

UNIVERSITÉ MOHAMMED V

Faculté des Sciences de Rabat



Master Ingénierie de Données et Développement Logiciel

Natural Language Processing

Réalisé par : - Fadwa Saoiabi

- Fatima Zahra Atourabi

Encadré par : Mr Abdelhak Mahmoudi

Année universitaire : 2020-2021

I- Introduction :

Durant ce semestre, nous avons été initiés au Machine Learning et au NLP ainsi que leurs différentes techniques. Nous avons assimilé plusieurs algorithmes et méthodes utiles pour réaliser tout projet concernant le Natural Language Processing. Afin de matérialiser ces informations acquises et avoir une meilleure idée de ce qu'est le Machine Learning, notre professeur nous a donné libre recours de choisir un projet à réaliser concernant le NLP.

Pour notre part, notre choix s'est vite porté sur le **Sentiment Analysis**.

L'objectif de notre **projet** est de réaliser une analyse sentimentale des commentaires politiques présents dans notre jeu de données, du fameux site marocain Hespress. Dans un second temps, le but sera de parvenir à classifier les sentiments de ces commentaires selon qu'ils soient plutôt positifs, neutres ou négatifs à l'aide des différents modèles disponibles en Python.

II- Natural Language Processing :

Le traitement naturel du langage, ou Natural Language Processing (**NLP**) en anglais, est une technologie d'intelligence artificielle visant à permettre aux ordinateurs de comprendre le langage humain.

L'objectif de cette technologie est de permettre aux machines de lire, de déchiffrer, de comprendre et de donner sens au langage humain. D'importants progrès ont été effectués dans ce domaine au fil des dernières années, et le traitement naturel du langage est aujourd'hui exploité pour une large variété de cas d'usage...

- Quelles sont les différentes techniques de NLP ?

Les deux principales techniques utilisées pour le traitement naturel du langage sont **l'analyse syntaxique** et **l'analyse sémantique**.

L'analyse syntaxique consiste à identifier les règles grammaticales dans une phrase afin d'en déchiffrer le sens. L'analyse sémantique, quant à elle, consiste à **déchiffrer directement le sens d'un texte** en utilisant des algorithmes pour analyser les mots et la structure des phrases.

III - Sentiment Analysis :

L'analyse de sentiments (Sentiment Analysis ou Opinion Mining) est l'interprétation et la **classification des émotions** (positives, négatives et neutres) dans les données textuelles à l'aide des techniques d'analyse de texte. Cette analyse permet aux entreprises d'identifier l'opinion des clients à l'égard des produits, des marques ou des services dans les conversations et les commentaires en ligne.

L'analyse de sentiment se concentre sur **la polarité** (positive, négative, neutre), les sentiments et émotions (colère, joie, tristesse, etc.), et même sur les intentions (par exemple, intéressé contre non intéressé). Cela permet donc différents types d'utilisation de cette méthode d'analyse.



III- Implémentation et outils :

1- Langage de programmation :

Python est connu depuis longtemps comme un langage de programmation simple à maîtriser, du point de vue de la syntaxe. Il possède une communauté active et un vaste choix de bibliothèques et de ressources. On dispose donc d'une plate-forme de programmation capable d'être utilisée avec les technologies émergentes telles que l'apprentissage automatique et la Data Science.

2- Notebook :

La Data Science est itérative : il faut souvent tenter plusieurs approches et étudier les résultats avant de décider de la bonne façon de traiter un problème. C'est la raison pour laquelle les notebooks sont parfaitement adaptés à cette particularité. Un **notebook** est une interface web dans laquelle on peut taper du code Python, l'exécuter et voir directement les résultats, y compris une visualisation à l'aide de graphiques. Dans notre cas, nous avons choisi de travailler avec **Jupyter Notebook**.

