

тсюда можно сделать вывод, что у разработанного ритма минимизирована работа оператора по ручной обработке результатов, т.е. хотя скорость обработки несколько меньше, но алгоритм позволяет существенно разгрузить операторов за счёт интеллектуальных систем принятия решений, чего не может предложить алгоритм прямого сравнения.

ри сравнении экономических характеристик разработанного программного обеспечения на основе выаемого алгоритма с процедурой прямого снения, для годового объёма идентификации в 000 физических лиц были получены следующие е: трудовые затраты на обработку информации году нечёткого сравнения по сравнению с методом прямого сравнения уменьшены в 6,7 раза, абсолютное снижение трудовых затрат составило 1446 , годовые затраты при использовании метода прямого сравнения уменьшились в 3 раза по сравнению с аналогичным периодом применения метода прямого сравнения, а годовой экономический эффект составил 580000 руб. Для наглядности некоторые стоимые показатели, формирующиеся при использовании разработанного и применявшегося до настоящего времени программного обеспечения, отображены на амме, приведенной на рис. 2. Величины затрат сены по оси ординат в рублях.

Выводы

обучающиеся системы позволяют освободить еческие ресурсы для выполнения творческих В этой области технология Data Mining предоставляет полный набор теоретических и практических в для выбора, разработки или использования актуальных компьютерных систем.

рассмотренную в статье процедуру идентификации рассматривать как часть системы поддержки гия решений (СППР). Процедура не требует ьства оператора, накапливает опыт и самостя в процессе работы, позволяя, тем самым, стью освободить специалистов от низкопродой, неэффективной, ручной работы напрямую с ми реквизитов физических лиц, хранящихся в данных. Данная процедура реализована на языке L СУБД Oracle 11g и успешно функционирует с 2007 г. в муниципальном учреждении ттинский городской информационный центр». ская структура разработанного алгоритма (рис. оляет реализовать его на любом популярном программном обеспечении.

ерспективе данный алгоритм обладает возмож- успешного внедрения в системы глобального нения хранилищ государственных или коммер- организаций для ведения единой базы данных ния любой страны мира. Масштабируемость

алгоритма позволяет применять программные проце- дуры на его основе как в малых организациях, так и в крупных корпорациях, везде, где ведётся и актуализируется реестр данных физических лиц. Возможные примеры использования: портал госуслуг, медицинские электронные системы, кадровые и бухгалтерские системы учёта служащих, банковские системы хранения данных о клиентах и т.п.

ЛИТЕРАТУРА

1. Чубукова И.А. «Data Mining»: учебный курс. // Издательство Интернет-университета информационных технологий (<http://www.intuit.ru/>), 2006.
2. Международный фонд автоматической идентификации. Технологии автоматической идентификации. (<http://www.fond-ai.ru/art1/art223.html>).
3. Подборка материалов о международном опыте законодательного регулирования использования систем идентификации личности (<http://www.kongord.ru/Index/Prison/SViP.htm>)
4. Отчёт о выполнении научно-исследовательской, опытно-конструкторской работы «Разработка механизмов однозначной идентификации данных о физических лицах и объектах недвижимости, хранящихся в различных информационных системах органов государственной власти и местного самоуправления» (http://www.nisse.ru/business/article/article_464.html)
5. Положение о персональном идентификационном номере граждан Российской Федерации, проживающих или пребывающих на территории Санкт-Петербурга. (<http://iac.spb.ru/shablon.asp?subpage=171&id=40&dir=0>).
6. Проект «Социальная карта москвича» (<http://www.soccard.ru>).
7. Сборник тезисов городской научно-практической конференции студентов, аспирантов, преподавателей вузов и специалистов муниципальных учреждений г.Тольятти «Информатизация в социальной сфере» (<http://it-exclusive.ru/idperson/docs/stat.doc>).
8. Урман С. «ORACLE 9i – Программирование на языке PL/SQL»: учебное пособие // Oracle Press – издательство «Лорд», 2004.

