



HAL
open science

Constitution d'un jeu de données social à partir des données twitter

Violaine Guichet, Mathilde Plard

► **To cite this version:**

Violaine Guichet, Mathilde Plard. Constitution d'un jeu de données social à partir des données twitter. 2018. hal-01943615

HAL Id: hal-01943615

<https://hal.science/hal-01943615v1>

Preprint submitted on 4 Dec 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

CONSTITUTION D'UN JEU DE DONNÉES SOCIAL PARTIR DES DONNÉES TWITTER

AUTEURES

Violaine GUICHET (IE), Mathilde PLARD (CNRS, UMR 6590 ESO)

RESUME

Le présent article présente les étapes de construction d'un jeu de données sociales sur un événement de course à pied. Les informations sources sont extraites du réseau social Twitter.

MOTS-CLES

Graphe du web, réseau social, Twitter, scraping, extraction, logiciel de visualisation et d'analyse réseaux



<https://running-datalab.com/>

CONTEXTE ET BESOIN

CONTEXTE INSTITUTIONNEL

Le programme de recherche CHALLENGE vise à documenter l'actuel engouement pour les événements sportifs associés à la course à pied à l'échelle internationale. Le Running DataLab (RDL) s'occupe d'une part de recenser et de qualifier ces événements pour alimenter la réflexion de CHALLENGE. D'autre part, le projet RDL analyse ces événements de courses à pied pour mieux comprendre leurs impacts sociaux et spatiaux.

BESOIN

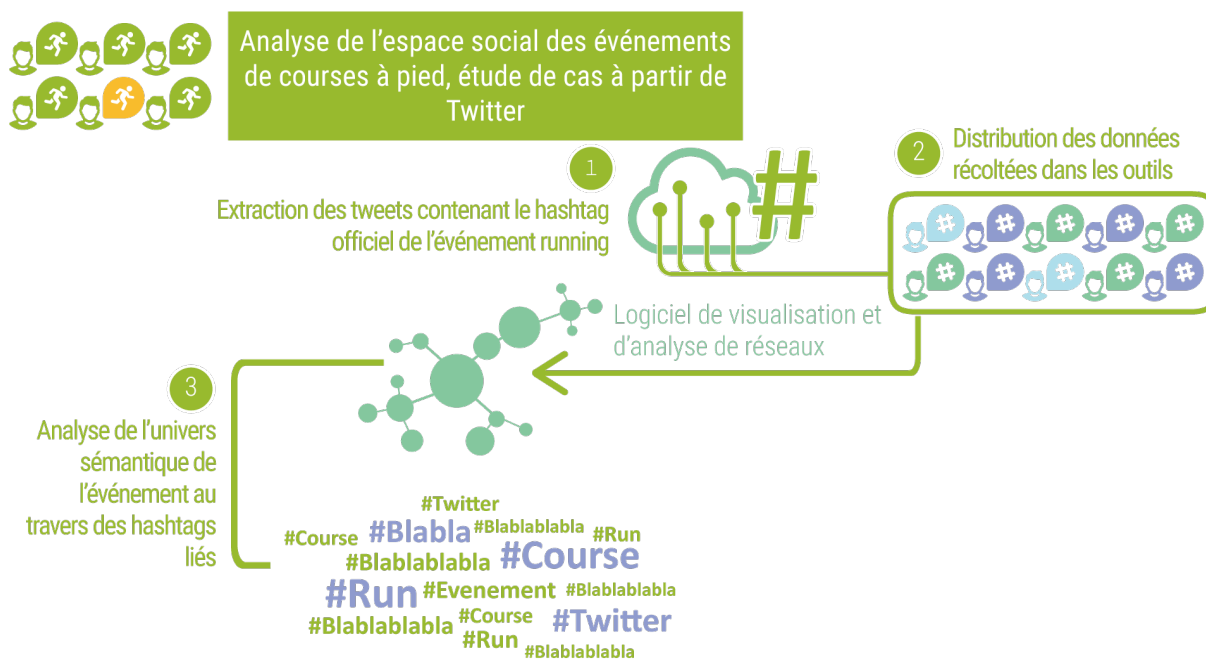
L'objectif est d'étudier l'espace social des événements de courses à pied. Plusieurs interrogations sont portées sur la résonance sociale ou territoriale d'un événement. Ici, le travail présenté est un travail exploratoire. Il s'agit à partir des données Twitter d'étudier les possibilités d'analyse d'un événement de courses à pied en fonction des variables extraites et de leurs potentielles exploitations.

