

Post 1

الموديل في الآخر يحتاج مني اني ادخله أرقام ، الأرقام ديه هو من خلالها يقدر يشوف ال pattern من خلالها بطريقة ما ، فلما نيجي ننقل لل NLP هلاقى الأرقام ديه او ال Features مبتقش متوفرة علطول عشان اديها لل model فهو يقدر يربط الدنيا ببعض ويفهم ، بقيت أنا محتاج اني دلوقت اعمل ال Features وديه ومش لاي ارقام لا لارقام هو فعلا يقدر يلاقى خلالها pattern فيفهم انا محتاج ققولة ايه من خلال الأرقام ديه فبقى عندي مرحلتين مهمين بالنسبة ليا قبل ما ال model بتاعى يشتغل ، المرحلة الأولى هي Preprocessing لل Text بتاعى ، سواء اني اعمله cleaning او اشيل ال Stop words او اعمل Lemmmization وغيره من الخطوات الى بظبط بيها ال Text بتاعى في Form معينة تتناسب مع ال Application بتاعى لان بعض الخطوات ممكن تكون متضمنة في Application ما واعر لا ، وكمان لأن بعض الخطوات ديه او معظمها بيكون ليه بعض الإيجابيات وبعض السلبيات ، بعد كده بتيجي المرحلة الى فيها تحول ال Text ده بقا لأرقام وهي ال Preparation ال يقدر من خلالها ادى ال Data بتاعى الى كانت مجرد Text من لغة ما لأرقام الموديل يقدر يشتغل معاها ويطلع نتائج ، المرحلة الأولى ديه فيها بعض الخطوات المهمة جدا والى من خلالها أصلا هقدر اني اتعامل مع معظم الخطوات الى بعدها وهي ال Tokenization وده نتكلم عنه بإذن الله اليوست الجاي

Post 2

ال Tokenization هو اني اعمل Segmentation للكلام الى عندي يعنى لو عندي paragraph احولى لمجموعة من الجمل ، والجمل ديه احولها للكلمات ، في الآخر انا هحتاج اني احول ل Basic level الى بيكون الجملة وهو الكلمة ، بعد كده بيتدى أشوف الكلمات ال Unique في الداتا الى عندي واستخدمها في اني احولها لأرقام بناء على بعض الطرق هنجلها بعددين والارقام ديه الى هتاخذها الكلمات هبتدى اعوض بيها مكان الكلمة في كل جملة وبرضه ده هنعرفه قدام بإذن الله ، لكن نرجع تاني لعملية ال Tokenization لانها بتختلف باختلاف اللغة وال Application الى بتعامل معاها ، مثلا هل كلمة زى dont هعتبر كلمة واحد ولا هحتاج اني احولها ل do and not طيب كلمة زى Ice cream هل هعتبرها ك Token واحد ولا اثنين ؟ طيب هل الكلمة ديه ليها Predefined rule في اللغة او معنى معين زى مثلا dont ديه بتعبر عندي عن حالة نفى ، كمان الكلمات نفسها و أشكالها المختلفة زى كمة end and ending هل كل كلمة من دول هعتبرها كلمة لوحدها ولا محتاج ارجع الكلمة لأصلها زى ما هنشوف في ال Steaming and lemmization ولكن نلاقى مثلا كلمة زى run and running لو هرجعها ل run ك أصل الكلمة محتاج اشيل ning بينما في ending محتاج اشيل ال ing فقط ولو كلمة زى bus هل ال s في الآخر هي suffix زى ما في كلمة words كل ديه cases وغيرها محتاج اعرفها وانا شغال وبتعامل معاها ازاي لانها هتتسبب في تكبير الكلمات بتاعى او ما يسمى Bag-of-words ، في الآخر ال Tokenization هو تقسيم الجملة او الكلام لمجموعة من الكلمات ولكن محتاج اعرف التقسيم ده هيكوّن بناء على ايه ، وهنا ال Libraires المختلفة منها الى بيستخدم rules من اللغة ومنها الى شغال عن طريق انه يعمل split عن طريق ال space وممكن انتا تعمل ال Tokenization الخاص ببيك بناء على Regular expression معين او rules معينة بإستخدام Regex or re libraires .

```
In [2]: sentence = """Thomas Jeferson began building Monticello at the age at 26."""
        sentence.split()

Out[2]: ['Thomas',
         'Jeferson',
         'began',
         'building',
         'Monticello',
         'at',
         'the',
         'age',
         'at',
         '26.']
```

Post 3

من اول stage ال nlp هلاقى فيه عندي challenges كثير لأن كل حاجه ممكن يكون ليها منفعة ما ممكن تضرنى بشكل اخر ، وهنا بحتاج اشوف كل خطوة وارتباطها بالابليكشن الى انا شغال فيه وازاي هتأثر فيها بشكل او باخر فمحتاج اشوف مثلا الفرق بين اني اعمل Steaming and lemmitazation واخد في اعتباري اني فيه بعض الكلمات بسببهم ممكن ميكونش ليها اي