3- Librairies utilisées :

Parmi les librairies Python que nous avons utilisées, on trouve :

- **NLTK** : C'est une librairie fondamentale pour la construction de programmes Python pour travailler avec des données de langage humain. Elle offre des interfaces faciles à utiliser sur des corpus ou ressources lexicales telles que WordNet, ainsi que des outils pour le traitement de texte, la classification, la tokenisation, le stemming, le balisage, l'analyse et le raisonnement sémantique.
- **SKLEARN** : C'est une librairie incontournable en Machine Learning et très bien documentée, destinée à l'**apprentissage automatique**. Elle comprend notamment des fonctions pour estimer des forêts aléatoires, des régressions logistiques, des algorithmes de classification ainsi que les machines à vecteurs de support.
- **Numpy** : C'est une librairie fondamentale pour effectuer des **calculs numériques** avec Python. Elle facilite grandement la **gestion des tableaux de données** et met à disposition également tout un **arsenal de fonctions pour effectuer des calculs mathématiques complexes** comme les fonctions trigonométriques ou encore les fonctions exponentielles et logarithmes.
- **Pandas** : C'est une **librairie très utilisée en data science** qui permet la manipulation et l'analyse des données. Elle propose en particulier des structures de données et des opérations de manipulation de tableaux numériques et de séries temporelles.



IV- Présentation du projet:

• Première analyse visuelle de notre Dataset :

Tout d'abord, voici un premier aperçu du jeu de données que nous avons à disposition. Comme nous pouvons l'apercevoir, cette data recensant au total 45857 de commentaires, contient deux colonnes intitulées "Comments" et "Score".

Comments	Score
فليجتمعوا لما فيه خير ليبيا وبنهوا هذا الصراع الذي عمر طويلا ودمر بلدا عزيزا علينا.	69
ايوى الرؤساء لن يحضروا يعني هناك ضغوط خارجية والجولة ستشغل بدونها كما يتمناه أعداء الليبيين الذين يشلون حركة التفاوض اما عن قريب أو عن بعيد	-18
الامبراطورية المغربية المعظمى لها تاريخ وشعب عريق واخواننا واشقاؤنا الليبيين يحبون ويتفوقون في المغرب .وعكس الجزائر أصبحت متورطة في استقبال القتلى وعائلته , والشعب الـ	29
الحل الوحيد هو تقسيم ليبيا الى الشرق والغرب هذا هو الحل .اما المفاوضات فهي مجرد مضجمة الوقت انا مع تقسيم ليبيا ان شاء الله ياذن الله	-87
تتمنى لاجوة الليبين ان يتوافقوا ويتفقا على ما فيه خير ليبيا والليبين . ليبيا المدنية الديمقراطية الموحدة دولة المؤسسات التي يسيرها ابتداء بعيدا عن الاجندات الخارجية التي تريد اـ	27
بعد الف جولة وجوله من المفاوضات ستبقى الحالة على ما عليها . ولكم في حماس وفتح لخير دليل على ذلك . في احد المحاولات العديدة لم اقتيادهم للتوقيع على اتفاق في رحاب	0
الكبريات اهون لديهم ان يتقاتل الليبيين الى الابد من ان ينسب للمغرب اي فضل في الاخاء بينهم	17
اللهم وفقهم واصلح ذات بينهم وهذا سينتفخ إن شاء الله إن إبتعد عنهم الشياطين خاصة الشياطين العرب الذين لا زالوا يظنون أن ليبيا التي يتحكم فيها القذافي لازالت قائمة .إنكم د	3
كان عليهم ان يجتمعوا هذه المرة بزعير	-1
Algeria is the tumor of the region	5
ليبيا للبين بالتوفيق انشاء الله	3
لقد كان من البديهي عدم حضور السادة عقيلة صالح وخالد المشري مع حفظ الاقالب و تلويح الاخير بالاستقالة لأن حضورهما يعني قبولهما بتصفيتها الجسدية من الجهات التي تا	2
صرحة الوضع معقد في ليبيا الشقيقة بحيث تتداخل فيه المصالح الأجنبية ودول الجوار ومع ذلك يمكن ان نتفاد بحيث اللقاء الاول ببرنامج كسب الثقة بين المتحاورين وجعلهم يتنا	0
تتمنى لاشقاؤنا الليبيين كل الخير و ان يجندوا حلا نهائيا لايتمهم في بلدهم الثاني المغرب وبمع السلام والازدهار بالتضامن والوحدة في المنطقة المغاربية والويل والمار لكل من له نية	2
المغرب لم يخذل أبدا إخوانه العرب!! صبر على تهوؤهم وطيشهم فكان حليما بقوة لفته في نفسه!!!! العرب لم يفهموا بعد أن الرزانة أم الحلول!!!!!! لهذا المغرب مستعد دائما لاستض	0
الوحي الصهيوني جند بعض الجيران للإيقاع على زعزعة الاستقرار في المنطقة لخدمت أجندتها التوسعية وقد التضح هذا من خلال بعض التصريحات التي أدلى بها بعض القادة	0
الاقوياء والرجال الصالح ،،،،سوف يتناقشون يوم الاثنين القادم تحت قيادة المانيا وبحضور المشاركين في مؤتمر برلين ..هذا هو المؤتمر الذي سوف يذكر التاريخ وضع المغرب المند	-3
تطلب من الله سبحانه وتعالى ان يبعد عن الأشقاء الليبيين بلا كبريات فرنسا. أمين يا رب العالمين	2
وزير الشغل في حكومتك ومن حزبك ووزير اخر من حكومتك وحزبك لم يسجلوا موقفهم ب cnss ولا شيء اتخذ ضدهم ومزال عندك الوجه تقول بحال هاد الكلام!!!!!!التيبي على من	167
الهجرة فالتلفازة حاجا اما الواقع حاجا!!!!!!اخرى اسي العماني بغينا التطبيق ماشي الشفوي بغينا حفا من التروة بغينا نتخلصوا مزيان ماشي تعطونا الشياطة ديال 2500 درهم اللي نتوه	140
حيثما نسمع كلمة إصلاح اعلم أنها بالنسبة لحقوق الشعب ومكاسبه ما هي إلا إفساد، من قبيل الزحف على حق النقاد وعلى القدرة الشرائية وعلى الوظيفة العمومية وعلى حقوق ا	58
خرجتو على الدرازي الى مغاوش افراو ماخلينوهوش التعلو الصنعة ما حدهم صغار يحكم انهم فاصرين حتى كبروا صنعة لا فريانة ابتلاوة بقرقوي الحشيش ابرد لهم لكاف اصبح	27
ما مصير مباراة المنتدبين القضائيين و التفتين المتخصصين فالكهرايا و الماء لقد سلمنا وطال انتظارا ؟؟؟؟؟	15

