



Project 5 Data Science / OpenClassrooms

Presented by:

Marc Felix DEGNI

Apprentice Data Scientist Consultant OC/CC

AutoTag Suggestion StackOverFlow



Data description

The screenshot displays the Stack Overflow interface for questions tagged with 'python'. The page title is 'Questions tagged [python]'. A search bar at the top contains '[python]'. The left sidebar shows navigation links: Home, PUBLIC, Stack Overflow, Tags, Users, Jobs, and Teams. The main content area lists questions. The first question is 'What does the "yield" keyword do?' with 8486 votes, 38 answers, and 1.8m views. It is annotated with a blue circle around the title, a red bracket around the body text, and a green box around the tags. The tags are 'python', 'iterator', 'generator', 'yield', and 'coroutine'. The right sidebar shows 'Related Tags' and 'Hot Network Questions'.





Annotations:

- Title:** What does the "yield" keyword do?
- Body:** What is the use of the yield keyword in Python? What does it do? For example, I'm trying to understand this code: `def _get_child_candidates(self, distance, min_dist, max_dist): if self...`
- Tags:** python, iterator, generator, yield, coroutine

Data cleansing

StackExchange

Search on Super User...



Log In

Sign Up

Home

Questions

Tags

Users

Unanswered

Ask a question


Title


What's your computer software or computer hardware question? Be specific.


Body


B


I





























LinksImagesStyling/HeadersListsBlockquotesCodeHTMLadvanced help »

Tags

e.g. (hard-drive windows-8.1 macos)

How to Ask

Is your question about computer software or computer hardware?

We prefer questions that can be *answered*, not just discussed.

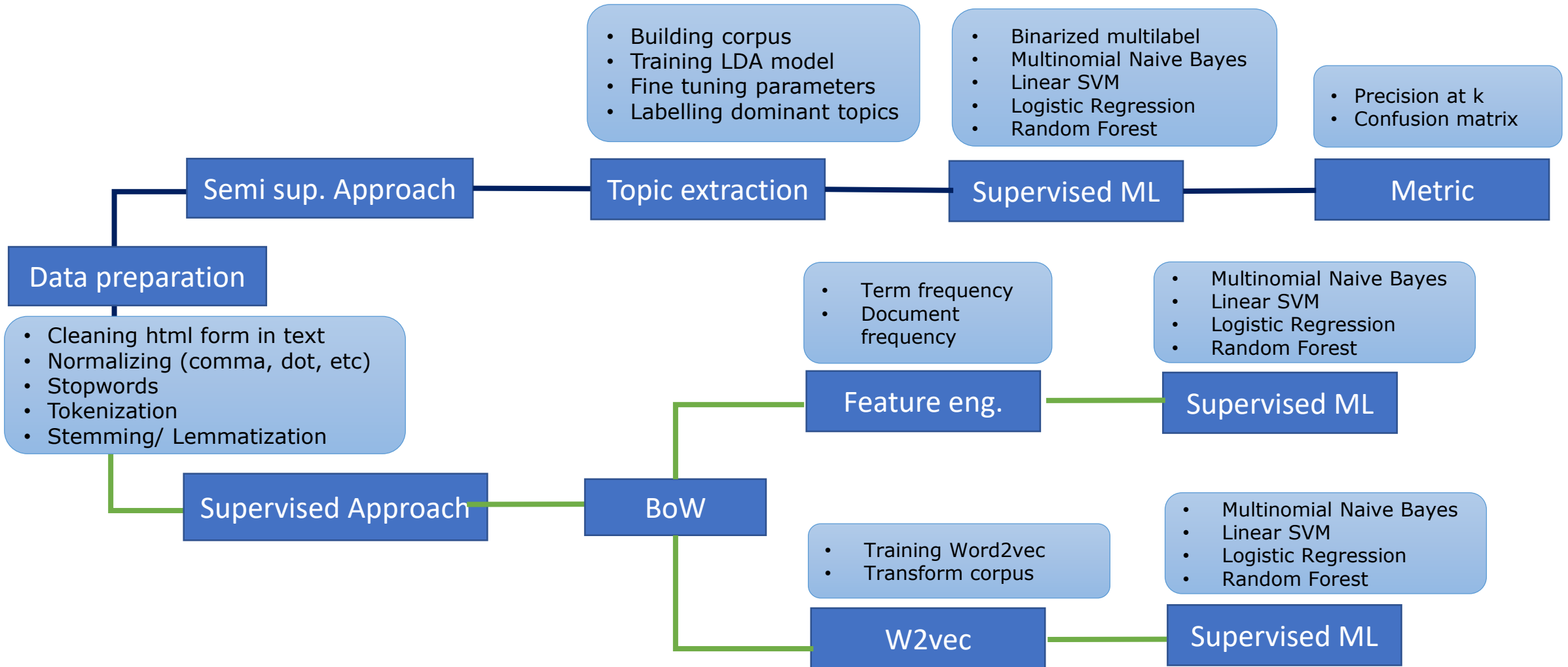
Provide details. Share your research.

