МЕТОД УЧЁТА СТРУКТУРЫ БИГРАММ В ТЕМАТИЧЕСКИХ МОДЕЛЯХ

М. А. Нокель

Московский государственный университет им. М. В. Ломоносова

Поступила в редакцию 24.10.2014 г.

Аннотация. В статье представлены результаты экспериментов по добавлению сходства между униграммами и биграммами в тематические модели. Вначале изучается возможность применения ассоциативных мер для выбора и последующего включения биграмм в тематические модели. Затем предлагается модификация оригинального алгоритма PLSA, учитывающая похожие униграммы и биграммы, начинающиеся с одних и тех же букв. И в конце статьи предлагается новый итеративный алгоритм без учителя, показывающий, как темы сами могут выбирать себе наиболее подходящие биграммы. В качестве текстовой коллекции была взята подборка статей из электронных банковских журналов на русском языке. Эксперименты показывают значительное улучшение качества тематических моделей по всем целевым метрикам..

Ключевые слова: тематические модели, биграммы, согласованность тем, перплексия, PLSA.

Annotation. The paper presents the results of experimental study of integrating word similarity and bigram collocations into topic models. First of all, we analyze a variety of word association measures in order to integrate top-ranked bigrams into topic models. Then we propose a modification of the original algorithm PLSA, which takes into account similar unigrams and bigrams that start with the same beginning. And at the end we present a novel unsupervised iterative algorithm demonstrating how topics can choose the most relevant bigrams. As a target text collection we took articles from various Russian electronic banking magazines. The experiments demonstrate significant improvement of topic models quality.

Keywords: topic models, bigrams, topic coherence, perplexity, PLSA.

1. ВВЕДЕНИЕ

Вероятностные тематические модели (далее просто тематические модели) – одно из современных приложений машинного обучения к анализу текстов. Тематические модели предназначены для описания текстов с точки зрения их тем. Они определяют, к каким темам относится каждый документ в текстовой коллекции и какие слова образуют каждую такую тему. При этом темы представляются в виде дискретных распределений на множестве слов, а документы – в виде дискретных распределений на множестве тем [1]. Пользова-

телям темы предоставляются, как правило, в виде некоторых списков часто встречающихся рядом друг с другом слов, упорядоченных по убыванию степени принадлежности им.

С момента своего появления тематические модели достигли значительных успехов в задачах информационного поиска, разрешении морфологической неоднозначности, многодокументного аннотирования, кластеризации и категоризации документов и других задачах. Самыми известными представителями являются латентное размещение Дирихле (LDA) [1], использующее априорное распределение Дирихле, и метод вероятностного латентного семантического анализа (PLSA) [2], не связанный ни с какими параметрическими априорными распределениями.

[©] Нокель М. А., 2014

Работа частично поддержана грантом РФФИ 14-07-00383

Одним из главных недостатков тематических моделей является использование модели «мешка слов», в которой каждый документ рассматривается как набор встречающихся в нем слов. Данная модель не учитывает порядок слов и основывается на гипотезе независимости появлений слов в документах друг от друга. На данный момент проведено множество исследований, посвященных изучению вопроса добавления словосочетаний и *n* -грамм в тематические модели. Однако часто это приводит к ухудшению качества модели в связи с увеличением размера словаря или к значительному усложнению модели [3], [4], [5].

В статье предлагается новый подход, позволяющий учесть взаимосвязь между похожими словами (в частности, однокоренными) в тематических моделях (такими, как банк – банковский – банкир). На основании данного метода в статье описывается и новый подход к добавлению биграмм в тематические модели, который рассматривает биграммы уже не как «черные ящики», а учитывает взаимосвязь между ними и униграммами, основанную на их внутренней структуре. Предлагаемый алгоритм улучшает качество тематических моделей по двум целевым метрикам: перплексии и согласованности тем [6].