معنى واصبحت اصلا ملهاش علاقة باللغة وبذات في حالة ال steaming لانه بيثيل ال suffix and prefix فكلمة مثلا زى ending هيشيل منها ing طب بالنسبة لكلمة زى sign لو شيلنا ال ing بقت مجرد حرف s وكلمة زى bus and words لما اشيل ال s اصبحت ليها معنى مع words لكن بالنسبة ل bus بقت bu وكلمة زى believe لما اعمل steaming هتكون believ برضه ملهاش اى معنى فى اللغة ، بينما ال Lemmization بيكون ليها علاقة بال Grammers بتاع اللغة وبياخذ في اعتبارة الأشكال المختلفة للكلمة بناء على السياق الى جت فيه ، وبعض المكتبات بتحدد فيها الأولوية لل verb or noun لما يكون عنده اكثر من اختيار للكلمة ، والمهم فى الخطوتين دول انهم ببساعدوا ال model بعد كده انه ميحصلش فيه over fitting لانه بدل ما كان بيدى attention لكل أشكال الكلمة لا دلوقت بقى عنده شكل او اثنين منها فى ال vocab الى انتا بتبنيها وده بيقلل ال over fitting لانه لو الكلمة كانت جت فى شكل معين مرة واحدة فى كل ال data هيديها اهمية لانك هتعتبرها من ضمن ال vocab فمش هيقدر يعمل generalization وبجانب ده كمان انتا بتقلل ال Vocab بتاعتك الى بدوره بيقلل معاك ال computation الى بتحتاجها ، طب ايه هي ال Vocab ديه وازاي ببنيها هنشوفها فى البوست الجاي بإذن الله

Post 4

بعد ما عملنا Tokenization لل Text وحددنا لو هنعمل Lemmtization or stemming بنبيتي ننقل من مرحلة ال Pre-processing text لمرحلة انى اعمل Pre-paration لل text وده هو انى احول الكلام ده لارقام الموديل يقدر يفهمها ، واحدة من الطرق ديه هي ال One-hot encoder وهو انى مجرد ما خلاص حددت ال Vocabulary (nuique words) الى هستخدمها هبتدى بقا انى اعمل matrix n*v ، لكل document عدد ال Token الى فيه فى عدد ال Vocabulary يعنى لو كان document من الداتا مثلا 100 كلمة وعدد الكلمات الى فى الداتا كلها 20000 فهنا محتاج matrix 100*20000 فقط ل document واحد، والحقيقة انى ده رغم انه مكلف جدا جدا وفى الحياه العملية مينفعش يشتغل لان اللغة نفسها natural ولو اعتبرنا كل كلمات اللغة وال symbols وغيره من الحاجات المتضمنه من اللغة هيكون مستحيل اعمل process وهتكون عملية مكلفة جدا وفى الآخر بيكون عندى sparse matrix كلها اصفار وغير ، الا ان ال one-hot من خلاله يقدر احافظ على ال context and grammar that words comes in فالموديل هيقدر ياخذ بالة من ال sequence بتاع الكلمات لكنه طبعا لو قلنا عندنا 3000 كتاب وكل كتاب فية حوالى 3500 جملة وكل جملة 15 كلمة ف احنا محتاجين شوف محتاجين بقا مساحة قد ايه غير انى ده كله هيكون مضروب فى عدد ال Vocabulary بتاعتك ، بعد كده ابتدينا نتجاهل ال order عشان ننقل لفكرة ال Bag of words وده فى البوست الجاي.

```

=====
Thomas Jefferson began building Monticello at the age at 26.
=====
it[4]: array([[0, 0, 0, 1, 0, 0, 0, 0, 0],
              [0, 1, 0, 0, 0, 0, 0, 0, 0],
              [0, 0, 0, 0, 0, 0, 1, 0, 0],
              [0, 0, 0, 0, 0, 0, 0, 1, 0],
              [0, 0, 1, 0, 0, 0, 0, 0, 0],
              [0, 0, 0, 0, 0, 1, 0, 0, 0],
              [0, 0, 0, 0, 0, 0, 0, 0, 1],
              [0, 0, 0, 0, 1, 0, 0, 0, 0],
              [0, 0, 0, 0, 0, 1, 0, 0, 0],
              [0, 0, 0, 0, 0, 1, 0, 0, 0],
              [1, 0, 0, 0, 0, 0, 0, 0, 0]])