If your question is about this website, [ask it on meta](#) instead.

[visit the help center »](#)
[asking help »](#)

Flight delays prediction

Data approach



Data cleansing

	Body	Title
0	<p>I want to use a track-bar to change a form'...	Convert Decimal to Double?
1	<p>I have an absolutely positioned <code>div</code>...	Percentage width child element in absolutely p...
3	<p>Given a <code>DateTime</code> representing ...	How do I calculate someone's age in C#?
4	<p>Given a specific <code>DateTime</code> valu...	Calculate relative time in C#
6	<p>Is there any standard way for a Web Server ...	Determine a User's Timezone

token_tag
[c#, floating-point, type-conversion, double, ...]
[html, css, css3, internet-explorer-7]
[c#, .net, datetime]
[c#, datetime, time, datediff, relative-time-s...
[javascript, html, browser, timezone, timezone...

Questions	token_questions	tokens_clean	tokens_clean_lemma
Convert Decimal to Double? <p>I want to use a track-bar to change a form's opacity. </p>\r\n\r\n<p>This is my code:</p>\r\n\r\n<pre><code>decimal trans = trackBar1.Value / 5000;\r\nthis.Opacity = trans;\r\n</code></pre>\r\n\r\n<p>When I build the application, it gives the following error:</p>\r\n\r\n<blockquote>\r\n<p>Cannot implicitly convert type <code>'decimal'</code> to <code>'double'</code>.</p>\r\n</blockquote>\r\n\r\n<p>I tried using <code>trans</code> and <code>double</code> but then the control doesn't work. This code worked fine in a past VB.NET project.</p>\r\n	['convert', 'decimal', 'to', 'double', 'i', 'want', 'to', 'use', 'a', 'track', 'bar', 'to', 'change', 'a', 'form', 's', 'opacity', 'this', 'is', 'my', 'code', 'when', 'i', 'build', 'the', 'application', 'it', 'gives', 'the', 'following', 'error', 'cannot', 'implicitly', 'convert', 'type', 'to', 'i', 'tried', 'using', 'and', 'but', 'then', 'the', 'control', 'doesn', 't', 'work', 'this', 'code', 'worked', 'fine', 'in', 'a', 'past', 'vb', 'net', 'project']	['convert', 'decimal', 'double', 'want', 'use', 'track', 'bar', 'change', 'form', 'opacity', 'code', 'build', 'application', 'gives', 'following', 'error', 'implicitly', 'convert', 'type', 'tried', 'using', 'control', 'doesn', 'work', 'code', 'worked', 'fine', 'past', 'vb', 'net', 'project']	['convert', 'decimal', 'double', 'want', 'use', 'track', 'bar', 'change', 'form', 'opacity', 'code', 'build', 'application', 'give', 'follow', 'error', 'implicitly', 'convert', 'type', 'try', 'use', 'control', 'does', 'work', 'code', 'work', 'fine', 'past', 'vb', 'net', 'project']