- **Importation des librairies nécessaires :**

Avant de commencer, nous devons impérativement importer les bibliothèques qui nous seront utiles, à savoir NLTK, Pandas, Numpy, SkLearn etc..

```
import numpy as np
import pandas as pd
import nltk
import string
import re
import seaborn as sns
import matplotlib.pyplot as plt

from nltk import word_tokenize
from nltk.corpus import stopwords
from nltk.tokenize import wordpunct_tokenize

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.pipeline import make_pipeline

from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
from sklearn.model_selection import train_test_split, GridSearchCV
```

Ensuite, nous avons importé notre dataset comme on le voit si bien sur la capture d'écran ci-dessous :

Entrée [2]: `data = pd.read_csv('comments.csv', sep=',', encoding='utf-8')`
`data.head()`

Out[2]:

	comment	score
0	... فليجتمعوا لما فيه خير ليبيا و ينهو هذا الصراع	69
1	...ايوى الرؤساء لن يحضروا يعنى هناك ضغوط خارجية و	-18
2	... لا امبراطورية المغربية العظمى لها تاريخ وشعب ع	29
3	...الحل الوحيد هو تقسيم ليبيا الى الشرق والغرب هذ	-87
4	...نتمنى لاخوة الليبيين ان يتوافقوا و يتفقوا على م	27

Entrée [3]: `data.shape`

Out[3]: (45857, 2)

- **Pre-Processing :**

La première étape à chaque fois que l'on fait du NLP est de faire le nettoyage de nos données, autrement dit faire le **cleaning et le traitement** de notre Dataset. On a commencé d'abord par enlever les stopwords, ponctuations, diacritics, les mentions etc.. afin de se débarrasser des éléments qui ne serviront pas à grand-chose lors de notre Sentiment Analysis.

