



بسم الله الرحمن الرحيم

Sudan University of Science and Technology

Collage of Computer Science and Information Technology

Department of Information Technology and Software Engineering

Graduation Project



Build Sentiment Analysis Model for the Sudanese Dialect for the Internet Service (Sudani - Case Study)

Preparation:

- Athar Abdalgauom (information technology)
- Rabab Khalifa (Software Engineering)
- Ryan Jafar (Software Engineering)

Supervising:

A. Muhammad Al-Fateh Othman

October 2022

مقدمة :

سنوضح في هذا المستند المنهجية التي اتبعت حتى الان بغرض التوضيح و التقييم بعرض الخطوات كل على حدى (التعليقات قبل و بعد، صور للاكواد و المكتبات المستخدمة ، عرض و شرح و صور للنتائج) مع مقارنة المنهجية بالمقالات عن المجال .

1- جمع التعليقات :

تم جمع التعليقات يدويا من صفحة سوداني للاتصالات لم يتم مراعاة اعتبار لسلبيه او ايجابية او حيادية التعليق او علاقتها مع المنشور

2- تنظيف البيانات :

تم تنظيف التعليقات لازالة الايموجي و التشكيل و الاحرف الانجليزية و الارقام (normalization) باستخدام مكتبة (aranorm)

الكود:

```
from aranorm import normalize_arabic_text
from nltk.corpus import stopwords
from nltk import word_tokenize
import pandas as pd

def normlizeComment(file1,file2):
    with open(file1, 'r') as file1:
        with open(file2, 'w',newline='') as file2:
            for row in file1:

                file2.write(normalize_arabic_text(row))
                file2.write('\n')
```

التعليقات :

117. ♥سوداني جوة وجداني
118. والله ماف كلام سوداني تماااااام
119. عشان بعدين تزيديوا السعر الله يغطس حجركم
120. احلي شي في الصفحة دي انو الادمن برد عليك والله دا في حد زاتو تواضع منكم موفقين يارب
121. الشبكة تعبانة خالص مع انو شريحتي خط لآكن النت ضعف جداً
122. بجد خدماتكم كلهااااا تبارك الله بس رجعوا ليئا ف غيابك دي
123. فعلا كلنا سوداني/ربنا يحفظكم ويبارك فيكم جميعا
124. انا سواني واقتخر
125. النت غالي يااااا أهلنا
126. نعمل 500ميغا بس اشرب شاوي واجي القاهم 30ميغا الله يخليكُم والي الامام ف الغش
127. والله صراحة أفضل أنترنت في العالم من غير منازل؟
128. سودان كل الحوووووب ياخ
129. النت أصبحت ثقيل ومرات لا يوجد المكان الابيض شمال كردفان مع العلم كان سريعاً جداً
130. نحن نعاني معاناته ما بعدها من معاناته من شبكة سوداني ضيفه جدا لحدي عدم الوجود
131. 😊♥شبكة في الخرطوم بس
132. ♥♥♥♥ لا يوجد شبكة في بعض الولايات
133. حقيقة سوداني وشيكه عجيبه ماشاء الله مزيدا من التقدم والنجاح ودائما يارب تكونو في القمة وسودانيه واقتخررر
134. موفقين يا سوداني .سمعتكم اصبحت فوق كل الشركات
135. باقاتكم مربحة والله وموفقين دايما
136. خدماتكم سمحة بس وفر ليئا الشبكة
137. مستقبل واحد لشركة كبيرة موفقين
138. دائما في المقدمة، إلى الأمام خطوات نجاح يسبقها توفيق من الله
139. ♥مع سوداني أيا منا اجمل

التعليقات بعد المعالجة:

- سوداني جوه وجداني
والله ماف كلام سوداني تماااااااااااا
عشان بعدين تزيدوا السر الله يغطس حجركم
احلي شي في الصغحه دي انو الادمين برد عليك والله دا في حد زاتو تواضع منكم موفقين يارب
الشبكة تعبانة خالص مع انو شريحتي خط لائن التت ضعف جدا
بجد خدماتكم كلهاااا تبارك الله بس رجعو لنا ف غيايلك دي
فعلا كلنا سوداني ربنا يحفظكم ويبارك فيكم جميعا
انا سواني واقتخر
التت عالي ياااا اهلنا
نعمل ميغا بس اشرب شاوي واجي القاهم ميغا الله يخليكم والي الامام ف الغش
والله صراحه افضل انترنت في العالم من غير منازع
سودان كل الحوووووب باخ
التت اصبحت ثقيل ومرات لا يوجد المكان الابيض شمال كردفان مع العلم كان سريعاً جداً
نحن نعاني معاناته ما بعدها من معاناته من شبكه سوداني ضيفه جدا لحدي عدم اللاوجود
شبكه في الخرطوم بس
لا يوجد شبكه في بعض الولايات
حقيقه سوداني وشيكه عجيبيه ماشاء الله مزيدا من التقدم والنجاح وداوما يارب تكونو في القمة وسودانيه واقتخرررر
موفقين يا سوداني سمعتكم اصبحت فوق كل الشركات
باقاتكم مريحه والله وموفقين دايما
خدماتكم سمحها بس وفر لنا الشبكه
مستقبل واحد لشركه كبيره موفقين
داوما في المقدمه الي الامام خطوات نجاح يسبقها توفيق من الله
مع سوداني ايامنا اجمل
سوداني جوه وجداني

ثم تم استخدام مكتبة (nltk) لازالة كلمات التوقف :

```
def removeStopWords(file1,file2):

    arabicStopWords= stopwords.words("arabic")

    with open(file1, 'r') as file1:
        with open(file2, 'w',newline='') as file2:

            for row in file1:
                tokenizedRow = word_tokenize(row)
                commentWithNoStopWords= ' '.join([i for i in tokenizedRow if i not in arabicStopWords])
                file2.write(commentWithNoStopWords)
                file2.write('\n')
```

التعليقات بعد ازالة كلمات التوقف :

سوداني جوه وجداني
والله ماف كلام سوداني تماااااام
عشان بعدين تزيديوا السعر الله يغطس حجركم
احلي شي الصغحه دي انو الادمين برد والله دا حد زاتو تواضع منكم موفقين بارب
الشبكة تعبانته خالص انو شريحتي خط لائنك انت ضعفت جدا
يجد خدماتكم كلهاااا تبارك الله رجعو لنا غيابك دي
فعلا كلنا سوداني ربنا يحفظكم ويبارك فيكم جميعا
انا سواني واقتخر
النت غالي ياااا اهلنا
نعمل ميغا اشرب شاوي واجي القاهم ميغا الله يخليك والى الامام الغش
والله صراحه افضل انترنت العالم منازع
سودان الحووورروب باخ
النت اصبحت ثقيل ومرات لا يوجد المكان الابيض كرفاق العلم سريعاً جداً
نعاتي معاناته بعدها شبكه سوداني ضيفه جدا لحدي عدم اللاوجود
شبكه الخرطوم
يوجد شبكة الولايات
حقيقه سوداني وشبكته عجيببيه ماشاء الله مزيدا التقدم والنجاح وداوماً بارب تكونو القمة وسودانية وافتررر
موفقين سوداني سمعتكم أصبحت الشركات
باقاتكم مربحة والله وموفقين دايماً
خدماتكم ممجده وفر لنا الشبكة
مستقبل واحد لشركه كبيره موفقين
داوماً المقدمه الي الامام خطوات نجاح يسبقها توفيق الله
سودانى ايامنا اجمل

