

# 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.** Twitter Trends<sup>1</sup> 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, i.e., tweets supposed to be read. Named entities are extracted using 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 location, with the recorded named entities.

## 1 Introduction

### 1.1 Twitter Trends

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

### 1.2 Google Chrome Extensions

Google Chrome extensions<sup>2</sup> 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 files into a single file that gets usually (but not necessarily) distributed through the Chrome Web Store. There are several types of extensions, for this paper we focus on extensions based on so-called content scripts. Content scripts are JavaScript programs

---

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

<sup>2</sup> Google Chrome Extensions: <http://code.google.com/chrome/extensions/index.html>. Text adapted from the description to be found there.

that run in the context of Web pages, similar to the Firefox Greasemonkey extension<sup>3</sup>. 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 hyperlinks to remove potential @target="\_blank" attributes, or increasing the font size.

### 1.3 Google Analytics

## 2 Twitter Swarm NLP Extension

With our Twitter Swarm NLP extension<sup>4</sup>, we inject JavaScript code via a content script into the Twitter.com homepage. The extension first checks if the user is logged in, and if so, retrieves the tweets of the logged-in user's timeline one-by-one, and performs NLP analysis via a remote NLP Web service on each of the tweets. The extracted entities are then displayed on the righthand-pane of the Twitter.com homepage, and sent to Google Analytics for further processing.

### 2.1 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 OpenCalais<sup>5</sup>, Zemanta<sup>6</sup>, AlchemyAPI<sup>7</sup>, and DBpedia Spotlight<sup>8</sup>.

### 2.2 Dealing With Extracted Entites On the Client Side

### 2.3 Dealing With Extracted Entites On the Google Analytics Side

## 3 Related Work

### 3.1 Linked Open Social Signals (TWARQL)

Previous work of Mendes et al. [1] have shown a possible implementation of real-time information both pushed and pulled from Twitter. TWARQL<sup>9</sup> is based on a distributed architectures which features:

- a client-side application which typically a Javascript-enabled web browser
- a "Social Sensor Server" to receive tweets and filter them according to the user's request. It is worth noting here that TWARQL filtering is based on web-semantic technologies: SPARQL, hash-tag resolution through glossaries and LOD cloud are used to extract the highest amount of information possible from the Twitter Streaming API.

<sup>3</sup> Firefox Greasemonkey extension: <http://www.greasespot.net/>

<sup>4</sup> <https://chrome.google.com/webstore/detail/dpbphenfakflfmdlanimlemacankjol>

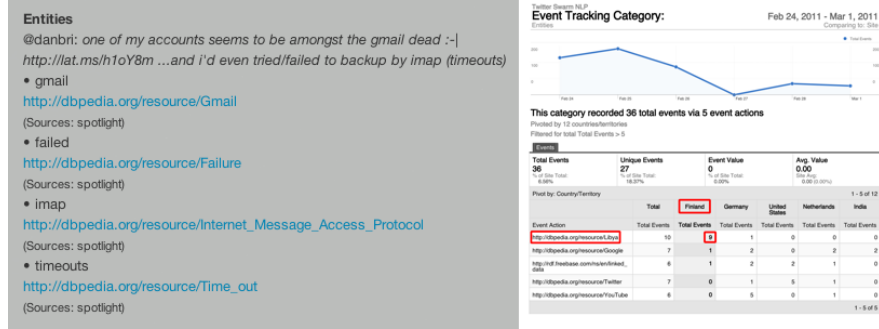
<sup>5</sup> <http://www.opencalais.com/>

<sup>6</sup> <http://www.opencalais.com/>

<sup>7</sup> <http://www.alchemyapi.com/>

<sup>8</sup> <http://dbpedia.org/spotlight>

<sup>9</sup> <http://wiki.knoesis.org/index.php/Twarql>



**Fig. 1.** Left: Screenshot of the extracted entites of a particular tweet as displayed by the Twitter Swarm NLP Extension. Right: Test.

- a number of distributed PuSH hubs which update clients as information flows (pushed-information model)
- another server – "Semantic Publisher" – registers user's interest and updates the hubs. The updated information is eventually displayed on the user's screen.

### 3.2 Twopular

Twopular<sup>10</sup> is a work of Martin Dudek. It aims at analysing current Twitter trends. Since March 2008, Twopular takes advantage of OpenCalais services to extract entities from tweets retrieved from the Twitter Streaming API. Semantic entities are then used to reflect Twitter's current "trends".

## 4 Conclusion

Contributions: time filters (via Analytics), geographical pivoting (via Analytics)

As seen in the Related Work section, semantic analysis of a (real-time) Twitter stream is not new and has been successfully exploited to analyse tweets produced by the Twitter community. What we propose here is an insight into tweets consumers' interests to provide a more accurate view of Twitter trends.

## References

1. 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.

<sup>10</sup> Twopular website: <http://twopular.com/>