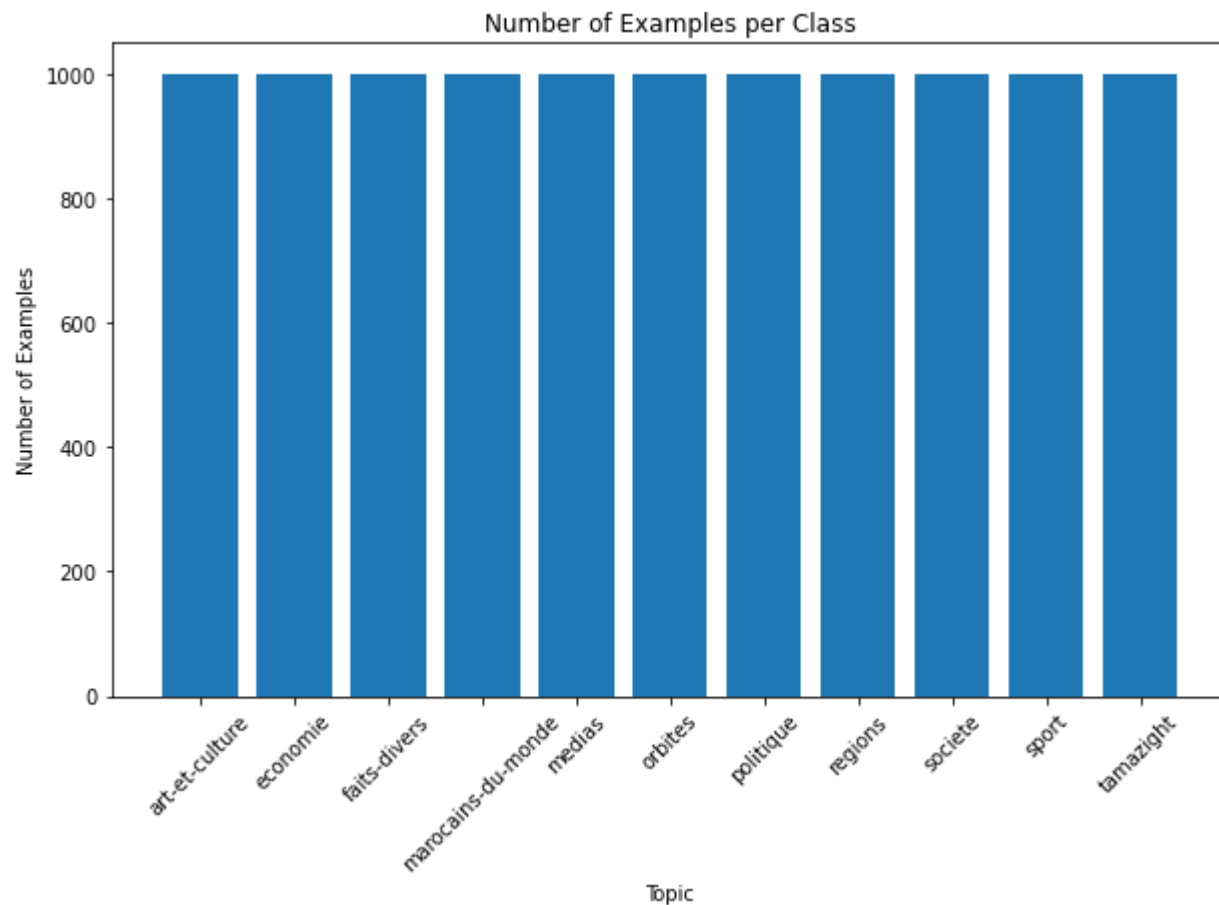
وجه "بيت الشعر في المغرب" إلى وزير الثقافة والشباب والرياضة رسالة موسومة بـ"لماذا تحولت أحلام بيضاء إلى كوابيس سوداء؟"، أشار من خلالها إلى أن "بيت الشعر تأسّس سنة 1996 بهدف تّ تحقيق جُملة من الأهداف التي تروم جميعها تعزيز مكانة الشّعر في المجتمع والحياة وترسيخ مكانته بين الناس كحاملٍ لقيم الحلم والخيال؛ فمن تلك الأهداف توطين الشعر المغربي في المقررات الدراسية وتّ شجيع التلاميذ والطلبة على قراءته وتذوّق جمالياته، خاصّة في اللحظة التي ينتصرُ فيها لكل ما هو مدهش وإنساني". وجاء في الرسالة، التي توصلت هسبريس بنسخة منها، أن "بيت الشعر نجح في المغّ رب خلال مسيرته الطويلة التي تمتدّ على مدى ربع قرن في إقناع المنظمة العالمية للتربية والثقافة والعلوم (اليونيسكو) بإحداث يوم عالمي للشّعر عن طريق المقترح الذي تقدّم به، والذي تبنته الحكومة ا لمغربية في عهد الراحل عبد الرحمان اليوسفي، كما نجح في ضمان مكانة عالمية لجانزته الشعرية المعروفة "الأركانّة" التي صار شعراء العالم يتطلعون للفوز بها، علاوة على انتظام منشوراته الشعرية ومجلته الرصينة "البيت"، واستدامة برامجها الشعرية بالتعاون مع عدد من الشركاء الذين آمنوا بجدية مشروع مؤسستنا وجودة مستواه فنيا وجماليا". وأشارت الرسالة ذاتها إلى أنه "سعيًا من بيت الشعر ف ي المغرب إلى تقديم خدمة جديدة للشعر المغربي، يربطه بالفضاء الإبداعي العام وخاصة الموسيقي والغنائي منه، واستعادة لحظات طيبة الذكر تعانق فيها شعُرنا المغربي (عبد الرفيع جواهري وإدريس الج اي والخمار الكنوني وحسن المفتي وأحمد الطيب لعلج وعلي الحداني...) مع أوتار وألحان (عبد السلام عامر وعبد النبي الجيراري وحسن القديري...)، تقدّم بمشروع فني تحت عنوان أحلام بيضا ء". وأوضح "بيت الشعر" أن المشروع "عبارة عن مقطوعات غنائية جعلت من الشعر المغربي متنا لها، بهدف استعادة ماضٍ جميل، أثرى خلاله الشعراء المغاربة السّجل الشعري للأغنية المغربية، فقد منحوها قصائد تحوّلت، من خلال أوتار ملّحنين مُقتردين، إلى أغانيّ تتردّد على ألسنة وشفاه الناس في المناسبات والأفراح والأعراس، بل إنّ بعضها سيحظّى باهتمام مطربين مشاركة، أعادوا أدائها وتسجيله ا من جديد بأصواتهم". وورد في الرسالة ذاتها أنه "نجاحٌ ما كان له أن يتحقّق لولا أنّ هؤلاء الشعراء كتبوا قصائدهم من داخل تجربتهم الإنسانية وفي أفق الرؤية التي امتلكوها تجاه اللغة والمجتمع والك ون، وليس تحت الطلب أو إرغامات سوق الغناء". وحسب الرسالة ذاتها، فإنّه "لأوّل مرّة في تاريخ برنامج الوزارة لدعم الأغنية المغربية تتقدّم حياة متخصصة كبيت الشعر في المغرب بمشروع ذي رؤية مندمجة، تربط الشعر المغربي بأفقه الفنّي من أجل أن يحظى بمساحةٍ أوسع للانتشار عبر الأغنية؛ حيث يصيرُ بمقدور الجمهور المغربي أن يتعرّف على شعرائنا المغاربة، ليس من خلال دواوينهم الشعرية أو عبر أمسياتهم الثقافية، بل من خلال الأغنية ك لحظة فنية تتجمّع وتنصهر فيها عدّة أبعاد شعرية، ولحنية موسيقية، وطربية غنائية...". وقد سبق للفنانة صباح زيداني أن جعلت من الشعر المغربي أفقا لتفكيرها واشتغالها، تضيف الرسالة ذاتها، "عندما قدمت بعضا من نصوصه الجميلة للشعراء بوجمعة العوفي وعبد الهادي السعيد، كما خاضت تجربة فنية مع الشاعر المغربي الكبير عبد الله زريقة بمعية ث لة من الشعراء والموسيقيين الأجانب؛ وهو ما يجعل وجودها ضمن هذا المشروع تميّنا لهذه الإرادة التي تلقّي فيها برغبة بيت الشعر في المغرب، في أن يكون شعُرنا المغربي حاضرا في مختلف الحوام ل التي تتيح له الذبوع والانتشار (مسرح، تشكيل، أغنية...)". وأكّد "بيت الشعر" أن "تشبيك الفنّون في ما بينها وتجسير الصلات بين مكوناتها أحد أهم مخرجات هذا المشروع، الذي يجمع شعراء مر موقين، بملّحنين مقتدرين وموزعين أكفاء وعازفين ماهرين، علاوة على الحضور الواعي والباذخ للفنانة صباح زيداني؛ جميع هؤلاء انخرطوا في هذا المشروع من أجل أنّ تعيد للكلمة الشعرية بريقها في ل حظة تعالقها مع اللحن والموسيقى". واختتم "بيت الشعر" رسالته بالتساؤل: "لماذا تحولت الأحلام البيضاء إلى كوابيس سوداء؟ ولماذا خيّبت اللجنة التي شكلتموها من أجل دعم الأغنية المغربية والارتقاء ب مستواها ظنكم وظن الشعب المغربي الذي يَؤمّل من ماله العام سياسة ثقافية فاشلة؟ وماذا ستفعلون لإنصاف بيت الشعر في المغرب وتحقيق المساواة وتكافؤ الفرص التي ينص عليها دستور المملكة؟ وكيف ستردون من خلال إجراءات عملية وتدابير إدارية عن الضجة- الفضيحة التي خلفتها نتائج هذه السنة في مجال دعم الأغنية؟"، مضيفا "إننا نطعن في نتائج هذه الدورة من برنامج دعم الموسيقى والأغني ءة"، وإننا ننتظر ما ستقومون به