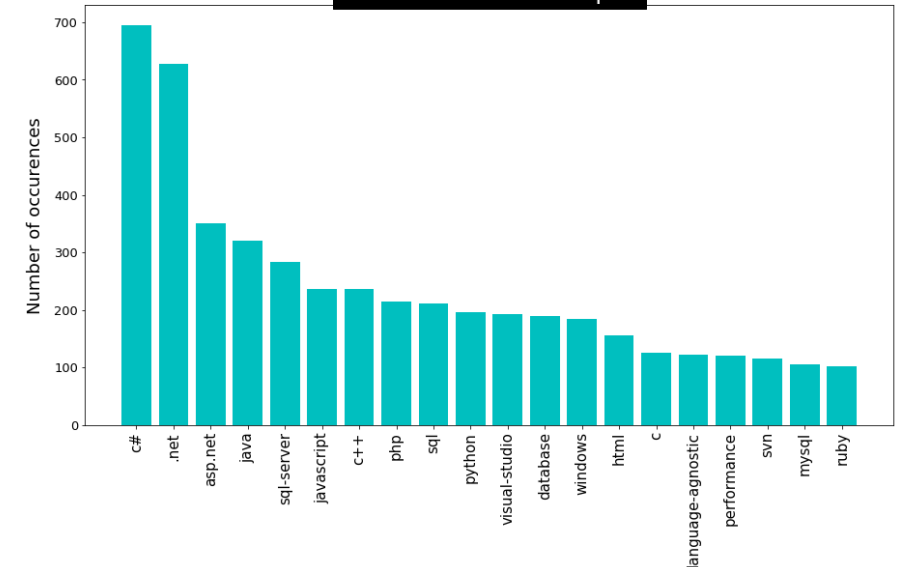
Cleaning

Tokenization

Stopword

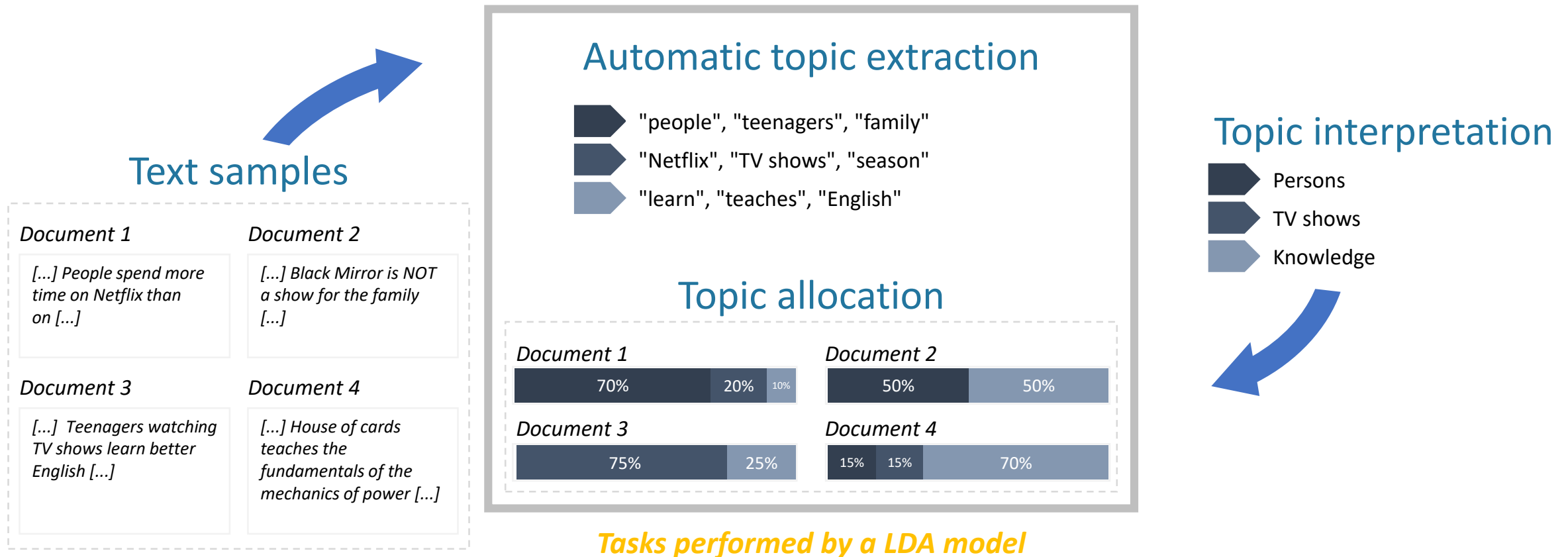
Lemmatization

