

Титульный лист

Содержание

1	Введение	3
2	Описание задачи	3
3	Анализ существующих подходов	4
3.1	NED	4
3.2	Выявление событий с помощью LDA	6
3.3	Information nuggets	7
4	Теоретическая часть предложенного решения	8
4.1	Нахождение экстремумов	9
4.2	Извлечение ключевых слов	10
4.3	Модифицированная модель HDP	11
4.3.1	Иерархический процесс Дирихле	11
4.3.2	Chinese Restaurant Franchise	12
4.3.3	Модификация частичной обучаемости	13
4.4	Анализ частотной характеристики сообщений темы	14
5	Реализация и тестирование решения	15
5.1	Описание данных	15
5.2	Извлечение событий	15
6	Заключение	18

1 Введение

В современном мире информация быстро устаревает, поэтому способы вовремя находить нужные данные — постоянный объект для исследований. Одним из направлений в этой области является извлечение данных из платформ микроблогов.

Платформы микроблогов стали очень популярным способом размещения данных в Сети. В них можно найти сообщения пользователей практически на любую тему, начиная стихийными бедствиями и заканчивая рейтингами музыкальных исполнителей. Правильная обработка доступной информации — нетривиальная задача, которая имеет множество областей применения. Последние несколько лет эта тема активно исследуется во многих университетах мира.

Отслеживание сообщений о стихийных бедствиях в реальном времени поможет вовремя организовать спасательные операции и сохранить жизни людей[1]. Руководствуясь сообщениями пользователей микроблогов можно судить о популярности товаров и вовремя принимать экономически целесообразные решения. Можно делать предположения о рейтингах политических деятелей и эффективности рекламы на основании информации в микроблогах. Помимо перечисленных способов применения доступной информации в микроблоггинговых платформах можно привести множество других.

В данной работе в дальнейшем будет рассматриваться сервис микроблогов Твиттер¹. В нем помимо текстовой информации можно публиковать фото, видео и геотеги, что так же может быть использовано при анализе, но в этой работе не рассматривается. В доступном наборе данных можно проводить анализ разных сущностей, эта работа посвящена выявлению событий среди потоков информации. Поскольку трактовка событий в сообщениях микроблогов может быть субъективной, выделим несколько свойств, которыми характеризуется событие.

Событие в первую очередь должно быть чем-то аномальным на фоне остальных данных. Оно определяется резким изменением частотных характеристик некоторых слов в сообщениях. События в микроблогах носят взрывной характер, в течении нескольких часов частоты релевантных слов возрастают в десятки раз и так же быстро опускаются до нормального уровня. Примером события могут быть: стихийное бедствие, выход законопроекта на резонансную тему, получение фильмом награды на кинофестивале.

Существуют разные подходы для решения описанной задачи. В следующих двух разделах будут даны формальные определения и рассмотрены некоторые подходы для решения задачи извлечения событий.

2 Описание задачи

Цель дипломной работы состоит из нескольких частей:

- исследовать существующие подходы по извлечению событий из сообщений пользователей, выделить возникающие проблемы и рассмотреть возможные методы их решения,
- исследовать возможность применения тематических моделей для решения задачи выявления событий,

¹<http://www.twitter.com>

- разработать метод для извлечения событий из сообщений пользователей сети Твиттер на основе иерархического процесса Дирихле,
- протестировать работу алгоритма на реальных данных.

Объект изучения этой работы — алгоритм, который по входным данным строит множество событий. В качестве данных для задач подобного рода служит корпус документов (сообщений):

$$\Omega = \{D_i \mid i \in \overline{1, n}\}. \quad (1)$$

Будем считать, что каждый документ D_i имеет временную метку t_i . В свою очередь документ D_i определяется как упорядоченный набор слов:

$$D_i = \{w_j \mid j \in \overline{1, l_i}\}, \quad (2)$$

при этом слова в документах принадлежат некоторому словарю V .

Событие — некоторая сущность, которая характеризуется временем возникновения и ключевыми словами. Оно вызывает резкий подъем частотных характеристик некоторых слов. Событием может быть футбольный матч и музыкальный концерт. В социальной сети Твиттер есть популярные темы, которые всегда генерируют много сообщений. Например это сообщения с ключевыми словами *iphone* и *ipad*. Но такие сообщения нельзя считать событиями. Также событиями нельзя считать еженедельные пятничные сообщения о конце рабочей недели[2].

Особенности социальной сети Твиттер состоят в следующем:

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

3 Анализ существующих подходов

Прежде чем составлять решение задачи были изучены существующие подходы. Все они комбинируют техники машинного обучения, обработки текстов, вероятностных графических моделей и других разделов науки. Разные методы лучше работают для одних данных и хуже для других, на их качество также влияет природа данных. Выделим несколько подходов, для того чтобы изучить общую схему подобных алгоритмов и чтобы обозначить идеи, которые могут быть использованы при составлении метода извлечения событий из сообщений сети Твиттер.

3.1 NED

Исторически первым подходом к извлечению событий принято считать NED (New Event Detection)[3]. NED предназначен для того, чтобы находить первый документ на тему, которая не встречалась раньше. Следующие документы на эту тему уже не будут новыми и не будут помечены алгоритмом. Для того, чтобы отвечать на вопрос, является ли документ новым, необходимо указать способ как определять степень сходства двух документов.

