

# Содержание

<b>1. Обзор</b>	<b>2</b>
1.1. Терминология . . . . .	2
1.2. Существующие подходы к решению задачи . . . . .	3
1.2.1. Определение схожести текстов на основе частотности употреб-	
ления слов . . . . .	4
1.2.2. Обобщённый метод, сопоставляющий новостной статье выска-	
зывания из социальных медиа . . . . .	5
1.2.3. Связывание твитов с новостями на основе словарей соответствий	5
1.2.4. Метод WTMF . . . . .	6
1.2.5. Метод WTMF-G . . . . .	8
1.3. Выбор подхода для решения задачи . . . . .	9
<b>2. Постановка задачи</b>	<b>11</b>
<b>3. Установления взаимосвязей между новостями и твитами</b>	<b>12</b>
3.1. Архитектура . . . . .	12
3.2. Обработка естественного языка . . . . .	15
3.3. Метод WTMF . . . . .	15
3.4. Метод WTMF-G . . . . .	16
3.5. Эффективная работа с матрицами . . . . .	17
<b>4. Эксперименты</b>	<b>19</b>
4.1. Методы оценки качества . . . . .	19
4.1.1. Метрика качества $MRR$ . . . . .	19
4.1.2. Метрика качества $TOP_I$ . . . . .	19
4.2. Оптимизация качества WTMF, путём варьирования параметров . . . .	20
4.3. Оптимизация качества WTMF-G, путём варьирования параметров . . .	23
4.4. Сравнительные результаты . . . . .	25

# 1. Обзор

Решение задачи по установлению взаимосвязи между твитами и новостными статьями в общем случае представляет собой решение задачи определения семантического сходства между короткими текстами. Методы естественной обработки языка не позволяют с высокой степенью точности определить семантическое сходство между короткими текстами, поэтому установление связей между твитами и новостями должно опираться на дополнительную информацию о предметной области.

## 1.1. Терминология

Решение задачи предполагает использование наработок различных дисциплин, таких как: обработка естественного языка, машинное обучение, информационный поиск, а основным источником данных выступают интернет ресурсы. Поэтому в работе используется специфичная терминология.

*Твиттер* — социальная сеть для публичного обмена короткими (до 140 символов) сообщениями при помощи веб-интерфейса, SMS, средств мгновенного обмена сообщениями или сторонних программ-клиентов.

*Твит* — термин сервиса микроблоггинга Твиттер, обозначающий сообщение, публикуемое пользователем в его твиттере. Особенностью твита является его длина, которая не может быть больше 140 знаков.

*Ретвит* — сообщение, целиком состоящее из цитирования сообщения одного пользователя Твиттера другим.

*Новость* — оперативное информационное сообщение, которое представляет политический, социальный или экономический интерес для аудитории в своей свежести, то есть сообщение о событиях произошедших недавно или происходящих в данный момент.

*Хэштег* — слово или фраза, которым предшествует символ #, используется в различных социальных сетях (Twitter, Facebook, Instagram) для объединения группы сообщений по теме или типу. Например: #искусство, #техника, #смешное, #анекдоты и т.д.

*URL* (от англ. Uniform Resource Locator — единый указатель ресурса) — единый образный определитель местонахождения ресурса. URL служит стандартизированным способом записи адреса ресурса в сети Интернет.

*Обработка естественного языка* (англ. Natural language processing) — направление математической лингвистики, которое изучает проблемы компьютерного анализа и синтеза естественных языков.

*Именованная сущность* — последовательность слов, являющаяся именем, названием, идентификатором, временным, денежным или процентным выражением.

*Аннотирование текста* — краткое представление содержания текста в виде аннотации (обзорного реферата).

*Информационный поиск* — процесс поиска неструктурированной документальной информации, удовлетворяющей информационные потребности, и наука об этом поиске.

*TF-IDF* (от англ. TF — term frequency, IDF — inverse document frequency) — статистическая мера, используемая для оценки важности слова в контексте документа, являющегося частью коллекции документов или корпуса. Вес некоторого слова пропорционален количеству употребления этого слова в документе, и обратно пропорционален частоте употребления слова в других документах коллекции.

*TF-IDF матрица* — матрица, строки которой соответствуют словам из корпуса, а столбцы текстам. Значение ячейки матрицы  $(i, j)$  равно значению метрики tf-idf для слова, соответствующего строчке  $i$ , и текста, соответствующего столбцу  $j$ .

*WTMF* — метод машинного обучения, применяемый для анализа схожести между короткими текстами [?].

*Тематическое моделирование* — способ построения модели коллекции текстовых документов, которая определяет, к каким темам относится каждый из документов.

*LDA* (от англ. Latent Dirichlet allocation — Латентное размещение Дирихле) — методов тематического моделирования, позволяющий объяснять результаты наблюдений с помощью неявных групп.

## 1.2. Существующие подходы к решению задачи

Задача автоматического установления связей между твитами и новостными статьями до сих пор не имеет устоявшегося решения. В рамках предварительного исследования были отобраны наиболее перспективные подходы к решению задачи, а именно:

- метод WTMT-G, представляющий собой доработку метода WTMF, которая позволила учитывать информацию о связях между текстами;
- обобщённый метод, позволяющий по новости находить относящиеся к ней высказывания из социальных медиа;
- связывание твитов с новостями на основе словарей соответствий;

Также рассматривается классическое решение задачи определения схожести текстов на основе частотности употребления слов. Ниже представлен краткий обзор выбранных методов.

Стоит также ввести несколько определений употребляемых в дальнейшем: под связью *текст-текст* подразумевается определение двух текстов как схожих на основе некоторой дополнительной информации; под связью *текст-слово* подразумевается определение двух текстов как схожих только на основе слов, из которых состоит текст.

### 1.2.1. Определение схожести текстов на основе частотности употребления слов

Наиболее простым и очевидным подходом к решению задачи связывания твитов с новостными статьями, является связывание текстов, наиболее близких по частотности употребления слов. Способ основан на сравнении значений метрики TF-IDF, поэтому в дальнейшем будем называть этот способ TF-IDF методом.