*Наталья Игоревна Лиманова,
д-р техн. наук, профессор.
Тольяттинский гос. университет*

*Максим Николаевич Седов,
инженер-программист I-ой категории.
Мэрия городского округа Тольятти.*

Ю.А. Ставенко, Д.А. Романов, А.И. Громов

ОПРЕДЕЛЕНИЕ ЗНАЧИМОСТИ ЛЕКСИКИ НА ОСНОВЕ СТАТИСТИЧЕСКОГО КОНТЕНТ-АНАЛИЗА

В статье приведен алгоритм контент-анализа корпоративных документов с целью выявления уровня экспертизы сотрудников в той или иной предметной области. В процессе контент-анализа были выявлены наиболее значимые слова, описывающие смысл документов, сгенерированных сотрудником за определенный период времени, а также был предложен способ отсеивания общеупотребительной лексики.

Ключевые слова: контент-анализ, значимость лексики, поиск экспертизы, корпоративный поиск, TF, IDF

Yu.A. Stavenko, D.A. Romanov, A.I. Gromov

ESTIMATION OF TERM INFORMATIVENESS BASED ON A STATISTICAL CONTENT-ANALYSIS

The article is dedicated to the content analysis algorithm of corporate documents in order to extract the level of expertise of each expert in a particular domain area. During the process of content analysis the most significant words were identified that describe the meaning of documents generated by the employee for a period of time, and a method of screening vernacular vocabulary was proposed.

Keywords: content analysis, significant words, expertise search, TF, IDF

Введение

Информация, представленная в виде текстов на естественном языке, является одним из самых распространенных видов информации. Соответственно, существует множество способов контент-анализа текстов (например, аннотирование, конспектирование, реферирование), позволяющих извлекать и оценивать значимую информацию (знания) из текстов. Философский смысл контент-анализа как исследовательского метода, состоит в восхождении от многообразия текстового материала к абстрактной модели содержания текста.

Прародителем контент-анализа можно считать проведенный анализ сборника церковных гимнов в Швеции в 18 веке. Эти гимны сначала прошли государственную цензуру и приобрели большую популярность, после чего были обвинены в несоответствии религиозным догматам. Для того чтобы доказать некорректность текстов с точки зрения официального учения церкви, в них был осуществлен подсчет религиозных символов, количество которых было затем сравнено с другими религиозными текстами, в том числе с теми, которые считались еретическими. Частота использования определенных, заранее собранных слов и тем позволяла судить о том, насколько корректен текст [1].

Что касается дальнейшего развития контент-анализа, то огромную популярность этот метод приобрел после появления различного рода политических партий в начале 20 века, например, он использовался Максом Вебером в 1910 году для анализа освещаемости прессой политических акций в Германии. Позднее, в 1937 году метод контент-анализа был использован в США

рациональных речей американских президентов, а затем в военных целях и в целях исследования направленной пропагандистской деятельности.

Огромный вклад в контент-анализ текстов был внесен французским журналистом Ж. Кайзером, который разработал методику статистического анализа публикаций СМИ, а затем Э. Морэн, которая ввела в научный оборот термин «единица информации» - блок, содержание которого определяет смысловую нагрузку текста.

Затем метод контент-анализа получил большое распространение в различных сферах политологии и социологии, и успешно используется в них до сих пор за счет своей эффективности. Эта эффективность достигается с помощью упрощения (редукции) текста к стандартному набору смысловых единиц (терминов), которые могут быть обработаны методами статистики, в результате чего на выходе будет получен значимый результат для целей исследования. Методами обработки наиболее часто являются частота упоминаний терминов в массиве текстов, на основании которой делается вывод об информативности термина и характеристиках текста.