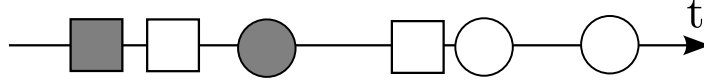


Рис. 1: Документы, соответствующие двум разным темам. Новые документы помечены серым цветом.

Для этого в алгоритме NED используется техника Incremental TF-IDF (Term Frequency – Inverse Document Frequency). TF-IDF — базовый метод выяснять насколько отдельно взятые слова характеризуют весь документ, другими словами, насколько большой вес имеет слово w в документе d . Пусть $f(d, w)$ — количество слов w в документе d . Определим значение $df_t(w)$ как количество документов, поступивших не позднее времени t , в которых встречается слово w . Используя введенные величины, можно записать значение веса определенного слова w в документе d . В момент времени t имеем:

$$\text{weight}_t(d, w) = \frac{1}{Z_t(d)} f(d, w) \cdot \log \frac{N_t}{df_t(w)}, \quad (3)$$

где N_t — общее количество документов, поступивших не позднее времени t , $Z_t(d)$ — нормализационное значение:

$$Z_t(d) = \sqrt{\sum_w \left[f(d, w) \cdot \log \frac{N_t}{df_t(w)} \right]^2}. \quad (4)$$

Теперь можно записать значение похожести двух документов, q и d :

$$\text{sim}_t(d, q) = \sum_w \text{weight}(w, d) \cdot \text{weight}(w, q). \quad (5)$$

Указанные формулы записаны в косинусной метрике, также могут быть использованы метрики Хеллингера, Кульбака–Лейблера и другие.

Для того, чтобы понять, является ли добавленный в момент времени t документ q новым, необходимо вычислить степень его похожести со всеми предыдущими документами. Пусть d^* — документ, максимально похожий на q :

$$d^* = \underset{d}{\operatorname{argmax}} \text{sim}_t(d, q). \quad (6)$$

Тогда значение

$$\text{score}_t(q) = 1 - \text{sim}_t(d^*, q) \quad (7)$$

может быть использовано для того, чтобы определить, является ли документ q новым. Новыми будем считать все документы q , у которых значение $\text{score}_t(q)$ больше, чем пороговое значение θ_s . В обратном случае считается что существует документ d^* , достаточно похожий на q и поэтому q не представляет собой сообщение на новую тему. Для того, чтобы определить подходящее значение θ_s , можно использовать размеченный корпус и посчитать значение $\text{sim}_t(d, q)$ среди документов, соответствующих одному и разным событиям.

Недостатком алгоритма NED является то, что он не учитывает отношения между словами. Например синонимичные слова в NED будут считаться совершенно разными. Для того, чтобы решить эту проблему могут быть применены тематические модели.

3.2 Выявление событий с помощью LDA

На следующем примере рассмотрим как тематические модели (topic models) могут быть использованы для задачи распознавания событий. Для этого опишем схему, по которой работает алгоритм, предложенный в статье [4].

Задача состоит в том, чтобы по данным к набору фотографий Flickr² распознать все события на определенную тему, проходившие в обозначенных городах. Есть несколько вариантов постановок задачи, они различаются тематикой событий и городами.

Фотографии содержат в себе описания, они будут являться документами в задаче распознавания. Также фотографии могут включать геотеги. Авторы статьи предлагают разбить решение задачи на пять составных частей:

1. предобработка текстовых данных,
2. определение в каком городе сделана фотография по ее описанию,
3. распознавание темы, к которой относится фотография,
4. распознавание события, к которому относится фотография,
5. оптимизация описания события, полученного на шаге 4.

В качестве предобработки, авторы предлагают выполнить следующее: удалить стоп-слова и html теги, провести стемминг слов³, перевести не английские слова на английский язык используя сервис Google Translate.

Так как задача упоминает отдельные города, необходимо разработать метод распознавания города по документу. Указанные в фотографии географические координаты были доступны авторам в 20% случаев, из этих координат были выявлены города используя сервис Google Tables. Используя технику TF-IDF, в описаниях фотографий с известными городами были извлечены ключевые слова. По этим ключевым словам появилась возможность назначить фотографиям без геотегов “ближайший” город с точки зрения похожести описаний. Для того чтобы определять похожесть документов, был использован подход, описанный в (3.1). В случаях, когда “ближайший” город выявить не удавалось, авторами использовались следующие предположения: считалось что один и тот же фотограф не мог побывать в один день более чем в двух разных городах, и что путешествие из одного города в другой занимает как минимум два часа. Эти предположения позволили улучшить классификатор, таким образом более 97% фотографий были привязаны к городу.

Далее необходимо для каждого города кластеризовать документы по темам и рассмотреть только те из них, которые заданы в описании задачи. Для распознавания тем была использована тематическая модель LDA (Latent Dirichlet Allocation)[5], для определения параметров которой применялось сэмплирование по Гиббсу[6]. LDA работает из предположения, что каждый документ D_i характеризуется случайным распределением над темами, в то время как каждая тема является мультиномиальным распределением над словами.

²<http://www.flickr.com>

³стемминг (stemming) — удаление окончаний у слов для их нормализации

Оставшаяся часть — извлечение событий и их оптимизация. Для того, чтобы алгоритм выявил событие, отвечающее теме k в день d , необходимо чтобы количество документов D_i по этой теме в день d превосходило некоторое пороговое значение θ . Оптимизация событий подразумевает под собой объединение событий на одну тему в последовательные дни и разделение событий в разных городах.

