

Санкт-Петербургский государственный университет
Факультет Прикладной математики - Процессов управления

Кафедра Технологии программирования

Пекша Любовь Станиславовна
Выпускная квалификационная работа бакалавра

Название

Направление 01.03.02
Прикладная математика и информатика

Научный руководитель:
доцент Блеканов И.С.

Санкт-Петербург
2019

Оглавление

Введение	3
Постановка задачи	4
Обзор литературы	5
1. Данные	6
1.1. Сбор данных	6
1.2. Формирование траекторий	8
1.3. Обзор данных	9
2. Модель тематических траекторий	10
2.1. Базовая модель	10
2.1.1. Описание модели	10
2.1.2. Достоинства и недостатки	13
2.2. Модернизация рекуррентного слоя	14
2.2.1. Идея	14
2.2.2. Эксперимент	14
2.3. Изменение функции ошибки	14
2.3.1. Идея	14
2.3.2. Эксперимент	14
Выводы	15
Заключение	16
Список литературы	17

Введение

Digital Humanities (цифровые гуманитарные науки) — стремительно развивающееся в наше время направление. На это есть множество причин, например успехи в области анализа данных, появление открытых данных, практическая польза, а также желание исследователей применить имеющийся математический аппарат в новой сфере.

Анализ текстов СМИ и анализ упоминаний персоналий в СМИ в частности — слабо затрагиваемая область Digital Humanities (цифровые гуманитарные науки), которая не только представляет собой не только интересную с научной точки зрения задачу, но и имеет практическое применение. Например, сравнение упоминаний персоналии в различных СМИ или получение общего представления о личности путем анализа СМИ. Также особый интерес представляет собой анализ персоналии во времени, то есть изучение того, как изменяется новостная повестка относительно персоналии с течением времени.

Причины по которым данная сфера не исследована кроются в следующем: отсутствие размеченных данных, отсутствие способа формально измерить качество работы алгоритма, сравнительно небольшое количество работ посвященных анализу текста во времени в целом

Настоящая работа посвящена анализу упоминаний персоналий в СМИ в долгосрочном периоде. Данная тема предоставляет широкий спектр возможных направлений для работы такие как выделение траекторий персоналий — ранжированного по времени набора статей о персоналии, выявление долгосрочных трендов траектории, а также анализ эмоциональной окраски траекторий.

Постановка задачи

Основной задачей проекта является построение модели осуществляющей анализ траекторий персоналий в долгосрочном периоде, ее реализация и оценивание.

Для достижение цели ставятся следующие задачи:

1. Формирование обучающей выборки
2. Построение модели, способной выделять долгосрочные тренды и тематики
3. Модернизация модели таким образом, чтобы она выделяла траектории, отражающие эмоциональную окраску упоминания персоналии в СМИ
4. Оценивание моделей

Обзор литературы

Нет подходов, реализующих анализ траекторий СМИ. Но похожие модели реализованы для художественной литературы[1, 2].

Особый интерес представляет собой работа Feuding Families and Former Friends: Unsupervised Learning for Dynamic Fictional Relationships[1]. В этой работе представляется модель RMN, которая совместно обучает набор глобальных дескрипторов отношений между героями из художественных текстов, а также сопоставляет каждому отрывку текста веса дескрипторов, которые отражают на сколько хорошо каждое из слов-дескрипторов описывает траекторию в текущий момент времени.

RMN превзошла модель HTMM[2], решающую сходную задачу, а также модели тематического моделирования LDA[3] и NUBBI[4], которые способны выделять тематики в неразмеченном тексте, но не предназначены для построения траекторий в долгосрочном периоде.

****Возможно что-то стянуть из описания базовой модели****

1. Данные

Для формирования обучающей выборки необходимо выполнить следующие шаги

1. Сбор данных с сайтов СМИ
2. Выделение траекторий персоналий - ранжированного по времени набора статей о персоналии
3. Обработка текста статей с целью выделить релевантную информацию

Для дальнейшей работы необходимо собрать данные с сайтов СМИ. В открытом доступе есть датасеты со статьями различных СМИ. Но в силу того, что такие датасеты не предоставляют информацию о дате выхода статьи, возникла необходимость собрать необходимые данные непосредственно с новостных сайтов.

