

Catégorisez Automatiquement des Questions

L.khellouf¹ Mentor: B.Beaufils²

¹²OpenClassRooms

Projet 21, mai 2021

- Problématique
- Introduction
- Récupération Des Données
- Data Profiling
- Nettoyage de la donnée
- Préprocessing
- Modelisation et Analyse
- Evaluation
- Api de recommandation de tags
- Conclusion

[Stack Overflow](#) est un site de questions-réponses liées au développement informatique. Pour poser une question sur ce site, il faut entrer plusieurs étiquettes liées à la question. L'objectif de ce projet est de développer un système de suggestion de tags pour une question posée sur le ce site. Le but est d'aider les membres et aussi les débutants du site à proposer une classifications préliminaire et automatique à leurs questions afin de gagner du temps et ainsi avoir des réponses potentiellement plus pertinentes.

Introduction

Dans un premier temps, je vais récupérer les données à partir d'une API du site [Stack Overflow](#), puis je vais les analyser et traiter en utilisant des méthodes propres au traitement du langage naturel afin d'en tirer tout leur potentiel.

Dans un second temps, je vais mettre en œuvre 2 approches différentes de recommandation de tags. La première, [supervisée](#), pour suggérer les étiquettes ou tags La seconde, [non supervisée](#), pour suggérer les mots clés ou keywords.

Le système de recommandation de tags mettant en œuvre les 2 approches sera intégré au travers d'une simple application web. Pour finir, des points d'améliorations seront proposées.

Récupération Des Données

- Récupération de 11935 questions pour 23 colonnes en utilisant l'API Stack exchange
- Aucun valeur manquante et aucune question en doublon.
- Fusion des colonnes Body et Title dans une colonne appelée Text

The screenshot displays the StackExchange Data Explorer interface. At the top, there's a navigation bar with 'StackExchange Data Explorer', 'Home', 'Queries', 'Users', and a 'Compose Query' button. Below this, the 'Viewing Query' section shows a query: `SELECT * FROM posts WHERE FavoriteCount > 50`. To the right of the query is a 'Database Schema' panel showing the structure of the 'Posts' table, including columns like 'Id', 'PostTypeId', 'AcceptedAnswerId', 'ParentId', 'CreationDate', 'DeletionDate', 'Score', 'ViewCount', 'Body', 'OwnerUserId', and 'OwnerDisplayName'. Below the query and schema, there's a 'Results' panel showing a table of query results. The table has columns: 'Title', 'Tags', 'AnswerCount', 'CommentCount', 'FavoriteCount', 'CreationDate', and 'CommunityOwnedDate'. The results show various questions and their associated statistics.

Title	Tags	AnswerCount	CommentCount	FavoriteCount	CreationDate	CommunityOwnedDate
Difference between global maximum and min...	javascript	1	3	66		
Is Mono ready for prime time?	java open-source mono	17	3	79	2015-08-03 15:54:59	
Why do we need entity classes?	sql database orm sqlalchemy	40	5	74	2016-06-10 19:17:52	2016-12-16 19:17:10
Quick way to map a RegEx to a list of strings?	python regex regular-expressions	17	3	90	2016-10-20 10:10:09	
FluentBuilder() fails for Programming?	scala builder-builder-pattern	114	3	436	2014-09-26 22:06:47	2016-09-16 22:26:21
Why should I learn Lua?	lua lua-programming lua	26	4	62	2012-08-10 12:10:53	2016-04-02 13:05:33
Learning Regular Expressions	regex	1	3	522	2011-11-28 10:50:46	2016-11-11 17:31:17
How does a "back number" error and how d...	javascript jquery/jquery	9	1	51		
FluentBuilder() fails for Programming?	scala builder-builder-pattern	8	1	64		
How can I play sound in Java?	java audio	10	1	64		
Change templates in Rails	ruby-on-rails rails-on-rails rails-on-rails	10	3	66		
SQL: addDate in MySQL alternative	sql sql	20	3	74		
Python module for converting PDF to text	python pdf pdf-to-text pdf-to-text	13	4	68	2015-03-25 13:26:02	
How to find a Java Memory Leak	java memory-memory-leak java	12	2	66		
How do you beta test an iPhone app?	ios ios-testing	9	2	336		
How many constructor arguments to too many?	parameters constructor constructor	15	1	65		