3- ازالة التكرار من الاحرف و الكلمات تم باستخدام مكتبة (re, itertools):

```
import stanza
nlp2 = stanza.Pipeline('ar')

def remove_duplicated_characters(text):
    result=[]
    text=''.join(i for i,_ in itertools.groupby(text))
    result.append(text)
    return ''.join(result)

def remove_duplicated_words(text):
    s1=[]
    text=' '.join(k for k,v in itertools.groupby(text.replace("</Sent>","").split()))
    s=re.sub(r'b(.+)(\s+\b)+',r'\1',text)
    return ''.join(s)

with open("commentsWithNoStopWords", 'r',encoding="utf8") as file1:
    with open("commentsWithNoDuplicated", 'w',newline='',encoding="utf8") as file2:
        for row in file1:
            phace1=remove_duplicated_characters(row)
            phace2=remove_duplicated_words(phace1)
            file2.write(phace2)
            file2.write('\n')
```

100. شكراً على خدمة العملاء
 101. سووووداني وافخرلو اكتب فيك طووول عمري يا سوداني م بوفيكي حقك ويعجز عنو لساني
 102. سوداني الاجمل
 103. شبكتكم رفت
 104. شكراً على خدمة العملاء
 105. 🍀💖 مدروشين بس باقاتكم م بنقدر نستغني عنها
 106. سرقة اموال وازمان
 107. 😊 باسعاركم الناشفه دي
 108. الكهرياء قاطعه والننت ضعيف
 109. الننت اخر ايام كان ما كويس نهائي
 110. خدماتكم سمحه بس وفر لينا الشبكة
 111. دائماً في المقدمة، إلى الامام خطوات نجاح يسبقها توفيق من الله
 112. مستقبل واعد لشركة كبيرة موفقين
 113. مع سوداني ايامنا اجمل
 114. مشاركة شنو انتو اسوأ شبكة في الإنترنت في العالم
 115. بالتوفيق والنجاح ان شاء الله
 116. "
 117. ❤️سوداني جوة وجداني
 118. والله مايف كلام سوداني تمام!!!!!!ام
 119. عشان بعيدين تزيدوا السعر الله يغطس حجركم
 120. احلى شي في الصفحة دي انو الادمن برد عليك والله دا في حد زاتو تواضع منكم موفقين يارب
 121. ...
 122. ...

التعليق بعد ازالة التكرار

107 شكرا علي خدمة العملاء
 108 سوداني وافخرلو اكتب فيك طول عمري سوداني بوفيكي حقك ويعجز عنو لساني
 109 سوداني الاجمل
 110 شبكتكم رفت
 111 شكرا علي خدمة العملاء
 112 مدروشين باقاتكم بنقدر نستغني عنها
 113 سرقة اموال وازمان
 114 باسعاركم الناشفه دي
 115 الكهرياء قاطعه والننت ضعيف
 116 الننت اخر ايام كويس نهائي
 117 خدماتكم سمحه وفر لينا الشبكة
 118 داءما المقدمة الي الامام خطوات نجاح يسبقها توفيق اله
 119 مستقبل واعد لشركة كبيره موفقين
 120 سوداني ايامنا اجمل
 121 مشاركته شنو انتو اسوا شبكة الانترنت العالم
 122 بالتوفيق والنجاح ان شاء اله

4-استخراج النصوص : تم استخدام ميزة استخراج النصوص في مكتبة (stanza) التي تعطي نتائج دقيقة باللغة الانجليزية

```
In [12]: import stanza
nlp2 = stanza.Pipeline('ar')
doc = nlp2('ولد')
```

```
print("-----")
for sentence in doc.sentences:
    print(sentence.dependencies)
```

```
2022-09-27 07:31:14 INFO: Loading: depparse
2022-09-27 07:31:17 INFO: Loading: ner
2022-09-27 07:31:43 INFO: Done loading processors!
```

```
-----
[({
  "id": 0,
  "text": "ROOT",
}, 'root', {
  "id": 1,
  "text": "ولد",
  "lemma": "وَلَدَ",
  "upos": "VERB",
  "xpos": "VP-P-3MS--",
  "feats": "Aspect=Perf | Gender=Masc | Number=Sing | Person=3 | Voice=Pass",
  "head": 0,
  "deprel": "root",
  "start_char": 0,
  "end_char": 3
})])
```

تم اخذ معلومات عن كل وحدة ونوعها اظهرت النتائج 297 كلمة من اصل حوالي 2500 كلمة مميزة كما اخطأت النتائج في اغلب انواع الوحدات .