Авторы статьи тестировали алгоритм на трех разных вариантах условия задачи, в таблице 1 приведены результаты по каждому из них.

Данные	Точность	Полнота	F-мера
#1	80.98	19.25	31.10
#2	91.21	77.85	84.00
#3	90.76	81.91	86.11

Таблица 1: Точность, полнота и F-мера алгоритма при разных условиях задачи

По результатам из таблицы можно видеть что в первом случае алгоритм справляется со своей задачей существенно хуже, чем в других. Авторы объясняют это тем, что в задаче #1 необходимо было находить научные конференции, а так как они проходят на разные темы, не получалось выделить конкретный набор ключевых слов. По этой причине полнота алгоритма в случае #1 относительно низкая. В задаче #2, #3 напротив удалось обозначить необходимые ключевые слова, о чем свидетельствуют результаты.

3.3 Information nuggets

Рассмотрим подход к извлечению событий, описанный в [1]. Авторы ставят перед собой задачу составить алгоритм описания подсобытий в социальной сети Твиттер. Были использованы данные, полученные во время торнадо Joplin в 2011 году. Данные представляют из себя сообщения пользователей, содержащие хэштег #joplin, собранные 22 мая 2011 года на протяжении нескольких часов, пока плотность сообщений не стала относительно низкой. Авторы ставили цель извлекать из потока сообщений так называемые золотые самородки информации — короткие и информативные сообщения, описывающие происходящие события. По этой причине подход назван information nuggets.

Авторы видели основную проблему в том, что даже при наличии большого количества сообщений об одном событии, ими трудно пользоваться, потому что они имеют разную природу. Например это может быть сообщение очевидца о происходящем стихийном бедствии или сообщение о перечислении правительством средств на восстановление разрушенных построек. Статья предлагает делить сообщения на пять разных категорий по степени информативности:

- *Персональное*: информация в сообщении может быть полезна только автору и его кругу общения. Она не является интересной для людей, которые не знают автора сообщения непосредственно.
- *Информативное (напрямую)*: сообщение может быть полезно людям вне круга общения автора, и эта информация написана прямым участником или очевидцем событий.

- *Информативное (косвенно)*: сообщение может быть полезно людям вне круга общения автора, при этом автор пишет о том, что он слышал по телевидению, радио или любому другому источнику информации.
- *Информативное (напрямую или косвенно)*: сообщение может быть полезно людям вне круга общения автора, но невозможно ответить на вопрос как автор связан с происходящими событиями.
- *Другое*: сообщение или не на английском языке, или не может быть классифицировано.

В дальнейшем рассматриваются только сообщения с информативным типом, так как они наиболее вероятно содержат полезные сведения. Затем сообщения разбиваются на подтипы по направленности информации. Всего авторами использовалось 32 разных подтипа, среди которых:

- *Предостережения и советы*: документ содержит предупреждение о возможном опасном происшествии.
- *Жертвы и разрушения*: в тексте сообщается о потерях, вызванных стихийным бедствием.
- *Сбор средств*: сообщения описывают пожертвования денег и предметов пострадавшим от чрезвычайного происшествия.

Для того, чтобы классифицировать сообщения по типам и подтипам, в алгоритме используется наивный байесовский классификатор, который предварительно тренируется на размеченных данных. В классификаторе используется большое количество свойств, бинарных, скалярных и текстовых. В таблице 2 приведены результаты работы алгоритма на некоторых подтипах. Для тестирования использовались размеченные вручную сообщения.

Подтип	Точность	Полнота	F-мера
Предостережения и советы	0.618	0.598	0.605
Жертвы и разрушения	0.578	0.645	0.610
Сбор средств	0.546	0.632	0.585

Таблица 2: Результаты работы алгоритма information nuggets.

4 Теоретическая часть предложенного решения

Эта работа ставит перед собой задачу разработать алгоритм нахождения событий в потоке сообщений пользователей сети Твиттер. В качестве данных были использованы сообщения с 4 июня 2013 года по 1 июля 2013 года, содержащие в себе хэштег

#texas. Всего корпус включает порядка 240 тысяч сообщений и 1.5 миллиона слов. Приведем псевдокод алгоритма, а затем рассмотрим каждый шаг более подробно.

Алгоритм 1: Извлечение событий из сообщений Твиттера

Исходные параметры: множество сообщений Ω , параметры алгоритма M , L , l_D , J , R параметры для HDP модели

Результат: для каждого события набор ключевых слов и частота сообщений по теме события

```

1 в частотной характеристике сообщений из  $\Omega$  найти  $K$  экстремумов, пусть они
  соответствуют временам  $t_1, t_2, \dots, t_K$ ;
2 для каждого  $t \in \{t_1, t_2, \dots, t_K\}$  выполнять
3   | извлечь  $M$  ключевых слов из сообщений в  $\Omega$  в некотором радиусе  $R$ 
   | момента времени  $t$ ;
4   | сформировать фиктивные документы  $D_0 = \{D_0^j\}_{j=1}^J$ , каждый из которых
   | имеет длину  $l_D$  и состоит из полученных на шаге 3 ключевых слов
   | пропорционально их весу;
5   | добавить  $D_0$  к множеству  $\Omega$ ;
6   | каждому сообщению из  $\Omega$  назначить тему, с помощью модифицированной
   | модели HDP (Hierarchical Dirichlet Process);
7   | пусть документам  $D_0$  была назначена тема  $k_0$ , удалить  $D_0$  из  $\Omega$ ;
8   | найти частотную характеристику сообщений в  $\Omega$  с темой  $k_0$ ;
9   | если частотная характеристика определяет новое событие тогда
10  |   | найти  $L$  ключевых слов темы  $k_0$  из HDP;
11  |   | добавить ключевые слова, частотную характеристику в результат;
12  | конец условия
13 конец цикла