Data Profiling

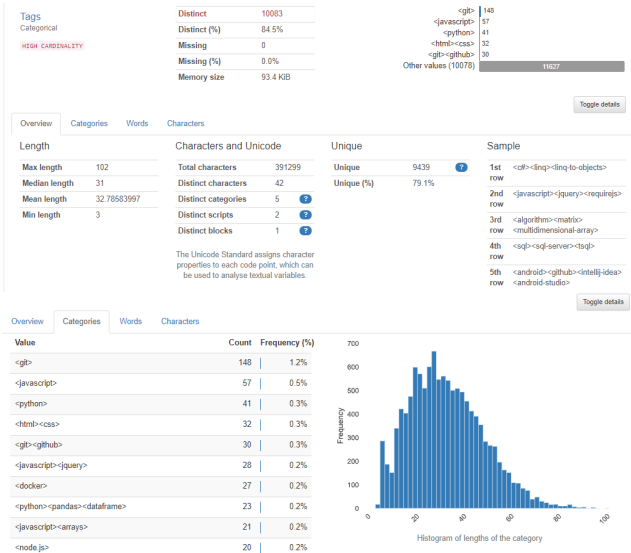
Analyse de la variable Body



Data Profiling

Analyse de la variable Tags

Il existe **5407** tags distincts, et nous avons conservé que **20** tags les plus fréquents.



Nettoyage de la donnée

- Dans la nouvelle dataset, nous avons gardé uniquement les variables [Body](#), [Title](#), [Tags](#)
- Suppression des balises [HTML](#) pour la variable Body avec le package [Beautiful Soup](#)
- Mise en minuscules de l'ensemble de champs avec la fonction [lower\(\)](#)
- Remplacer les formes contractées par leur forme longue par l'utilisation d'un dictionnaire de correspondance.
- Suppression des stop words et ponctuation avec [Spacy](#), usage des attributs ['is_stop'](#) et ['is_punst'](#) sur les tokens

- Suppression des chiffres y compris à l'intérieur des chaînes.
- Suppression des espaces multiples et des mots de tailles faibles (inférieures à 2)
- Lemmatisation pour réduire les mots à leur forme neutre canonique

Conclusion après NLP :

- Tous les documents possèdent au moins un tag.
- Il y a 8 763 questions

Transformation de la variable **Tags** à l'aide d'un **MultilabelBinarizer** pour la modélisation supervisée.

Tags						
Doc 1	[c#]	MultilabelBinarizer				
Doc 2	[Python]					
Doc 3	[javascript, jquery]					

	C#	python	javascript	jquery
Doc 1	1	0	0	0
Doc 2	0	1	0	0
Doc 3	0	0	1	1

Préprocessing

Bag of Words

Actuellement, nous avons les messages sous forme de listes de tokens (également appelés lemmes) et nous devons maintenant convertir chacun de ces messages en un vecteur avec lequel les modèles d'algorithme de **SciKit Learn** peuvent fonctionner.

Bag of words (BoW)

Very good drama although it appeared to have a few blank areas leaving the viewers to fill in the action for themselves. I can imagine life being this way for someone who can neither read nor write. This film simply smacked of the real world: the wife who is suddenly the sole supporter, the live-in relatives and their quarrels, the troubled child who gets knocked up and then, typically, drops out of school, a jackass husband who takes the nest egg and buys beer with it. 2 thumbs up.... very very very good movie.



(the', 8),
(', 5),
(very', 4),
(', 4),
(who', 4),
(and', 3),
(good', 2),
(it', 2),
(to', 2),
(a', 2),
(for', 2),
(can', 2),
(was', 2),
(of', 2),
(drama', 1),
(although', 1),
(appeared', 1),
(have', 1),
(few', 1),
(blank', 1)

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 vectorizer = CountVectorizer()
3 corpus = [
4     "This is the first document.",
5     "This is the second second document.",
6     "and the third one.",
7 ]
8
9 X = vectorizer.fit_transform(corpus)
10 print(X)
```