Для работы были выбраны следующие источники: Lenta.ru[5], Tvrain[6], Meduza[7], РИА Новости[8]. Данные СМИ имеют внушительный архив данных, а также ведут свою деятельность продолжительный период времени.

1.1. Сбор данных

Сбор данных осуществлялся с помощью Scrapy[9]. Это написанная на Python платформа, которая нацелена на простой, быстрый и автоматизированный обход (краулинг) веб-страниц, имеющий большую популярность. Одним из главных преимуществ Scrapy является то, что он построен поверх Twisted, асинхронного сетевого фреймворка. Асинхронность означает, что не нужно ждать завершения запроса, прежде чем сделать еще один, это позволяет добиться высокого уровня производительности. Тот факт, что Scrapy реализован

с использованием неблокирующего (асинхронного) кода для параллелизма, делает его одним из самых эффективных фреймворков для краулинга.

В ходе работы по сбору данных возникло несколько проблем.

Первая из них — отсутствие архива статей на сайте СМИ Meduza[7]. Данная проблема была решена с помощью группы ВКонтакте данного интернет-издания. Ссылки на новостные статьи были собраны при помощи VK API — интерфейса, который позволяет получать информацию из базы данных vk.com. Для удобства работы с VK API была использована библиотека vk для Python. Эта библиотека предоставляет удобный интерфейс для работы с VK API, а также не требует авторизации.

Еще одна проблема возникла из-за того, что сайт телеканала Дождь блокирует запросы от краулера (программы, осуществляющей сбор данных), если эти запросы поступают от него слишком часто. Из-за этого краулер без дополнительных изменений будет скачивать лишь часть доступной информации. Выход из этой ситуации — использовать разные useragent для сетевых запросов. Useragent — это клиентское приложение, использующее определённый сетевой протокол. При посещении веб-сайта клиентское приложение обычно посылает веб-серверу информацию о себе. Это текстовая строка, являющаяся частью HTTP запроса, обычно включающая такую информацию, как название и версию приложения, операционную систему компьютера и язык. Библиотека fake_useragent для Python позволяет создавать случайные useragent для каждого запроса на сайт.

Для хранения полученных данных используется база данных на основе SQLite. Для работы с базой данных используется SQLAlchemy. Это набор инструментов SQL с открытым исходным кодом и ORM (технология программирования, которая связывает базы данных с концепциями объектно-ориентированных языков программирования)

для языка программирования Python.

1.2. Формирование траекторий

Для формирования траекторий необходимо уметь извлекать из текста имена упоминаемых в нем людей. Для этого необходимо использовать программу, решающую задачу извлечения именованных сущностей. Извлечение именованных сущностей — это класс подзадач извлечения информации, цель которой найти и классифицировать упоминания именованных сущностей в неструктурированном тексте по заранее определенным категориям, таким как имена людей, организации, адреса, даты и т. д. Библиотека DeepPavlov предоставляет модель, которая решает задачу извлечения имен для текстов на русском языке.

Далее необходимо сопоставить извлеченные имена из разных статей друг с другом чтобы определять, что разные статьи относятся к одному и тому же человеку. Для этого необходимо привести имя человека к нормальной форме. Для этой цели используется библиотека `rumorphy2`. Затем полученные слова сортируются в алфавитном порядке. Разделение на имена и фамилии не используется, так как задача определения фамилий работает недостаточно хорошо, особенно для иностранных фамилий. К тому же, благодаря деловому стилю написания новостных статей, при первом упоминании человека как правило используется его полное имя и фамилия. Более того, упоминание в новостной статье фамилии без имени обычно используется в устойчивых словосочетаниях, например "пакет Яровой". Таким образом появляется возможность в траекториях отделить упоминания самого человека от упоминания связанного с ним устойчивого словосочетания.

Для более точной работы будущей модели, учитываются имена, извлеченные из первых двух предложений и только при условии, что

в этом же предложении нет других упоминаний имен. Это практически гарантирует то, что главным фигурантом новостной статьи будет именно тот человек, имя которого извлечено.