```

Post 5

ال Bow of words هو فكرة بتعتمد على ال Vocabulary الى فى الداتا ال unique words وليه طرق مختلفة من binary to counts to frequency and lastly the TF-idf ، بعد فكرة ال one-hot الى كانت بتحافظ على ال sequence بتاعى لكنها كانت استحالة نقدر نتعامل بيها لانها مكلفة جدا جدا ومحتاجه resource كبيرة لداتا ممكن تكون قليلة جدا بالنسبة للداتا الحالية الى بنتعامل معاها ، فقلنا لو تجاهلنا ال order بتاع الكلمات وروحنا عملنا sum لل columns هنلاقى بقا عندنا one vector بعدد ال vocabulary عبارة عن 0 و 1 لو كانت الكلمة ديه جت هنلاقى نتيجة ال column ده 1 لو مكنتش جت فى ال document الى انا فيه هنلاقى 0 ولكن هنا مبقاش فيه order وممكن علاجة حاجة زى كده هي تقسيم ال

document نفسه ل sentence وال ml فى التعلم هيقدر برضه يحصل على معلومات لان الجملة فى الاخر عبارة عن كام كلمة لو حتى ال order بتاعها اتلغبط احنا نفسنا بنقدر نفهم بنتكلم عن ايه ونرتبها كمان بظبط زى ما بتيجى تتعلم لغة جديدة ، بعد كده بدل ال binary representation ده بقينا نعبر بال counts عشان يكون فيه different weights والموديل يقدر ياخد باله من الكلمات الى عدد تكرارها اكثر لكنه بقا فيه مشكلة وهو انى الكلمات الى استخدامها شائع جدا فى document هتاخذ weights اكبر ، ومشكلة اخر فى ال document الكبيرة هنلاقى بعض الكلمات خدعت counts كبير جدا والبعض الاخر واخذ counts قليل جدا وبقا فيه ranges مختلفة فرجعنا لعملية ال normalization وهو انى اعمل normalize للكلمات بالنسبة لل Vocabulry فبقيت فى range من 0 ل 1 ولكن ما زال ده كل مجرد counts ومش بيعبر عن لا معانى الكلمات او حتى اهميتها وبالنسبة لل rare words الكلمات الاقل شيوعا الى بتكون بالنسبة document ما هى الاثاثة هنلاقى ال weights بتاعتها لا تذكر وبالتالى ال model مش هيقدر ياخد باله منها وده كان ال frequency بتاعت الكلمات، ولكن بعد كده ابنتت ظهور ال TF-idf وده الى هنتكلم عنه المرة الجاية بإذن الله .

Out[10]:

	Thomas	Jeferson	began	building	Monticello	at	the	age	26.	Construction	...	South	Pavilion	In	1770.	Turning	a	neoclassical	masterpiece
sent_1	1	1	1	1	1	1	1	1	1	1	0	...	0	0	0	0	0	0	0
sent_2	0	0	0	0	0	0	0	0	0	0	1	...	0	0	0	0	0	0	0
sent_3	0	0	0	0	0	0	0	1	0	0	0	...	1	1	1	1	0	0	0
sent_4	0	0	0	0	1	0	0	0	0	0	0	...	0	0	0	0	1	1	1

4 rows x 31 columns

Post 6

ال Tokenization بتاعى دوما بيبقا محتاج improvement وبيختلف باختلاف ال application الى انا شغال فيه ، ساعات بحتاج اعتبر ال punctuations ديه tokens وفى اوقات تانى لا فمثلا لو بنتكلم فى tweets فال symbols معظمها بيكون مهم بالنسبة ليا لانها بتؤدى لمعنى ما زى ال emotions وفى نفس الوقت ممكن ال application نفسه يكون مبنى على tweets لكن مش مهم فيه ال emotions مثلا انى اعتبرها tokens لو انا مثلا شغال Named Entity Recognition ، والتعامل مع ال Tokenization بيختلف بإستخدام ال library الى بتتعامل معاها وال Tokenizer نفسه مثلا NLTK tree bank tokenizer بيشغل بناء على Grammar rules زى spacy لكن spacy بترجع spacy token ومحتاج اخذ بالى من حاجه زى كده واحولها ل strings عشان لما اجى اتعامل معاها زى ال sparse matrix مع tf-idf بحتاج احولة ل array لما اجى اديه للموديل ، وفيه casual tokenizer وده فى حالة انك بتتعامل مع Tweets فبياخد باله من معظم ال symbols ، وفى نفس الوقت انتا ممكن بإستخدام ال Regular expression تعمل implementation لل Tokenizer بناء على regs expression معين وطبعا كل حاجه من ديه بتختلف فى دقتها وسرعتها وده مثلا فرق بين nltk and spacy ل 10 تويئات، كمان ساعات بحتاج يكون ال token نفسة يكون عبارة عن كلمتين او ثلاثه لانه بيمثل على بعض expression او intuition ما زى ice cream ولكن ده تاثيره بيكون سىء لو اعتبرت كل ال 2-grams لان معظم الكلمات مش بتيجى فى sequence مرتبط ببعضه زى ال ice cream وكل ما ال n-grams بتاعتك بتكبر بيكون ال statics بتاعت حدوث ال sequence ده من الكلمات أقل ، وكمان ده بيساعدنا اننا نحصل على جزء حتى لو بسيط من ال sequence ال order الى تجاهلنا من بعد ال one-hot ، المشكلة هنا زى ما قلنا ال statics بتاعت حدوث ال n-grams ده مع بعض بتاثر على الموديل بتاعى وبيحصل overfitting لان الموديل هيبندى ياخذ باله من sequence هو اصلا بيحدث نادرا فلما يجى فى ال testing هنلاقى نتيجة سيئة ،

```
The number of Tokens are: 188
#####
The number of Vocabulary are: 133
vocabulary
0 12398
1 59
2 اتفا فيه
3 اجنه
4 استاذ
5 اصا به
6 اصل
7 اكتشاف
8 الاتصالات
9 الاطفال
```

```

The number of Tokens are: 192
#####
The number of Vocabulary are: 134
vocabulary
0
1      12398
2      59
3      اتفقيه
4      احنه
5      استاذ
6      اما به
7      اصل
8      اكتشاف
9      الاتصالات