Решение задачи связывания твитов с новостными статьями на основе частотности употребления слов можно представить в виде небольшого алгоритма:

1. объединить тексты всех твитов и тексты всех новостей — (для новости текст это конкатенация заголовка и краткого изложения);
2. в качестве корпуса использовать объединение начальных форм всех слов, используемых в текстах, за вычетом списка стоп-слов (под списком стоп-слов подразумевается набор часто употребляемых слов языка, которые вне контекста не несут смысловой нагрузки, к примеру, предлоги);
3. по множеству текстов и построенному корпусу построить TF-IDF матрицу;
4. каждому тексту сопоставить столбец TF-IDF матрицы, соответствующий тексту (вектор для сравнения);
5. рассматривая вектор для сравнения в качестве координат в метрическом пространстве, для каждого твита найти список наиболее похожих на него новостей.

В работе в качестве меры близости в метрическом пространстве используется косинусная мера близости — мера численно равная косинусу угла между векторами. В дальнейшем каждый раз, когда говорится о схожести или близости двух векторов, подразумевается близость согласно косинусной мере.

### 1.2.2. Обобщённый метод, сопоставляющий новостной статье высказывания из социальных медиа

В рамках метода решается следующая задача: по новости необходимо найти высказывания в социальных сетях, которые на неё неявно ссылаются. Метод был предложен в статье Linking Online News and Social Media [?]. Поставленная задача решается в три этапа:

1. по заданной новостной статье формируется несколько моделей запросов, которые создаются как на основе структуры статьи, так и на основе явно связанных со статьёй высказываний из социальных медиа.
2. построенные модели используются для получения высказываний из индекса целевого социального медиа, результатом являются несколько ранжированных списков;
3. полученные списки объединяются с использованием особой техники слияния данных.

Авторы также предлагают способ, созданный для борьбы с дрейфом запроса (порождение менее подходящего запроса), который возникает при большом объёме используемого текста. Способ основан на выборе дополнительных отличительных условий.

Для экспериментальной оценки используются данные из различных медиа, таких как Twitter, Digg, Delicious, the New York Times Community, Wikipedia, а также из блогов.

В результате работы показано, что модели запросов, основанные на различных источниках данных, повышают точность выявления высказываний из социальных медиа; методы слияния ранжированных списков приводят к значительному повышению производительности в сравнении с другими подходами.

### 1.2.3. Связывание твитов с новостями на основе словарей соответствий

Метод связывания твитов с новостными статьями, основанный на словарях соответствий, предложен в статье Bridging Vocabularies to Link Tweets and News [?]. *Словарь соответствий* — множество слов, которые встречаются только в твитах и, соответственно, не встречаются в новостях. Авторы предложили способ автоматического установление связи между множеством твитов и множеством новостей определённой темы. Темы извлекаются из новостей на основе методов тематического моделирования.

Значительную сложность при решении проблемы связывания твитов с новостями вызывают малый размер твита и различия в словарях: в твитах используются аббревиатуры, неформальный язык, сленг; в новостях, напротив, используется литературный язык. В частности, твиты могут вообще не нести смысловой нагрузки.

Твиттер предлагает хэштеги, как механизм для категоризации твитов. Но этот подход обладает рядом недостатков, таких как: не все записи содержат хэштеги, хэштег не содержит информацию о событии, хэштег сформулирован в слишком общей форме, твит содержит несколько хэштегов. Следовательно использование одних хэштегов приведёт к низкому качеству связывания твитов с новостями.

Для решения задачи и преодоления описанных выше проблем, авторами работы предлагается следующий подход:

1. С помощью метода LDA из множества новостей извлекается набор тем. Тема характеризуется распределением частот слов, характерных для этой темы.
2. Каждой полученной теме сопоставляется множество наиболее близких к ней твитов.
3. Из полученных твитов извлекаются слова, которые дополняют характеристику рассматриваемой темы.
4. Полученные слова образуют словарь соответствий и служат «мостом» к другим твитам.

В результате работы продемонстрирован способ установления связей между множеством твитов и множеством новостей с использованием словарей соответствий.

#### 1.2.4. Метод WTMF

Метод WTMF предназначен для определения семантической близости коротких текстов. Этот метод учитывает отсутствующие в тексте слова в виде признаков короткого текста. Под отсутствующими словами подразумеваются все слова из корпуса, составленного из всех текстов, за исключением слов из рассматриваемого короткого текста, то есть отсутствующие слова можно трактовать как негативный сигнал.

Работа метода WTMF основана на разложении TF-IDF матрицы  $X$  в произведение двух матриц  $P$  и  $Q$ :

$$X \sim P^T Q.$$

На рисунке 1 показано как матрица  $X$  приближается произведением двух матриц  $P^T$  размера  $M \times K$  и  $Q$  размера  $K \times N$ .

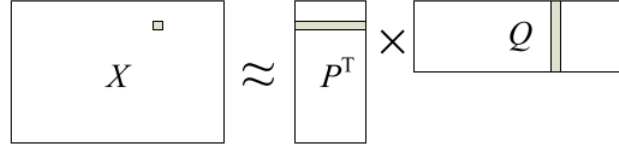


Рисунок 1 — Разложение TF-IDF матрицы ( $X$ ) на произведение матриц  $P$  и  $Q$

Каждый текст  $s_j$  представлен в виде вектора  $Q_{\cdot,j}$  размерности  $K$ , каждое слово  $w_i$  представлено в виде вектор  $P_{i,\cdot}$ . Если  $X_{ij} = (P_{i,\cdot}, Q_{\cdot,j})$  близко к нулю, то это трактуется как отсутствующее слово.

Задачей метода является минимизация целевой функции ( $\lambda$  - регуляризирующий член, матрица  $W$  определяет вес элементов матрицы  $X$ ):

$$\sum_i \sum_j W_{ij} (P_{i,\cdot} \cdot Q_{\cdot,j} - X_{ij})^2 + \lambda \|P\|_2^2 + \lambda \|Q\|_2^2.$$

Для получения векторов  $P_{i,\cdot}$  и  $Q_{\cdot,j}$  используется алгоритм описанный в статье [?]. Сначала  $P$  и  $Q$  инициализируются случайными числами. Затем запускается итеративный пересчёт  $P$  и  $Q$  по следующим формулам (эффективный способ расчёта описан в [?]):

$$P_{i,\cdot} = (QW'_iQ^T + \lambda I)^{-1}QW'_iX_{i,\cdot}^T,$$

$$Q_{\cdot,j} = (PW'_jP^T + \lambda I)^{-1}PW'_jX_{\cdot,j}.$$

Здесь  $W'_i = \text{diag}(W_{i,\cdot})$  - диагональная матрица полученная из  $i$ -ой строки матрицы  $W$ , аналогично  $W'_j = \text{diag}(W_{\cdot,j})$  - диагональная матрица полученная из  $j$ -ого столбца матрицы  $W$ . Матрица  $W$  определяется следующим образом:

$$W_{ij} = \begin{cases} 1, & \text{if } X_{ij} \neq 0, \\ w_m, & \text{otherwise.} \end{cases},$$

где  $w_m$  положительно и  $w_m \ll 1$ .