# exploratory data analysis

```
In [19]: # Create a custom-sized figure
plt.figure(figsize=(10, 6)) # Change the values (width, height) as needed

# Plot number of examples for each class
class_counts = combined_data['topic'].value_counts()
plt.bar(class_counts.index, class_counts.values)
plt.xlabel('Topic')
plt.ylabel('Number of Examples')
plt.title('Number of Examples per Class')
plt.xticks(rotation=45)

# Display the plot
plt.show()
```



```

In [11]: def preprocess_text(text):
# Tokenize the text into words
tokens = word_tokenize(text)
# Define a list of Arabic punctuation characters
arabic_punctuation = ['`', '-', ',', ';', '"', '.', ':']
# Combine Arabic and ASCII punctuation characters
all_punctuation = arabic_punctuation + list(string.punctuation)
# Remove unwanted punctuation marks
cleaned_tokens = [token for token in tokens if token not in all_punctuation]

return cleaned_tokens

def get_ngrams(text, n):
cleaned_tokens = preprocess_text(text)
return list(ngrams(cleaned_tokens, n))

def plot_top_ngrams(data, n, top_n=10):
ngram_counter = Counter()

for text in data:
ngrams_list = get_ngrams(text, n)
ngram_counter.update(ngrams_list)

top_ngrams = ngram_counter.most_common(top_n)

ngram_labels, ngram_counts = zip(*top_ngrams)
print("ngram: ", ngram_labels)

# Convert ngram_counts tuple to list
ngram_counts = list(ngram_counts)

# Plotting the n-grams
plt.barh(range(len(ngram_labels)), ngram_counts)
plt.yticks(range(len(ngram_labels)), ngram_labels)
plt.xlabel('Frequency')
plt.ylabel('N-gram')
plt.title('Top N-grams')
plt.gca().invert_yaxis()
plt.tight_layout()
plt.show()