```

4.1 Нахождение экстремумов

Понятие экстремума в этом разделе отличается от привычного определения в математическом анализе. Назовем экстремумом те точки, в которых плотность сообщений позволяет утверждать что в этот момент произошло событие. В первую очередь, понятие экстремума для задачи выявления событий должно быть локальным. С другой стороны не все локальные максимумы можно корректно обозначить как события. Задача в общем виде получила развитие в контексте обработки сигналов [7]. Существует множество методов, каждый из которых будет больше подходить к определенному типу данных. Алгоритм нахождения экстремумов может меняться вне зависимости от других составных частей исходного алгоритма. Несколько возможных подходов нахождения экстремумов для функции $f(x)$ с сеточной областью определения X перечислены ниже:

- Алгоритм помечает как экстремумы все точки, в которых функция достигает максимума в некотором окне, и при этом ее значение больше некоторого заранее определенного θ . Опционально θ может зависеть от среднего значения функции и других статистических характеристик.

- Для любых двух значений аргумента x и y , где $x < y$ определим функции $T(x, y)$ (travel) и $R(x, y)$ (rise):

$$\begin{aligned} T(x, y) &= \sum_{x \leq k < y} |f(k+1) - f(k)|, \\ R(x, y) &= f(y) - f(x) + \epsilon, \end{aligned} \quad (8)$$

где ϵ — некоторое небольшое значение. Тогда значения $T(x, y)/R(x, y) > 1$ соответствуют пику функции между точками x и y . Алгоритм может находить все экстремумы, где значение T/R больше некоторого порога.

- Подход заключается в том, чтобы выделить в функции стандартную пиковую подфункцию, например функцию Гаусса. Для этого применяются согласованные фильтры. Пусть $g(x)$ — функция Гаусса с некоторыми μ и σ :

$$g(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}. \quad (9)$$

Тогда можно определить степень похожести исходной функции на функцию $g(x)$:

$$\text{sim}(f, g) = \frac{(f, g)}{\|f\| \cdot \|g\|}. \quad (10)$$

В косинусной метрике выражение примет следующий вид:

$$\text{sim}(f, g) = \frac{\sum_{x \in X} f(x) \cdot g(x)}{\sqrt{\sum_{x \in X} f^2(x)} \cdot \sqrt{\sum_{x \in X} g^2(x)}}. \quad (11)$$

Аналогично предыдущим случаям, будем считать что функции “похожи” если значение $\text{sim}(f, g)$ не меньше некоторого θ .

- Один из способов найти экстремумы — применить смешанную модель (mixture model) к входной функции. Компонентами модели могут быть нормальные распределения. Затем нужно найти параметры модели используя ЕМ (expectation maximization) алгоритм или сэмплирование по Гиббсу. Этот метод будет корректно обрабатывать случай, когда экстремумы накладываются друг на друга или расположены очень близко.

Каждый из перечисленных методов можно комбинировать со сглаживающими фильтрами если шум мешает выделить события. Для исходной задачи был выбран первый метод из списка. Данные показали что экстремумы сильно выделяются на фоне средней плотности сообщений и расположены на значительном расстоянии. Таким образом выбранный метод экспериментально обоснован для имеющихся данных.

4.2 Извлечение ключевых слов

Для того, чтобы использовать на шаге 6 модель HDP с частичным обучением, необходимо знать ключевые слова события, которое мы хотим проследить с помощью тематической модели. Их извлечение происходит на шаге 3. В реализации алгоритма учитывая особенности сети Твиттер был реализован наиболее простой способ

извлечь M ключевых слов из набора сообщений, а именно отсортировать слова по невозрастанию частоты и взять первых M слов.

Этот способ показал хорошие результаты потому что события в сети Твиттер носят "взрывной" характер и частота релевантных слов во время события в десятки и сотни раз превосходит нормальную частоту этих слов. В случае когда "взрывного" характера событий не наблюдается, при извлечении ключевых слов можно использовать технику TF-IDF. Также в этом случае можно рассматривать биграммы и применять различные методы NLP⁴ к анализируемому тексту.

4.3 Модифицированная модель HDP

Для извлечения всех сообщений, принадлежащих событию, в алгоритме используется модифицированная модель HDP, в которую добавлена возможность частичного обучения.

