ETUDES NUMERIQUES ET INNOVATION



Topic Modeling

Bilel Benbouzid - Alexis Perrier - Décembre 2017

Slack: <u>upem-numi.slack.com</u>

Invitation: bit.ly/2yyUF83

Alexis Perrier - Data Scientist



- Learning analytics @berklee.edu
- Signal processing @splice.com
- NLP, healthcare @docenthealth.com

- Twitter: <a>@alexip
- Linkedin.com/in/alexisperrier
- alexis.perrier@gmail.com

Post-Discharge Call



Plan de la semaine



2 corpus:

- Est républicain 1999
- God Emperor Trump https://www.facebook.com/GodEmperorTrump



3 Labs:

- Lundi: Gensim python Est Républicain
- Mardi: STM R Est Républicain
- MJV (3j): God Emperor Trump



Outils de travail





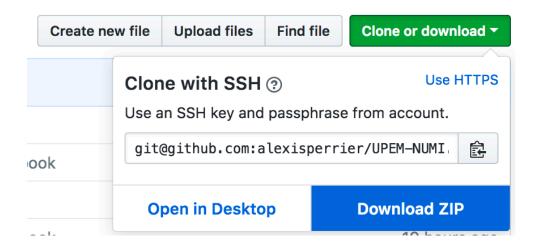


• Slack: <u>upem-numi.slack.com</u>

• Invitation: bit.ly/2yyUF83

• Notebooks: <u>35.196.112.120:5000</u>, jupyter.org

• Code, slides, sur Github: https://github.com/alexisperrier/UPEM-NUMI



Topic Modeling, pourquoi?



Vous avez des milliers de documents, comment savoir ce qu'ils contiennent ?



- Exemples
 - Article de presse sur une longue période
 - Réseaux sociaux, commentaires, tweets, ...
 - Articles scientifiques, textes légaux
 - Débats retranscrit,
 - Littérature, fiction, théâtre, ...

Top Down vs Bottom Up



Approche supervisée, top down:

- Vous définissez les sujets <u>a priori</u> comme autant de catégories.
 Vous classez les documents dans ces catégories.
 - => Classification de documents

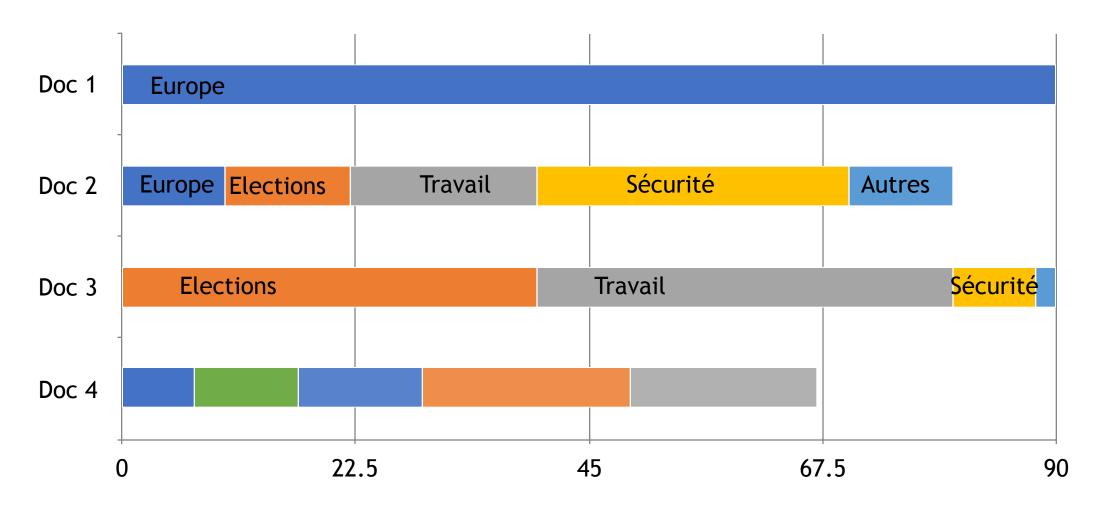
Approche non supervisée, bottom up:

- Aucun a priori préalable sur le contenu des documents Inférer les topics
 - => Topic Modeling

Topic Modeling



Chaque document peut contenir aucun, un ou plusieurs topics Chaque topic est composé de plusieurs mots



Ensuite



- Croisement avec des variables externes
 - Evolution dans le temps
 - Auteur ou Locuteur, parti politique, context,
- Utiliser les topics pour classer les documents sous topics, analyse de sentiment, vocabulaire, ...
- Réduction de dimension et modèles prédictifs

Exemples

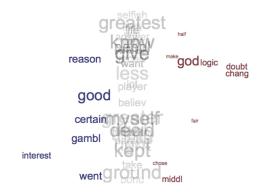


- <u>Followers sur twitter</u> (gensim, LDA, LSA, python)
- <u>Debats présidentiels américains</u> (stm, LDA, R)
- Analyse des notes en milieu hospitalier pour prédire patient satisfaction et identifier les problèmes