Столбцы построенной матрицы  $Q$  представляют собой вектора для сравнения текстов между собой. Тексту, на основе которого построена  $i$ -я строка TF-IDF матрицы  $X$ , в соответствие ставится  $i$ -й столбец матрицы  $Q$ .

В статье предложен подход для поиска семантической близости текстов, который на небольших текстах работает лучше чем подход, основанный на частотности слов.

### 1.2.5. Метод WTMF-G

Метод WTMF-G решает задачу установления связей между твитами и новостными статьями, путём построения модели, которая учитывает неявные связи между текстами. Метод был предложен в статье *Linking Tweets to News: A Framework to Enrich Short Text Data in Social Media* [?].

Метод WTMF-G (WTMF on Graphs) представляет собой доработанный метод WTMF, позволяющий хорошо моделировать семантику коротких текстов, но не учитывающий некоторые специфичные для твитов и новостей характеристики, которыми обладает исходная выборка и которые взаимосвязаны с семантической близостью текстов:

1. хештеги, которые являются прямым указанием на смысл твита;
2. именованные сущности, которые с высокой точностью можно извлекать из новостей;
3. информацию о времени публикации твитов и новостей.

Метод WTMF-G расширяет возможности метода WTMF, путём учёта взаимосвязи текстов на основе специфичных для твитов и новостных статей характеристик, то есть позволяет учесть информацию о взаимосвязи текст-текст.

Для решения задачи необходимо иметь эталонный набор данных, на котором будет производиться оценка качества полученного решения. Сначала за общий период времени собираются твиты и новости. Для твита помимо текста хранится информация о времени публикации и авторе работы. Для новости хранится время публикации, заголовки, краткое изложение и URL.

На основе собранной информации строится набор данных, который состоит из трёх частей:

1. множество новостей — все собранные новости;
2. множество связей твит-новость, под связью подразумевается явное указание URL новости в тексте твита;
3. множество твитов — все твиты, имеющие связь с одной из собранных новостей.

К построенному набору данных применяется метод WTMF-G — то есть метод WTMF, расширенный путём добавления связей текст-текст. Добавление связей



текст-текст происходит путём модификации регуляризующего члена  $lambda$ . Для каждой пары связанных текстов  $j_1$  и  $j_2$ :

$$\lambda = \delta \cdot \left( \frac{Q_{\cdot,j_1} \cdot Q_{\cdot,j_2}}{|Q_{\cdot,j_1}| |Q_{\cdot,j_2}|} - 1 \right)^2,$$

коэффициент  $\delta$  задаёт степень влияния связей текст-текст.

Так как новый регуляризующий член  $lambda$  зависит от  $|Q_{\cdot,j}|$ , который меняется во время итерации, вводим упрощение: длина вектора  $Q_{\cdot,j}$  не изменяется во время итерации. Также необходимо модифицировать итеративный процесс построения матриц  $P$  и  $Q$  следующим образом:

$$P_{i,\cdot} = (QW_i'Q^T + \lambda I)^{-1}QW_i'X_{i,\cdot}^T,$$

$$Q_{\cdot,j} = (PW_j'P^T + \lambda I + \delta L_j^2 Q_{\cdot,n(j)} \text{diag}(L_{n(j)}^2) Q_{\cdot,n(j)}^T)^{-1} (PW_j'X_{j,\cdot} + \delta L_j Q_{\cdot,n(j)} L_{n(j)}).$$

В этих формулах  $n(j)$  — список связанных текстов с текстом  $j$ .  $Q_{\cdot,n(j)}$  — матрица, состоящая из связанных векторов для  $Q_{\cdot,j}$ .  $L_j$  — длина вектора  $Q_j$  на начало итерации,  $L_n(j)$  — вектор длин векторов связанных с  $j$ , то есть  $Q_{\cdot,n(j)}$ , полученный на начало итерации.

В статье показано, что добавление информации о взаимосвязи текст-текст позволяет повысить качество установления связей между твитами и новостными статьями. Качество метода WTMF-G, измеренное с использованием метрики MRR (метрика описана в главе 4.1.1), в сравнении с такими популярными подходами как: TF-IDF, LDA, WTMF, показано в таблице 1.

Таблица 1: Значение метрики MRR для алгоритма WTMF-G в сравнении с другими подходами.

Алгоритм	TF-IDF	LDA	WTMF	WTMF-G
Значение MRR	0.4602	0.1313	0.4531	0.4791

В таблице 1 показано, что алгоритм WTMF-G даёт лучшее качество, чем прочие подходы.

### 1.3. Выбор подхода для решения задачи

В качестве основного подхода, на основе которого строится решение задачи по установлению связей между твитами и новостями, был выбран WTMF-G. Основной причиной подобного выбора является то, что большинство подходов учитывают

только статистические зависимости вида текст-слово; метод WTMF-G, напротив, не ограничивается зависимостями текст-слово, а позволяет учесть взаимосвязь текст-текст, что, как ожидается, даст прирост качества в решении задачи установления связей.

Также, в рамках работы задача по установления связей между твитами и новостями решена классическим подходом для установления связей между текстами — определение схожести текстов на основе частотности употребления слов. Этот подход даёт хорошие результаты на больших текстах. Результаты этого метода помогут оценить влияние связей вида текст-текст в методе WTMF-G на качество полученного решения.

## 2. Постановка задачи

В разделе 1 проведено исследование существующих методов автоматического установления связей между сообщениями твиттера и новостными статьями, выбраны методы, на основе которых необходимо реализовать программный комплекс, позволяющий устанавливать связи между твитами и новостными статьями.

Для установления связей должен быть собран эталонный набор данных, которой состоит из множества твитов, новостей и связей между ними. Выбранный алгоритм WTMF-G накладывает ограничение на формат эталонного набора: для каждого твита существует связь с единственной новостью.

