

Извлечение описаний событий предопределённых типов из потока сообщений пользователей микроблогов

Ефим Пышнограев

Филиал МГУ в городе Севастополе

22 апреля 2014

Цель работы:

- исследовать существующие подходы по извлечению описаний событий из сообщений пользователей, выделить возникающие проблемы и рассмотреть возможные методы их решения,
- исследовать возможность применения тематических моделей для решения задачи выявления событий,
- разработать метод для извлечения описаний событий из сообщений пользователей сети Твиттер на основе иерархического процесса Дирихле,
- протестировать работу алгоритма на реальных данных.

Постановка задачи

Данные:

- множество документов (сообщений) $\Omega = \{D_i \mid i \in \overline{1, n}\}$,
- документ — множество слов и временная метка t_i
 $D_i = \{w_j \mid j \in \overline{1, l_i}\}$.

События:

- резкое увеличение частоты ключевых слов, затем спад до нормального уровня,
- должно быть вызвано реальным событием,
- не носит периодический характер.

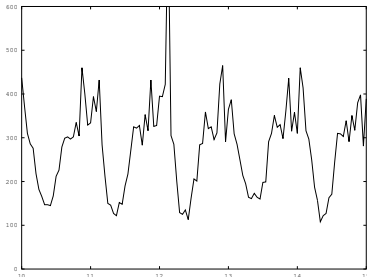
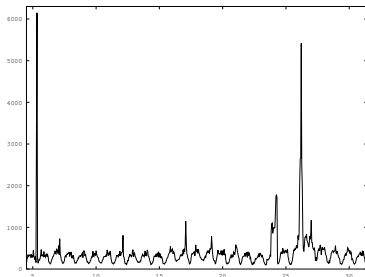
Особенности сети Твиттер:

- короткие сообщения (до 140 символов),
- наличие шума и ошибок,
- большая плотность сообщений,
- “взрывной” характер событий.

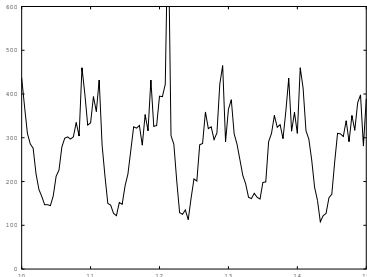
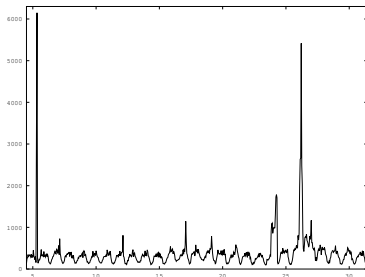
Составные части:

- ➊ нахождение максимума в частотной функции, который будет соответствовать событию,
- ➋ извлечение ключевых слов, характерных этому максимуму,
- ➌ применение модели HDP с частичным обучением для того чтобы выделить все сообщения этой темы,
- ➍ проверка насколько полученный результат соответствует новому событию.

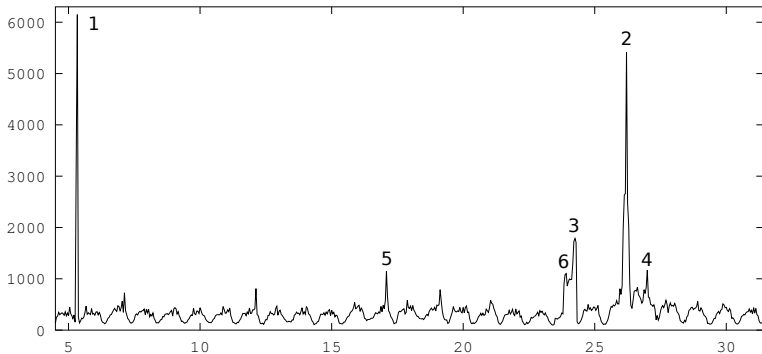
Данные: 240 тысяч сообщений с 4 июня по 1 июля 2013 года, содержащие хэштег #texas.



Данные: 240 тысяч сообщений с 4 июня по 1 июля 2013 года, содержащие хэштег #texas.



Шаг 1: выделение экстремумов.



#1		year, read, prison, dragon, amazon,
#2		sb5, standwithwendy, talk, made, histori
#3		houston, california, england, artist, rt
#4		execut, 500th, mccarthy, kimberli, 1976
#5		florida, z1, america, usa, bahrain
#6		houston, california, rt, watch, england

