

A Tweet Consumer's Look At Twitter Through Linked Data Goggles Via Google Analytics

Thomas Steiner and Arnaud Brousseau

Google Germany GmbH, ABC-Str. 19, 20354 Hamburg, Germany
{tomac, arnaudb}@google.com

Abstract. The Twitter Trends feature allows for a global or local view on “what’s happening in my world right now” from a tweet producers’ point of view. In this paper, we discuss the possibility to complete the functionality provided by Twitter Trends by having a closer look at the other side: the tweet consumers’ – i.e., readers’ – point of view. While Twitter Trends works by analyzing the frequency of terms and their velocity of appearance in tweets being written, our approach is based on the popularity of extracted named entities (in the sense of Linked Data) in tweets being read. Our experimentation architecture takes advantage of the possibility to use a client-side browser extension to harvest and dissect tweets from users’ timelines, search result pages, or profile pages, i.e., tweets supposed to be read. Named entities are extracted via several third-party Natural Language Processing (NLP) Web services in parallel, and are then reported to Google Analytics, which is used to store, analyze, and compute trends by pivoting Analytics data, e.g., users’ geographic locations, with the reported named entities.

1 Introduction

The Twitter Trends feature was introduced by Twitter in September 2008 with the objective to reflect “what’s happening in my world right now” on the micro-blogging platform (in the beginning globally¹, since January 2010 also locally²). To compute trends, Twitter uses an algorithm which analyzes the words and hashtags in tweets. Although the concrete implementation details are kept secret, it is known³ that the algorithm considers three main characteristics of a term/hashtag to determine if it is part of a “trend”:

- The **quantity**, i.e, the absolute number of appearances of a given term/hashtag among all users’ tweets.
- The **velocity**, i.e, the frequency of appearance of that term/hashtag. The higher the frequency, the more popular the term/hashtag.

¹ <http://blog.twitter.com/2008/09/twitter-trends-tip.html>

² <http://blog.twitter.com/2010/01/now-trending-local-trends.html>

³ <http://blog.twitter.com/2010/12/to-trend-or-not-to-trend.html>

- The **newness**, e.g., the (at the time of writing brand new) hashtag `#ipad2` will be given priority over an old and generic tag like `#awesome` to be featured as “trend”. The freshness of trends is also assured by the fact that the frequency is computed over a recent set of tweets, although we have no further details of the exact computation period – Twitter keeps this information secret.

Twitter Trends analysis is a fascinating field to many, with both an academic or an industry (or even social media hobbyist) background (e.g., [4] by Kannan et al. has a good overview of Twitter Trends analysis and visualization efforts). What all these approaches have in common is that they focus on the tweet producers’ point of view. Either they are based on the output from Twitter’s Trends API⁴ (i.e., use Twitter’s pre-calculated trends data), or use the Twitter Streaming API⁵ (i.e., calculate trends data on their own). What to the best of our knowledge has been missing so far is the tweet consumers’ point of view. If many Twitter users tweet about a topic (e.g., a hashtag like `#typeofsex` that seems to persist in Twitter Trends over a period of several weeks), this does not necessarily mean that tweet consumers also read these tweets. Currently the closest approximation to estimating whether a tweet has been read is to check whether it has ever appeared on someone’s timeline, search result page, or profile page (the tweet producer having followers is not a sufficient condition).

In this paper, we suggest a solution for enabling tweet tracking combined with named entity extraction (NEE) for the Twitter.com page. This solution allows to combine data from classic Web tracking (such as user location) with our interpretation of Twitter trends, which is based on named entities. Others⁶ have created visualizations of user location combined with trends, however, user location so far has been an approximation of what Twitter users expose on their Twitter profile pages, which is not necessarily the same as their current physical location, but oftentimes their hometown.

The remainder of the paper is structured as follows: Section 2 introduces the two background technologies Google Chrome extensions and Google Analytics required for our experimentations. Section 3 introduces our Twitter Swarm NLP Google Chrome extension. Section 5 gives an overview on related work. Section 4 contains an evaluation of our experimentation results so far. Section 6 finalizes the paper with an outlook on future work and a conclusion.