В дальнейшем для работы модели используется текст нескольких первых предложений, если в них не упоминается какая-либо другая личность. Суммарная длина этих предложений не должна превышать 200 слов (в среднем это 10 предложений). Это обоснованно тем, что суть новости обычно укладывается в эти 200 слов. Последующий текст чаще всего является уточнением или справкой о каких-то событиях или организациях. Также к этим предложениям добавляются те, в которых упоминается извлеченная в начале личность, если такие предложения есть в дальнейшем тексте.

1.3. Обзор данных

****Единый формат, написать даты, дать оценку****

Итоговый датасет представляет собой набор новостных статей сгруппированных по траекториям. Представленные в датасете СМИ: Лента, Дождь, Медуза, РИА Новости. СМИ Лента соответствует 450 траекторий состоящих из 31645 статей в сумме, средняя длина траектории 69.4. Данные СМИ Дождь это 405 траектория из 14136 статей средней длины 34.9. СМИ Медуза соответствует 63 траектории из 1618 статей, средняя длина 25.7. РИА Новости - 189 траекторий, 37058 статей, средняя длина траектории 196.1 Всего 1113 траекторий, 84457 статей, средняя длина траектории 75.9. Максимальная длина траектории - 500 статей.

2. Модель тематических траекторий

В данной главе рассматривается модель, которая сопоставляет дескрипторам долгосрочны

1.1 Выделение интерпретируемых дескрипторов — тематического базиса для составления траекторий

2.2 Создание модели, способной сопоставлять траектории распределение на дескрипторах (веса дескрипторов), которое отображает долгосрочные тренды в текстах статей траектории а также их изменение

2.1. Базовая модель

Базовая модель создана на основе модели RMN, представленной в статье Feuding Families and Former Friends: Unsupervised Learning for Dynamic Fictional Relationships[1]

Перед RMN стоит две задачи. Первая задача — построить матрицу R размерности $K \times dim$, где K — задаваемое количество дескрипторов, dim — размерность векторного представления слов, состоящую из строк векторных представлений слов, которые и будут дескрипторами модели. Вторая — сопоставить каждой статье с номером t из траектории веса дескрипторов (распределение на дескрипторах), представленные вектором d_t размерности K .

вставить схему, чтобы было понятнее

2.1.1. Описание модели

Модель получает на вход векторное представление \hat{u}_t новостной статьи

$$\hat{u}_t = \frac{1}{|\hat{S}_t|} \sum_{w \in \hat{S}_t} u_w \quad (1)$$

где t — номер текущей статьи в траектории, S_t — множество всех слов траектории, $\hat{S}_t = \{w_i \in S_t \text{ with probability } p\}$ — множество слов текущей статьи, в которое каждое слово попадает с вероятностью p — параметр word dropout, u_w — векторное представление слова.

Word dropout используется для того, чтобы сделать модель более устойчивой. Также при использовании word dropout модель лучше выделяет долгосрочные тренды при использовании рекуррентного слоя.

Далее модель вычисляет вектор скрытого состояния

$$h_t = \text{ReLU}(W_h \cdot u_t) \quad (2)$$

где $\text{ReLU}(x) = \max(0, x)$ — функция активации, W_h — матрица весов слоя.

Затем с помощью рекуррентного слоя RNN с использованием результата предыдущего шага вычисляется текущее распределение на дескрипторах d_t

$$\begin{aligned} d_t &= \alpha \cdot \text{softmax}(W_d \cdot [h_t; d_{t-1}]) + (1 - \alpha) \cdot d_{t-1} & t > 1 \\ d_t &= \text{softmax}(W_d \cdot [h_t; d_{t-1}]) & t = 1 \end{aligned} \quad (3)$$

где $\alpha \in (0, 1]$ — параметр сглаживания, $\text{softmax}(X)_i = \frac{\exp(x_i)}{\sum_{x_j \in X} \exp(x_j)}$.

Таким образом сумма всех компонент дескриптора равна единице.

Далее модель вычисляет вектор-реконструктор

$$r_t = R^T d_t \quad (4)$$

Матрица R нормированная и обучаемая. Задача вектора-реконструктора — приближать начальный вектор $u_t = \frac{1}{|S_t|} \sum_{w \in S_t} u_w$ новостной статьи, который был вычислен без использования word dropout.