Etape 2: Text pre-processing

```
#List of arabic and english punctuations:
punctuations = '""÷x-“”!|+|~{}',.°"/-][%^&*()_<>:'' + string.punctuation

# Arabic stop words with NLTK :
stop_words = stopwords.words()

arabic_diacritics = re.compile("""
    | # Shadda
    | # Fatha
    | # Tanwin Fath
    | # Damma
    | # Tanwin Damm
    | # Kasra
    | # Tanwin Kasr
    | # Sukun
    - # Tatwil/Kashida
""", re.VERBOSE)

def preprocess(text):

# Remove punctuations :

    translator = str.maketrans(' ', '', punctuations)
    text = text.translate(translator)

# Remove diacritics :
|
    text = re.sub(arabic_diacritics, '', text)

# Remove_mention :

    text = re.sub(r'@\S+', '', text)

# Allow_only_ar :

    text = re.sub(r'^[\u0600-\u06ff\u0750-\u077f\u0fb50-\ufbc1\u0bd3-\ufd3f\u0fd50-\ufd8f\u0fe70-\ufefc\uFDF0-\uFDFD]+',

# Remove Longation :

    text = re.sub("[\u0640\u0641\u0642\u0643\u0644\u0645\u0646\u0647\u0648\u0649\u0650\u0651\u0652\u0653\u0654\u0655\u0656\u0657\u0658\u0659\u0660\u0661\u0662\u0663\u0664\u0665\u0666\u0667\u0668\u0669\u0670\u0671\u0672\u0673\u0674\u0675\u0676\u0677\u0678\u0679\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb\u06gc\u06gd\u06ge\u06gf\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb\u06gc\u06gd\u06ge\u06gf\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb\u06gc\u06gd\u06ge\u06gf\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb\u06gc\u06gd\u06ge\u06gf\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb\u06gc\u06gd\u06ge\u06gf\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb\u06gc\u06gd\u06ge\u06gf\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb\u06gc\u06gd\u06ge\u06gf\u0680\u0681\u0682\u0683\u0684\u0685\u0686\u0687\u0688\u0689\u0690\u0691\u0692\u0693\u0694\u0695\u0696\u0697\u0698\u0699\u06a0\u06a1\u06a2\u06a3\u06a4\u06a5\u06a6\u06a7\u06a8\u06a9\u06aa\u06ab\u06ac\u06ad\u06ae\u06af\u06b0\u06b1\u06b2\u06b3\u06b4\u06b5\u06b6\u06b7\u06b8\u06b9\u06c0\u06c1\u06c2\u06c3\u06c4\u06c5\u06c6\u06c7\u06c8\u06c9\u06d0\u06d1\u06d2\u06d3\u06d4\u06d5\u06d6\u06d7\u06d8\u06d9\u06e0\u06e1\u06e2\u06e3\u06e4\u06e5\u06e6\u06e7\u06e8\u06e9\u06f0\u06f1\u06f2\u06f3\u06f4\u06f5\u06f6\u06f7\u06f8\u06f9\u06ga\u06gb
```

Ainsi, nos données seront représentées comme suit :

	comment	score
0	69	فليجتمعوا خير ليبيا و ينهوا الصراع عمر طويلا و ...
1	18-	ايوي الرءساء يحضروا يعني ضغوط خارجيه والجوله س... ا
2	29	لاامبراطوريه المغربيه العظمي تاريخ وشعب عريق ...
3	87-	الحل الوحيد تقسيم ليبيا الي الشرق والغرب الحل ...
4	27	نتمنى لاخوه الليبين ان يتوافقوا و يتفقوا على خ...

- **Drawing a WordCloud :**

Nous avons par la suite tracé un nuage de mots dans lequel on peut apercevoir **les mots les plus utilisés dans les différents commentaires des utilisateurs**. Plus la taille de la police d'un mot est grande, plus le nombre d'occurrences de celui-ci est important.

Etape 3: Drawing a WordCloud

```
import os
import codecs
from wordcloud import WordCloud
import arabic_reshaper
from bidi.algorithm import get_display

text_data = str(data['comment'])

# Make text readable for a non-Arabic Library like wordcloud
text = arabic_reshaper.reshape(text_data)
text = get_display(text)

# Generate a word cloud image
wordcloud = WordCloud(font_path='NotoNaskhArabic-Regular.ttf').generate(text)

# Export to an image
img = wordcloud.to_file("arabic_example.png")
plt.imshow(img, interpolation='bilinear')
plt.axis("off")

(-0.5, 399.5, 199.5, -0.5)
```



- **Sentiment Analysis with different techniques :**

Concernant cette étape, nous avons commencé par classer nos commentaires en se basant sur le score de ces derniers, ce qui nous a permis de déterminer si un commentaire est positif ou négatif.