2 Required Technologies

We begin with an overview of the used technologies for this experiment, namely we introduce Google Chrome extensions for the Google Chrome browser, and give a brief summary of the Web analysis solution Google Analytics.

⁴ <http://dev.twitter.com/doc/get/trends/>

⁵ http://dev.twitter.com/pages/streaming_api

⁶ E.g., <http://trendsmap.com/>

2.1 Google Chrome Extensions

Google Chrome extensions⁷ are small software programs that can be installed to enrich the browsing experience with the Google Chrome browser. They are written using a combination of standard Web technologies, such as HTML, JavaScript, and CSS. Chrome extensions bundle all their resources into a single file that gets usually (but not necessarily) distributed through the Chrome Web Store⁸. There are several types of extensions, for this paper we focus on extensions based on so-called content scripts. Content scripts are JavaScript programs that run in the context of Web pages via dynamic code injection, similar to the Firefox Greasemonkey extension⁹. By using the standard Document Object Model (DOM), they can read or modify details of the Web pages a user visits. Examples of such modifications are, e.g., changing the behavior of hyperlinks by removing potential `@target="_blank"` attributes in order to open all links on the same browser tab, or increasing the font size for better legibility. Content scripts can run on any website, or be limited to just certain websites. This can be controlled by the extension developer with a so-called manifest file in JSON format. When a user installs an extension, the access rights of the extension are displayed and must be acknowledged by the user.

2.2 Google Analytics

Google Analytics¹⁰ is Google's free Web analysis solution allowing for detailed statistics about the visitors of a website. The software is implemented by including the Google Analytics Tracking Code (GATC), an invisible snippet of JavaScript code that a webmaster needs to add onto every of the to-be-tracked pages of a website. This code collects visitor data by requesting a specific 1 x 1 pixel image on Google's servers, during which the page and user data is encoded in the query part of the image's URL. In addition to that, the snippet sets a first party cookie on visitors' computers in order to store anonymous information such as the timestamp of the current visit, whether the visitor is a new or returning visitor, and the referrer of the website that the visitor came from. Part of the shared visitor information is the IP address, which allows for IP-based geolocation of visitors, i.e., given an IP address, a mapping to a geographical location is possible using lookup tables. Depending on the region, the accuracy of IP-based geolocation can be down to urban district level.

3 Twitter Swarm NLP Extension

With our Twitter Swarm NLP extension¹¹, we inject JavaScript code via a content script into the Twitter.com homepage. By installing this extension, users

⁷ Google Chrome Extensions: <http://code.google.com/chrome/extensions/index.html>. Text partly adapted from the description to be found there.

⁸ <https://chrome.google.com/webstore/>

⁹ Firefox Greasemonkey extension: <http://www.greasespot.net/>

¹⁰ <http://www.google.com/analytics/>

¹¹ <https://chrome.google.com/webstore/detail/dpbphenfakflfmdlanimlemacankjol>

explicitly opt-in to their data as a Twitter.com visitor being tracked by Google Analytics. The extension first checks if the user is logged in to Twitter.com, and if so, retrieves the tweets of the logged-in user’s timeline (<http://twitter.com/#>), or search result page (e.g., <http://twitter.com/#!/search/%23semweb>), or profile page (e.g., http://twitter.com/#!/timberners_lee) on a one-by-one basis, and performs Named Entity Extraction (NEE) via Natural Language Processing (NLP) using a remote NLP Web service (see Section 3.1) on each of the tweets. The extracted entities are then displayed on the righthand-column of the Twitter.com homepage (see Figure 1), and sent to Google Analytics for further processing (see Figure 2).

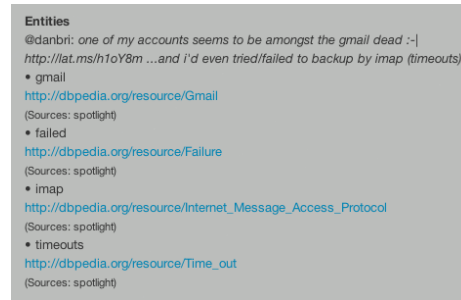


Fig. 1. Screenshot of the extracted entites of a particular tweet as displayed by the Twitter Swarm NLP Extension.