В данной статье будет сделана попытка применения метода контент-анализа к корпоративным текстам, созданным сотрудниками организации. Выбор данного массива текстов не случаен. В организациях любого уровня каждый день встают вопросы поиска сотрудника, имеющего достаточный уровень компетенции для решения вопросов конкретной предметной области. Поиск необходимой компетенции - одна из самых приоритетных задач управления знаниями, которая, к сожалению, во многих организа-

ручного профилирования. Основная проблема заключается в том, что неявные знания, присутствующие у сотрудников, слабо структурированы и должны быть извлечены из корпоративной документации. Соответственно, задача поиска неявных знаний сильно взаимосвязана с задачей контент-анализа текстов документов.

Цель и актуальность исследования

Одной из задач извлечения информации из текста является выделение ключевых терминов, с определённой степенью достоверности отражающих тематическую направленность документа. Автоматическое извлечение ключевых терминов можно определить как автоматическое выделение важных тематических терминов в документе. Оно является одной из подзадач более общей задачи – автоматической генерации ключевых терминов, для которой выделенные ключевые термины не обязательно должны присутствовать в данном документе [2].

В последние годы было создано множество подходов, позволяющих проводить анализ наборов документов различного размера и извлекать ключевые термины, состоящие из одного, двух и более слов. Самыми распространёнными схемами для расчёта весов терминов являются *TF-IDF* и различные его модификации. В общем-то, индекс *TF-IDF* уже можно использовать для ранжирования результатов поиска – чем выше индекс, тем больше релевантность данного термина в данном документе. Применительно к задаче поиска экспертов метрика *IDF* может быть переопределена как логарифм отношения числа авторов, хотя бы раз употребивших термин, к общему числу авторов. Но поскольку число авторов в корпоративной среде определяется числом сотрудников компании и фактически равно константе (или меняется достаточно медленно), метрика *IDF* зависит только от числа авторов, употребивших данный термин. Тем не менее, для задачи поиска компетенции этого всё ещё недостаточно, так как при сравнении вектора *TF-IDF* поискового запроса и вектора *TF-IDF* множества документов автора незначимые, часто встречаемые термины вносят слишком большую вариацию в результат. Поэтому представляется важным различать общеупотребительность и частоту встречаемости термина. Поскольку для решения задачи поиска экспертов применяется, и будет применяться ограниченный по объёму корпус текстов, накопленный за конечное (и, часто, весьма небольшое) время, некоторые вполне общеупотребительные, но редко встречающиеся слова могут просто «не успеть» проявиться в лексике всех авторов. В этом случае метрика *IDF* для таких терминов окажется

завышенной, хотя сами по себе термины не будут являться значимыми.

Целью исследования является разработка алгоритма для программного сервиса внутрикорпоративного поиска экспертизы на основании контент-анализа. С математической точки зрения необходимо разработать алгоритм нахождения наиболее релевантных экспертов с помощью представления запроса пользователя и контента, генерируемого сотрудниками организации в виде набора значимых терминов $\{t_k\}$ и $\{t_i\}$ соответственно. Каждому термину t_i должна соответствовать мера $Weight(t_i; p_i)$, определяющая значимость данного термина t_i для эксперта p_i .

Методологическая база. Выделение значимых терминов

Извлечение значимых терминов является базисным этапом для выделения экспертов. Значимыми терминами являются ключевые термины, которые могут дать интерпретацию смысла документов, созданных сотрудниками. Это задача обработки естественного языка, которая решается универсальными методами с помощью подсчета частотных характеристик употребления терминов, например, в поисковых системах.

Методологической базой для алгоритмов, основанных на выявлении статистических закономерностей распределения частоты появления различных слов в текстах, принято считать закон Ципфа (Zipf's Law) [3]. Для его интерпретации необходимо посчитать частоту появления слов в тексте *TF* (term frequency). Если расположить слова текста в порядке убывания частоты их появления *TF*, начиная с наиболее часто встречающихся, то произведение частоты слова (*TF*)_{*i*} на порядковый номер частоты будет постоянным для любого данного слова t_i :

$$TF_i \times r_i = C, \quad (1)$$

где *C* – некоторая константа, r_i – порядковый номер (ранг) частоты слова.