MÉTHODOLOGIE GLOBALE

METHODE

Dans un premier temps, un événement de course à pied est sélectionné. Une fois sélectionnés, tous les tweets contenant l'hashtag officiel de l'événement sont extraits. Les données extraites sont ensuite insérées dans un logiciel de visualisation et d'analyse réseaux afin de réaliser les graphes.

Figure 1. Méthodologie globale



Conception/Réalisation : Violaine Guichet, Ingénieure d'études (2018)
Projet Running DataLab, sous la direction de Mathilde Plard, chercheuse CNRS

CIBLAGE DE L'ÉVÈNEMENT

ÉVÈNEMENT DE COURSES A PIED RETENU

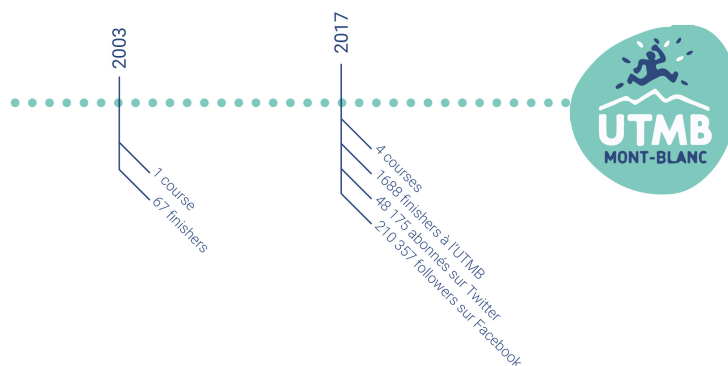
L'événement de courses à pied retenu pour cette manipulation est l'UltraTrail du Mont-Blanc (UTMB).

PRESENTATION DE L'UTMB

L'UTMB a été créé en 2003. Au début, l'événement était constitué d'une course. En quinze ans, l'UTMB s'est développé en proposant plusieurs courses (4 en 2017 et 5 en 2018) qui permettent de cibler un plus large public. Le nombre de finishers (coureurs finissant la course) est passé de 67 à 1 688.

L'événement rassemble plus de 250 000 abonnés sur les différents réseaux sociaux.

Figure 2. Frise chronologique simplifiée de l'UTMB



Conception/Réalisation : Violaine Guichet, Ingénieure d'études (2018)
Projet Running DataLab, sous la direction de Mathilde Plard, chercheuse CNRS

EXTRACTION DES DONNEES

RENSEIGNEMENT SUR L'ÉVÉNEMENT

Lorsqu'un événement est ciblé, l'objectif est d'extraire les tweets dont son hashtag officiel est mentionné. L'hashtag officiel d'un événement est renseigné sur sa page Twitter.

Figure 3. Page Twitter de l'UTMB



Source : Twitter

RECOLTE DES DONNEES

EXTRACTION

Pour récolter les données, la librairie GetOldTweets est utilisée. Cette librairie écrite en python, permet de récolter les tweets selon un certain nombre de critères (nom de l'utilisateur, date, hashtag, etc.).

Concernant notre étude, les tweets mentionnant l'hashtag officiel de l'événement sont extraits.

Figure 4. Script d'extraction des tweets

```
Python Exporter.py -- querysearch « #UTMB »
```

RESULTAT DE L'EXTRACTION

Le script extrait les informations recherchées dans un csv. Concernant l'UTMB, il y a 49 593 tweets extraits qui ont été publiés entre le 14/02/2009 au 3/12/2017 (jour de l'extraction des données).

Pour chaque tweet est renseigné :

- username : l'utilisateur ayant envoyé le tweet
- date : date d'envoi
- retweets : nombre de retweets (= RT) du tweet
- favorites : nombre de favoris pour l'utilisateur
- text : message du tweet
- geo : localisation lors de l'envoi du tweet (— mention : autre(s) utilisateur(s) mentionné(s))
- hashtags : hashtag(s) contenu(s) dans le tweet
- id : id du tweet
- permalink : page sur laquelle se trouve le tweet