3.1 Twitter Swarm NLP Web Service

We have created a wrapper NLP Web service that merges results from existing third-party NLP Web services, namely from OpenCalais¹², Zemanta¹³, AlchemyAPI¹⁴, and DBpedia Spotlight¹⁵. All NLP Web services return named entities together with their types and potential subtypes, names, relevance, and URIs that link into the Linked Open Data (LOD) cloud¹⁶. The problem is that each service has implemented its own type system (e.g., the type “city” is represented by <http://dbpedia.org/ontology/City> in DBpedia vs. <http://rdf.freebase.com/ns/location.citytown> in Freebase), and providing mappings for all of them would be a rather time-consuming task. However, as all services offer links into the LOD cloud, the type information can be pulled from there if need be, in a true Linked Data manner. In order to clarify what our wrapper service provides, the merged results for the sample query “Google Translate” are depicted below. For the sake of clarity, we just show one entity with two URIs, while the original result at

¹² <http://www.opencalais.com/>

¹³ <http://www.opencalais.com/>

¹⁴ <http://www.alchemyapi.com/>

¹⁵ <http://dbpedia.org/spotlight>

¹⁶ <http://lod-cloud.net/>

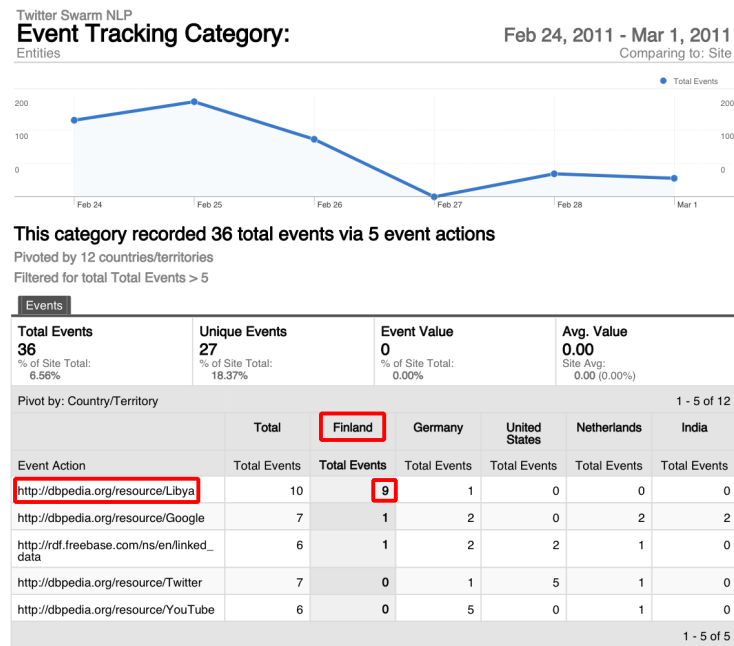


Fig. 2. Named entities pivoted by countries/territories. The named entity represented by the URL <http://dbpedia.org/resource/Lybia> appeared in nine tweets read by users located in Finland (red borders in the screenshot).

the time of writing contained seven entities, among which six were directly relevant, and one was closely related (the complete result can be seen at the URL <http://tomayac.no.de/entity-extraction/combined/Google%20Translate>).

```
[
  {
    "name": "Google Translate",
    "relevance": 0.7128319999999999,
    "uris": [
      {
        "uri": "http://dbpedia.org/resource/Google_Translate",
        "source": "alchemyapi"
      },
      {
        "uri": "http://rdf.freebase.com/ns/en/google_translate",
        "source": "zemanta"
      }
    ],
    "source": "alchemyapi,zemanta"
  }
]
```

3.2 Technical Implementation Of the Twitter Swarm NLP Extension

Twitter.com is an Ajax-dependent website, which makes use of so-called hashbang URIs¹⁷. When the page Twitter.com gets loaded, the static parts of the page can be cached, while the dynamic parts – controlled by the content of the hashbang URI – get pulled in via JavaScript. Using a so-called escaped fragment in the query part of the URL (`?_escaped_fragment_`), the state of the Twitter.com page is still accessible to Web crawlers, as specified in [7]. Currently the extension is implemented in a way to run once upon page load, however, not upon Ajax refresh events on the page.