На рис. 1 представлен график закономерности «ранг-частота» для текста научной статьи [3].

Ципф выявил, что для закономерностей ранг-частота всегда существуют три явно различаемые зоны ранговых распределений: зона ядра рангового распределения (наиболее часто употребляемые слова языка: «the», «a», «that»), центральная зона и зона усечения. В центральной зоне находятся слова, максимально характеризующие данный текст и выражающие его специфичность и тематику, то есть ключевые термины: «zipf», «low». В зоне усечения находятся слова, не несущие основной смысловой нагрузки. От ширины центральной зоны зависит

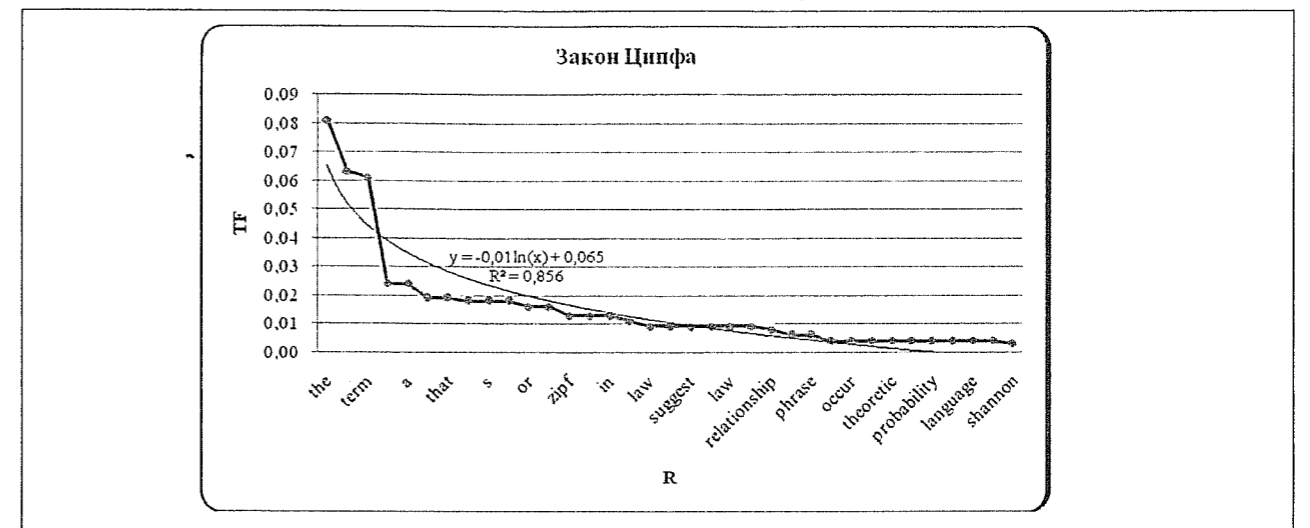


Рис. 1. Закон Ципфа

качество выделяемых ключевых терминов. На этапе определения ширины центральной зоны первостепенной задачей является отсечение информационного шума – нейтральных стоп-слов. Традиционный способ решения данной задачи заключается в применении специализированных словарей, содержащих наиболее часто встречающиеся языковые единицы (местоимения, предлоги, союзы и т.п.).

В настоящем исследовании также стоит вопрос выделения значимых терминов в текстах, написанных сотрудниками организации и выделения их профессиональной экспертизы. Воспользовавшись следствием из закона Ципфа, на первом этапе предварительной обработки из массива выделенных слов из текстов сотрудников были удалены все стоп-слова.

Следующим этапом необходимо выделить ключевые термины, характеризующие профиль сотрудника (его предметную область). При этом подобная задача должна решаться в соответствии со спецификой каждой организации (социальной среды).