Все эксперименты, описанные в статье, проведены на основе алгоритма PLSA и его модификаций на коллекции текстов банковской тематики на русском языке, взятых из электронных журналов.

2. БЛИЗКИЕ РАБОТЫ

2.1. Тематические модели

На сегодняшний день разработано достаточно много различных тематических моделей. Исторически одними из первых появились модели, основанные на традиционных методах кластеризации текстов. При этом после окончания работы алгоритма кластеризации каждый получившийся кластер рассматривается как отдельная тема для вычисления вероятностей входящих в него слов по следующей формуле:

$$P(w|t) = \frac{f(w|t)}{\sum_{w} f(w|t)},$$

где f(w|t) – частотность слова w в теме t.

Естественным ограничением таких моделей является отнесение каждого документа лишь к одной теме.

В последнее время появились вероятностные механизмы нахождения тем в документах, рассматривающие каждый документ в виде смеси тем, а каждую тему в виде некоторого вероятностного распределения над словами. Вероятностные модели порождают слова по следующему правилу:

$$P(w \mid d) = \sum_{t} P(w \mid t) P(t \mid d),$$

где $P(t \mid d)$ и $P(w \mid t)$ – распределение тем по документам и слов по темам, а $P(w \mid d)$ – наблюдаемое распределение слов по документам.

Самыми известными представителями данной категории являются метод вероятностного латентного семантического анализа (PLSA) [2] и латентное размещение Дирихле (LDA) [1].

2.2. Словосочетания в тематических моделях

Все описанные в прошлом разделе алгоритмы работают только со словами, основываясь на гипотезе о независимости слов друг от друга – модели «мешка слов». Идея же использования словосочетаний в тематических моделях сама по себе не нова. На данный момент существуют 2 подхода к решению данной проблемы: создание унифицированной вероятностной модели и предварительное извлечение словосочетаний и *п*-грамм для их последующего добавления в тематические модели.

Большинство исследований на данный момент посвящено первому подходу. Так, первая попытка выйти за пределы модели «мешка слов» была предпринята в работе [3], где была представлена Биграммная Тематическая Модель. В этой модели вероятности слов зависят от вероятностей непосредствен-

но предшествующих им слов. Модель словосочетаний LDA [4] расширяет Биграммную Тематическую Модель за счет введения дополнительных переменных, способных генерировать и униграммы, и биграммы. В работе [5] представлена Тематическая *N*-граммная Модель, усложняющая предыдущие для обеспечения возможности формирования биграмм в зависимости от контекста.

Несмотря на то, что все описанные выше модели имеют теоретически элегантное обоснование, у них очень большое число параметров для настройки, что ведёт к неприменимости на реальных данных. Так, например, число параметров Биграммной Тематической Модели равно W^2T , в то время как для LDA оно равно WT, для PLSA – WT+DT, где W – размер словаря, D – количество документов в коллекции и T – число тем. Поэтому такие модели представляют в основном чисто теоретический интерес.

Алгоритм, предложенный в работе [7], относится ко второму типу методов, добавляющих словосочетания в тематические модели. На этапе предобработки авторы извлекают биграммы с помощью t-теста и заменяют отдельные униграммы лучшими по данной мере биграммами. При этом используются 2 метрики оценивания качества полученных тем: перплексия и согласованность тем [6]. В статье показано, что добавление биграмм в тематические модели приводит к ухудшению перплексии и к улучшению согласованности тем.

Данная работа также относится ко второму типу методов и отличается от работы [7] в том, что описываемый здесь подход учитывает внутреннюю структуру биграмм и взаимосвязь между ними и составляющими их униграммами, что приводит к улучшению обоих показателей: и перплексии, и согласованности тем.