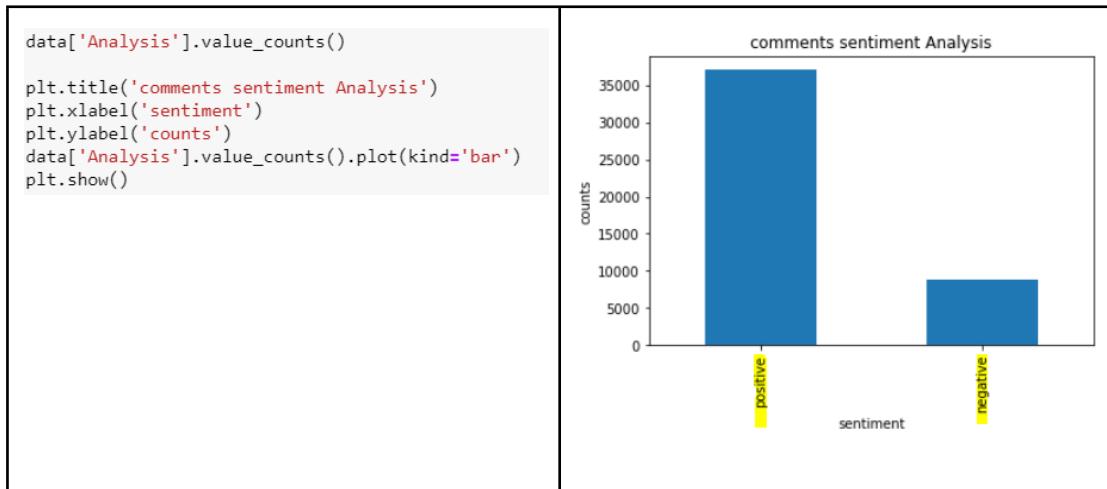
Sentiment Analysis with different techniques

```
#using score column
def getAnalysis(score):
    if score>=0:
        return 'positive'
    elif score<0:
        return 'negative'

data['Analysis'] = data['score'].apply(getAnalysis)
data
```

	comment	score	Analysis
0	... فليجتصعوا خبز ليبيا و ينهو الصراع عمر طويلا و	69	positive
1	...ايوي الرعاء يحضروا يحيي صغوط خارجيه والجوله س	-18	negative
2	... الامبراطوريه المغريه العظمي تاريخ وشعب عريق	29	positive
3	... الحل الوحيد تقسيم ليبيا الي الشرق والغرب الحل	-87	negative
4	...لثمني لآخوه الليبيين ان يتوافقوا و يتفقوا علي خ	27	positive
...
45852	...المغرب يخطو بالزيادة الرعايه سياسيا اجتماعيا ا	-4	negative
45853	...وزير الخارجيه الامريكي زار المغرب الاستقلال اق	3	positive
45854	...انا ارد علي سيد ابو وليد تغليفه خصوصا تبخيس مس	3	positive
45855	...ردا علي السيد فءاء تكلم علي الضباط السامون للقي	3	positive
45856	...ردا علي ابو صفاء علي علمي ان الرجل رغم انه دكت	1	positive

- Nous avons tracé par la suite un graphe qui montre clairement que notre Dataset contient + des avis positifs que négatifs.



• Multinomial Naïve Bayes :

L'algorithme **Multinomial Naive Bayes** est une méthode d'apprentissage probabiliste principalement utilisée dans le traitement du langage naturel (NLP).