4.3.1 Иерархический процесс Дирихле

Пусть (Θ, \mathcal{B}) — измеримое пространство, с вероятностной мерой G_0 . Пусть α_0 — положительное действительное число. Определим процесс Дирихле $DP(\alpha_0, G_0)$ как случайное распределение вероятностной меры G над (Θ, \mathcal{B}) , такое, что для любого конечного измеримого разбиения Θ (A_1, A_2, \dots, A_r) , случайный вектор $(G(A_1), G(A_2), \dots, G(A_r))$ будет распределен согласно конечномерному распределению Дирихле с параметрами $(\alpha_0 G_0(A_1), \alpha_0 G_0(A_2), \dots, \alpha_0 G_0(A_r))$:

$$(G(A_1), G(A_2), \dots, G(A_r)) \sim \text{Dir}(\alpha_0 G_0(A_1), \alpha_0 G_0(A_2), \dots, \alpha_0 G_0(A_r)). \quad (12)$$

Иерархический процесс Дирихле это распределение над множеством вероятностных мер над (Θ, \mathcal{B}) . Процесс определяет множество вероятностных мер G_j , по одному на каждую группу, и глобальную меру G_0 . Глобальная мера G_0 распределена согласно процессу Дирихле с параметрами γ и базовым распределением H :

$$G_0 \mid \gamma, H \sim DP(\gamma, H). \quad (13)$$

Групповые меры G_j условно независимы при заданном G_0 и имеют следующее распределение:

$$G_j \mid \alpha_0, G_0 \sim DP(\alpha_0, G_0). \quad (14)$$

Гиперпараметрами HDP являются базовое распределение H и концентрационные параметры γ и α_0 . Иерархический процесс Дирихле может быть использован как априорное распределение над параметрами групповых данных. Для каждого j пусть $\theta_{j1}, \theta_{j2}, \dots$ — независимые распределенные по закону G_j случайные величины. Каждая величина θ_{ji} будет соответствовать наблюдаемому значению x_{ji} . Правдоподобие может быть записано в следующем виде:

$$\begin{aligned} \theta_{ji} \mid G_j &\sim G_j, \\ x_{ji} \mid \theta_{ji} &\sim F(\theta_{ji}). \end{aligned} \quad (15)$$

Процесс, описанные выше носит название смешанная модель для иерархического процесса Дирихле (hierarchical Dirichlet process mixture model)[8].

4.3.2 Chinese Restaurant Franchise

HDP — непараметрическая тематическая модель, процесс генерации данных по ней может быть описан в терминах процесса "chinese restaurant franchise" (CRF). В этом процессе существует сеть ресторана, в каждом из которых одинаковое меню. В ресторане находится неограниченное число столов, за каждым столом неограниченное число мест; каждый посетитель, заходя в ресторан, садится либо за уже занятый стол, либо за новый стол. На каждом непустом столе в ресторане заказано одно блюдо из меню, которое едят все посетители, сидящие за этим столом. Блюдо на стол заказывает первый человек, севший за этот стол.

В CRF группам соответствуют рестораны, а параметры θ_{ji} считаются посетителями. Переменные ϕ_1, \dots, ϕ_K составляют единое меню для сети ресторанов, эти переменные распределены согласно базовому распределению H . Пусть ψ_{jt} — переменная, связанная со столом t в ресторане i и обозначающее блюдо, которое заказано за этим столом.

Каждый посетителю θ_{ji} соответствует стол ψ_{jt} , в то время как каждому столу ψ_{jt} соответствует одно блюдо ϕ_k . Для того, чтобы указать это соответствие, используются индексы:

- t_{ji} — индекс, связывающий посетителя θ_{ji} со столом ψ_{jt} , за которым сидит этот посетитель.
- k_{jt} — индекс, связывающий стол ψ_{jt} с блюдом ϕ_k , которое заказано для этого стола.

Для описания процесса, необходимо ввести следующие обозначения:

- n_{jtk} — количество посетителей в ресторане j за столом t , едящих блюдо k ;
- n_{jt} — количество посетителей в ресторане j за столом t ;
- $n_{j \cdot k}$ — количество посетителей в ресторане j , едящих блюдо k ;
- m_{jk} — количество столов в ресторане j , на которых заказано блюдо k ;
- $m_{j \cdot}$ — количество столов в ресторане j ;
- $m_{\cdot k}$ — количество столов, на которых заказано блюдо k ;
- $m_{\cdot \cdot}$ — общее количество столов.

Запишем условное распределение θ_{ji} , где G_j исключены интегрированием:

$$\theta_{ji} \mid \theta_{j1}, \theta_{j2}, \dots, \theta_{j,i-1}, \alpha_0, G_0 \sim \sum_{i=1}^{m_j} \frac{n_{jt}}{i-1 + \alpha_0} \delta_{\psi_{jt}} + \frac{\alpha_0}{i-1 + \alpha_0} G_0. \quad (16)$$

Аналогично, интегрируя по G_0 , получим:

$$\psi_{jt} \mid \psi_{11}, \psi_{12}, \dots, \psi_{21}, \dots, \psi_{j,t-1}, \gamma, H \sim \sum_{k=1}^K \frac{m_{\cdot k}}{m_{\cdot \cdot} + \gamma} \delta_{\phi_k} + \frac{\gamma}{m_{\cdot \cdot} + \gamma} H. \quad (17)$$

⁴Natural Language Processing — обработка текстов на естественном языке

Уравнение (16) описывает выбор стола посетителем, в то время как уравнение (17) определяет выбор блюда.

Если перевести используемые термины в область тематических моделей, посетители будут соответствовать словам, рестораны документам, столы внутренним темам документов, блюда внешними темами всего корпуса[8].

4.3.3 Модификация частичной обучаемости

HDP является моделью без учителя, она не требует размеченных данных для того чтобы назначить документам соответствующие темы. Для того, чтобы контролировать процесс определения тем, необходимо соответствующе изменить модель. В алгоритме 1 необходимо добиться от HDP того, чтобы она выделила ключевые слова из некоторого множества \mathcal{K} в отдельную тему k_0 . Есть несколько способов добиться требуемого результата, один из них описан ниже [10].