Решение задачи установления связей между твитами и новостными статьями в общем случае неоднозначно: как твиту может соответствовать несколько новостей, так и новостной статье может соответствовать несколько твитов. Отталкиваясь от существующего ограничения: твит связан с единственной новостью, получаем что для оценки качества установления связей хорошо подходят метрики, принятые в информационном поиске. Для использования подобных метрик будем рассматривать твит как запрос, в терминологии информационного поиска, а список новостей как ответ нашей системы установления связей. То есть для каждого твита мы получаем список новостей, ранжированный по мере убывания их схожести, в дальнейшем будем называть подобный список рекомендацией, а процесс установления связей построением рекомендаций.

Целью работы является создание программного комплекса, который реализует такие методы машинного обучения как WTMF, WTMF-G, TF-IDF, позволяет для произвольного твита построить рекомендацию новостей с использованием любого из предложенных методов машинного обучения, а также способен оценить качество используемого метода машинного обучения.

### 3. Установления взаимосвязей между новостями и твитами

Задача автоматического установления связей между твитами и новостями решена посредством написания программного комплекса, который обладает следующими возможностями:

1. сбор необходимой для решения задачи информации;
2. построение наборов данных;
3. применение к наборам данных методов машинного обучения;
4. получение рекомендаций новостей для произвольных твитов;
5. вариативность в выборе метода для построения рекомендаций;
6. возможность получить информацию о качестве используемого метода.

Программный комплекс реализован с использованием языка программирования Python версии 2.7.

Ниже приводится описание архитектуры программного комплекса, а также разбор отдельных моментов.

#### 3.1. Архитектура

Программный комплекс состоит из набора подсистем, которые выполняют следующий набор функций:

1. получение данных из твиттера;
2. получение данных из новостной rss-ленты;
3. расшифровка коротких URL;
4. автоматическое построение набора данных;
5. построение набора данных на основе вручную построенных заготовок;
6. построение моделей для методов WTMF и WTMF-G;
7. построение рекомендаций для методов WTMF, WTMF-G и поиска схожести на основе частотности употребления слов (TF-IDF);

8. оценка качества рекомендаций;
9. получение результатов рекомендаций в пригодном для чтения формате;