Cet algo est basé sur le théorème de Bayes, il prédit la balise d'un texte tel qu'un e-mail ou un article de journal et calcule la probabilité de chaque étiquette pour un échantillon donné, puis donne en sortie l'étiquette avec la probabilité la plus élevée.

```
X = CountVectorizer(analyzer = preprocess, dtype="uint8").fit_transform(data["comment"]).toarray()
y = data["score"].apply(lambda x: 'positive' if x >= 0 else 'negative')

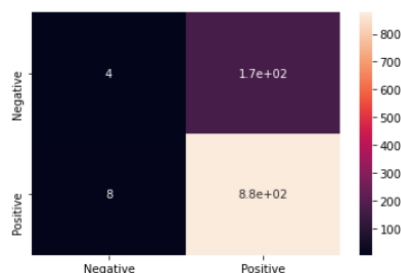
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.023)

NB_classifier = MultinomialNB()
NB_classifier.fit(X_train, y_train)

MultinomialNB()

y_predict_test = NB_classifier.predict(X_test)
cm = confusion_matrix(y_test, y_predict_test)
sns.heatmap(cm, annot=True, xticklabels=['Negative', 'Positive'], yticklabels=['Negative', 'Positive'])
print(classification_report(y_test, y_predict_test))
```

	precision	recall	f1-score	support
negative	0.33	0.02	0.04	170
positive	0.84	0.99	0.91	885
accuracy			0.84	1055
macro avg	0.59	0.51	0.48	1055
weighted avg	0.76	0.84	0.77	1055



- **Support Vector Machine :**

SVM est l'un des algorithmes d'apprentissage supervisé les plus populaires. Il est principalement utilisé pour les problèmes de classification ainsi que de régression. C'est même l'un des meilleurs modèles de classification.

Support Vector Machine (SVM)

```
#SVC :  
  
clf = svm.SVC()  
clf.fit(X_train, y_train)  
pred=clf.predict(X_test)  
print("SVC accuracy_score",accuracy_score(y_test, pred)*100)  
  
cm = confusion_matrix(y_test,pred)  
sns.heatmap(cm, annot=True, xticklabels = ['Negative', 'Positive'],yticklabels=['Negative','Positive'])  
  
print(classification_report(y_test,pred))
```

SVC accuracy_score 83.88625592417061

	precision	recall	f1-score	support
negative	0.00	0.00	0.00	170
positive	0.84	1.00	0.91	885
accuracy			0.84	1055
macro avg	0.42	0.50	0.46	1055
weighted avg	0.70	0.84	0.77	1055



- **Logistic Regression :**

La régression logistique est une méthode statistique pour effectuer des **classifications binaires**. Elle prend en entrée des **variables prédictives qualitatives et/ou ordinales** et mesure la probabilité de la valeur de sortie en utilisant la **fonction sigmoïd**.

```
from sklearn.linear_model import LogisticRegression  
logisticRegr = LogisticRegression()  
logisticRegr.fit(X_train, y_train)  
pred=logisticRegr.predict(X_test)  
print("LogisticRegression accuracy_score",accuracy_score(y_test, pred)*100)  
  
cm = confusion_matrix(y_test,pred)  
sns.heatmap(cm, annot=True, xticklabels = ['Negative', 'Positive'],yticklabels=['Negative','Positive'])  
  
print(classification_report(y_test,pred))
```

```

LogisticRegression accuracy_score 83.88625592417061
              precision    recall  f1-score   support

   negative      0.00      0.00      0.00      170
   positive      0.84      1.00      0.91      885

 accuracy
macro avg      0.42      0.50      0.46      1055
weighted avg      0.70      0.84      0.77      1055

```



V- Conclusion :

Le Machine Learning, nouvel outil pour utiliser de la donnée, n'est pas encore déployé à son plein potentiel. Néanmoins, l'avancée technologique de cette méthode permet de découvrir de nouveaux cas d'études et de nouvelles opportunités.

Une fois comprise par le plus grand nombre d'entreprises, **l'analyse de sentiments** permettra de mieux comprendre les clients ainsi que de donner de nouvelles perspectives aux équipes afin d'obtenir un travail plus productif et de meilleure qualité.

Aujourd'hui, l'image qu'a une entreprise est très importante. Les répercussions d'une mauvaise image peuvent se manifester très rapidement, notamment avec les réseaux sociaux. C'est la raison pour laquelle il est recommandé d'utiliser l'analyse de sentiments, parce que ça peut aider à surveiller et à contrôler l'image d'un client ou d'une entité.