We overload Google Analytics event tracking¹⁸ for the purpose of tracking named entities. We interpret a Google Analytics event as the occurrence of a named entity in a never seen before tweet. In its original sense, event tracking is meant to track on-page events, like, e.g., playing or stopping an embedded video on a page. The anatomy of an event is as follows: an event is defined by its `category` (req.), `action` (req.), `label` (opt.), and `value` (opt.). In our concrete case, the `category` is always “Entities”, the `action` is any of the particular named entity’s URIs (as they all represent the same thing, we simply choose the first URI of the list). The `label` is the particular entity’s name and the sources of that named entity in parenthesis (as outlined in Section 3.1), and finally for the `value` we pass the timestamp of the moment where the entity extraction happened (which is closest to what we assume is the time where the tweet was read). A concrete call to log an occurrence of a named entity as injected onto Twitter.com by the content script might look like this (the `_gaq` array is the Google Analytics queue for asynchronous tracking):

```
_gaq.push(['twitter_swarm_nlp._trackEvent',
  'Entities',
  'http://dbpedia.org/resource/JSON',
  'json (zemanta,opencalais)',
  1299776578]);
```

3.3 Visual Presentation Of Named Entities

Currently the extension displays all extracted named entities grouped by originating tweets and roughly in the tweet order at the righthand-column. In order not to flood the Twitter Swarm NLP Web service (and in turn the underlying third-party NLP Web services), the requests to perform NEE on tweets get sent with a 1s pause in between two tweets. Dependent on the response time for each particular request, the extracted named entities get added to the righthand-column. Each named entity gets displayed with its label, one representing URI (out of potentially several URIs), and the sources (as outlined in Section 3.1). See Figure reffig:danbri for a screenshot.

¹⁷ See <http://www.jenitennison.com/blog/node/154> for a detailed analysis of Twitter’s use of hashbang URIs

¹⁸ <http://code.google.com/apis/analytics/docs/tracking/eventTrackerGuide.html>

3.4 Avoiding Duplicate Tracking

Each tweet gets sent one-by-one to the Twitter Swarm NLP Web service, as described in Section 3.1. While the extracted named entities of all tweets always get displayed to the user (as shown in Figure 1), only named entities from tweets never seen before get sent to Google Analytics. This is to ensure that duplicate data from the same user does not falsify the overall statistics. It is to be noted that we do allow the same tweet to be tracked more than once (because one tweet can be read by more than one user), however, we want to exclude the case that the same tweet gets tracked more than once from the same user (we assume each tweet gets read at most one time). Twitter has introduced a new system to generate unique tweet IDs, called snowflake¹⁹, which guarantees k-sorted IDs (see [1]) with at least a 1s bound. Hence this allows us to store the ID of the latest tweet in the extension’s local storage, and check whether the current tweet’s ID is greater (i.e., the tweet is younger than the previous latest tweet and has thus not been seen before), and only if so, send its named entities to Analytics. While this strict limitation to only consider new tweets leaves out many old tweets from being tracked in Analytics, it is the only alternative to storing the ID of each and all tweets in local storage to assure duplication-free tracking.

4 Evaluation And Discussion

In this section we first provide an evaluation of our experimentations and discuss the approach in the second part.

5 Related Work

As stated in Section 1, Twitter Trends is both an active field of academic and industry research, but also a playground for social media hobbyists. In the following we present four exemplary cases from both categories.