Words occurrence : top20



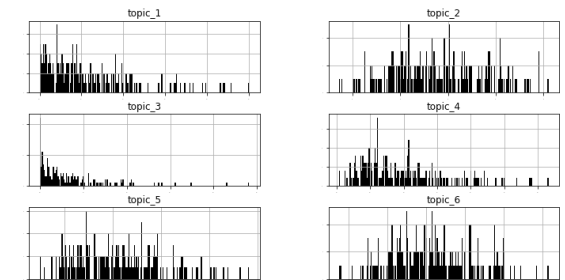
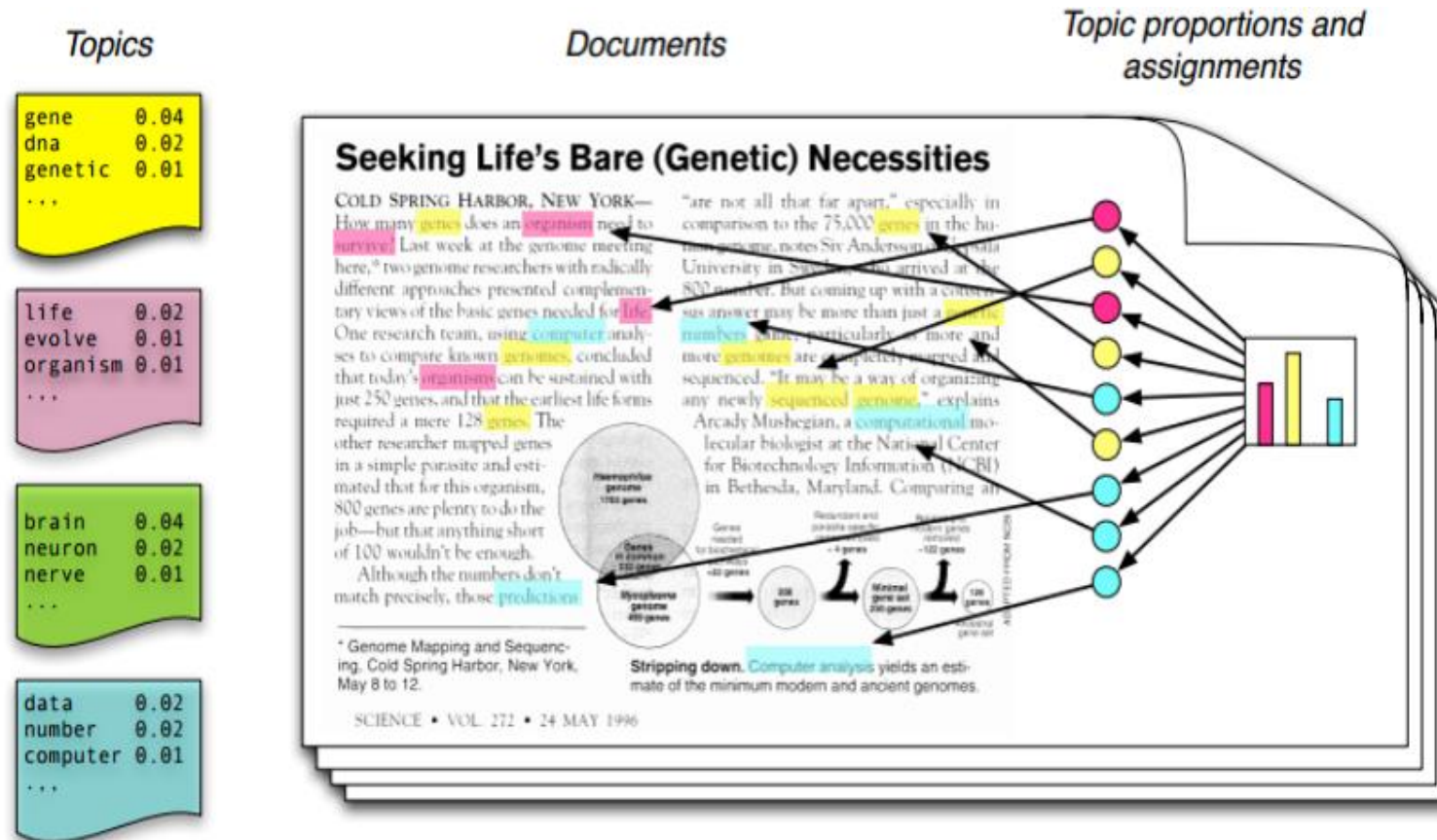
Topic extraction (Part I)