Таким образом матрица R состоит из строк, которые можно ин-

терпретировать как векторные представления слов. Эти слова можно выявить, найдя ближайшие к строкам матрицы векторные представления слов используя косинусное расстояние. Вектор d_t отражает на сколько каждое из этих слов описывает исходную новостную статью. За счет параметра α обеспечивается гладкость распределения на дескрипторах во времени

Функция ошибок

$$L(\Theta) = J(\Theta) + \lambda X(\Theta) \quad (5)$$

состоит из двух слагаемых.

Первое слагаемое

$$J(\Theta) = \sum_{n \in N} \max(0, 1 - r_t \cdot u_t + r_t \cdot u_n) \quad (6)$$

минимизирует косинусное расстояние между вектором новостной статьи и вектором-реконструктором. Также эта функция максимизирует это расстояние между вектором-реконструктором и случайно выбранными N векторами u_n (negative sampling). Таким образом модель принуждается обучать более уникальные дескрипторы для каждой статьи.

Второе слагаемое

$$X(\Theta) = \|RR^T - I\| \quad (7)$$

отвечает за ортогональность матрицы дескрипторов R . Это означает, что вектора дескрипторов должны быть как можно более отдаленными друг от друга в смысле косинусного расстояния. λ - параметр обучения, ((в базовой конфигурации равен $10e-4$)).

В базовой конфигурации используется **где-то написать про конфигурацию и параметры обучения**

2.1.2. Достоинства и недостатки

Достоинством модели является то, что она не только выделяет тематики траекторий во времени, но и находит дескрипторы, которые должным образом кластеризируют блоки траекторий. Более того, эта модель не нуждается в обучающем множестве.

Несмотря на свои достоинства, модель имеет недостатки, связанные с ее архитектурой.

Задача параметра α - сглаживать траекторию, то есть обеспечивать незначительное отклонение текущего распределения на дескрипторах от предыдущего. Это обеспечивает выделение в качестве дескрипторов более долгосрочных тематик. Данный способ кажется слишком "жестким". К тому же, есть все основания полагать, что функция ошибки будет меньше при α близком к 1, но в этом случае распределение на дескрипторах будет недостаточно гладким. Возникает желание сконструировать такую сеть, в которой уменьшение значения функции ошибок будет однозначно отражать то, что модель лучше справляется с поставленной задачей.

Рекуррентный слой является устаревшим и имеет ряд недостатков (можно сослаться на работу про проблемы rnn). К настоящему времени существуют архитектуры, которые решают проблемы rnn и работают сравнительно лучше.

2.2. Модернизация рекуррентного слоя

2.2.1. Идея

2.2.2. Эксперимент

2.3. Изменение функции ошибки

2.3.1. Идея

2.3.2. Эксперимент

Выводы

Выводы

Заклучение

Заклучение

Список литературы

- [1] Feuding Families and Former Friends: Unsupervised Learning for Dynamic Fictional Relationships / Mohit Iyyer, Anupam Guha, Snigdha Chaturvedi [и др.]
- [2] Amit Gruber, Yair Weiss, and Michal Rosen-Zvi. 2007. Hidden topic markov models. In Proceedings of Artificial Intelligence and Statistics
- [3] David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. / Latent dirichlet allocation. Journal of Machine Learning Research, 3
- [4] Jonathan Chang, Jordan Boyd-Graber, and David M Blei. 2009a. / Connections between the lines: augmenting social networks with text. In Knowledge Discovery and Data Mining
- [5] Lenta.ru [Электронный ресурс]: URL: <https://lenta.ru/> (дата обращения:).
- [6] Tvrain [Электронный ресурс]: URL: <https://tvrain.ru/> (дата обращения:).
- [7] Meduza [Электронный ресурс]: URL: <https://meduza.io/> (дата обращения:).
- [8] РИА Новости [Электронный ресурс]: URL: <https://ria.ru/> (дата обращения:).
- [9] scrapy [Электронный ресурс]: URL: <https://scrapy.org/> (дата обращения:).