5.1 Linked Open Social Signals (TWARQL)

In the Linked Open Social Signals project [5] Mendes et al. investigate the representation of microposts as Linked Open Data (the authors call the opinions, observations, and suggestions contained in microposts “social signals”, hence the project name Linked Open Social Signals). Mendes et al. address the problem of information overload caused by the sheer amount of microposts. While the micropost community has come up with hashtags in order to categorize microposts, these hashtags are ambiguous and have to be explicitly added to the micropost by the author. Given the typical length limitations of microposts (often 140 characters), sometimes hashtags are left out in favor of more text. The project’s

¹⁹ <https://github.com/twitter/snowflake#readme>

main goal is thus to enable collective analysis of social signals for sensemaking by making use of Linked Open Data principles in combination with realtime push models. The approach consists of the following steps:²⁰

- Extract content (entity mentions, hashtags and URLs) from microposts
- Encode content in RDF format using common vocabularies
- Enable SPARQL querying of microposts
- Enable subscription to micropost streams that match a given query
- Enable scalable real-time delivery of streaming data via SparqlPuSH

The authors maintain a client-side JavaScript application²¹ allowing for users to search for tweets matching a customizable SPARQL query or to subscribe to tweet streams in a realtime “push” way, filtered according to the user’s request (so-called concept feeds).

5.2 Twitris 2.0

With Twitris 2.0 [3], Jadhav et al. present an application to find out what is being said about an event and when, to detect how topics of discussion are changing over a period of time, and finally to check whether there are regional differences in the opinions on a given topic. The approach consists of picking trending hashtags from Twitter, which are then expanded by data from Google Insights for Search²². Using this set of search terms, find related hashtags by performing a Twitter search in order to detect topic drifts. In order to locate the origin of a tweet, Twitris uses the approximation of the location given in the user’s Twitter profile. This approach works well if there is geocodable data (like “Austin, Texas”), however, fails if there is generic data (like “somewhere under the rainbow”). Tweets about a given set of topics can then be examined on a map view, enriched by relevant photo and video content. For each tweet location hotspot a tag cloud with related tags is displayed, and the data can be sliced by time.

5.3 Semantic-MicrOBlogging (SMOB)

SMOB [6] by Passant et al. is a Semantic MicroBlogging framework that enables a distributed, open and semantic microblogging experience based on Semantic Web and Linked Data technologies. Microposts get annotated with common vocabularies like FOAF²³ and SIOC²⁴. SMOB relies on distributed autonomous hubs that communicate with each other to exchange microblog posts and subscriptions, which can also be cross-posted to Twitter. The authors suggest people to use meaningful hashtags such as `#dbp:Eiffel.Tower`, or `#geo:Paris_France`,

²⁰ Steps adapted from <http://wiki.knoesis.org/index.php/Twarql>

²¹ <http://knoesis1.wright.edu/twarql/query.html>

²² <http://www.google.com/insights/search/>

²³ <http://www.foaf-project.org/>

²⁴ <http://sioc-project.org/>

in the style of RDF prefixes for DBpedia [2] and GeoNames. SMOB allows for manually annotating hashtags with URIs and RDF, with the objective of making microposts accessible for, e.g., lookup services such as Sindice [8].

5.4 Twopular

Twopular²⁵ is an experiment by Martin Dudek with the objective of analyzing current Twitter trends. Therefore, for a given set of at the particular moment current trends the most recent tweets are obtained, and in an interval of five minutes run through the OpenCalais Web service in order find tags. By having tags for Twitter trends, another way of searching for trends – and more importantly the possibility to interrelate trends based on tag similarity – gets enabled. The author sees the feature more like a “linguistic experiment”²⁶, however states that the first results seem to be promising. In our tests of the service we could affirm the author's self-assessment, e.g., the Twitter trend (at the time of writing, March 7, 2011, 2PM CET) “Prince Andrew” was mapped to the OpenCalais tag “Prince Andrew, Duke of York”, and related tags were, among others, “Sex offender registration”, and “X-Offender”, where the story behind the trend was that people were tweeting about Prince Andrew of England, whose close friend was found to be a pedophile.