```

```

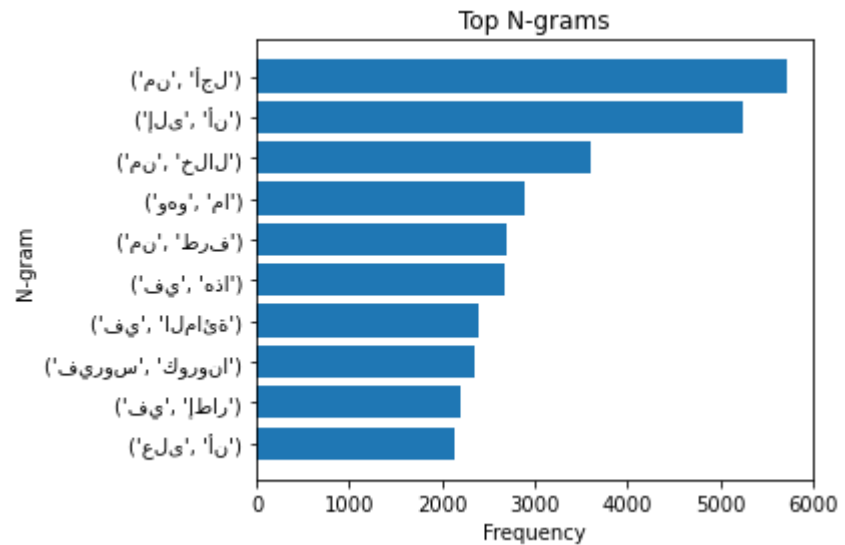
In [12]: # plot top frequent n-grams generally
n = 2
top_n = 10
print("-"*20, "top frequent n-grams generally for all classes", "-"*20)
print('*'*100)
plot_top_ngrams(combined_data['story'], n, top_n)