3. ТЕКСТОВАЯ КОЛЛЕКЦИЯ И МЕТОДЫ ОЦЕНИВАНИЯ КАЧЕСТВА ТЕМАТИЧЕСКИХ МОДЕЛЕЙ

3.1. Текстовая коллекция и предобработка

В экспериментах, описанных в данной статье, использовалась текстовая коллекция из

10422 статей на русском языке, взятых из некоторых электронных банковских журналов (таких, как Аудитор, РБК, Банковский журнал и др.). В данных документах содержится почти 15.5 млн слов.

На этапе предобработки был проведен морфологический анализ документов. В экспериментах рассматривались только существительные, прилагательные, глаголы и наречия, поскольку служебные слова не играют значительной роли в определении тем. Кроме того, из рассмотрения исключались слова, встретившиеся менее 5 раз во всей текстовой коллекции.

На этапе предобработки из документов также извлекались биграммы в формах сущ. + сущ. в родительном падеже и прил. + сущ. В экспериментах рассматривались только такие биграммы, поскольку темы, как правило, задаются именными группами.

3.2. Методы оценивания качества тематических моделей

Для оценивания качества полученных тем в статье рассматриваются две метрики.

Во-первых, использовалась **перплексия**, являющаяся стандартным критерием качества тематических моделей. Эта мера несоответствия модели p(w | d) словам w, наблюдаемым в документах коллекции, определяется через логарифм правдоподобия:

Perplexity(D) = exp
$$\left(-\frac{\sum_{d \in D} \sum_{w \in d} n_{dw} \ln p(w \mid d)}{n}\right)$$

где n – число всех рассматриваемых слов в текстовой коллекции, D – множество всех документов в коллекции, n_{dw} – частота слова w в документе d, p(w|d) – вероятность появления слова w в документе d.

Чем меньше значение перплексии, тем лучше модель предсказывает появление слов w в документах коллекции D. Поскольку известно, что перплексия, вычисленная на той же самой обучающей коллекции документов, склонна к переобучению и может давать оптимистически заниженные значения [1], в данной статье используется стандартный ме-

тод вычисления контрольной перплексии. Коллекция документов изначально разбивалась на 2 части: обучающую D, по которой строилась модель, и контрольную D', по которой вычислялась данная метрика. Хотя на данный момент существует множество исследований, утверждающих, что перплексию нельзя применять для оценивания качества тематических моделей, данная метрика по-прежнему широко используется для сравнения различных тематических моделей.

В то же время неоднократно предпринимались попытки предложить способ автоматического оценивания качества тематических моделей, никак не связанного с перплексией и коррелирующего с мнениями экспертов. Данная постановка задачи является очень сложной, поскольку эксперты могут достаточно сильно расходиться во мнениях. Однако в недавних работах [6], [8] было показано, что возможно автоматически оценивать согласованность тем, основываясь на семантике слов с точностью, почти совпадающей с экспертами. Предложенная метрика измеряет интерпретируемость тем, основываясь на способах оценивания экспертом [6]. Поскольку темы, как правило, предоставляются экспертам для проверки в виде первых топ-N слов, согласованность тем оценивает то, насколько данные слова соответствуют рассматриваемой теме. Newman в работе [6] предложил использовать автоматический способ вычисления данной метрики исходя из меры взаимной информации:

$$TC - PMI(t) = \sum_{j=2}^{10} \sum_{i=1}^{j-1} \ln \frac{P(w_j, w_i)}{P(w_i)P(w_j)},$$

где $(w_1, w_2, \ldots, w_{10})$ – топ-10 слов в рассматриваемой теме t, $P(w_i)$ и $P(w_j)$ – вероятности униграмм w_i и w_j соответственно, а $P(w_j, w_i)$ – вероятность биграммы (w_j, w_i) . Итоговая мера согласованности тем вычисляется усреднением TC-PMI(t) по всем темам t.

Данная метрика показывает очень высокую корреляцию с оценками экспертов [6]. Предложенная метрика рассматривает только первые топ-10 слов в каждой теме, поскольку они, как правило, предоставляют достаточно информации для формирования

предмета темы и отличительных черт одной темы от другой.