Figure 5. Aperçu du fichier csv en extraction

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	username	date	retweets	favorites	text	geo	mentions	hashtags	id	permalink							
2	Koponen	03/12/17 23:21	0	0	Pikane kisamatka #saintely #saintely	9,4E+17				https://twitter.com/KoponenTuomas/status/937446817909936130							
3	valmente	03/12/17 10:48	0	1	Ah enfin quelqu'un qui le #Temple	9,4E+17				https://twitter.com/valmente/status/937257216591192064							
4	shaunste	02/12/17 23:07	0	5	No @wser lottery l @wser @ #UTMB	9,4E+17				https://twitter.com/shaunstewie13/status/937080878176107939							
5	Snapshuf	02/12/17 21:10	0	0	What a day for frayed nerves! #wser #h	9,4E+17				https://twitter.com/Snapshuffler/status/937051280266203136							
6	whoisaar	02/12/17 20:13	0	1	Me too until January & the #UTMB	9,4E+17				https://twitter.com/whoisaaron/status/93703691399205888							
7	SWilliam	02/12/17 01:37	4	11	Congrats to Tamer @UroDroi #UTMB #	9,4E+17				https://twitter.com/SWilliams_MD/status/936756127584473088							
8	Longevill	02/12/17 00:52	0	1	28 janvier 2018 Trail de l'aqu #trail #ac	9,4E+17				https://twitter.com/Longeville8/status/936744967556489217							
9	lainmski	01/12/17 22:25	0	4	The Barkley Marathons: The #UTMB #	9,4E+17				https://twitter.com/lainmski/status/936707754907095040							
10	runner_k	01/12/17 16:41	0	16	Well you got a pretty sweet #UTMB	9,4E+17				https://twitter.com/runner_kc/status/936621184766107649							
11	Kingsleyj	30/11/17 23:05	0	1	Getting a #UTMB # @UTMB! #UTMB #	9,4E+17				https://twitter.com/Kingsleyjones/status/936355565101568001							
12	BimTaks	30/11/17 20:32	0	2	#TBT to #Chamonix! @ajthijss #TBT #	9,4E+17				https://twitter.com/BimTaks/status/936316909804556288							
13	scottrunr	30/11/17 16:55	0	5	Ultra Running has i @andyysy #NOSH	9,4E+17				https://twitter.com/scottrunning/status/936262343402622976							
14	scottspor	30/11/17 16:55	0	2	Ultra Running has i @andyysy #NOSH	9,4E+17				https://twitter.com/scottsports/status/936262342609915909							
15	lclclcha	30/11/17 12:43	1	2	With the #UTMB each August #UTMB #	9,4E+17				https://twitter.com/lclclcha/status/93619888983660032							
16	DavidFtC	30/11/17 09:14	0	0	Punts aconseguits, ara a crev #utmb #	9,4E+17				https://twitter.com/DavidFOliver/status/936146393047207936							
17	DavidFtC	30/11/17 08:59	0	0	Punts aconseguits, ara a crev #utmb #	9,4E+17				https://twitter.com/DavidFOliver/status/936142522996396032							
18	UTMBM	29/11/17 19:44	3	5	L'À@quipe #UTMB s'agrandi #UTMB	9,4E+17				https://twitter.com/UTMBMontBlanc/status/935942436941070336							
19	techArac	29/11/17 11:50	0	2	Via @lequipe - @fr @lequip #UTMB #	9,4E+17				https://twitter.com/techArace/status/935823150910017536							
20	UsamaYa	28/11/17 19:17	7	17	Plasmacytoid urothelial carc #GUpath	9,4E+17				https://twitter.com/UsamaYassi/status/935573467877007361							
21	pep909	28/11/17 11:40	0	0	#GGUT 28.7 #UTMB 31.8 Ein #GGUT #	9,4E+17				https://twitter.com/pep909/status/93545835539251200							
22	pure_TRA	28/11/17 10:28	0	1	@Nariokotomeboy @Nariok #tds #ccc	9,4E+17				https://twitter.com/pure_TRAIL/status/935440207079452672							
23	pep909	28/11/17 09:01	0	2	erst mal sehen ob mich der T #UTMB #	9,4E+17				https://twitter.com/pep909/status/935418439765807105							
24	melinpop	27/11/17 13:43	0	0	#Rafa #UTMB #Chamonix ht #Rafa #U	9,4E+17				https://twitter.com/melinpop/status/935126871146037248							
25	3diasTrai	26/11/17 10:34	0	6	#4dayoff #countdown #3dia #4dayoff	9,3E+17				https://twitter.com/3diasTrailbiza/status/934717016107307008							
26	CarrerasT	26/11/17 08:58	0	3	Teresa Nimes: #çæLo m'Àis #UTMB	9,3E+17				https://twitter.com/CarrerasTrail/status/934692920158052352							
27	akiniaiki	25/11/17 02:40	0	0	#GowesBdgCibodas #gedeps #GowesB	9,3E+17				https://twitter.com/akiniaiki/status/934235335843774464							
28	brunoma	24/11/17 21:43	0	0	Faces of UTMB #UTMB #tra #UTMB #	9,3E+17				https://twitter.com/brunomagnien/status/934160578741981184							
29	Run_and	24/11/17 16:46	0	1	Interview with the @UTMB! #buff #b	9,3E+17				https://twitter.com/Run_and_travel/status/934085824412635136							
30	nrgiao	24/11/17 08:39	0	1	Bom dia Over the top 2016 #UTMB #	9,3E+17				https://twitter.com/nrgiao/status/933963389988473216							
31	mikewari	24/11/17 03:04	0	4	We were thinking i @utmbn #thanksg	9,3E+17				https://twitter.com/mikewardian/status/933879043870547969							
32	MundoEn	23/11/17 20:29	1	3	I had a dream #UTMB #utmb #UTMB #	9,3E+17				https://twitter.com/MundoEnzo/status/93377965098732066							

