

UNIVERSITÉ CLAUDE BERNARD LYON I

MASTER DATA SCIENCE

PROJET DATA MINING

Détection d'événements sur Twitter

Auteurs :

Gregory HOWARD 11207726

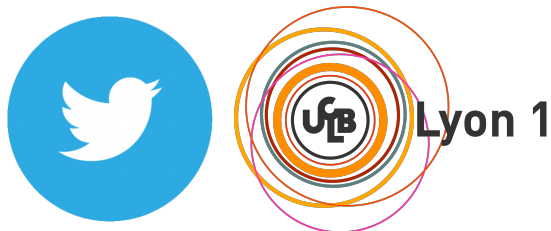
Marine RUIZ 11208141

Jules SAUVINET p1412086

Professeur :

Marc PLANTEVIT

20 janvier 2017



Résumé

L'objectif est de détecter des évènements à partir de tweets en exploitant des informations spatio-temporelles et textuelles.

Table des matières

1	Introduction	2
2	Description de nos données	2
3	L'objectif	2
4	L'approche	3
5	Description des éléments utiles à la détection d'événements	4
6	Algorithme de la détection d'événements	4
7	Les options sur le clustering	4
8	La visualisation du résultat	5
9	Interprétation de nos résultats	5
10	Notre essai sur DBScan	5
11	Conclusion	5

1 Introduction

L'objectif est de détecter des événements à partir de tweets en exploitant des informations spatio-temporelles et textuelles. Notre approche utilise la technique de détection d'événements multi-échelles dans les réseaux sociaux introduite dans l'article de recherche [1] et développé par AHMED ANES BENDIMERAD et AIMENE BELFODIL. Nous avons adapté cette version à nos données à laquelle nous avons ajouté des mécanismes permettant de filtrer les tweets envoyés par des robots en amont du clustering permettant la détection d'événements. Nous avons également développé une approche de détection spatio-temporelle d'événements avec l'algorithme de clustering DBScan puis nous avons comparé les résultats des deux types de clustering.

2 Description de nos données

Source du fichier

Nous avons effectué notre détection d'évènement à partir d'un fichier de 1,7 Go comportant 10982005 tweets récupérés sur Twitter. Chaque tweet contient les informations suivantes :

- un identifiant
- la date et l'heure à laquelle le tweet a été posté
- la position de la personne quand le tweet a été posté (latitude et longitude)
- l'identifiant du lieu géographique où a été posté le tweet
- le lieu du tweet au moment où il a été posté
- l'ensemble des références et hashtags du tweet
- l'identifiant de la personne qui a posté le tweet
- le pseudo de la personne qui a posté le tweet
- le nom complet de la personne qui a posté le tweet
- l'identifiant du lieu de création du compte de la personne qui a posté le tweet
- l'information de la certification ou non du compte de l'utilisateur ayant posté le tweet
- le nombre de personnes qui suivent la personne qui a posté le tweet
- le nombre d'amis de la personne qui a posté le tweet
- le nombre de tweet déjà posté par la personne qui a posté le tweet

Pour la détection d'évènement, nous nous sommes servi de l'identifiant du tweet, du nom de la personne qui a posté le tweet, des hashtags, de la date et de la position d'envoi du tweet.

3 L'objectif

Nous avons pour objectif de détecter des événements journaliers à l'aide des tweets de notre fichier de données. Nous allons tester deux algorithmes différents puis comparer les résultats afin de tester les performances de chacun et de garder le meilleur.

A détailler

4 L'approche

Pour notre premier algorithme, nous faisons des sous-ensembles de tweets en découpant nos données par date. On a ainsi un ensemble de tweets par jour. On applique sur chaque ensemble journalier de tweets, l'algorithme de clustering qui calcule la matrice de similarité et construit les clusters.

Pour notre second algorithme, nous supprimons avant tout traitement, les comptes robots que nous sommes susceptibles d'avoir dans nos données. Il est fort probable qu'ils faussent notre clustering.

De plus, chacun de nos événements est décrit par 20 hashtags pertinents. Ces hashtags pertinents sont définis comme étant les plus récurrents dans les tweets contenus dans l'événement. Cependant, certains hashtags sont présents dans la description de beaucoup d'autres événements, ce qui les rend moins pertinents. Nous voulons les détecter et ne pas les prendre en compte lors du clustering afin d'avoir une vraie description de l'événement.

5 Description des éléments utiles à la détection d'événements

Le tweet

On décrit un tweet par :

- son identifiant
- le nom du Tweetos
- les hashtags
- la date et l'heure à laquelle le Tweetos a tweeté
- sa position en latitude/longitude au moment du tweet
- La date

On découpe nos données par date c'est à dire par journée. Une date est caractérisée par une journée, un mois et une année.

La matrice de similarité

La matrice de similarité est une matrice triangulaire supérieure de taille (nombre de tweets) * (nombre de tweets). Décrire construction Citer l'article en parlant rapidement de la compression de Haar

Le clustering

Le clustering se fait à partir de la matrice de similarité. Décrire mieux l'algo

Citer le jar et expliquer rapidement ce qu'on maximise / minimise dans l'algorithme

L'événement

On décrit un événement par :

- une heure de début et une heure de fin qui permettent de calculer une durée
- une position centrale (longitude, latitude) calculée à partir de la moyenne de toutes les positions (moyenne des longitudes, moyenne des latitudes)
- le nombre de personnes différentes ayant posté des tweets
- la liste des 20 hashtags les plus importants permettant de le décrire

6 Algorithme de la détection d'événements

on utilise Knime pour filtrer les robots évidents — ceux qui tweet beaucoup trop par rapport aux autres

on traite ensuite tous les tweets d'une même journée

on enlève ceux qui sont régi par une loi Géométrique (a voir)

on fait du clustering sur les tweets et on renvoie les clusters pertinents sous forme d'événements

7 Les options sur le clustering

L'élasticité

On suppose qu'un hashtag est fréquent pour un ensemble de tweets donné s'il apparaît plus de 20 fois.

Ce nombre est statique, ce qui ne permettrait pas de détecter de petits événements dans une journée où il y a eu peu de tweets. A l'inverse, on détecte beaucoup d'événements sur une journée où il y a eu beaucoup de tweets. Cela peut être intéressant, mais il est parfois plus judicieux de d'adapter au nombre de tweets postés dans une même journée. Ainsi, on définit un hashtag fréquent s'il apparaît dans plus de 20. Pour coder cette option nous avons défini un booléen, que l'on met à Vrai pour activer l'élasticité et à Faux pour ne pas l'activer. Les valeurs 20 et 20 module d'exécution.

La géolocalisation

La loi de Poisson géographique

8 La visualisation du résultat

L'utilité de la visualisation

Décrire

La source de données Décrire

L'organisation de la visualisation Décrire

9 Interprétation de nos résultats

10 Notre essai sur DBScan

11 Conclusion

Les combinaisons intéressantes Quand, comment

Acknowledgments

Nos remerciements vont à Marc Plantevit pour l'enseignement de son cours "Data Mining" à l'Université Claude Bernard Lyon I et son accompagnement durant ce projet.

Références

- [1] X. Dong, M. D., Calabrese, and P. Frossard. Multiscale event detection in social media. *Data Min Knowl Disc*, June 2015. 2