В соответствии с подходом, изложенным в работе [8], в данной статье вероятности униграмм и биграмм вычисляются путем деления количества документов, в которых встретилась та или иная униграмма или биграмма, на число всех документов в коллекции. Другой вариант вычисления меры согласованности тем на основе логарифма от условной вероятности (TC-LCP), предложенный в работе [8], не рассматривается, поскольку в работе [7] было показано, что этот вариант работает хуже, чем TC-PMI.

4. ДОБАВЛЕНИЕ БИГРАММ В ТЕМАТИЧЕСКИЕ МОДЕЛИ

На первом этапе экспериментов исследовалось, может ли улучшиться качество тематической модели путем добавления в неё биграмм в качестве отдельных элементов словаря. Для этой цели были извлечены все биграммы, встретившиеся в коллекции, с частотностью не меньше 5. Для последующего упорядочения извлечённых биграмм применялись ассоциативные меры - математические критерии, определяющие силу связи между составными частями фраз, основываясь на частотах встречаемости отдельных слов и словосочетаний целиком. В экспериментах были использованы следующие 15 ассоциативных мер: Взаимная Информация (MI), Дополненная MI, Кубическая MI, Нормализованная MI, Настоящая MI, Коэффициент Dice (DC), Модифицированный DC, T-Score, Симметричная Условная Вероятность, Коэффициент Простого Соответствия, Коэффициент Кульчинского, Коэффициент Юла, Хи-Квадрат, Отношение логарифмического правдоподобия и Лексическая Связность.

В соответствии с результатами [7] в тематические модели добавлялись топ-1000 биграмм для каждой ассоциативной меры. Так, в каждом эксперименте к словарю в качестве отдельных элементов добавлялись топ-1000 биграмм, и в каждом документе, содержащем любые из добавляемых словосочетаний, из частот образующих их униграмм вычитались

частоты биграмм, а сами словосочетания добавлялись в его разреженное представление. Отдельно следует отметить, что во всех экспериментах число тем фиксировалось рав-Таблица 1

Результаты добавления биграмм

Ассоциативная	Перплексия	TC-PMI
мера		
Оригинальный	1694	86.4
PLSA		
MI	1683	79.2
Настоящая MI	2162	110.7
Кубическая МІ	2000	95
DC	1777	89.6
Модифицирован-	2134	94.1
ный DC		
T-Score	2189	104.9
Лексическая	1928	101.3
связность		
Хи-квадрат	1763	89.6

ным 100.

Хотя эксперименты были проведены для всех 15 упомянутых выше ассоциативных мер, в табл. 1 представлены только наиболее характерные результаты добавления топ-1000 биграмм наряду с результатом оригинального алгоритма PLSA без добавления биграмм.

Как видно, добавление топ-1000 биграмм, упорядоченных по той или иной ассоциативной мере, как правило, приводит к увеличению размера словаря и, следовательно, ухудшению перплексии, в то время как согласованность тем становится лучше. Эти выводы полностью согласуются с результатами, описанными в работе [7]. Однако, используя некоторые ассоциативные меры (например, Взаимную Информацию), можно получить немного лучше перплексию, но чуть хуже согласованность тем, что обусловлено добавлением нестандартных и низкочастотных биграмм.

5. ДОБАВЛЕНИЕ СХОЖИХ УНИГРАММ И БИГРАММ В ТЕМАТИЧЕСКИЕ МОДЕЛИ

5.1. Добавление схожих униграмм в тематические модели

Оригинальные тематические модели (PLSA и LDA) используют модель «мешка слов», предполагающую независимость слов друг от друга. Однако в документах есть много слов, связанных между собой по смыслу – в частности, однокоренные слова, например: банк – банковский – банкир. Поэтому на следующем этапе экспериментов исследовалась возможность учета в тематических моделях подобных похожих слов – а именно, слов, начинающихся с одних и тех же букв.