подробное описание архитектуры системы приведённой на flowchart

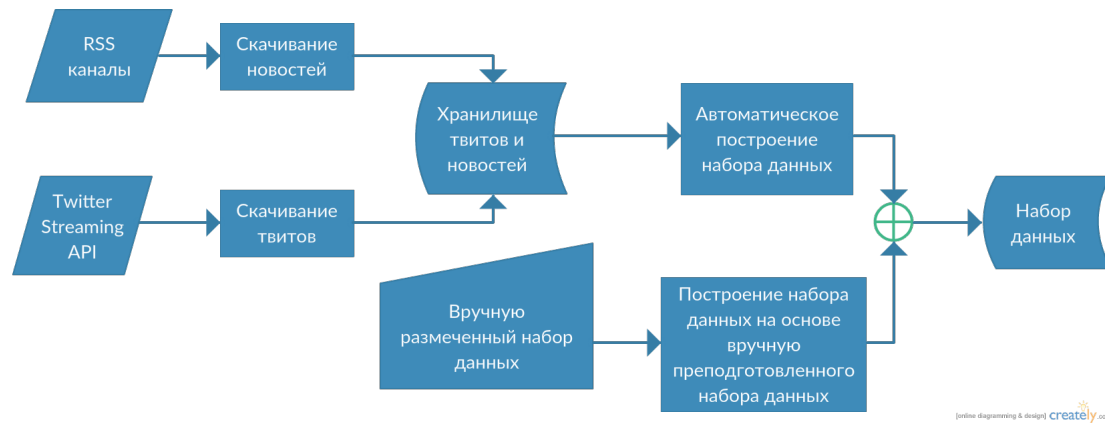


Рисунок 2 — flow chart

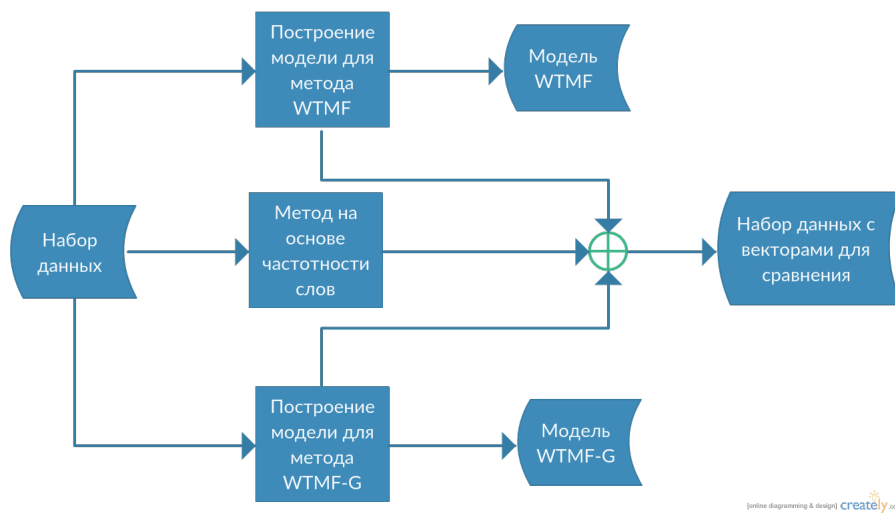


Рисунок 3 — flow chart

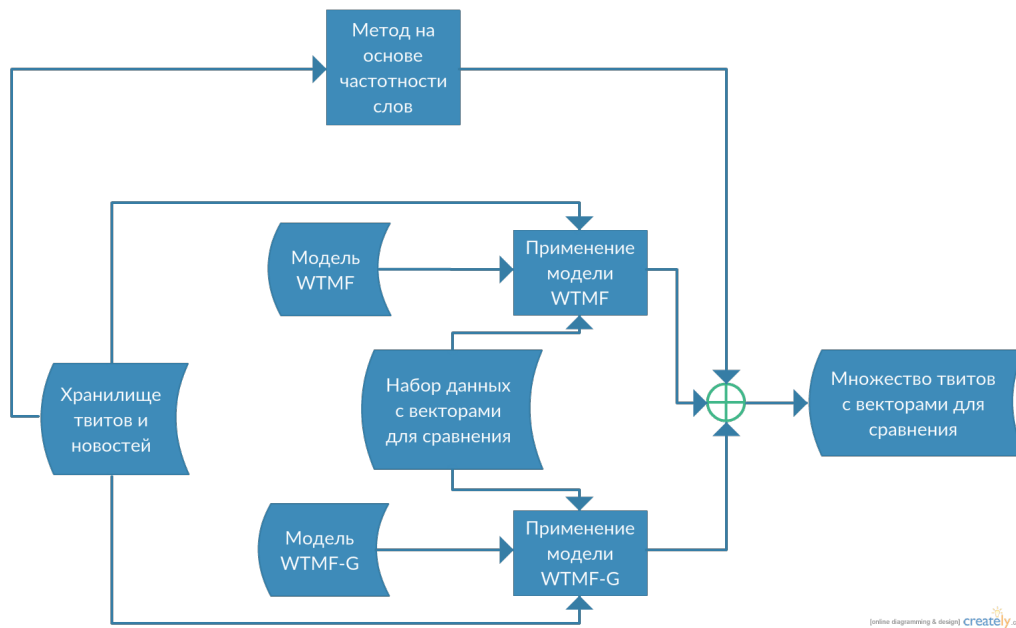


Рисунок 4 — flow chart

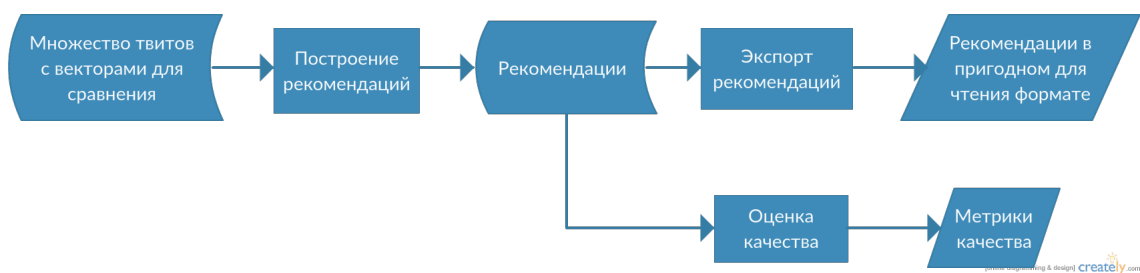


Рисунок 5 — flow chart

### 3.2. Обработка естественного языка

Работа посвящена поиску семантической близости текстов, поэтому в ней имеет место использование решений таких задач обработки естественного языка, как:

1. токенизация — разбиение предложения на слова;
2. лемматизация — процесс приведения словоформы к лемме;
3. извлечение именованных сущностей.

Описанные выше задачи решены с использованием набора сторонних библиотек для языка Python, а именно:

1. nltk — платформа, для написания приложений на языке Python, обрабатывающих естественный язык;
2. pymorphy2 — морфологический анализатор;
3. polyglot — библиотека, позволяющая извлекать именованные сущности из текстов на разных языках.

Для решения задачи токенизации используется стандартный токенизатор, реализованный в nltk. Задача лемматизации решается в случае русского языка с помощью морфологического анализатора pymorphy2, в случае английского языка с помощью морфологического анализатора WordNet, реализованного в nltk.

Извлечение именованных сущностей происходит с помощью библиотеки polyglot. В используемой библиотеке реализуется выявление именованных сущностей на основе заранее сформированного и размеченного корпуса именованных сущностей. Корпус формируется на основе данных из Википедии.

### 3.3. Метод WTMF

Модель для метода WTMF построена на основе заранее подготовленного набора данных. В контексте работы набор данных состоит из множества новостей и твитов, из которых в процессе работы извлекается набор текстов (для твита — текст твита, для новости — конкатенация заголовка и краткого изложения статьи).

По множеству текстов, которые получены из набора данных, построена модель, пригодная для сериализации, состоящая из матрицы  $P$  (здесь и далее используются обозначения введенные в главе 1.2.4). Построение модели зависит от четырёх констант:

1.  $K$  — размерность вектора, по которому производится сравнение (если TF-IDF матрица  $X$  была размера  $M \times N$ , то по завершении работы алгоритма будут получены две матрицы  $P$  размера  $K \times M$  и  $Q$  размера  $K \times N$ );
2.  $I$  — число итераций алгоритма построения модели;
3.  $w_M$  — коэффициент, задающий вес негативного сигнала при построении матрицы весов  $W$ ;
4.  $\lambda$  — регуляризующий член.

Применение полученной модели на множество твитов представляет собой следующий процесс: сначала строится TF-IDF матрица  $X$  для новостей из набора данных и множества твитов, затем на основе новой матрицы  $X$  строится весовая матрица  $W$ , и наконец на основе построенных матриц  $X$  и  $W$  и посчитанной на этапе обучения матрицы  $P$  выполняется половина итерации алгоритма обучения, а именно получение матрицы  $Q$  по матрице  $P$ :

$$Q_{:,j} = (PW_j'P^T + \lambda I)^{-1}PW_j'X_{j,:}$$

В результате получаем вектора для сравнения твитов из заданного множества.

### 3.4. Метод WTMF-G

Построение модели для метода WTMF-G основывается на построение модели метода WTMF. Набор данных состоит из множества новостей и твитов и связей вида текст-текст, из которых, в процессе работы извлекается набор текстов. (для твита — текст твита, для новости — конкатенация заголовка и краткого изложения статьи).

По множеству текстов, которые получены из набора данных, построена пригодная для сериализации модель, представляющая собой матрицу  $P$ . Построение модели зависит от четырёх констант:

1.  $K$  — размерность вектора, по которому производится сравнение (если TF-IDF матрица  $X$  была размера  $M \times N$ , то по завершении работы алгоритма будут получены две матрицы  $P$  размера  $K \times M$  и  $Q$  размера  $K \times N$ );
2.  $I$  — число итераций алгоритма построения модели;
3.  $w_M$  — коэффициент, задающий вес негативного сигнала при построении матрицы весов  $W$ ;



4.  $\delta$  — коэффициент, задающий степень влияния связей вида текст-текст.