Составим из ключевых слов в \mathcal{K} документ длины M , в котором количество каждого ключевого слова пропорционально его весу. Перед тем как находить параметры HDP для корпуса документов, добавим J сгенерированных документов в корпус.

Для того, чтобы извлечь параметры из модели HDP в алгоритме применяется сэмплирование по Гиббсу. В этом методе изначальные документы каким-то образом распределяются по темам, а затем итеративно каждое слово удаляется из состояния и назначается на новое место, согласно уравнениям (16) и (17). Краткий псевдокод сэмплирования по Гиббсу для модели HDP приведен ниже:

Алгоритм 2: Сэмплирование по Гиббсу

Исходные параметры: множество документов Ω , параметры для HDP
модели α_0, γ

```

1  назначить всем словам тему 0;
2  повторять
3      для каждого документа  $d \in \Omega$  выполнять
4          для каждого слова  $w \in d$  выполнять
5              удалить  $w$  из состояния модели;
6              выбрать внутреннюю тему  $t$  согласно уравнению (16);
7              если выбрана новая внутренняя тема  $t$  тогда
8                  выбрать внешнюю тему  $k$  согласно (17);
9                  назначить внутренней теме  $t$  внешнюю тему  $k$ ;
10             иначе
11                 извлечь назначенную ранее внешнюю тему  $k$  для внутренней
12                 темы  $t$ ;
13             конец условия
14         назначить слову  $w$  внутреннюю тему  $t$  и внешнюю тему  $k$ ;
15     конец цикла
16 до тех пор, пока не прошло  $L$  итераций;
```

Для того, чтобы внести в модель частичную обучаемость, необходимо по-разному обрабатывать оригинальные документы и дополнительные. Для первоначальных документов должен сохраниться родной алгоритм, а словам в дополнительном доку-

менте всегда будем назначать внешнюю тему k_0 . Приведем псевдокод для модифицированной HDP:

Алгоритм 3: Сэмплирование по Гиббсу с частичным обучением

Исходные параметры: множество документов Ω , параметры для HDP модели α_0, γ, k_0

```

1 назначить всем словам тему 0;
2 повторять
3   для каждого документа  $d \in \Omega$  выполнять
4     для каждого слова  $w \in d$  выполнять
5       если документ  $d$  дополнительный тогда
6         | назначить слову  $w$  внутреннюю тему 0 и внешнюю тему  $k_0$ ;
7       иначе
8         | удалить  $w$  из состояния модели;
9         | выбрать внутреннюю тему  $t$  согласно уравнению (16);
10        если выбрана новая внутренняя тема  $t$  тогда
11          | выбрать внешнюю тему  $k$  согласно (17);
12          | назначить внутренней теме  $t$  внешнюю тему  $k$ ;
13        иначе
14          | извлечь назначенную ранее внешнюю тему  $k$  для внутренней
15          | темы  $t$ ;
16        конец условия
17        назначить слову  $w$  внутреннюю тему  $t$  и внешнюю тему  $k$ ;
18      конец условия
19    конец цикла
20 до тех пор, пока не прошло  $L$  итераций;
```

Изменяя параметры M и J можно определять насколько сильной будет привязка к конкретным ключевым словам для искомой темы. Экспериментально установлено что в качестве параметра k_0 можно взять значение 1.

4.4 Анализ частотной характеристики сообщений темы

Возможны случаи, когда сообщения, найденные в экстремуме не являются новыми событиями. Это возможно в следующих случаях:

1. уже было найдено такое же событие,
2. график события не является одним максимумом.

Первый случай можно выявить, найдя нормализованное скалярное произведение двух функций частоты в косинусной метрике. Если значение будет близко к единице, события можно считать тождественными. Чтобы различать второй случай, предлагается посчитать императивное стандартное отклонение. Если оно будет достаточно велико, событие нельзя считать подходящим.

Полученный метод не будет работать в случае когда событие состоит из двух достаточно удаленных пиков. В этом случае можно использовать смешанную модель, как это было описано выше в методах нахождения экстремума. При тестировании подобных событий не было найдено [11].

5 Реализация и тестирование решения

Для реализации программы, описанной алгоритмом 1, был выбран язык Java. Выбор мотивирован тем, что в качестве свободной реализации для модели HDP использовалась реализация на языке Java, написанная Arnim Bleier et al [9]. Java объектно-ориентированный язык высокого уровня, что оказалось очень удобным для реализации алгоритма 1.

5.1 Описание данных

В качестве данных для задачи были использованы сообщения пользователей сети Твиттер с 4 июня по 1 июля 2013 года, содержащие хэштег `#texas`. Всего в корпусе находится около 240 тысяч сообщений. Цель алгоритма — найти события, показать их на графике частоты сообщений и указать ключевые слова к событиям.

На рисунке 5.1 изображена плотность сообщений во входных данных. На графике есть несколько пиков, которые обозначают события.