Latent Dirichlet Allocation:
an unsupervised algorithm for topic modelling

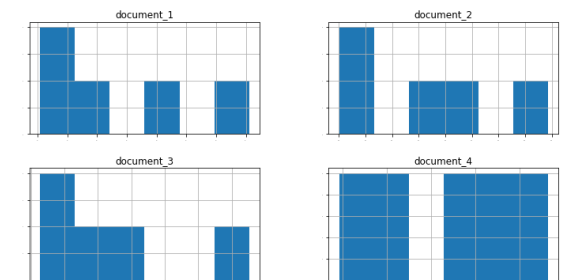


Topic extraction (Part 2)

LDA optimizes the probability that documents occur knowing topics are distributed among them



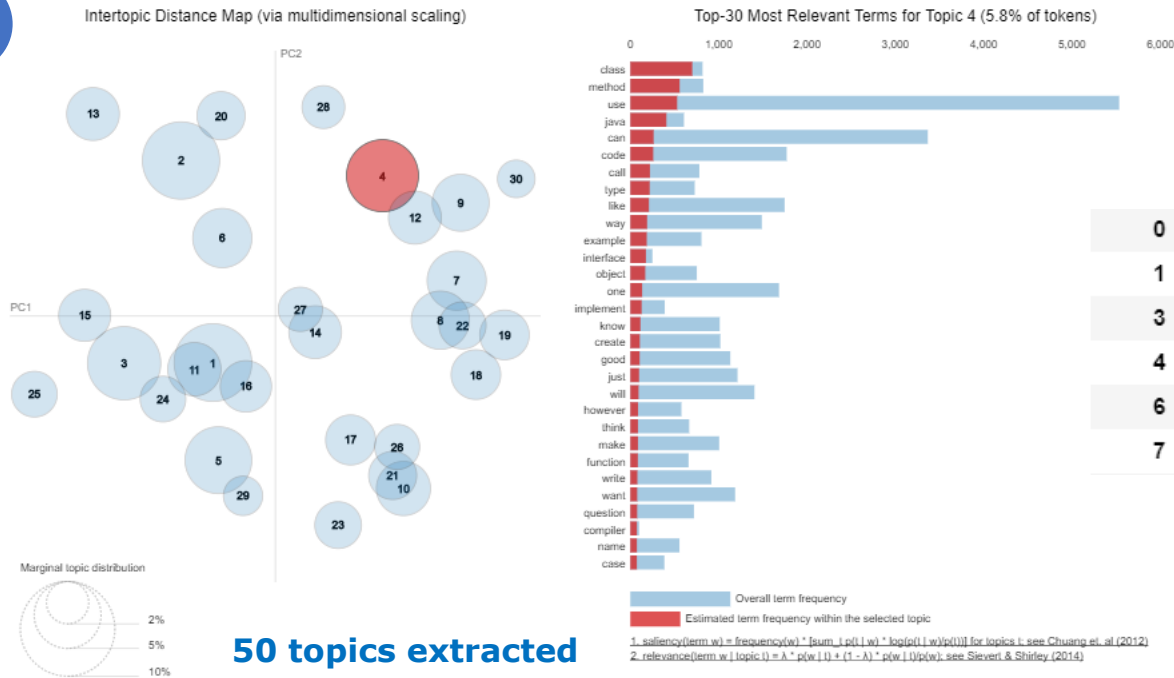
Topics are represented by different word probabilities



Similarly, documents are represented by different topic probabilities

Results approach and performance

1



2

Dominant topic labelled

Dominant_Topic	token_tag
0	4.0 [c#, floating-point, type-conversion, double, decimal]
1	1.0 [html, css, css3, internet-explorer-7]
2	2.0 [c#, .net, datetime]
3	27.0 [c#, datetime, time, datediff, relative-time-span]
4	20.0 [javascript, html, browser, timezone, timezoneoffset]

3

How to deal with multilabel ? Binarizing !!!