```

----- top frequent n-grams generally for all classes -----

\*\*\*\*\*

ngram: (('من', 'أجل'), ('إلى', 'أن'), ('من', 'خلال'), ('ما', 'وهو'), ('من', 'طرف'), ('في', 'هذا'), ('في', 'المائة'), ('في', 'كورونا'), ('في', 'فيروس'), ('على', 'أن'), ('ي', 'إطار'))



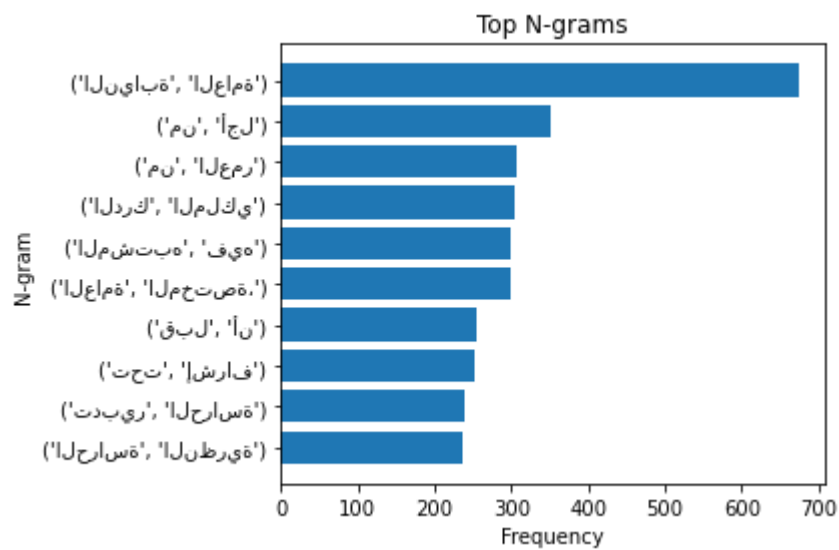
```
In [13]: # plot top frequent n-grams for each class separately
n = 2
top_n = 10
for topic in combined_data['topic'].unique():
    print("-"*30, f"topic: {topic}", "-"*30)
    print('*'*100)
    plot_top_ngrams(combined_data[combined_data['topic'] == topic]['story'], n, top_n)
    print('*'*100)
```

----- topic: art-et-culture -----

\*\*\*\*\*

ngram: (('من', 'خلال'), ('من', 'أجل'), ('إلى', 'أن'), ('في', 'هذا'), ('مجموعة', 'من'), ('العديد', 'من'), ('ما', 'وهو'), ('عدد', 'من'), ('في', 'هذه'), ('إطار', 'في'))