Одним из традиционных методов выделения значимых терминов является расчет показателя *IDF* (inverse document frequency) – инверсия частоты, с которой некоторый термин встречается в документах коллекции, то есть значение *IDF* тем меньше, чем чаще встречается слово. По сути *IDF* оценивает количество информации, свойственной слову. Применительно к задаче составления профиля сотрудника показатель *IDF* будет рассчитываться следующим образом для термина t_i :

$$IDF(t_i) = \frac{\lg N}{nt}, \quad (2)$$

где *N* – общее число экспертов; *nt* – число экспертов, употребивших термин t_i , хотя бы один раз.

Ключевыми в данном случае расчета значимости будут являться термины, набравшие наибольший вес *IDF*. Таким образом, для слов, которые употребляются большим числом экспертов, *IDF* будет близок к нулю.

Смысл применения подобного алгоритма для поиска экспертов состоит в доказательстве гипотезы, что значимость термина пропорциональна количеству встречаемости термина в документах автора и обратно пропорциональна количеству встречаемости термина во всей коллекции документов. К сожалению, эта гипотеза может быть опровергнута, так как метрика *IDF* не может учитывать зависимость появления терминов во времени, что критично для анализа динамически изменяющегося информационного поля организации. Поскольку для решения задачи поиска экспертов применяется и будет применяться ограниченный по объёму корпус текстов, накопленный за конечное (и, часто, весьма небольшое) время, некоторые вполне общеупотребительные, но редко встречающиеся слова могут просто «не успеть» проявиться в лексике всех авторов. В этом случае метрика *IDF* и, соответственно, значимость для таких терминов окажется завышенной, хотя сами по себе они не будут являться значимыми. Зависимость от организационно-социальной среды и зависимость от времени обуславливают необходимость динамического определения значимости термина.

Гипотеза исследования

Для разрешения вопроса отделения значимых терминов от общеупотребительных, описанного выше, можно сформулировать гипотезу, что характеристиками рангового распределения обладают не только зависимость ранга от частоты употребления слова в