```
(0, 8) 1
(0, 3) 1
(0, 6) 1
(0, 2) 1
(0, 1) 1
(1, 8) 1
(1, 3) 1
(1, 6) 1
(1, 2) 1
(1, 5) 2
(2, 6) 1
(2, 0) 1
(2, 7) 1
(2, 4) 1
```

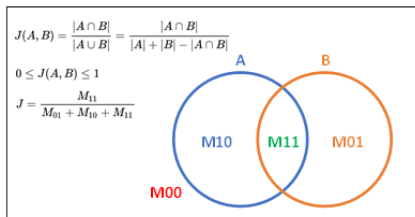
Cette méthode consiste à compter le nombre d'occurrences des tokens présents dans le corpus pour chaque texte, que l'on divise ensuite par le nombre d'occurrences total de ces même tokens dans tout le corpus.

Words	TF (for A)	TF (for B)	IDF	TFIDF (A)	TFIDF (B)
Jupiter	1/5	0	$\ln(2/1) = 0.69$	0.138	0
Is	1/5	1/8	$\ln(2/2) = 0$	0	0
The	1/5	2/8	$\ln(2/2) = 0$	0	0
largest	1/5	0	$\ln(2/1) = 0.69$	0.138	0
Planet	1/5	1/8	$\ln(2/2) = 0$	0.138	0
Mars	0	1/8	$\ln(2/1) = 0.69$	0	0.086
Fourth	0	1/8	$\ln(2/1) = 0.69$	0	0.086
From	0	1/8	$\ln(2/1) = 0.69$	0	0.086
Sun	0	1/8	$\ln(2/1) = 0.69$	0	0.086

Le but de cette méthode est réduire l'importance relative de certains mots très présents mais peu sinifiatifs.

Métrique utilisée est Indice de Jaccard:

- comparer un ensemble de labels prédits à un échantillon aux labels réels correspondant.
- Indice de Jaccard pondéré détermine la moyenne des métriques calculée pour chaque label, pondérée par leur distribution réelle observée.

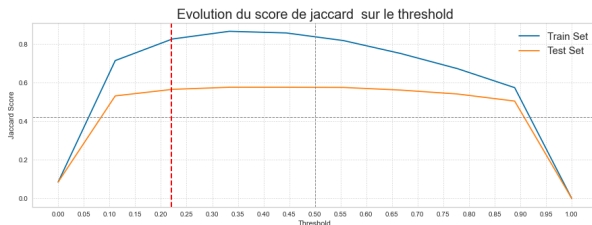


- Mise en place d'une optimisation d'hyperparamètres sur la représentation **Bag_of_words** et **TF_IDF** Pour la représentation vectorielles de mots.
- L'algorithme utilisé est **Regression logistic** et il dispose dans sa documentation les meilleurs paramètres à optimiser.

	Les hyperparamètres à optimisé sur Bag_of_words et TF_IDF	Best hyperparamètres sur TF_IDF
✓ OneVsRestClassifier	✓ Params= {'estimator__penalty': ['l1', 'l2'], 'estimator__C': [0.1, 1, 10, 100, 1000]}	✓ {'estimator C': 10, 'estimator_penalty': 'l1'}
✓ KFold		
✓ GridSearch	✓ KFold= 5	✓ Jaccard Score= 63.94
✓ LogisticRegression		

Approche Supervisée

Optimisation du threshold

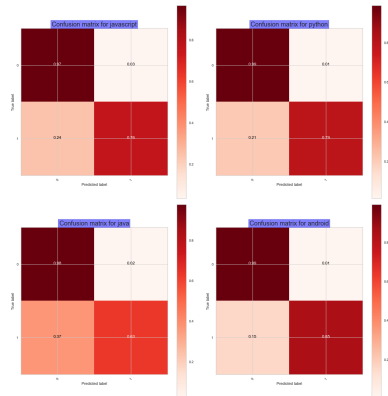
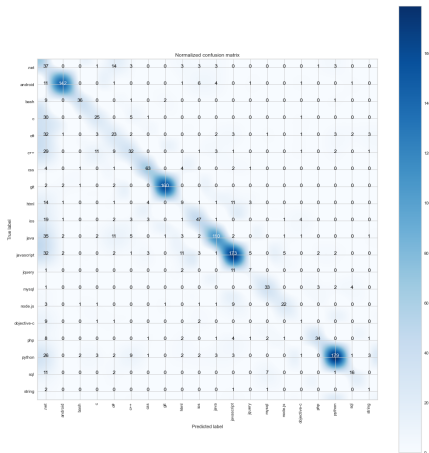


	y_true	y_pred_0.5	y_pred_0.22
20	(css, git, html)	(css, html, javascript)	(css, html, javascript)
21	(java,)	(java,)	(java,)
22	(c, c++)	(c,)	(c, javascript)
23	(git,)	(git,)	(git,)
24	(mysql, python)	(python,)	(python,)
25	(java,)	(java,)	(c#, java)
26	(git,)	(git,)	(git,)
27	(css,)	(css,)	(css, html)
28	(objective-c,)	()	()
29	(mysql, php)	(mysql,)	(mysql,)
30	(android,)	(android,)	(android,)
31	(python,)	(python,)	(python,)
32	(git,)	(git,)	(git,)
33	(ios,)	(objective-c,)	(ios, objective-c, python)
34	(.net, c#)	(c#, string)	(c#, string)
35	(python,)	(python,)	(python,)
36	(ios, objective-c)	()	()
37	(git,)	(git,)	(git,)
38	(css,)	(css, html)	(css, html)
39	(git,)	(git,)	(git,)

- Modifier le seuil par défaut ($=0.5$) n'améliore pas la capacité du modèle à mieux généraliser.
- Réduire le seuil à 0.22 a pour effet de favoriser la capacité du modèle à fournir une prédiction sans en dégrader la performance.

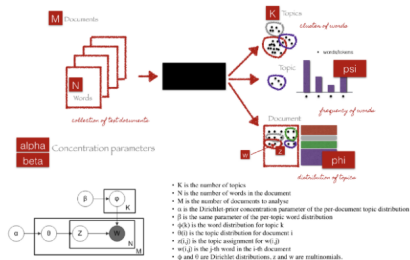
Approche Supervisée

Matrice de Confusion



Approche Non Supervisée

Nous avons utilisé l'algorithme **LDA(Latent Dirichlet Allocation)** qui est une méthode non supervisée générative pour les mots clés.



<http://chdoig.github.io/pytexas2015-topic-modeling/#3/4>

Nous pouvons décrire ce processus génératif comme, étant donné le nombre M de documents, le nombre N de mots et le nombre K préalable de sujets, le modèle s'entraîne à produire:

- ψ , la distribution des mots pour chaque sujet K
- ϕ , la répartition des sujets pour chaque document i .

Approche non supervisée

- Inférence : déterminer les thèmes, les distributions de chaque mot sur les thèmes, la fréquence d'apparition de chaque thème sur le corpus.
- Pour effectuer l'inférence de ce modèle, nous avons utilisé la librairie LDA de gensim.