Применение полученной модели на множество твитов производится аналогично применению модели для метода WTMF за исключением двух моментов: во-первых, необходимо на основе новостей из набора данных и множества твитов пере-строить связи текст-текст, во-вторых получение матрицы  $Q$  происходит по следующей формуле:

$$Q_{:,j} = (PW_j'P^T + \lambda I + \delta L_j^2 Q_{:,n(j)} \text{diag}(L_{n(j)}^2) Q_{:,n(j)}^T)^{-1} (PW_j'X_{j,\cdot} + \delta L_j Q_{:,n(j)} L_{n(j)}).$$

В результате получаем вектора для сравнения твитов из заданного множества.

### 3.5. Эффективная работа с матрицами

Построение и применение моделей WTMF и WTMF-G требует большого количества операций над матрицами, что на практике занимает продолжительное время. Поэтому актуальна задача по повышению эффективности работы с матрицами.

Для эффективной работы с матрицами используются программные библиотеки для языка Python `numru` и `scipy` (базируется на библиотеке `numru` и расширяет её функционал).

Повышение производительности при работе с матрицами производится на примере оптимизации времени расчёта формулы получения строк матрицы  $P$ , которая используется при построении моделей WTMF и WTMF-G. На каждой итерации построения модели происходит многократное выполнение формулы (число выполнений порядка  $10^4$ , зависит от размера корпуса):

$$P_{i,\cdot} = (QW_i'Q^T + \lambda I)^{-1} QW_i'X_{i,\cdot}^T.$$

В начале была написана наивная реализация алгоритма, которая показала производительность, не приемлемую в рамках решения задачи. Затем наивная реализация оптимизировалась следующим образом:

1. переход к перемножению матриц с использованием высокопроизводительной библиотеки для языка C `OpenBlass` (в библиотеке `numru` существует возможность перейти к использованию для работы с матрицами некоторых библиотек, написанных на языке C [?]);
2. сохранение в отдельной переменной переиспользуемых результатов вычислений над матрицами;

3. переписывание кода для работы с разреженными матрицами;
4. удаление лишних приведений матриц к формату python list и обратно.

Результаты оптимизации приведены в таблице 2.

Таблица 2: Оптимизация работы с матрицами

Добавленная оптимизация	Время за 100 итераций (с)	Прирост производительности (раз)
Наивная реализация	205	1
Перемножение с помощью OpenBlass	55	3.73
Переиспользование результатов	15.15	3.63
Работа с разреженными матрицами	0.75	20.2
Сокращение количества приведений типов	0.63	1.21

Получили, что оптимизированное решение работает в 325 раз быстрее наивной реализации. Дальнейшая оптимизация не производилась, так как получено решение работающее за приемлемое время.

## 4. Эксперименты

Главное целью проведения экспериментов является сравнение двух реализованных методов автоматического установления связей между твитами и новостными статьями: метод основанный на частотности употребления слов и WTMF-G. Для исследования влияния на качество добавления информации о взаимосвязях вида текст-текст также производится сравнительное тестирование методов WTMF и WTMF-G.

Ввиду малого числа твитов в наборах данных тестирование производится на тех же выборках, на которых производится обучение.

### 4.1. Методы оценки качества

Решение задачи установления связей между твитами и новостными статьями неоднозначно. Как твиту может соответствовать несколько новостей, так и новостной статье может соответствовать несколько твитов. Но в эталонном наборе данных для каждого твита существует связь с единственной новостью. В данном случае для оценки качества применимы метрики принятые в информационном поиске.

Мы рассматривает твит как запрос в терминологии информационного поиска, а список новостей как ответом. То есть для каждого твита мы получаем список новостей, ранжированный по мере убывания их схожести.

#### 4.1.1. Метрика качества $MRR$

$MRR$  (от англ. Mean reciprocal rank) — статистическая метрика, используемая для измерения качества алгоритмов информационного поиска. Пусть  $rank_i$  — позиция первого правильного ответа в  $i$ -м запросе,  $n$  — общее количество запросов. Тогда значение  $MRR$  можно получить по формуле:

$$MRR = \frac{1}{n} \sum_{i=1}^n \frac{1}{rank_i}.$$

#### 4.1.2. Метрика качества $TOP_I$

$TOP_I$  — группа метрик, используемых для оценки качества алгоритмов информационного поиска. Значение метрики  $TOP_I$  численно равно проценту запросов с правильным ответом, входящим в первые  $I$  ответов. Пусть  $n$  — общее количество запросов,  $Q_I(i)$  — равно 1, если правильный ответ на  $i$ -й запрос входит в первые  $I$  предложенных ответов, 0 — в противном случае. Тогда значение  $TOP_I$  можно полу-

читать по формуле:

$$TOP_I = \frac{1}{n} \sum_{i=1}^n Q_I(i).$$

В дальнейшем будут рассматриваться следующие три метрики из группы метрик  $TOP_I$ :  $TOP_1$ ,  $TOP_3$ ,  $TOP_{10}$ .

## 4.2. Оптимизация качества WTMF, путём варьирования параметров

Оптимизация параметров модели для метода WTMF будет производиться на наборе данных cutted, используя метрику MRR. Модель WTMF зависит от четырёх параметров:  $K$ ,  $I$ ,  $\lambda$ ,  $w_m$ . Параметры  $K$  и  $I$  влияют на время построения модели, а параметры  $\lambda$  и  $w_m$  не влияют на время построения модели.

В качестве начального приближения берутся значения параметров, которое использовали авторы работы [?], а именно:  $K = 30$ ,  $I = 3$ ,  $\lambda = 20$ ,  $w_m = 0.1$ .

Оптимизируются параметры, не влияющие на время работы алгоритма:  $\lambda$  и  $w_m$ . Для этого фиксируются остальные параметры:  $I = 1$ ,  $K = 30$ . Для начала находится оптимальный порядок значений начального приближения. Результаты занесены в таблицу 3.

Таблица 3: Качество работы алгоритма WTMF для различных значений  $\lambda$  и  $w_m$  при фиксированных значениях  $I = 1$ ,  $K = 30$ .