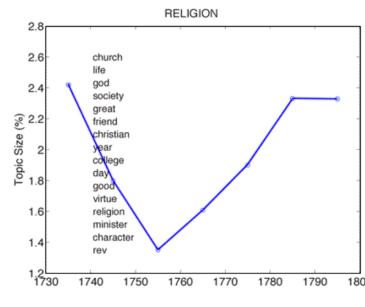
Autres travaux:

- Gazette de Pensylvanie http://www.common-place-archives.org/vol-06/no-02/tales/
- 3,346 works of 19th-century British, Irish, and American fiction evolution dans le temps
- Analyse graph / réseau des topics dans Proust (Stanford, Mallet, Java)
- <u>Structural Topic Models for Open-Ended Survey Responses</u> (Molly Roberts, STM)
- Topic modeling of <u>French crime fiction novels</u>

FIGURE 15 Intuitive Topic Allowing for Different Vocabularies Based on Gender



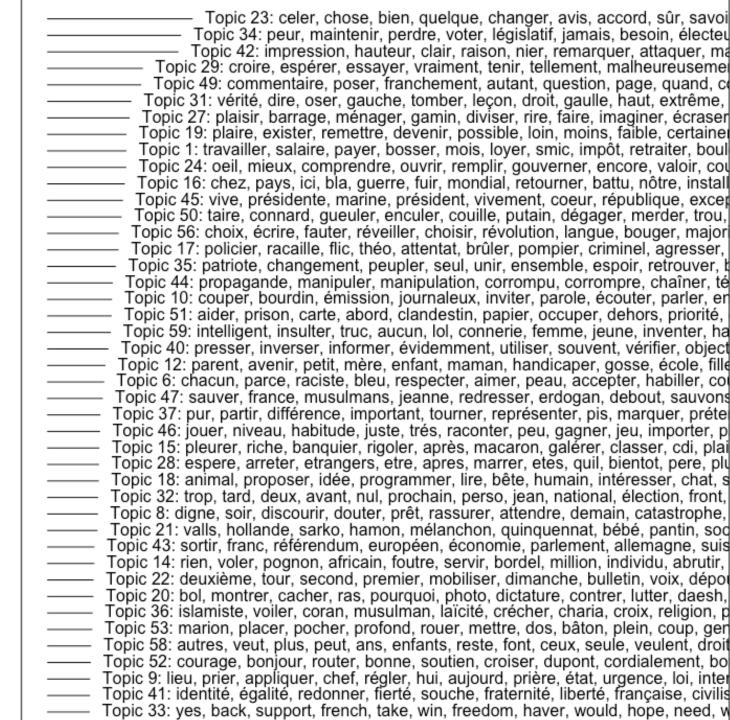
Male Female



Autre exemple

#Datapol

Facebook FN - Commentaires 60 topics



Etapes



Définition

- 1. Corpus brut
- 2. Unité d'analyse: tweet(s), commentaire(s), paragraphe(s), article(s), blog, livre, ...
- 3. Technique utilisée: LSA, LDA, pLSA, HLDA, DTM, Topic2Vec
- 4. Variable externes?

Post-Discharge Call

Etapes:

- 1. Pre-processing:
 - 1. Enlever le bruit, garbage-in garbage-out
 - 2. Adapter, transformer le contenu
- 2. Topic Modeling: faire tourner les modèles
 - 1. Combien de topics?
- 3. Interprétation des résultats
 - 1. Nommer les topics
 - 2. Mesurer leur qualité
- 4. Visualisation
- Wordcloud!

experience nurses experienceq stay Shelp great

Corpus



- Longueur: du tweet (280) à l'article
- Langue: anglais, français, franglais
- Niveau de langue: classique, argot, smiley
- Contenu: images, urls, html, texte,
- Source: sms, tweets, forum, presse, OCR, transcript

« Nous avons des devoirs envers notre pays.

Nous sommes les héritiers d'une grande histoire et du grand message humaniste adressé au monde.

Nous devons les transmettre d'abord à nos enfants, mais plus important encore, il faut les porter vers l'avenir et leur donner une sève nouvelle. »

Emmanuel Macron, 08 mai 2017

Emmanuel Macron, 08 mai 2017





Replying to @lchexo

j'avais mm pas vu j'sais pas ou g foutu mes écouteurs mais wsh crache tte ta haine du monde jtecouterai tjrs

Influence la transformation du texte

Translate from French

Bag of Words



L'information prise en compte est intégralement contenue dans la liste de mots issue du texte, indépendamment de leur position ou de leur fonction dans la phrase

Le texte original:

« Il sera une page dans un livre de dix mille pages que l'on mettra dans une bibliothèque qui aura un million de livres, une bibliothèque parmi un million de bibliothèques. » Ionesco, Le roi se meurt 1962

devient:

aura bibliothèque bibliothèque bibliothèques dans de de de dix il on l livre livres mettra mille million million page pages parmi que qui sera un un une une une

Vectorizer le texte





Mots

Documents

	Arbre	Gazon	Été
Doc 1	3	0	0
Doc 2	1	2	1
•••			
Doc N	0	2	4

Matrice Mots - Documents

Mots - Documents - TF-IDF



Comment évaluer la fréquence des mots dans un ensemble de documents?