الكود:

```
import stanza
nlp2 = stanza.Pipeline('ar')

with open("commentsWithNoStopWords", 'r',encoding="utf8") as file1:
    with open("commentTextExtraction", 'w',newline='',encoding="utf8") as file2:
        for row in file1:
            doc = nlp2(row)
            for sent in ([f'entity: {ent.text}\\ttype: {ent.type}'
                          for sent in doc.sentences for ent in sent.ents]):
                file2.write(sent)
                file2.write('\\n')
```

نتيجة استخراج النصوص من التعليقات :

1	entity: مائصال»type: PER
2	entity: دي»type: MISC
3	entity: دي»type: MISC
4	entity: شنو»type: PER
5	entity: دي»type: MISC
6	entity: ليك»type: LOC
7	entity: دي»type: MISC
8	entity: غيغا شويه»type: MISC
9	entity: والميقابايت»type: MISC
10	entity: دي»type: MISC
11	entity: دانا»type: MISC
12	entity: دي»type: MISC
13	entity: اناسوداني»type: PER
14	entity: هبي هبي»type: PER
15	entity: دي نت»type: PER

● نلاحظ النتائج الخاطئة بالنسبة للغة العربية

5- استخراج الجذور (lemmatization) تم باستخدام مكتبة (stanza) التي اعطت نتائج جيدة

```
import stanza
nlp2 = stanza.Pipeline('an')

with open("commentsWithNoStopWords", 'r', encoding="utf8") as file1:
    with open("commentsLemmatizaion3", 'w', newline='', encoding="utf8") as file2:
        for row in file1:
            doc = nlp2(row)
            for sent in doc.sentences:
                for word in sent.words:
                    file2.write(f'word: {word.text+" "}\tlemma: {word.lemma+" "}\tupos: {word.upos+" "}\txpos: {word.xpos+" "}')
                    file2.write('\n')
```

جذور الكلمات :

1	entity:	نَفْسِي\lemma:
2	entity:	طبيب\lemma:
3	entity:	قَلَّتِي\lemma:
4	entity:	رَمَز\lemma:
5	entity:	إِسْتَرَاكَ\lemma:
6	entity:	دي\lemma:
7	entity:	حَاوَل\lemma:
8	entity:	تَكْرَا\lemma:
9	entity:	مِيعَر\lemma:
10	entity:	تَم\lemma:
11	entity:	نَشَطَلت\lemma:
12	entity:	بِ\lemma:
13	entity:	عِشَان\lemma:
14	entity:	سُودَاتِي\lemma:
15	entity:	سِعر\lemma:
16	entity:	خِدْمَات\lemma:
17	entity:	طَيع\lemma:
18	entity:	تَبَاكَ\lemma:
19	entity:	أَنَا\lemma:
20	entity:	رَجَا\lemma:

التعويضات باستخدام الجذور عوضاً عن الكلمات الأصلية:

[illegible]