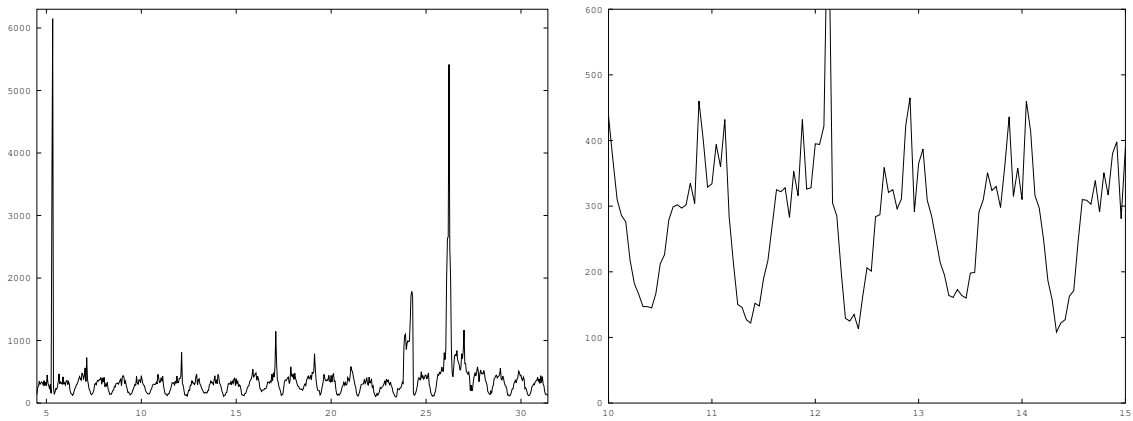


Рис. 2: Зависимость количества сообщений в час от времени. На втором графике видна зависимость от времени суток.

5.2 Извлечение событий

На рисунке 5.2 можно видеть экстремумы, отмеченные алгоритмом 1 на шаге 1. Путем выбора метода и параметров, можно пометить больше или меньше максимумов.

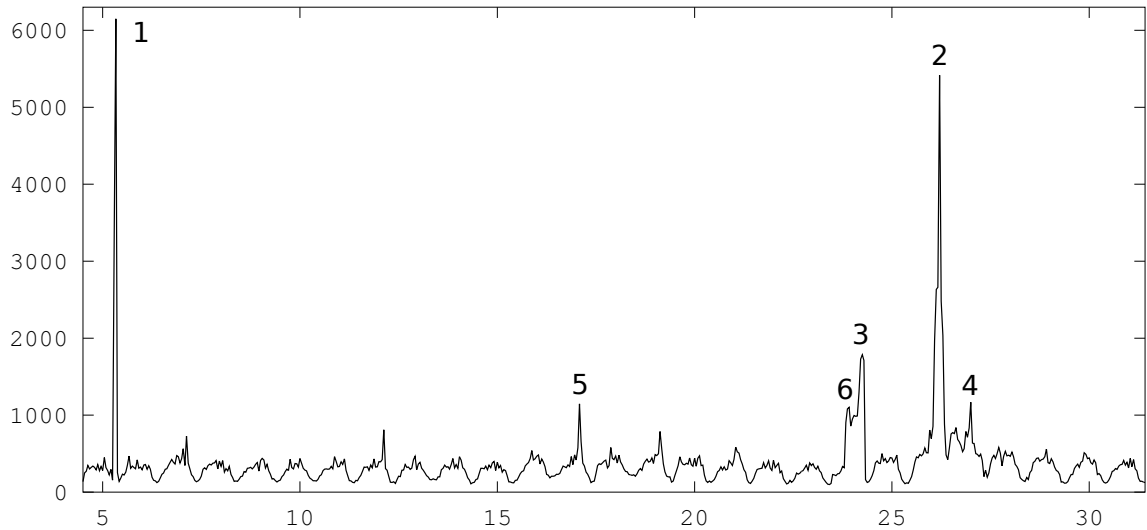


Рис. 3: График частоты сообщений с отмеченными экстремумами.

Следующие таблицы показывают параметры частотных характеристик сообщений, полученных на шаге 8 основного алгоритма. На основании этих данных те или иные события войдут в финальный результат или будут отброшены.

	σ^2
#1	7853
#2	8993
#3	3500
#4	34628
#5	30232
#6	4261

Таблица 3: Жирным шрифтом помечены события с $\sigma^2 > \sigma_0^2$.

#1	1					
#2	0.004	1				
#3	0.001	0.009	1			
#4	0.057	0.109	0.131	1		
#5	0.038	0.079	0.120	0.935	1	
#6	0.001	0.006	1.000	0.132	0.121	1
	#1	#2	#3	#4	#5	#6

Таблица 4: Значения нормализованного скалярного произведения функций частот в косинусной метрике. Жирным шрифтом помечены значения, позволяющие считать события тождественными.