Source : projet Running DataLab

ANALYSE DES CROISEMENTS POSSIBLES

Au vu de ces variables, différentes analyses apparaissent exploitables dans le cadre de l'étude des événements de courses à pied.

LE JEU D'ACTEURS

Le jeu d'acteurs peut être étudié de différente manière grâce à trois variables que sont username, retweets et favorites, nous permettent d'étudier la résonance sociale de l'événement de courses à pied : qui publie ? De qui s'agit-il ? Quelle portée ont ces publications ?

LA RESONNANCE TERRITORIALE

La résonance territoriale, soit la manière dont un événement fait écho sur d'autres territoires, pourrait être détaillée à l'aide de trois variables : geo, hashtag et text. Dans ces trois éléments peuvent être mentionnés des territoires. Seulement, la variable geo est souvent erronée, elle renvoie à des éléments qui ne font pas référence à un lieu.

Pour cette étude il faudra alors réaliser des tests à partir des hashtags et du text.

UNIVERS SEMANTIQUE

Enfin, l'univers sémantique ou ce à quoi est associé l'événement de courses à pied peut être exploré à travers une analyse plus fine des hashtags.

CHOIX POUR REALISER UN PREMIER ESSAI

Ici le choix est fait d'étudier dans un premier temps le jeu d'acteurs et sa résonance.

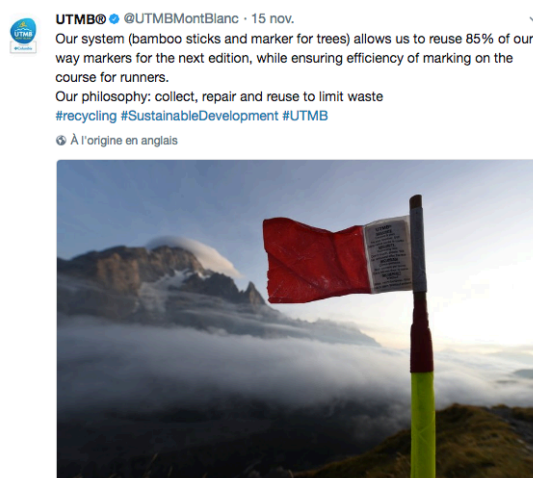
TRAITEMENT DANS GEPHY

RETRANSCRIPTION D'UN TWEET EN GRAPHE

Un graphe du web est structuré en nœud et en lien. En fonction de leur relation, les nœuds sont connectés les uns aux autres.