6- احدى خطوات ال (text extraction) هي استخراج ال (Upos , Xpos) تم ذلك باستخدام مكتبة (stanza)

```

7680 word: الاقوي — lemma: أَقْوَى — upos: ADJ — xpos: A-----MS2D
7681 word: سوداني — lemma: سُودَانِيّ — upos: ADJ — xpos: A-----MS1R
7682 word: المزاي — lemma: مَزَيَّة — upos: NOUN — xpos: N-----P2D
7683 word: الفريده — lemma: فَرِيد — upos: ADJ — xpos: A-----FS2D
7684 word: سوداني — lemma: سُودَانِيّ — upos: ADJ — xpos: A-----MS1R
7685 word: شيكتنا — lemma: شَيْكْتَنَّا — upos: X — xpos: U-----
7686 word: حاجه — lemma: حَاجَه — upos: NOUN — xpos: N-----S1R
7687 word: تخليك — lemma: تَخْلِي — upos: VERB — xpos: VIIA-2MS--
7688 word: تفخر — lemma: إِفْتَحَر — upos: VERB — xpos: VIIA-2MS--
7689 word: أن — lemma: أَنْ — upos: SCONJ — xpos: C-----
7690 word: ك — lemma: هُوَ — upos: PRON — xpos: SP---2MS4-
7691 word: عاءله — lemma: عاءله — upos: X — xpos: U-----
7692 word: سوداني — lemma: سُودَانِيّ — upos: ADJ — xpos: A-----MS1I
7693 word: بالتوفيق — lemma: تَوْفِيق — upos: NOUN — xpos: N-----S1D
7694 word: ليكم — lemma: لِيَكُم — upos: X — xpos: U-----

```

التي اعطت نتائج خاطئة للغة العربية

7- الترميز (tokenization) تم باستخدام مكتبة (nltk) استخراج الترميز الاحادي للتعليقات و الثنائي و الثلاثي كل في ملف منفصل.

الكود:

```

from nltk.tokenize import sent_tokenize, word_tokenize
with open("commentsWithNoStopWords", 'r', encoding="utf8") as file1:
    with open("commentsWithtokenization", 'w', newline='', encoding="utf8") as file2:
        for row in file1:
            file2.write(str(word_tokenize(row)))
            file2.write('\n')

```

الترميز الاحادي:

```

1  ['نفسى' و 'يوم' و 'تكتبوا' و 'السعر' و 'بدون' و 'مانسال']
2  ['طبيب' و 'تشرحو' و 'طريقه' و 'الاستراك' و 'الباقه' و 'دي']
3  ['قللتو' و 'الرسائل' و 'دي' و 'واديتونا' و 'ليها' و 'مقيقات' و 'يكون' و 'اجمل']
4  ['رمز' و 'الاستراك' و 'شغو']
5  []
6  ['الاستراك']
7  ['دي' و 'خط' و 'يتمشى' و 'المركز' و 'عشان' و 'تكتشط' و 'ليك']
8  ['حاولو' و 'اعملو' و 'كود' و 'للخدمات' و 'دي' و 'مش' و 'الزول' و 'مره' و 'يمشى' و 'يحول' و 'شريحته' و 'خدمه']
9  ['شكرا' و 'التوضيح' و 'مفيد' و 'اكرر' و 'الشكر' و 'سوداني' و 'الايداع' و 'والتميز']
10 ['السعر']
11 ['تم']
12 ['نشطت' و 'غيغا' و 'سويه' و 'ونزلت' و 'ملف' و 'الباقه' و 'خلصت' و 'عايزكم' و 'تفيدوني' و 'ليه' و 'خلصت' و 'الباقه' و 'وشكرا']
13 ['يتيالى' و 'نشطه' و 'قيفا' و 'نزلت' و 'ملف' و 'ديرها' و 'تنتهي' و 'كمان']
14 ['ساريا' و 'انت' و 'شكرا']
15 ['سوداني' و 'جميل']
16 ['سعر' و 'قيفا' و 'شريحه' و 'خط']
17 ['خدماتكم' و 'كلها' و 'ضابطه' و 'انا' و 'زعانته' و 'حاجه' و 'وحده' و 'شلتو' و 'منى' و 'شاره' و 'ابرز' و 'المعجبين']
18 ['طبعاً' و 'خجالتين' و 'تكتبوا' و 'السعر' و 'الزيادات' و 'الاخيره' و 'صح']

```


الترميز الثنائي :

