

**МОСКОВСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ  
УНИВЕРСИТЕТ ИМЕНИ Н. Э. БАУМАНА**  
**Факультет информатики и систем управления**  
**Кафедра теоретической информатики и компьютерных  
технологий**

Наброски дипломного проекта

«Автоматическое установление связей между сообщениями  
твиттера и новостными статьями»

Выполнил:

студент ИУ9-101

Выборнов А. И.

Руководитель:

Лукашевич Н.В.

Москва 2016

# Содержание

<b>Введение</b>	<b>3</b>
<b>1. Обзор литературы</b>	<b>5</b>
1.1. Linking Tweets to News: A Framework to Enrich Short Text Data in Social Media . . . . .	5
1.1.1. Перевод аннотации . . . . .	5
1.1.2. Идея статьи . . . . .	6
1.2. Linking Online News and Social Media . . . . .	7
1.2.1. Перевод аннотации . . . . .	7
1.3. Bridging Vocabularies to Link Tweets and News . . . . .	8
1.3.1. Основная идея . . . . .	8
<b>Список литературы</b>	<b>9</b>

# Введение

В современном мире всё больший вес приобретают социальные медиа (преимущественно социальные сети). Их главное отличие от традиционных медиа (газеты, тв) заключается в том, что контент порождается тысячами и миллионами людей. Социальные медиа не заменяют традиционные новостные источники, а дополняют их. Они могут служить полезным социальным датчиком того, насколько популярна история (тема) и как долго. Часто обсуждения в социальных медиа основаны на событиях из новостей и, наоборот, социальные медиа влияют на новостные события.

Одной из самых популярных социальных сетей является Twitter - социальная сеть для публичного обмена сообщениями. Главной особенностью Twitter является малый размер сообщений (140 символов), называемых твитами. Часто твиты являют собой описание, происходящего прямо сейчас события, отклик на него.

...

Выявление связи между сообщениями твиттера (твитов) и новостями позволит как расширить информативность твитов, так и обогатить новости.

...

Преимущества расширения новости с помощью твитов: определение отношения аудитории к новости, дополнительные признаки для тематической классификации новостей, дополнительная информация для аннотирования новостей.

Современные методы обработки естественного языка хорошо работают, используя большой массив текста в качестве входных данных, однако, они становятся неэффективными, когда применяются на коротких текстах, таких как твиты. Существенным преимуществом расширения твита с помощью новости является появляющаяся возможность использования большого количества методов обработки естественного языка

(Natural Language Processing).

...

Данная работа ставит целью исследование и разработку методов автоматического установления связей между сообщениями твиттера и новостными статьями.

Не существует стандартных решений, и есть считанное количество статей. На основе этих статей будет сделана попытка построить pipeline для получения подобной взаимосвязи.

# 1. Обзор литературы

В рамках предварительного исследования были разобраны несколько статей [1] [2] [3]. Ниже приводится краткое изложение основных идей, описанных в выбранных статьях.

## 1.1. Linking Tweets to News: A Framework to Enrich Short Text Data in Social Media

### 1.1.1. Перевод аннотации

Многие современные методы обработки естественного языка (NLP<sup>1</sup>) хорошо работают с большой массив текста в качестве входных данных. Однако они очень неэффективны при работе с короткими текстами (к примеру твиты). Преодоление этой проблемы мы видим в нахождении соответствующего твиту новостного документа. Решение этой задачи требует хорошего моделирования семантики коротких текстов.

Основной вклад статьи двойной:

1. представлено решение задачи нахождения взаимосвязи между твитами и новостями, из этого могут извлечь выгоду многие NLP задачи;
2. в отличие от предыдущих исследований, которые фокусируются на лексических особенностях коротких текстов (информация о связи текст-слово), мы предлагаем взаимосвязь, основанную на модели скрытой переменной, которая моделирует корреляцию между короткими текстами (информация о связи текст-текст). Необходимость этого обоснована наблюдением: твит обычно покрывает только один аспект события.

---

<sup>1</sup>Natural Language Processing

Мы покажем, что с помощью особенных признаков твита (хэштегов) и особых признаков новостей (именованных сущностей<sup>1</sup>) а также временных ограничений, мы можем получить взаимосвязь текст-текст, и, таким образом, дополнить семантическую картину короткого текста. Наши эксперименты показывают значительное преимущество нашей новой модели над baseline<sup>2</sup>.

### 1.1.2. Идея статьи

Современные методы обработки естественного языка плохо работают с короткими текстами. Для преодоления этого к твитам привязываются соответствующие новости.

Для формирования обучающей выборки, были выбраны твиты, которые имели ссылки на новости, опубликованные новостными агентствами (CNN или NYT) в тот же период.

Как показано в статье [5], добавление к твиту содержимого веб-страницы, ссылка на которую включена в этот твит, повышает **purity score** их кластеризации с 0.280 до 0.392.

Модели со скрытой переменной хорошо подходят для отображения коротких текстов в плотный малоразмерный вектор. В рамках решения задачи была применена модель со скрытой переменной, которая называется WTMF (Weighted Textual Matrix Factorization, подробное описание[6]), к твитам и к новостям. Модель была протестирована на двух схожих наборах данных из небольших сообщений. Как результат - используемая модель с большим запасом превзошла и LSA (Latent Semantic Analysis) и LDA (Latent Dirichlet Allocation). Эта модель позволила добавить информацию об отсутствующих словах в твит (модель WTMF добавляет более 1000 фичей к твиту, LDA лишь 14). Недостатком WTMF является то, что порождается только связь текст-слово, без

---

<sup>1</sup>Какой-то кривой перевод, найдо найти получше. In data mining, a named entity is a phrase that clearly identifies one item from a set of other items that have similar attributes.