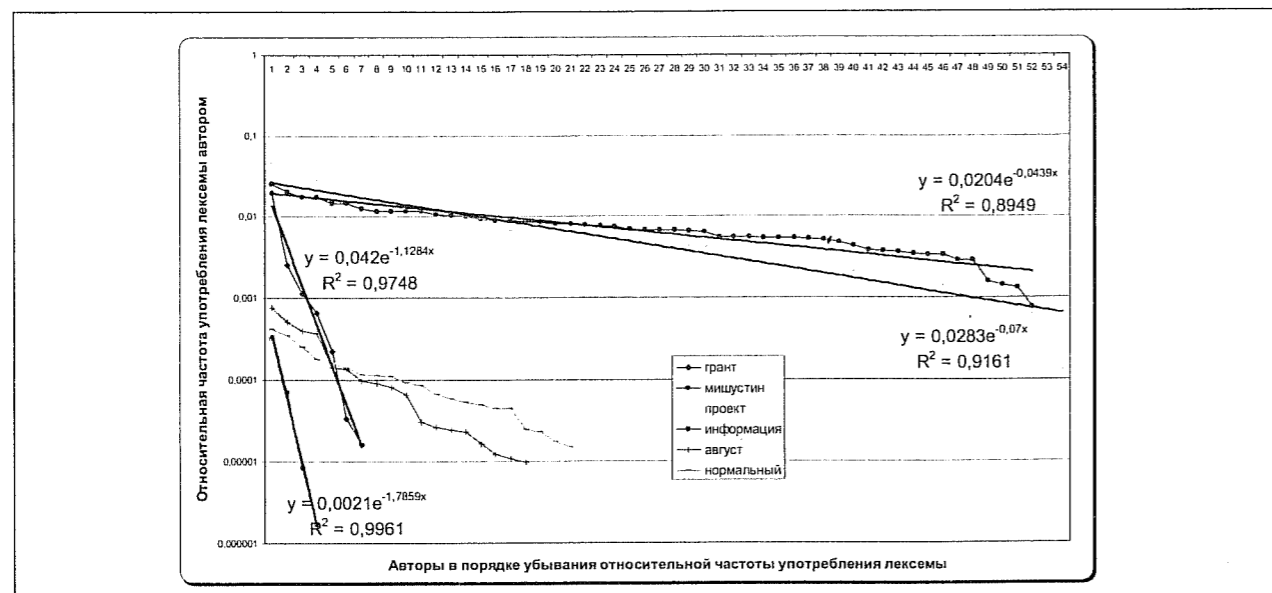


Рис. 2. Зависимость относительной частоты употребления термина автором от ранга в би-логарифмической шкале

тексте (закон Ципфа), но и зависимость ранга от относительной частоты употребления термина автором. Для этого необходимо подсчитать статистику *относительной частоты употребления термина* t_i для всех текстов, написанных конкретным сотрудником p_j :

$$TF(t_i, p_j) = \frac{m(t_i, p_j)}{\sum_k m_k} \quad (3)$$

где $m(t_i, p_j)$ – число употреблений термина t_i персонай p_j , а в знаменателе – общее число употреблений всех терминов персонай p_j .

Можно предположить, что значимые термины должны иметь сильно неравномерное распределение относительной частоты употребления среди сотрудников, а общеупотребительные – примерно одинаковую относительную частоту употребления. Построим подобные зависимости в двойной логарифмической шкале для различных терминов (см. рис. 2), выбранных из текстов сотрудников ИТ-компаний¹:

Интуитивно понятно, что такие термины как «грант» и «мишустин» значимы, так как «грант» – специфический термин предметной области, «мишустин» – имя собственное. Термины «проект», «информация» являются общеупотребительными в выбранной предметной области. Анализируя график видно, что для значимых слов и общеупотребительных отличается характер распределения, а именно дисперсия в наблюдаемом распределении относительных частот употребления терминов. В результате эксперимента можно сделать вывод, что именно

¹ Эмпирическую базу исследования составили документы электронной почты компании «АйТи».

более высокие значения дисперсии относительной частоты употребления термина являются тем самым критерием, который позволит отличить значимые термины от общеупотребительных.

Для подтверждения гипотезы в проведенном исследовании было проведено сравнение двух алгоритмов поиска экспертов на основании определения значимости терминов:

- алгоритма, основанного на подсчете IDF и используемого в векторной модели TF*IDF;
- алгоритма, основанного на вычислении дисперсии распределения относительной частоты встречаемости терминов.

Лингвистический анализ текста. Подготовка к работе алгоритмов

В проведенном исследовании поиск экспертов осуществлялся на основе обработки текстового контента, образующегося в корпоративной сети (включая архив электронной почты и массивы других видов текстовых сообщений). Лингвистический анализ текстов реализуется в виде лингвистического процессора, задачей которого является преобразование текста на естественном языке в некоторый набор структур, являющихся формальным представлением контента исходного текста.

Для выделения терминов из текстов применяется парсинг документов, созданных сотрудниками организации, цель которого – провести токенизацию (графематический анализ, лексический анализ), на основании которой происходит выделение токенов – лексических конструкций, состоящих из последова-

тельности символов одного типа (слов, цифровых записей).

После этого происходит морфологический анализ с помощью морфологического анализатора, на вход которого поступают токены. С помощью него происходит приведение выделенных токенов к словарной форме, а также установление грамматических характеристик словоформ (род, число, падеж), которые используются для выделения словосочетаний. Результаты морфологического анализа используются на этапе синтаксического анализа, в частности. Модуль синтаксического анализа позволяет частично разрешить грамматическую омонимию как для идентифицированных, так и для отсутствующих в словаре слов. На следующем этапе происходит синтаксический анализ – выделение не только слов, но и именных групп – это выражение, центральным элементом которого является существительное.

В итоге на выходе получается сформированный список терминов, привязанных к автору – сотруднику организации.