$\lambda \backslash w_m$	<b>0.001</b>	<b>0.01</b>	<b>0.1</b>	<b>1</b>	<b>10</b>	<b>100</b>
<b>0.2</b>	0.6855	0.6877	0.7482	0.3651	0.1526	0.1485
<b>2</b>	0.7000	0.7015	0.7173	0.7525	0.3707	0.1605
<b>20</b>	0.6964	0.7081	0.7149	0.7308	0.7507	0.3784
<b>200</b>	0.7075	0.6991	0.7010	0.7016	0.7146	0.7448
<b>2000</b>	0.6970	0.7070	0.6991	0.7114	0.6994	0.7044

Как видно из таблицы 3 в целом получена достаточно однородная картина для всех порядков  $\lambda$  и  $w_m$ . Заметное снижение качества происходит при большом порядке  $w_m$  и малом порядке  $\lambda$ . Максимальное значение метрики достигнуто при  $\lambda = 2$  и  $w_m = 1$ . Для уточнения значения коэффициентов, производится исследование качества работы алгоритма в окрестностях максимального значения метрики. Результаты приведены в таблице 4.

Из таблицы 4 получаем оптимальные значения коэффициентов  $\lambda = 0.95$  и  $w_m = 1.95$ .

Таблица 4: Качество работы алгоритма WMTF для различных значений  $\lambda$  и  $w_m$  при фиксированных значениях  $I = 1$ ,  $K = 30$ .

$\lambda \backslash w_m$	<b>0.9</b>	<b>0.95</b>	<b>1</b>	<b>1.1</b>	<b>1.2</b>
<b>1.9</b>	0.7442	0.7451	0.7536	0.7542	0.7544
<b>1.95</b>	0.7447	0.7554	0.7452	0.7439	0.7504
<b>2</b>	0.7507	0.7528	0.7504	0.7515	0.7566
<b>2.05</b>	0.7413	0.7505	0.7424	0.7525	0.7479
<b>2.1</b>	0.7405	0.7484	0.7485	0.7502	0.7501

Оптимизируются параметры, влияющие на время работы алгоритма:  $K$  и  $I$ . Для этого фиксируются остальные параметры:  $\lambda = 0.95$ ,  $w_m = 1.95$ . Для начала находится примерное значение коэффициента  $K$  и оптимальное значение  $I$ . Результаты занесены в таблицу 5.

Таблица 5: Качество работы алгоритма WMTF для различных значений  $K$  и  $I$  при фиксированных значениях  $\lambda = 0.95$ ,  $w_m = 1.95$ .

$K \backslash I$	<b>1</b>	<b>2</b>	<b>3</b>
<b>5</b>	0.1232	0.1593	0.1838
<b>10</b>	0.3521	0.4102	0.4437
<b>30</b>	0.7426	0.7422	0.7158
<b>60</b>	0.8326	0.8117	0.7620

Как видно из таблицы 5 увеличение  $K$  приводит к значительному улучшению качества работы алгоритма, увеличении  $I$  приводит к улучшению качества алгоритма только при малых значениях параметра  $K$ , при больших значениях  $K$  увеличение параметра  $I$  приводит к ухудшению качества. Максимальное значение метрики достигнуто при  $K = 60$  и  $I = 1$ . Для уточнения значения коэффициента  $K$ , производится исследование качества работы алгоритма при фиксированном значении коэффициента  $I$ . Результаты приведены в таблице 6.

Из таблицы 6 получаем оптимальные значения коэффициента  $K = 90$

В итоге оптимизации качества рекомендаций на основе алгоритма WMTF были получены оптимальные параметры:  $K = 90$ ,  $I = 1$ ,  $\lambda = 0.95$ ,  $w_m = 1.95$ .

Таблица 6: Качество работы алгоритма WMTF для различных значений  $K$  при фиксированных значениях  $I = 1$ ,  $\lambda = 0.95$ ,  $w_m = 1.95$ .

<b>K</b>	<b>Значение метрики RR</b>
<b>10</b>	0.3595
<b>20</b>	0.6460
<b>30</b>	0.7496
<b>40</b>	0.8003
<b>50</b>	0.8220
<b>60</b>	0.8424
<b>70</b>	0.8472
<b>80</b>	0.8535
<b>82</b>	0.8549
<b>84</b>	0.8597
<b>86</b>	0.8592
<b>88</b>	0.8572
<b>90</b>	0.8675
<b>92</b>	0.8580
<b>94</b>	0.8604
<b>96</b>	0.8612
<b>98</b>	0.8644
<b>100</b>	0.8655
<b>110</b>	0.8627

### 4.3. Оптимизация качества WTMF-G, путём варьирования параметров

Датасет cutted\_0.0

Начальное приближение  $\delta = 0.1$ ,  $w_m = 1.95$ ,  $K = 90$ ,  $I = 1$ ,  $\lambda = 0.95$ .

Таблица 7: Качество работы алгоритма WTMF-G для различных значений  $\lambda$  и  $\delta$  при фиксированных значениях  $I = 1$ ,  $K = 30$ ,  $w_m = 1.95$ .

$\lambda \backslash \delta$	<b>0.001</b>	<b>0.01</b>	<b>0.1</b>	<b>1</b>	<b>10</b>
<b>0.01</b>	0.3889	0.3842	0.3924	0.3900	0.3895
<b>0.1</b>	0.4895	0.4875	0.4886	0.4850	0.4847
<b>1</b>	0.8227	0.8256	0.8242	0.8225	0.8212
<b>10</b>	0.8477	0.8440	0.8496	0.8454	0.8495
<b>100</b>	0.8294	0.8318	0.8283	0.8240	0.8243

максимум при  $\lambda = 10$ ,  $\delta = 0.1$  рассмотрим окрестности.

Таблица 8: Качество работы алгоритма WTMF-G для различных значений  $\lambda$  и  $\delta$  при фиксированных значениях  $I = 1$ ,  $K = 30$ ,  $w_m = 1.95$ .