```

Post 7

بعد موضوع ال n-grams هنلاقي حصل زي conflict بين ال n-grams عموما وكل ما تكبر في ال n ديه بيحصل اني نادر ما ال sequence ده فعلا بيجي مع بعض ، وفي نفس الوقت لو شوفنا ال one-gram هنلاقي اني فيه كلمات rare يعني ذكرها قليل جدا لكنها مهمه جدا جدا بالنسبة لل document بتاعي ، بينما ال rare sequence في ال n-grams هو كده كده مش مهم بالنسبة ليا ، وال n-grams ده كل ما كان قريب من عدد ال docs او اكثر من ال documents of dataset بيحصل overfitting وحل حاجه زي كده بيكون عن طريق اني اظبط مثلا threshold معين لل grams الى هعتبرها زي بظبط موضوع ال stopwords بانى لو كلمة تعدت نسبة معينة هيتدى اتاجهله وده بيكون بناء على نسبتها بالنسبة لل document الى هي فيه مع ال documents الثانية وهو ده الى بيعمله ال Tf-idf انه بيحافظ على ال rare words وبيديها weights كبير بينما بالنسبة لل stopwords بنلاقيها واخده weights قليلة جدا ، ورغم اني عدد ال stopwords مش بيكون مهم اشيله او اسويه لانه بالنسبة لل vocabulary size هنلاقي الفرق بسيط وكمان ساعات بيكون ليه اهمية بذات في كلمات زي as or and from والكلمات الى زي كده مع حاجات ال three-grams زي work at home هنلاقي كلمة at هنا ليه تاثير كبير ، وبرضة من الحاجات المهمة هو ال normalization وهو اني اخلي ال text يكون ليه form واحدة وبينفع مع بعض ال application برضه على عكس application تانية زي موضوع التشكيل بالنسبة للغة زي العربى ساعات كتير بيغير المعنى لكنه بيوفر computation power كبير جدا بسبب اني الكلمة الواحدة ليه تشكيلات مختلفة لكن ممكن في classification problem ده ميكونش مهم او في حاجة زي Information Retrievers. كل الحاجات ديه وغيرها بتساعد في تقليل ال vocabulary الى عندي ، وبذات حاجه زي ال stemming and lemmatization وده الى هنجيلة باذن الله البوست الجاي .

Post 8

الكلمات في اللغة بيكون ليه inflections مختلفة أشكال مختلفة سواء فيه prefix or suffix او سواء تشكيل الكلمات زي العربى والتعامل مع ال inflections ديه بيكون عن طريق stemming or lemmatization ال stemming بيرجع الكلمة لل root بتاعها وبيشيل ال prefix وال suffix للكلمات زي كلمات ending, ends be end لكن لو كلمة زي running هحتاج اشيل ال ning وكلمة زي sing لو شيلت ال ing هيتبقى s وغيره لذلك بلاقي كلمات كتيرة مبقاش ليه معنى في اللغة اصلا ، على عكس ال lemmatization بيتعامل مع قواعد اللغة والمحتوى الى الكلمة جت فيه + اعتبار الشكل المناسب مع السياق من المعانى المختلفة الى الكلمة ممكن تيجي فيه من فعل او اسم او غيره ، لذلك فيه parameter مهم وهو ال Part of speech موجود لما بنستخدم nltk lemmatizer وهو pos عشان احدد لما يكون عندي اكثر من معنى للكلمة الافضل تكون في استخدام الفعل ولا الاسم ولا الصفة وغيره ، ورغم اني ال stemming بيؤدى ساعات لكلمات مش من اللغة الا انه اسرع من ال lemmatization .