```
1 no_top_words = 10
2 n_topics = 5
3 def display_topics(model, no_top_words, num_topics):
4     for idx, topic in model.show_topics(formatted=False, num_topics=num_topics, num_words=no_top_words):
5         print("-----")
6         print("Topic %d:" % (idx))
7         print(" ".join([w for w in topic]))
8         print("-----")
9
10 display_topics(lda_model, no_top_words, n_topics)
```

Topic 0:
java public int class new void string return static method

Topic 1:
file git use branch commit instal run error work line

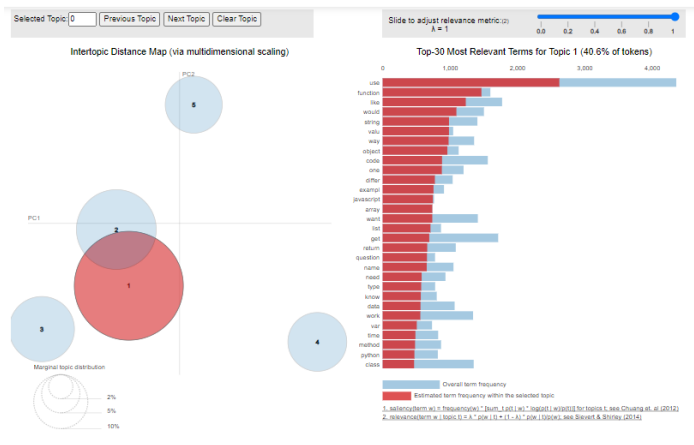
Topic 2:
use function like would string valu way object code one

Topic 3:
android view div layout imag button text height width content

Topic 4:
js node request server modul user use error get data

- Perplexity: -7.499418620423826
- Coherence Score: 0.5206279823824616

Approche non Supervisée



Le model est un bon topic car les bulles sont assez grosses et non superposées dispersées dans le graphique.

Api de recommandation de tags

- **API** propose une liste de tags [Stack Overflow](#) relatifs à une question saisie traitant de sujets informatiques.

Recommandation Tags Application

khellouf leila project

Enter your question Here

I'm kind of going nuts here. I need to start writing a stand alone non web application in JavaScript.

My problem is that I can run (using nodejs) the script just fine, but I can't split the code that I need to write among multiple files because then I don't know how to include them.

Unsupervised Tags Recommendation

and is of it string array int is git file commit div branch convert imag android css string json format

Supervised Tags Recommendation

```
@ : [  
  0 : "javascript"  
  1 : "jquery"  
  2 : "node.js"  
]
```

- Le text saisi passe par toutes les étapes de préporocessing NLP puis est transformé en matrices TF-IDF avant application respectivement des modèles supervisée et non supervisée.

- **Data set**: entrainer les modèles sur une grande base de donnée.
- Utilisation de plusieurs modèles supervisée.
- **le prétraitement**: Manipulation des URL, emoji par exemple,..etc
- **Apprentissage des features**: Utiliser les dernières techniques issues du deep learning comme BERT ou ELMo.
- Trouver de nouvelle **métrique d'évaluation**

The End