$\lambda \backslash \delta$	<b>0.06</b>	<b>0.08</b>	<b>0.1</b>	<b>0.12</b>	<b>0.14</b>
<b>6</b>	0.8589	0.8524	0.8511	0.8580	0.8493
<b>8</b>	0.8483	0.8528	0.8539	0.8439	0.8498
<b>10</b>	0.8504	0.8455	0.8416	0.8453	0.8408
<b>12</b>	0.8453	0.8398	0.8472	0.8376	0.8415
<b>14</b>	0.8462	0.8456	0.8387	0.8398	0.8377

рассмотрим максимум при  $\lambda = 6$ ,  $\delta = 0.06$

варьируем  $w_m$

0.01 0.8283 0.05 0.8296 0.1 0.8285 0.5 0.8359 1 0.8442 5 0.8639 10 0.8391 50 0.6094  
100 0.5035  
1.5 0.8474 2.0 0.8507 2.5 0.8563 3.0 0.8592 3.5 0.8585 4.0 0.8594 4.5 0.8603 5.0  
0.8597 5.5 0.8591 6.0 0.8586 6.5 0.8574 7.5 0.8536

берём  $w_m=5$ , варьируем  $K/I$

Оптимизация параметров ещё не завершена, существующая и очень, очень грубая оценка приведена ниже

Оптимизация параметров модели для метода WTMF-G будет производиться на наборе данных auto\_cleared, используя метрику MRR. Модель WTMF зависит от четырёх параметров:  $K$ ,  $I$ ,  $\delta$ ,  $w_m$ . Параметры  $K$  и  $I$  влияют на время построения модели, а параметры  $\lambda$  и  $w_m$  не влияют на время построения модели.

Таблица 9: Качество работы алгоритма WMTF-G для различных значений  $K$  и  $I$  при фиксированных значениях  $w_m = 5$ ,  $\lambda = 6$ ,  $\delta = 0.06$ .

for delta in [0.06, 0.08, 0.1, 0.12, 0.14]: for lmbd in [6, 8, 10, 12, 14]:

$K \backslash I$	1	2	3	4	5
30	0.7529	0.7927	0.7577	0.6736	0.5794
40	0.7992	0.8194	0.7695	0.6813	0.5828
50	0.8269	0.8349	0.7834	0.6830	0.5801
60	0.8419	0.8450	0.7984	0.7056	0.6006
70	0.8557	0.8466	0.7977	0.7036	0.6002
80	0.8614	0.8511	0.7990	0.7032	0.5957
90	0.8606	0.8522	0.8039	0.7088	0.6038
100	0.8606	0.8527	0.8022	0.7089	0.6021
110	0.8686	0.8553	0.8074	0.7123	0.6065
120	0.8693	0.8579	0.8097	0.7174	0.6085
130	0.8725	0.8588	0.8160	0.7264	0.6206
140	0.8740	0.8597	0.8157	0.7248	0.6241
150	0.8763	0.8620	0.8171	0.7263	0.6200

В качестве начального приближения параметров взяты оптимальные параметры для метода WTMF, а именно  $K = 90$ ,  $I = 1$ ,  $w_m = 1.95$ . В качестве начального приближения параметра  $\delta$  мы берем значение 0.1

Оптимизируется параметр  $\delta$ . Для этого фиксируются остальные параметры:  $K = 90$ ,  $I = 1$ ,  $w_m = 1.95$ . Для начала находится оптимальный порядок значений начального приближения. Результаты занесены в таблицу 10. Как видно из таблицы 10

Таблица 10: Качество работы алгоритма WMTF-G для различных значений  $\delta$  при фиксированных значениях  $K = 90$ ,  $I = 1$ ,  $w_m = 1.95$ .

$\delta$	Значение метрики RR
<b>0.001</b>	0.5508
<b>0.01</b>	0.5307
<b>0.1</b>	0.5695
<b>1</b>	0.5311
<b>10</b>	0.5303
<b>100</b>	0.5203

максимальное значение метрики получено при  $\delta = 0.1$ . Для уточнения значения коэффициента  $\delta$ , производится исследование качества работы алгоритма в окрестностях максимального значения метрики. Результаты приведены в таблице 11. Из таблицы 11 получаем оптимальные значения коэффициента  $\delta = 0.1$ .

В итоге оптимизации качества рекомендаций на основе алгоритма WMTF-G



Таблица 11: Качество работы алгоритма WMTF-G для различных значений  $\delta$  при фиксированных значениях  $K = 90$ ,  $I = 1$ ,  $w_m = 1.95$ .

$\delta$	Значение метрики RR
<b>0.05</b>	0.5340
<b>0.1</b>	0.5695
<b>0.15</b>	0.5380
<b>0.25</b>	0.5533
<b>0.3</b>	0.5195
<b>0.35</b>	0.5329

были получены оптимальные параметры:  $K = 90$ ,  $I = 1$ ,  $\delta = 0.1$ ,  $w_m = 1.95$ .

#### 4.4. Сравнительные результаты

Для выявления влияния добавления связей текст-текст на результаты работы метода WMTF-G производится сравнительное тестирование алгоритма WTMF и WTMF-G. Результаты тестирования приведены в таблице 12.

Таблица 12: Сравнительное тестирование алгоритмов WTMF и WTMF-G.

Набор данных	Метрика MRR		Метрика $TOP_1$		Метрика $TOP_3$	
	WTMF	WTMF-G	WTMF	WTMF-G	WTMF	WTMF-G
manual	0.7293	0.	0.	0.	0.	0.
auto	0.8640	0.	0.	0.	0.	0.
total	0.8196	0.	0.	0.	0.	0.
cutted	0.8630	0.	0.	0.	0.	0.
manual_nt	0.6194	0.	0.	0.	0.	0.
auto_nt	0.5297	0.5695	0.	0.	0.	0.
total_nt	0.5729	0.	0.	0.	0.	0.
cutted_nt	0.6495	0.	0.	0.	0.	0.

Как видно из таблицы 12 ...

Сравним метод основанный на частотности употребления слов и WTMF-G.

Метод основанный на частотности употребления слов обозначим как TF-IDF. Результаты тестирования приведены в таблице 13.

Как видно из таблицы 13 ...

объяснение влияния различных датасетов, специфики русского языка и сравнение с результатами статьи.

Таблица 13: Сравнительное тестирование алгоритмов TF-IDF и WTMF-G.

Набор данных	Метрика MRR		Метрика $TOP_1$		Метрика $TOP_3$	
	TF-IDF	WTMF-G	TF-IDF	WTMF-G	TF-IDF	WTMF-G
manual	0.8336	0.	0.	0.	0.	0.
auto	0.8817	0.	0.	0.	0.	0.
total	0.8610	0.	0.	0.	0.	0.
cutted	0.9075	0.	0.	0.	0.	0.
manual_nt	0.7565	0.	0.	0.	0.	0.
auto_nt	0.6048	0.5695	0.	0.	0.	0.
total_nt	0.6914	0.	0.	0.	0.	0.
cutted_nt	0.7485	0.	0.	0.	0.	0.