<sup>2</sup>Как перевести?

учёта взаимосвязи между короткими текстами.

Ввиду разреженности исходных данных, возникает ещё одна проблема: твит обычно отражает, только один аспект события.

Полученный подход не учитывает следующих характеристик, которым обладает исходная выборка:

1. Хэштеги, которые являются прямым указанием на смысл твита.
2. **Named entities** новостей. Из новостей можно с высокой точностью извлекать **named entities**, используя инструменты для NER (Named Entity Recognition). Если несколько текстов содержат схожие **named entities** они наверняка описывают одно и то же событие.
3. Информация о времени публикации для твитов и новостей. Если несколько текстов опубликованы примерно в одно и то же время, то велик шанс, что они описывают одно и то же событие

В статье описывается решение проблемы поиска взаимосвязи между текстами, с использованием описанных выше характеристик. Два связанных текста, должны иметь схожий скрытый вектор (семантическая модель твита достраивается из схожих твитов).

Это дополнительная информация была добавлена в модель WTMF. Было также показано различное влияние на связь текст-текст жанра твита и жанра новости. Был получен на порядок более лучший результат чем при использовании исходной WTMF модель.

## 1.2. Linking Online News and Social Media

### 1.2.1. Перевод аннотации

Многое из того, что обсуждается в социальных медиа вдохновлено событиями, описанными в новостях и, наоборот, социальные медиа предоставляют механизм, позволяющий влиять на новостные события.

Мы обращаемся к следующей задаче: по новости, найти в социальных сетях высказывания, которые неявно на неё ссылаются. Используется трехступенчатый подход: сначала получаются несколько моделей запросов по исходной статье, затем модели используются для получения высказываний из индекса целевого социального медиа, результатом являются несколько ранжированных списков, которые объединяются с использованием особой техники слияния данных. Модель запроса создаётся как на основе структуры статьи, так и на основе явно связанных со статьей высказываний из социальных медиа. Для борьбы с дрейфом запроса<sup>1</sup> при большого объёме используемого текста (либо в новости, либо в явно связанных высказываниях из социальных медиа), предлагается основанный на графике метод для выбора отличительных условий.

В нашей экспериментальной оценки для порождения моделей запросов, использованы данные из Twitter, Digg, Delicious<sup>2</sup>, the New York Times Community, Wikipedia и блогосферы. Показано, что другие модели запросов, основанные на различных источниках данных, не только обеспечивают дополнительную информацию, но и влияют на получение различных высказываний из социальных медиа по нашему целевому индексу. Как следствие, методы слияния данных приводят к значительному повышению производительности в сравнении с индивидуальными подходами. Показано, что основанный на графике метод выделения условий помог улучшить как эффективность, так и продуктивность.

## 1.3. Bridging Vocabularies to Link Tweets and News

### 1.3.1. Основная идея

Значительную сложность при решении проблемы связывания твитов с новостями преимущественно вызывают малый размер твита и различия в словарях: в твитах используются аббревиатуры, неформальный

---

<sup>1</sup>Порождение менее подходящего запроса.

<sup>2</sup>Веб-сайт, бесплатно дающий зарегистрированным пользователям услугу хранения и публикации закладок на страницы Всемирной сети.



язык, сленг, в новостях, напротив, используется литературный язык. Также твиты очень зашумлены и не содержат полезного содержания.

Твиттер предлагает хештэги, как механизм для категоризации твитов. Но этот подход далеко не совершенен, так как не только далеко не все записи содержат хештеги, но и записи содержащие хештеги обладают рядом проблем. Таковыми как: хештег не содержит информацию о событии, хештег сформулирован в слишком общей форме, твит содержит несколько хештегов. Из этого делается вывод, что использование только хештегов приведёт к низкому качеству связывания твитов с новостями.

Предлагается следующий подход: Используется LDA для построения моделей тем поверх новостей. Затем среди твитов ищутся наиболее близкие к конкретному топику. Из полученных твитов извлекаются слова, которые служат “мостом” к другим твитами.

## Список литературы

- [1] W. Guo, H. Li, H. Ji, and M. T. Diab. Linking tweets to news: A framework to enrich short text data in social media. - ACL, pages 239–249, 2013.
  - [2] Manos Tsagkias, Maarten de Rijke, Wouter Weerkamp. Linking Online News and Social Media. - ISLA, University of Amsterdam.
  - [3] T. Hoang-Vu, A. Bessa, L. Barbosa and J. Freire. Bridging Vocabularies to Link Tweets and News. - International Workshop on the Web and Databases (WebDB 2014), Snowbird, Utah, US, 2014.
  - [4] J. Sankaranarayanan, H. Samet, B. Teitler, M. Lieberman, J. Sperling. TwitterStand: news in tweets. - 17th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, 2009, Seattle, Washington.
  - [5] Ou Jin, Nathan N. Liu, Kai Zhao, Yong Yu, and Qiang Yang. 2011. Transferring topical knowledge from auxiliary long texts for short text clustering. In Proceedings of the 20th ACM international conference on Information and knowledge management.
  - [6] Weiwei Guo and Mona Diab. 2012a. Modeling sentences in the latent space. In Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics.
-