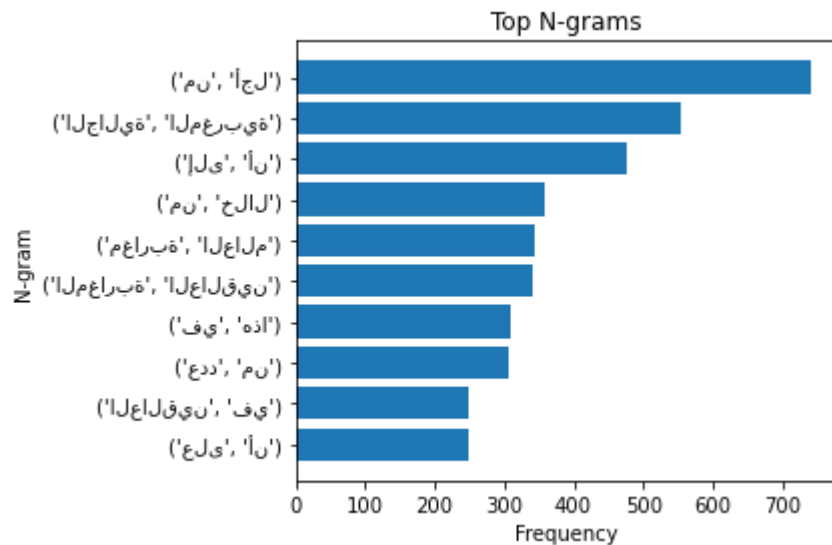


\*\*\*\*\*

----- topic: marocains-du-monde -----

\*\*\*\*\*

ngram: (( 'أجل' , 'من' ), ('الجالية' , 'المغربية' ), ('إلى' , 'أن' ), ('من' , 'خلال' ), ('مغاربة' , 'العالم' ), ('المغاربة' , 'العالمين' ), ('في' , 'هذا' ), ('م' , 'عدد' , 'على' , 'أن' ), ('العالمين' , 'في' ), ('ن' , 'ن' ))

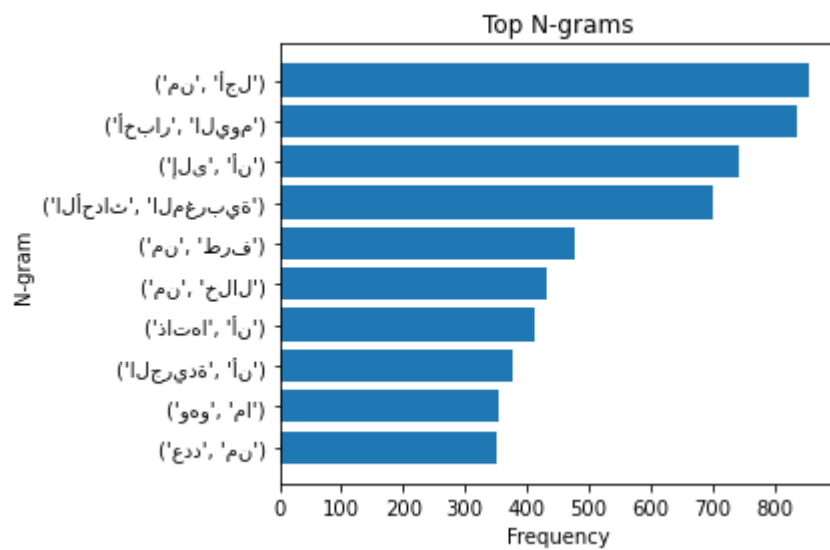


\*\*\*\*\*

----- topic: medias -----

\*\*\*\*\*

ngram: (( 'أجل' , 'من' ), ('اليوم' , 'أخبار' ), ('إلى' , 'أن' ), ('الأحداث' , 'المغربية' ), ('من' , 'طرف' ), ('من' , 'خلال' ), ('ذاتها' , 'أن' ), ('الجريدة' , 'أن' ), ('عدد' , 'من' ), ('وهو' , 'ما' ))

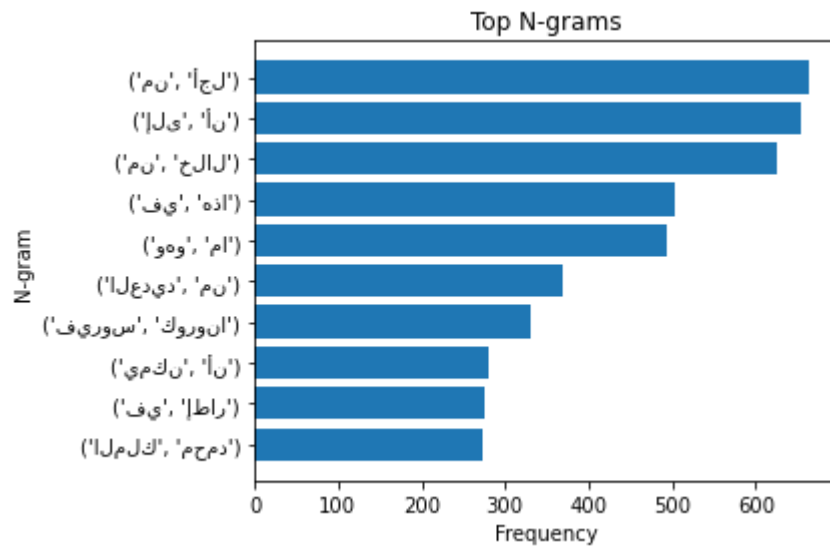


\*\*\*\*\*

----- topic: orbites -----

\*\*\*\*\*

ngram: (('أجل', 'من'), ('إلى', 'أن'), ('من', 'خلال'), ('في', 'هذا'), ('وهو', 'ما'), ('العديد', 'من'), ('يمكن', 'أن'), ('ف', 'ف'), ('ي', 'إطار'), ('الملك', 'محمد'))