Алгоритм поиска экспертов

Для каждого термина t_i и сотрудника p_j считается статистическая информация по количеству употребления термина сотрудником (см. формулу 3) – относительная частота употребления $TF(t_i, p_j)$. В итоге для каждого термина можно составить выборку относительной частоты употребления этого термина авторами $TF(t_i, p_j)$:

$$\begin{matrix} TF(t_1, p_1) & \dots & TF(t_1, p_N) \\ \vdots & & \vdots \\ TF(t_n, p_1) & \dots & TF(t_n, p_N) \end{matrix}, \quad (4)$$

где n – количество терминов в информационном поле организации, N – количество авторов.

Можно предложить несколько вариантов метрик для определения значимости термина, основанных на идее нелинейного распределения вероятности. Простейшие метрики для степени изменения значений выборки определяются путем вычисления их дисперсии $D(t_i)$, то есть дисперсии относительной частоты употребления термина t_i и стандартного отклонения $\sigma(t_i)$:

$$\begin{aligned} D(t_i) &= ([TF(t_i, p_1) - M]^2 + \\ &+ [TF(t_i, p_2) - M]^2 + \\ &+ [TF(t_i, p_N) - M]^2) / (N - 1) = \\ &= 1 / (N - 1) \sum_{k=1}^N [TF(t_i, p_k) - M]^2 \end{aligned} \quad (5)$$

где M – оценка математического ожидания (выборочное среднее) относительной частоты употребления термина, рассчитываемая по формуле:

$$\begin{aligned} \bar{M}(t_i) &= \\ &= \frac{TF(t_i, p_1) + TF(t_i, p_2) + \dots + TF(t_i, p_N)}{N} = \\ &= \frac{1}{N} \sum_{k=1}^N TF(t_i, p_k) \end{aligned} \quad (6)$$

Среднеквадратичное отклонение относительной частоты употребления термина – корень из дисперсии:

$$\sigma(t_i) = \sqrt{D(t_i)} \quad (7)$$

Тогда значимость термина для конкретного автора будет являться превышением над математическим ожиданием $\bar{M}(t_i)$, измеренным в среднеквадратических отклонениях $\sigma(t_i)$:

$$Impact(t_i; p_j) = \frac{(TF(t_i, p_j) - \bar{M}(t_i))}{\sigma(t_i)} \quad (8)$$

По формуле (8) получается, что отрицательные значения $Impact(t_i; p_j)$ будут у авторов, которые редко употребляли термины, то есть не экспертов по этому термину.

Для нормировки веса $Impact(t_i; p_j)$ возьмем арктангенс, тогда чем выше значение арктангенса, тем слово значимее для автора:

$$FinalImpact(t_i; p_j) = \arctg Impact(t_i; p_j) \quad (9)$$

Мерой близости компетенции эксперта к поисковому запросу (тексту) будет являться сумма весов термов, употребленных в тексте для автора:

$$Weight(p_j) = \sum_i^l FinalImpact(t_i; p_j) \quad (10)$$

где t_i – термины в запросе пользователя; l – количество терминов в запросе пользователя.

Произведя ранжирование сотрудников по мере $Weight(p_j)$ получим релевантный список экспертов.

Алгоритм поиска экспертов на основании векторной модели

Для определения эффективности предложенного алгоритма будем сравнивать его со стандартным алгоритмом TF*IDF, который используется при информационном поиске. Главной идеей этого алгоритма является рассмотрение документов и запросов как векторов в пространстве слов, а релевантность документа и запроса – как расстояние между ними. В качестве меры подобия используется «мера косинуса», которая считается как косинус угла между векторами. Она хороша тем, что принимает значения от 0 до 1, и равна 1 только при полном подобии векторов. В исследовании была использована аналогичная модель для определения релевантности эксперта запросу пользователя.

На первом шаге алгоритма происходит построение матрицы относительной частоты употребления тер-