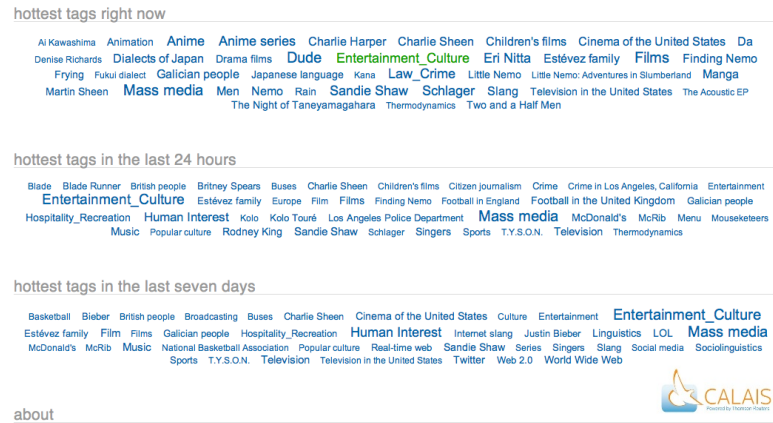


Fig. 3. Screenshot of the Twopular trends page.

6 Future Work And Conclusion

In this paper, we have suggested a solution for obtaining a consumers' point of view on Twitter Trends. Rather than measuring the “trendiness” of tweets

²⁵ <http://twopular.com/>

²⁶ <http://twopular.com/blog/?p=308>

being produced, we measure the “trendiness” of tweets being consumed. This is not only interesting *per se*, but together with the “by-product” of classical Web analysis allows for even richer insights into, e.g., the location of users interested in a certain trend, where a trend in our understanding is a named entity that is popular on a freely configurable period of time. Contributions: time filters (via Analytics), geographical pivoting (via Analytics), global ranking (sliceable per time unit)

References

1. T. Altman and Y. Igarashi. Roughly sorting: sequential and parallel approach. *J. Inf. Process.*, 12:154–158, January 1989.
2. S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. G. Ives. DBpedia: A Nucleus for a Web of Open Data. In *6th International Semantic Web Conference (ISWC’07)*, pages 722–735, Busan, Korea, 2007.
3. A. Jadhav, H. Purohit, P. Kapanipathia, P. Ananthram, A. Ranabahu, V. Nguyen, P. N. Mendes, A. G. Smith, M. Cooney, and A. Sheth. Twitris 2.0 : Semantically empowered system for understanding perceptions from social data. Semantic Web Challenge at the 9th International Semantic Web Conference (ISWC2010), Shanghai, China, 2010. http://www.cs.vu.nl/~pmika/swc/submissions/swc2010_submission_8.pdf.
4. A. Kannan, J. Patzer, and B. Avital. Trendtracker: A system for visualizing trending topics on twitter. University of California, Berkeley, Visualization Course CS294-10, Visualization Wiki, 2010. <http://vis.berkeley.edu/courses/cs294-10-sp10/wiki/images/archive/d/d4/20100511073322!FinalPaper.pdf>.
5. P. N. Mendes, A. Passant, P. Kapanipathi, and A. P. Sheth. Linked open social signals. *Web Intelligence and Intelligent Agent Technology, IEEE/WIC/ACM International Conference on*, 1:224–231, 2010.
6. A. Passant, T. Hastrup, U. Bojars, and J. Breslin. Microblogging: A semantic and distributed approach. In *Proceedings of the 4th Workshop on Scripting for the Semantic Web, Tenerife, Spain, June 02, 2008, CEUR Workshop Proceedings*, 2008. <http://CEUR-WS.org/Vol-368/paper11.pdf>.
7. K. Probst, B. Johnson, A. Mukherjee, E. van der Poel, L. Xiao, and J. Mueller. Making ajax applications crawlable. Google, 2009. <http://code.google.com/web/ajaxcrawling/docs/getting-started.html>.
8. G. Tummarello, R. Delbru, and E. Oren. Sindice.com: weaving the open linked data. In *6th International Semantic Web Conference (ISWC’07)*, pages 552–565, Busan, Korea, 2007.