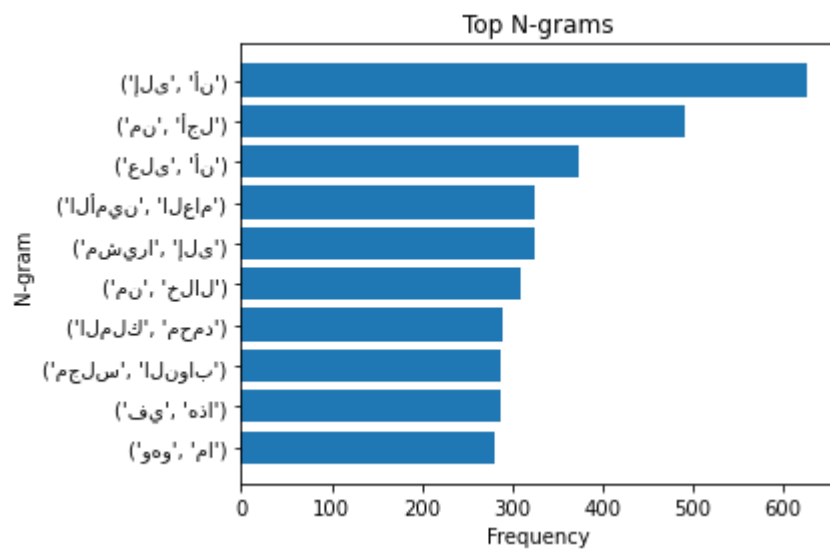
\*\*\*\*\*

----- topic: politique -----

\*\*\*\*\*

ngram: (('إلى', 'أن'), ('من', 'أجل'), ('على', 'أن'), ('الأمين', 'العام'), ('مشيرا', 'إلى'), ('من', 'خلال'), ('الملك', 'محمد'), ('مجلس', 'النواب'), ('وهو', 'ما'), ('في', 'هذا'))



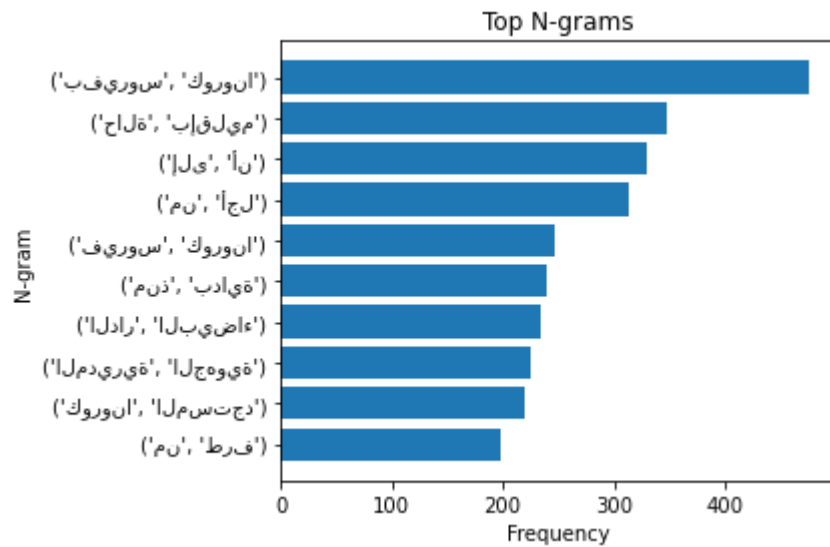


\*\*\*\*\*

----- topic: regions -----

\*\*\*\*\*

ngram: (( 'المديرية', ' '), ('الدار', ' البيضاء'), ('منذ', ' بداية'), ('فيروس', ' كورونا'), ('من', ' أجل'), ('إلى', ' أن'), ('حالة', ' بإقليم'), ('بفيروس', ' كورونا'), ('الجهوية'))