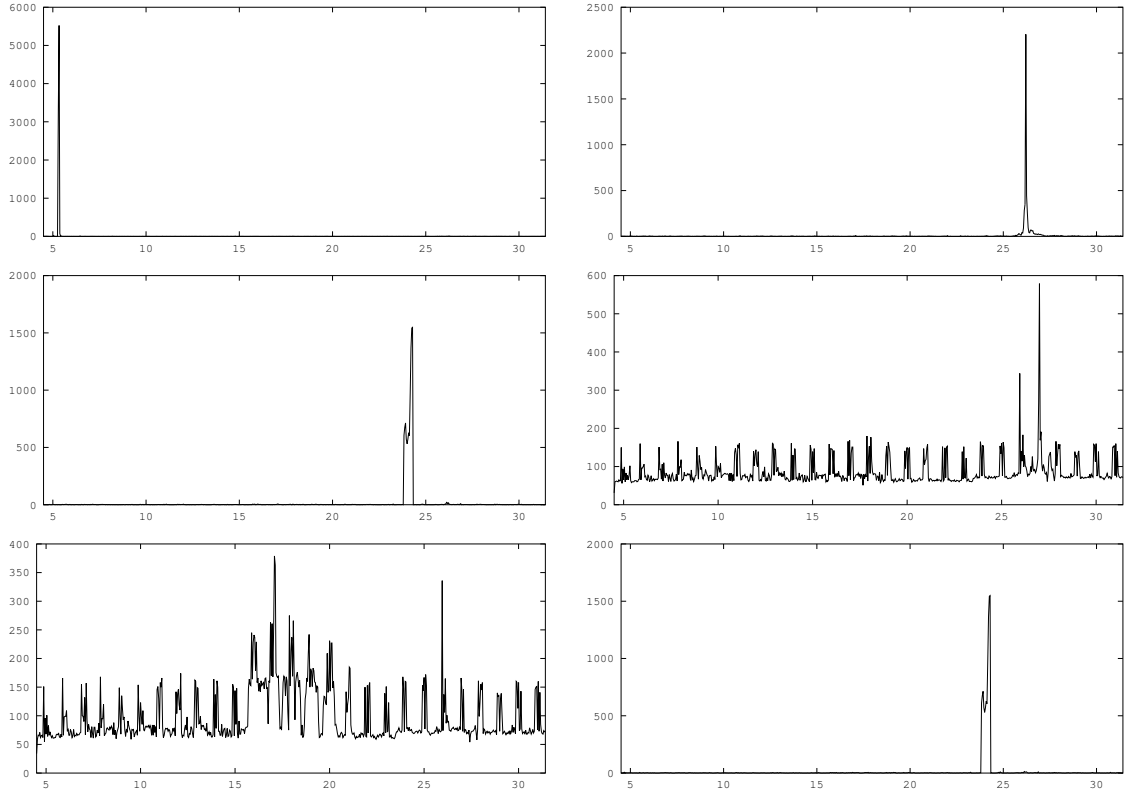


Рис. 4: Частоты сообщений для "событий" #1 — #6.

Из таблиц 3 и 4 можно видеть, что алгоритм выдаст в конечном счете только события #1, #2, #3. События #4 и #5 будут отброшены как имеющие большое среднеквадратичное отклонение. Событие #6 не войдет в финальный результат потому, что значение похожести функций частоты для событий #3, #6 больше допустимой нормы:

$$\text{sim}^{\cos}(f_{\#3}, f_{\#6}) = \frac{(f_{\#3}, f_{\#6})}{\|f_{\#3}\| \cdot \|f_{\#6}\|} > \theta_{\text{sim}}. \quad (18)$$

Приведем ниже таблицу 5 с описаниями событий #1, #2, #3:

	Ключевые слова	Вызвавшее событие
#1	year, read, prison, dragon, amazon, written, bestseller	Написанная в техасской тюрьме книга The Sword and the Dragon является бестселлером на Amazon более трех лет.
#2	sb5, standwithwendi, abort, senat	Обсуждение сенатом закона по запрету аборт, который носит название Senate Bill 5. Wendi Davis — политик, которая боролась с принятием этого закона.
#3	houston, california, england, artist, rt, watch, defjam	Хип хоп исполнитель из Хьюстон Devin the Dude объявил, что его восьмой студийный альбом будет называться One for the Road и будет выпущен в продажу в сентябре 2013 года. Альбом выпущен лейблом Def Jam Recordings.

Таблица 5: Таблица с описаниями событий

6 Заключение

В ходе выполнения дипломной работы было достигнуто:

- изучены существующие методы извлечения и описания событий из коротких пользовательских сообщений,
- изучена возможность применения тематических моделей для этой задачи,
- разработан и проанализирован алгоритм извлечения событий из социальной сети Твиттер на основе иерархического процесса Дирихле,
- алгоритм был реализован, параметры подобраны экспериментальным путем,
- тестирование показано, что алгоритм справляется с поставленной задачей.

Список литературы

- [1] Imran, Elbassuoni, Castillo, Diaz and Meier. Extracting Information Nuggets from Disaster-Related Messages in Social Media. 2013.
- [2] Xun Wang, Feida Zhu, Jing Jiang, Sujian Li. Real Time Event Detection in Twitter. 2011.
- [3] Thorsten Brants, Francine Chen, Ayman Farahat. A System for New Event Detection. 2003.
- [4] Konstantinos N. Vavliakis, Fani A. Tzima, and Pericles A. Mitkas. Event Detection via LDA for the MediaEval2012 SED Task. 2012.
- [5] David M. Blei, Andrew Y. Ng, Michael I. Jordan. Latent Dirichlet Allocation. 2003.
- [6] Tom Griffiths. Gibbs sampling in the generative model of Latent Dirichlet Allocation.

- [7] Girish Keshav Palshikar. Simple Algorithms for Peak Detection in Time-Series.
- [8] Yee Whye Teh, Michael I. Jordan, Matthew J. Beal, David M. Blei. Hierarchical Dirichlet Processes. 2005.
- [9] Yee Whye Teh. Hierarchical Bayesian Nonparametric Models with Applications. 2009.
- [10] Ramnath Balasubramanyan, William W. Cohen, Matthew Hurst. Modeling corpora of timestamped documents using semisupervised nonparametric topic models.
- [11] David Blei. COS 424: Interacting with Data. 2008.