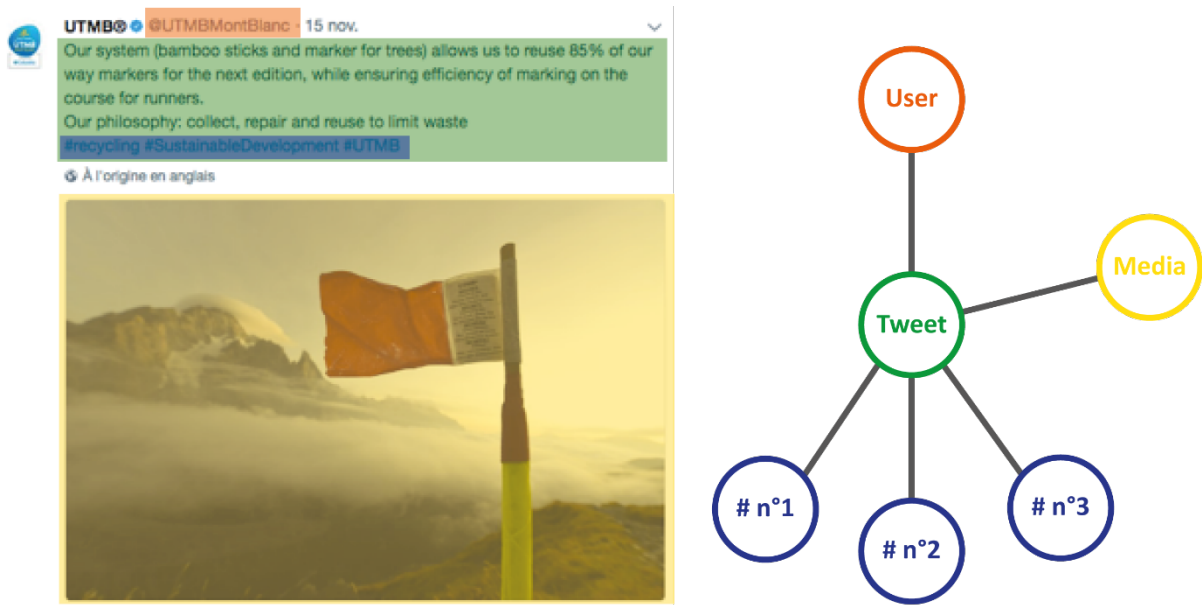
Un tweet est constitué des différents éléments, il y a : l'utilisateur qui envoie le tweet, le tweet (soit le message écrit), la langue du tweet, la date et l'heure d'envoi. Dans certains cas, l'utilisateur mentionne un ou plusieurs utilisateur(s) (@user2, @user3, etc.), hashtags (#hashtags1, #hashtags2, etc.), médias (photographie, vidéo, etc.), liens ou la localisation de l'envoi du tweet.

Pour exemple, le tweet ci-dessous contient : un tweet, l'utilisateur qui a envoyé le tweet, la date et l'heure d'envoi, la langue du tweet, un médiateur (photographie) et 3 hashtags.



Source : twitter

Figure 1. Retranscription d'un tweet en graphe



Conception/Réalisation : Violaine Guichet, Ingénieure d'études (2018)
 Projet Running DataLab, sous la direction de Mathilde Plard, chercheuse CNRS

ADAPTATION DES DONNEES TABLEUR AU LOGICIEL GEPHY

Afin d'insérer les données dans le logiciel de visualisation et d'analyse réseau, il faut structurer les données d'une certaine manière. Pour que le logiciel puisse traiter la donnée, il faut stipuler une source et une cible (= »target«). Chaque source et chaque cible sont un nœud, la cible est reliée à sa source. Plus une source ou cible a de liens avec d'autres nœuds, plus elle est importante.

Concernant les acteurs, les champs mobilisés sont :

- users (= target)
- mentions (= source)
- favorites (= attributs)