1 ['نفسى يوم' , 'يوم تكتبو' , 'تكتبو السع' , 'السع بدون' , 'بدون مائسل']
2 ['طبيب تشرحو' , 'تشرحو طريقه' , 'طريقه الاشتراك' , 'الاشتراك الباقه' , 'الباقه دي']
3 ['قللتو الرسائل دي' , 'الرسائل دي' , 'وايديتونا' , 'وايديتونا ليها' , 'ليها ميقات' , 'ميقات يكون' , 'يكون اجمل']
4 ['رمر الاشتراك' , 'الاشتراك شنو']
5 []
6 []
7 ['دي خط' , 'خط بتمشي' , 'بتمشي المركز' , 'المركز عشان' , 'عشان تنكشط' , 'تنكشط ليك']
8 ['حاولو اعملو' , 'اعملو كود' , 'كود للخدمات' , 'للخدمات دي' , 'دي مش' , 'مش الزول' , 'الزول مره' , 'مره يمشي' , 'يمشي يحول' , 'يحول شريحته' , 'شريحته خدمه']
9 ['شكرا التوضيح' , 'التوضيح مفيد' , 'مفيد اكرر' , 'اكرر الشكر' , 'الشكر سوداني' , 'سوداني الابداع' , 'الابداع والتميز']
10 []
11 []
12 نشطت غيغا , 'غيغا شويه' , 'شويه ونزلت' , 'ونزلت ملف' , 'ملف والياقه' , 'والياقه خلصت' , 'خلصت عايزكم' , 'عايزكم تقيدونى' , 'تقيدونى ليه' , 'ليه خلصت' , 'خلصت الباقه' , 'و']
13 ['يتبقى نشطه' , 'نشطه قيفا' , 'قيفا نزلت' , 'نزلت ملف' , 'ملف ديرها' , 'ديرها تنتهي' , 'تنتهي كمان']
14 ['عشان انت' , 'انت ساريا']
15 ['سوداني جميل']
16 ['سعر قيفا' , 'قيفا شريحه' , 'شريحه خط']
17 ['خدماتكم كلها' , 'كلها ضابطه' , 'ضابطه انا' , 'انا زعلانه' , 'زعلانه حاجه' , 'حاجه وحده' , 'وحده شلتو' , 'شلتو منى' , 'منى شاره' , 'شاره ابرز' , 'ابرز المعجبين']
18 ['طبعما خجلانين' , 'خجلانين تكتبوا' , 'تكتبوا السع' , 'السع الزيادات' , 'الزيادات الاخيره' , 'الاخيره صح']
19 الشيكه تعبانه , 'تعبانه شديد' , 'شديد حاولوا' , 'حاولوا حلو' , 'حلو المشكله' , 'المشكله عشان' , 'عشان نشترك' , 'نشترك واشتراكنا' , 'واشتراكنا بروح' , 'بروح على' , 'على']
20 الفاضى , 'الفاضى ساكت' , 'ساکت تحياتي' , 'تحياتي ليكم' , 'ليكم منطقه' , 'منطقه ابو' , 'ابو فسيه' , 'فسيه حله' , 'حله الدرجمان' , 'الدرجمان ابو' , 'ابو طباعه' , 'طباعه الفجع' , 'الفجع الجماعه']
21 ['انا مشترك' , 'مشترك عشره' , 'عشره قيفا' , 'قيفا جاتني' , 'جاتني الفاتوره' , 'الفاتوره ورفعتها' , 'ورفعتها ومازلت' , 'ومازلت قيفا' , 'قيفا اول' , 'اول الشهر']

الترميز الثلاثي:

1 ['نفسى يوم تكتبو' , 'يوم تكتبو السع' , 'تكتبو السع بدون' , 'السع بدون مائسل']
2 ['طبيب تشرحو طريقه' , 'تشرحو طريقه الاشتراك' , 'طريقه الاشتراك الباقه' , 'الاشتراك الباقه دي']
3 ['قللتو الرسائل دي' , 'الرسائل دي' , 'وايديتونا' , 'وايديتونا ليها' , 'ليها ميقات' , 'ميقات يكون' , 'يكون اجمل']
4 ['رمر الاشتراك' , 'الاشتراك شنو']
5 []
6 []
7 ['دي خط بتمشي' , 'خط بتمشي المركز' , 'بتمشي المركز عشان' , 'عشان تنكشط' , 'تنكشط ليك']
8 ['حاولو اعملو كود' , 'اعملو كود للخدمات' , 'كود للخدمات دي' , 'دي مش الزول' , 'الزول مره' , 'مره يمشي يحول' , 'يمشي']
9 ['شكرا التوضيح مفيد' , 'التوضيح مفيد اكرر' , 'مفيد اكرر الشكر' , 'اكرر الشكر سوداني' , 'الشكر سوداني الابداع' , 'سوداني الابداع والتميز']
10 []
11 []
12 نشطت غيغا شويه , 'غيغا شويه ونزلت' , 'ونزلت ملف' , 'ملف والياقه' , 'والياقه خلصت' , 'خلصت عايزكم' , 'عايزكم تقيدونى' , 'تقيدونى ليه' , 'ليه خلصت' , 'خلصت الباقه' , 'و']
13 ['يتبقى نشطه قيفا' , 'نشطه قيفا نزلت' , 'نزلت ملف' , 'ملف ديرها' , 'ديرها تنتهي' , 'تنتهي كمان']
14 ['عشان انت ساريا']
15 []
16 ['سعر قيفا شريحه' , 'قيفا شريحه خط']
17 ['كلها ضابطه' , 'كلها ضابطه انا' , 'انا زعلانه حاجه' , 'حاجه وحده' , 'وحده شلتو' , 'شلتو منى' , 'منى شاره' , 'شاره ابرز' , 'ابرز المعجبين']
18 ['طبعما خجلانين تكتبوا' , 'خجلانين تكتبوا السع' , 'تكتبوا السع الزيادات' , 'الزيادات الاخيره' , 'الاخيره صح']
19 تعبانه شديد , 'تعبانه شديد حاولوا' , 'حاولوا حلو' , 'حلو المشكله' , 'المشكله عشان نشترك' , 'عشان نشترك واشتراكنا' , 'واشتراكنا بروح' , 'بروح على' , 'على الفاضى' , 'الفاضى ساكت' , 'ساکت تحياتي' , 'تحياتي ليكم' , 'ليكم منطقه' , 'منطقه ابو' , 'ابو فسيه' , 'فسيه حله' , 'حله الدرجمان' , 'الدرجمان ابو' , 'ابو طباعه' , 'طباعه الفجع' , 'الفجع الجماعه']
20 مشترك عشره , 'مشترك عشره قيفا' , 'عشره قيفا جاتني' , 'جاتني الفاتوره' , 'الفاتوره ورفعتها' , 'ورفعتها ومازلت' , 'ومازلت قيفا' , 'قيفا اول' , 'اول الشهر']

8- الاشتقاق (stemming) تم باستخدام مكتبة (nltk) التي تعطي نتائج صحيحة للغة الانجليزية و لكن النتائج خاطئة للغة العربية حيث انها ترجع نفس الكلمة

```

from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize

ps = PorterStemmer()

with open("commentsWithNoStopWords", 'r',encoding="utf8") as file1:
    with open("commentsWithstemming", 'w',newline='',encoding="utf8") as file2:
        for row in file1:
            tokens = word_tokenize(row)

            stemmed = []

            for token in tokens:
                stemmed_word = ps.stem(token)
                file2.write(f'{token} : {stemmed_word}'+ "\n")
                file2.write('\n')

```

النتيجة:

1	تفسي : تفسي
2	
3	يوم : يوم
4	
5	تكتبو : تكتبو
6	
7	السعر : السعر
8	
9	يدون : يدون
10	
11	ماتسال : ماتسال
12	
13	طوب : طوب
14	
15	تشرحو : تشرحو
16	
17	طريقه : طريقه
18	
19	الاشتراك : الاشتراك
20	
21	تكتبو : تكتبو