Multilabel (Targets)

0	[c#, floating-point, type-conversion, double, decimal]
1	[html, css, css3, internet-explorer-7]
3	[c#, .net, datetime]
4	[c#, datetime, time, datediff, relative-time-span]
6	[javascript, html, browser, timezone, timezoneoffset]
7	[.net, math]

Labels Binarized

```
array([[0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0],
       [1, 0, 0, ..., 0, 0, 0],
       [0, 0, 0, ..., 0, 0, 0]])
```

4

Supervised method (LR/ MNB)

Probabilities		Probability	
c#	0.196276	c#	0.196276
.net	0.120634	.net	0.120634
c++	0.111925	c++	0.111925
php	0.074453	php	0.074453
c	0.070478	c	0.070478

Total supervised approach

Bag of word : TF-IDF

Bag of words

- Consider a corpus of documents, each document is a text
- Each document contains a certain number of words for which we can do tokenization
- Each word will be a token
- Bag-of-word: representing without any order the words in a document

Example:

- "It is super cool to be a Data Scientist"
- - "super" "it" "scientist" "be" "to" "cool" "is" "a" "data"

Term Frequency principle

- After creating a bag of words, it is necessary to think about the importance of each word in the text
- $tf(t, d)$ can be read as the term frequency of the term t in document d

- Formula:
$$tf(t, d) = \frac{n_{t,d}}{\sum_{i \in d} n_{i,d}}$$

- With $n_{i,d}$ the number of appearance of word i in document d
- Thus $\sum_{i \in d} n_{i,d}$ is the total number of words in the document d

Example:

For the next slides we will use the following examples:

- Doc 1: "It is super cool to be a Data Scientist"
- Doc 2: "Data Science is done with data by Data Scientists"
- Doc 3: "Poker is a cool game"

By using stopwords & lemmatization, only the following words remain: "super", "cool", "data", "science", "do", "poker", "game"

Bag of word : TF-IDF

Principle of TF-IDF

- Continuity of the previous slide
- It is more likely to have an insight on words which are really important
- The weights of the words which are extremely frequent like “the” or “of” are diminished
- It deals with how much information the word gives for the query
- Formula of Inverse Document Frequency:

$$idf(t) = \log \left| \frac{n}{\{d \in D: t \in d\}} \right|$$

- With n the total number of documents, D the set of all documents, d a given document in D , t is a given term. Here, the natural logarithm (base e) is used
- Thus, $\{d \in D: t \in d\}$ is the number of documents in which the term t appears
- Once Term Frequency and Inverse Document Frequency has been found, it is possible to compute TF-IDF:

$$tfidf(t, d) = tf(t, d) \times idf(t)$$

TF-IDF on data

	aa	aad	aardvark	ab	abandon	abap	abbreviate	abbreviation	abrrviate	abc	...	éâè	ôäü	šturc	caught
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0
5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0