Tableau 1. Tableau de base lors de l'extraction des données

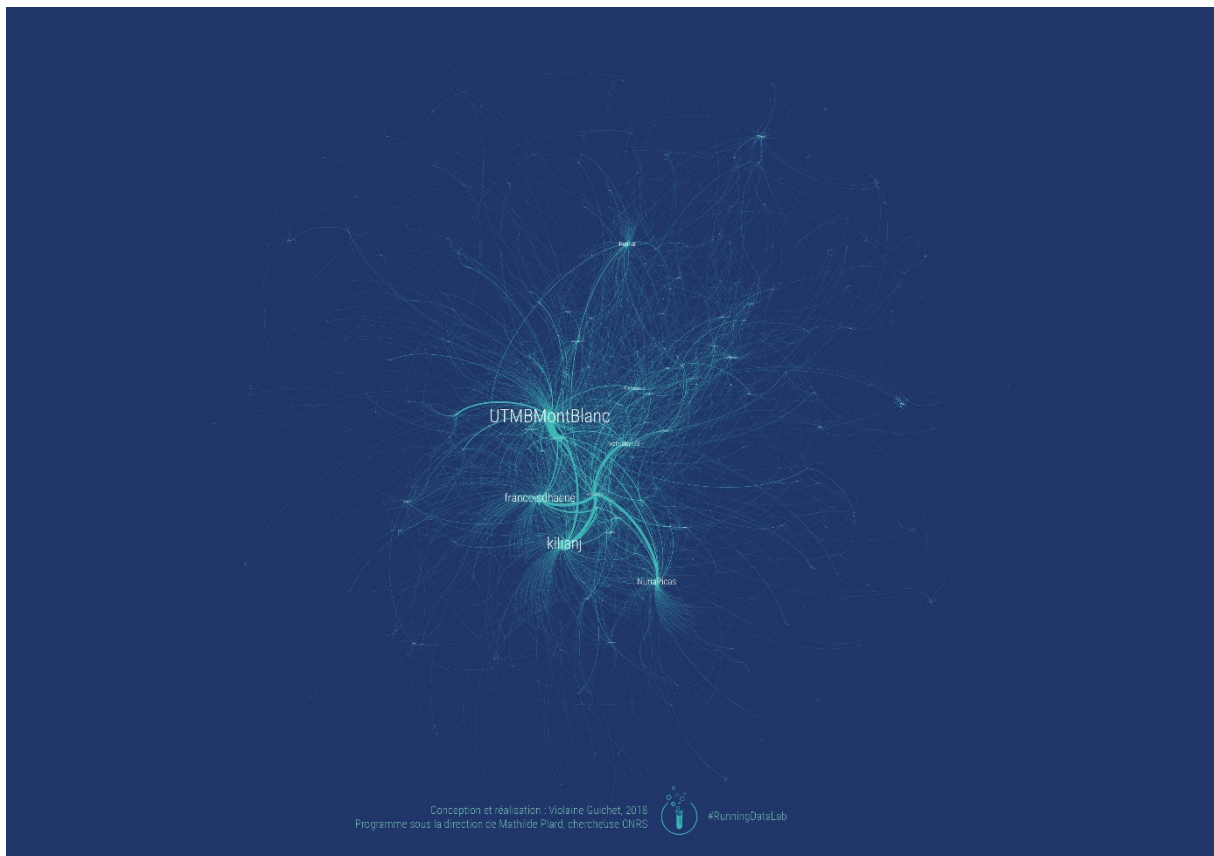
Tableau de base			
username	mentions	retweets	favorite
User1		0	0
User2	@User3, @User5	23	3
User4	@User1	44	132
User3	@User2	3	1
User1	@User4	4	0
User1		11	0
User5	@User3, @User2	6	4

Tableau 2. Restructuration des données pour import dans le logiciel de visualisation et d'analyse réseau

Tableau pour Gephi			
source	target	Retweets_total	favorite
User1		15	0
User2	User1	40	3
User2	User5	40	3
User4	User1	44	132
User3	User2	3	1
User1	User4	15	0
User1		15	0
User5	User3	6	4
User5	User2	6	4

RESULTAT

Le graphe



Première catégorisation des acteurs-Twitter

Un utilisateur peut tweeter souvent par rapport aux autres (ex : User1), être actif sur twitter, mais pour autant ne pas avoir de résonance, puisque non retweeté.

En revanche, d'autres utilisateurs sont moins actifs, mais lorsqu'ils tweetent sont beaucoup repris, il s'agit des influenceurs.

Puis il y a les groupes d'amis ou communautés qui échangent entre eux et sont plus imperméables aux autres utilisateurs.

Piste d'explorations

Ce type de constitution de jeu de données permet d'analyser :

- les acteurs qui tweetent sur l'événement
- le poids de ces acteurs en fonction des RT
- les groupes d'acteurs en fonction des conversations

LIMITES, AMELIORATIONS & PERSPECTIVES

QUALIFIER LES ACTEURS

Pour avoir une analyse plus fine des acteurs, il faudrait télécharger la liste des comptes certifiés Twitter (<https://twitter.com/verified/following>) qui permettrait de différencier les particuliers, des institutionnels ou personnalités.

QUALIFIER LES HASHTAGS

La catégorisation des hashtags dans différents domaines (économie, acteur, territoire, etc.) permettrait une analyse plus fine des événements de courses à pied. De plus, cela permettrait d'évaluer en fonction de chaque événement le poids de chacun des domaines, et ainsi avoir une « photographie » de l'univers sémantique de chacune des courses. D'autre part, cela permettrait de mettre en lumière la résonance territoriale.