9-تم استخراج ال(TF-IDF) باستخدام مكتبة (sklearn)

```

import pandas as pd
with open("commentsWithNoStopWords", 'r',encoding="utf8") as file1:

    vectorizer = TfidfVectorizer()
    vectors = vectorizer.fit_transform(file1)
    feature_names = vectorizer.get_feature_names_out()
    dense = vectors.todense()
    denselist = dense.tolist()

    with open("TFIDFDENSE", 'w',encoding="utf8") as file2:
        for i in range(len(denselist)):

            file2.write(str(i)+" - "+ str(denselist[0][i])+ "\t"+ str(denselist[1][i]))
            file2.write("\n")
    file2.close()

df = pd.DataFrame(denselist, columns=feature_names)

```

النتيجة:

```

180 179 - 0.0 → 0.0
181 180 - 0.0 → 0.0
182 181 - 0.0 → 0.0
183 182 - 0.0 → 0.0
184 183 - 0.0 → 0.0
185 184 - 0.0 → 0.0
186 185 - 0.0 → 0.0
187 186 - 0.0 → 0.0
188 187 - 0.0 → 0.0
189 188 - 0.0 → 0.0
190 189 - 0.0 → 0.3877088831278344
191 190 - 0.0 → 0.0
192 191 - 0.0 → 0.0
193 192 - 0.0 → 0.0
194 193 - 0.0 → 0.0
195 194 - 0.0 → 0.0
196 195 - 0.0 → 0.0
197 196 - 0.0 → 0.0
198 197 - 0.0 → 0.0
199 198 - 0.0 → 0.0
200 199 - 0.0 → 0.0

```

كما تم استخدام (tfidfvectorizer) لاستخراج اسماء المفردات المميزة حيث بلغت 2566 مفردة

2540	يقال
2541	يقدركم
2542	يقفل
2543	يقول
2544	يقولوا
2545	يقوم
2546	يكرمكم
2547	يكفي
2548	يكمل
2549	يكون
2550	يلا
2551	يمشي
2552	ينتقم
2553	ينتهي
2554	ينجحنا
2555	ينزل
2556	ينزلوا
2557	ينغمكم
2558	يهديكم
2559	يو
2560	يوجد
2561	يوديها!!!!
2562	يوفقكم
2563	يوم
2564	يومي
2565	يومين
2566	يوميين

10- تم استخراج (BOW) باستخدام مكتبة (sklearn)

الكود:

```
import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer

wordset = [

]

with open("commentsWithNoStopWords", 'r', encoding="utf8") as file1:
    for row in file1:
        wordset.append(row)

vectorizer = CountVectorizer()
X = vectorizer.fit_transform(wordset)

df_bow_sklearn = pd.DataFrame(X.toarray(), columns=vectorizer.get_feature_names_out())
df_bow_sklearn.head()

df_bow_sklearn.to_csv('BOW3.csv', encoding='utf-8')

file2.close()
print(X.shape)
```

[illegible]

كما تم استخراج الكلمات المفردة المميزة عن طريق (countVectoriser)

```
import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer

wordset = [

]

with open("commentsWithNoStopWords", 'r', encoding="utf8") as file1:
    for row in file1:
        wordset.append(row)

vectorizer = CountVectorizer()
X = vectorizer.fit_transform(wordset)
features=vectorizer.get_feature_names_out()
with open("featureWithBOW.csv", 'w', encoding="utf8") as file2:
    for i in features:
        file2.write(str(i+" \n"))

file2.close()
print(X.shape)
```

2548	يَكمُلْ
2549	يَكونْ
2550	يَلا
2551	يَمتَني
2552	يَنقُصْ
2553	يَنقُصِ
2554	يَنجَحنا
2555	يَنزِلْ
2556	يَنزِلُوْ
2557	يَنفَعكمْ
2558	يَهْدِيكمْ
2559	يُوْ
2560	يُوجدْ
2561	يُودوهاااااا
2562	يُوفِقكمْ
2563	يُومْ
2564	يُومي
2565	يُومينْ
2566	يُوميينْ
---	---