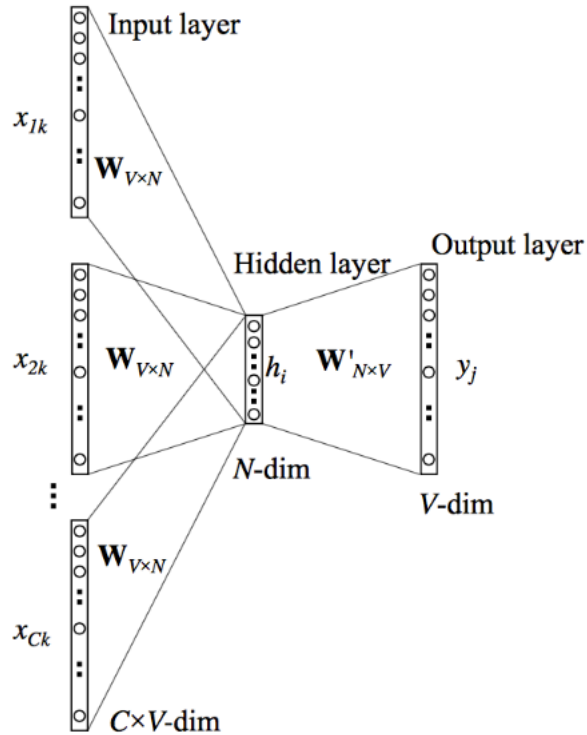
6 rows × 10757 columns

Supervized model

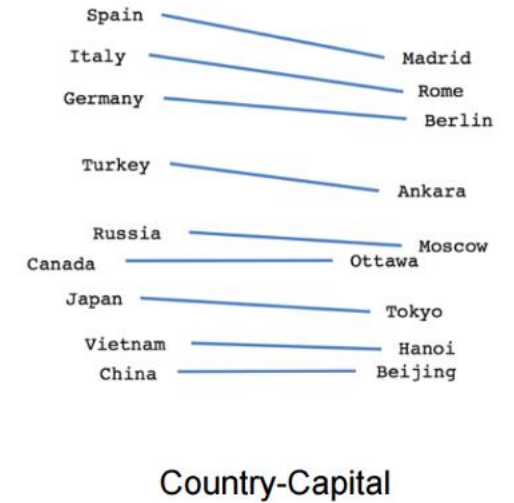
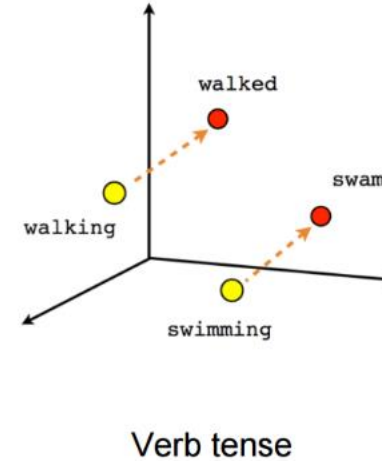
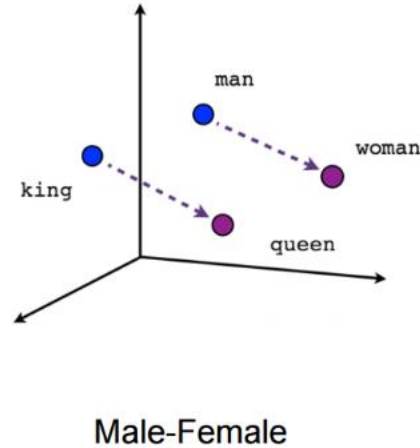
Model	BoW	Metrics
Multinomial NB	TF-IDF	27,12%
Linear SVM	TF-IDF	24,37%

Bag of words : Word Embedding (Word2vect)

Intuition behind Word2Vec



Représentation du modèle CBOW



Example of word similar

```
model.most_similar('convert')
executed in 16ms, finished 14:54:44 2018-09-17

[('hashcode', 0.9891046285629272),
 ('hash', 0.9856103658676147),
 ('iterate', 0.9854134321212769),
 ('aq', 0.9842157959938049),
 ('timezone', 0.9826910495758057),
 ('chicago', 0.9822856783866882),
 ('represent', 0.9810150861740112),
 ('decrypt', 0.9808613657951355),
 ('concatenate', 0.980566143989563),
 ('representation', 0.9798986315727234)]
```

Example of word vectorized

Supervised model

Model	BoW	Metrics
Multinomial NB	word2vec	8,90%
Linear SVM	word2vec	2,62%

DEPLOY MODEL ON HEROKU CLOUD

- Creating API with microframework (FLASK)
- Deploy on heroku cloud

Link for demonstration

<https://autotag-suggestion-nlp.herokuapp.com/>

CONCLUSION

- ✚ **The reduction of the sparsity** of the TF-IDF matrix by a principal component analysis or other dimension reduction method. The algorithms used will then have less computational problems and can focus their intelligence on the words having weight in at least one document.
- ✚ Another improvement would be **the deletion of words with a weight below a fixed threshold**, because these features will be common to each document and therefore will not bring any value.
- ✚ With the word2vec, instead of using a MeanEmbedding which aggregates the values of the weights of the words present in the questions, one could combine **word2vec and TF-IDF method**.