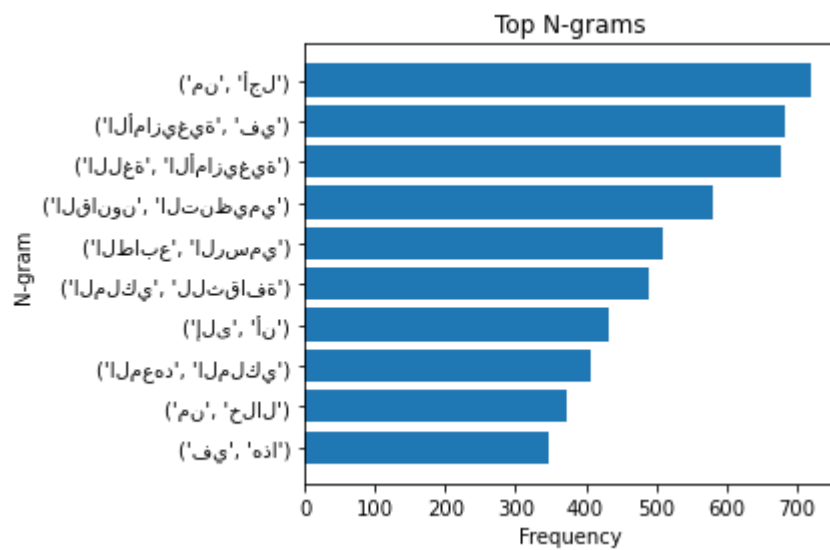
\*\*\*\*\*

----- topic: societe -----

\*\*\*\*\*

ngram: (( 'فيروس', ' كورونز'), ('في', ' تصريح'), ('وزارة', ' الصحة'), ('التربية', ' الوطنية'), ('في', ' المائة'), ('و', ' هو', ' ما'), ('من', ' أجل'), ('إلى', ' أن'), ('بفيروس', ' كورونا'), ('أ', ' على', ' مستوى'))





\*\*\*\*\*

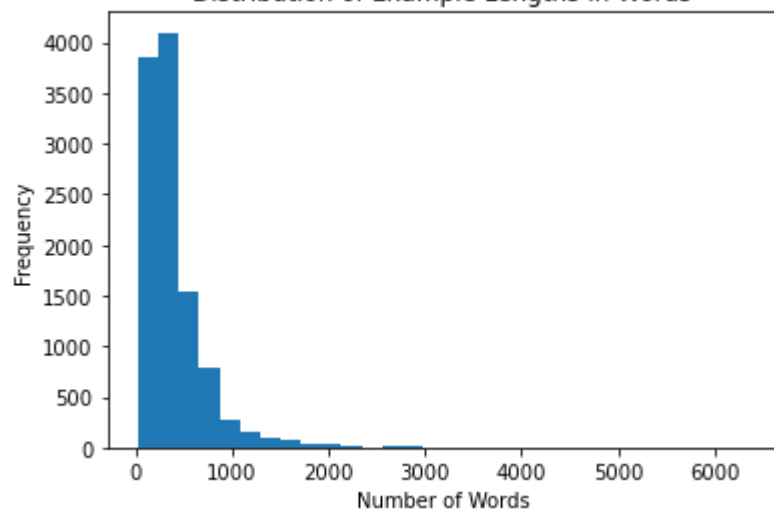
```
In [14]: def plot_length_distribution(data, measure):
    if measure == 'words':
        lengths = data.apply(lambda x: len(word_tokenize(x)))
        plt.xlabel('Number of Words')
        plt.title('Distribution of Example Lengths in Words')
    elif measure == 'letters':
        lengths = data.apply(len)
        plt.xlabel('Number of Letters')
        plt.title('Distribution of Example Lengths in Letters')
    else:
        raise ValueError("Invalid measure. Use 'words' or 'letters'.")

    plt.hist(lengths, bins=30)
    plt.ylabel('Frequency')
    plt.show()

    # Plot distribution of example lengths in words
    plot_length_distribution(combined_data['story'], measure='words')

    # Plot distribution of example lengths in Letters
    plot_length_distribution(combined_data['story'], measure='letters')
```

Distribution of Example Lengths in Words



Distribution of Example Lengths in Letters