TF-IDF: fréquence du mot dans un document, normalisée par sa présence dans l'ensemble du corpus

Pour un document donné:

- TF: Term Frequency: combien de fois le terme est dans le document
- IDF: Inverse document term frequency:
 # de documents / # documents contenant le terme

On obtient une matrice terme - document plus représentative que un simple comptage des mots

Voir:

- https://fr.wikipedia.org/wiki/TF-IDF
- TfidfVectorizer de scikit-learn

Latent Semantique Analysis

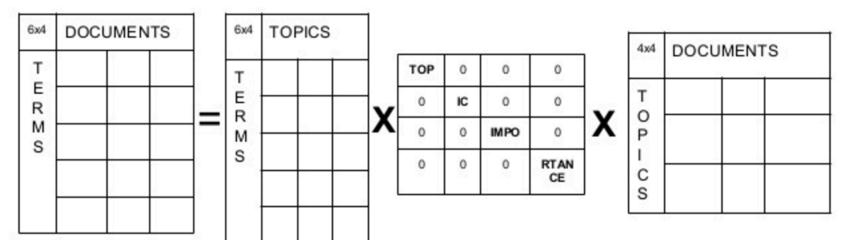


Aproche Deterministe

LSA

Nothing more than a singular value decomposition (SVD) of document-term matrix:

Find three matrices U, Σ and V so that: $X = U\Sigma V^t$



For example with 5 topics, 1000 documents and 1000 word vocabulary:

Original matrix: 1000 x 1000 = 106

LSA representation: 5x1000 + 5 + 5x1000 ~ 104

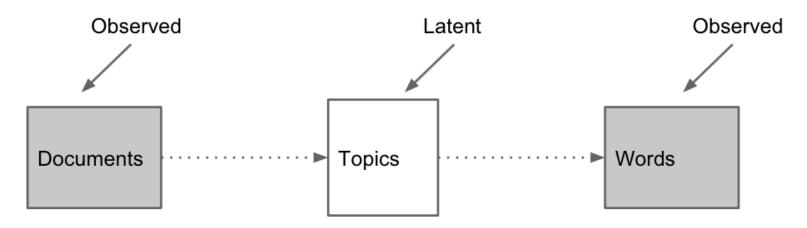
-> 100 times less space!

Latent Dirichlet Allocation



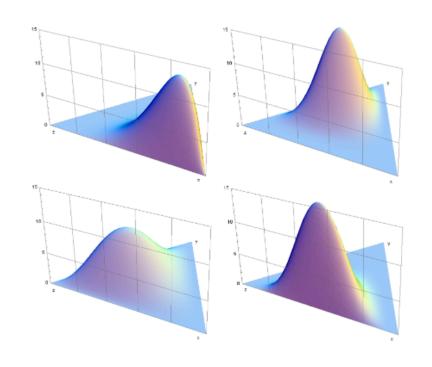
Blei 2002:

- Les topics sont distribués aléatoirement dans les documents
- Les mots sont distribués aléatoirement dans les topics
- L'aléatoire suit une distribution de Dirichlet
- K: Nombre de topics
- α: Nombre de topics par document
- B: Nombre de mots par topic



Distribution de Dirichlet

Sorte de gaussienne généralisée



The Dirichlet Distribution

Let
$$\Theta = \{\theta_1, \theta_2, \dots, \theta_m\}$$

We write:

$$\Theta \sim \text{Dirichlet}(\alpha_1, \alpha_2, \dots, \alpha_m)$$

$$P(\theta_1, \theta_2, \dots, \theta_m) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^m \theta_k^{\alpha_k - 1}$$

 Samples from the distribution lie in the m-1 dimensional probability simplex

https://fr.wikipedia.org/wiki/Loi de Dirichlet

Librairies



Python

- Gensim https://radimrehurek.com/gensim/
- LDA Python library: https://pypi.python.org/pypi/lda
- Scikit-learn: http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.LatentDirichletAllocation.html

R

- LDA, LSA
- Topicmodels
- STM package http://structuraltopicmodel.com/

Visualisation

- Gensim: LDAviz
- R STM: stmBrowser

Lab 01 - Est Républicain 1999



35.196.112.120:5000