Для данной цели был модифицирован оригинальный алгоритм PLSA. При описании проведённой модификации будет использоваться описание алгоритма PLSA, представленное в работе [9], и следующие обозначения:

- D коллекция документов;
- T множество полученных тем;
- $\bullet \ W$ словарь (множество уникальных слов в коллекции документов D);
- $\Phi = \{ \varphi_{wt} = p(w \mid t) \}$ распределение слов w по темам t;
- $\Theta = \{\theta_{td} = p(t \mid d)\}$ распределение тем t по документам d;
- $S = \{S_w\}$ множество похожих слов, где S_w множество слов, похожих на слово w;
- n_{dw} и n_{ds} частотности слов w и s в документе d;
 - n_{wt} оценка частотности слова w в теме t;
- n_{td} оценка частотности темы t в документе d;
- n_t оценка частотности темы t в коллекции документов D.

Псевдокод алгоритма PLSA-SIM представлен в Алгоритме 1. Единственная модификация оригинального алгоритма PLSA касается строчки 6, где в рассмотрение добавляются предварительно вычисленные множества похожих слов (в оригинальном алгоритме данная строчка отсутствует, а в строчке 9 вместо f_{dw} используется n_{dw}). Тем самым вес подоб-

Вход: коллекция текстов D, число тем |T|, множества похожих слов S**Выход**: распределения Φ и Θ

While не выполнится критерий остановки do

ппе не выполнится критерии остановки **do** For
$$d \in D$$
, $w \in W$, $t \in T$ **do** $n_{wt} = n_{td} = n_t = 0$ For $d \in D$, $w \in d$ **do** $Z = \sum_{t \in T} \varphi_{wt} \theta_{td}$, $f_{dw} = n_{dw} + \sum_{s \in S_w} n_{ds}$ For $t \in T$ **do** $\delta = \frac{f_{dw} \varphi_{wt} \theta_{td}}{Z}$, $n_{wt} = n_{wt} + \delta$, $n_{td} = n_{td} + \delta$, $n_t = n_t + \delta$ For $d \in D$, $w \in W$, $t \in T$ **do** $\varphi_{wt} = \frac{n_{wt}}{n_t}$, $\theta_{td} = \frac{n_{td}}{n_t}$

ных слов увеличивается в каждом документе коллекции.

Поскольку в русском языке достаточно богатая морфология, а темы в основном задаются именными группами, в качестве кандидатов в похожие слова рассматривались только сущ. и прил. В табл. 2 представлены результаты добавления похожих слов в тематические модели наряду с оригинальным алгоритмом PLSA.

Таблица 2 Результаты добавления похожих слов

1 сзультиты обоивления похожих слов			
Число одинаковых букв	Перплексия	TC-PMI	
0 букв (PLSA)	1694	86.4	
2 буквы	1852	187.2	
3 буквы	1565	432.9	
4 буквы	1434	2432.3	
5 букв	1620	2445.3	
6 букв	1610	1310.9	

Как видно, наилучшие результаты показывает модель, рассматривающая в качестве похожих слова, начинающиеся с 4 одинаковых букв. Однако в русском языке есть множество приставок длины в 4 буквы и больше. Учитывая это, был составлен список из 43 наиболее широко использующихся таких приставок (анти-, гипер-, пере- и др.) и введён дополнительный критерий: если слова начинаются на одну и ту же приставку, то они считаются похожими, если следующая буква после приставки также совпадает. Данный критерий позволил еще больше снизить перплексию до 1376 и оставить согласованность тем примерно на лучшем уровне - 2250. В дальнейших экспериментах, описываемых в данной статье, было решено использовать именно эти 2 критерия.

Следует отметить, что в результате добавления знаний о похожести слов в тематические модели такие слова с большей вероятностью окажутся в топ-10 в полученных темах. Тем самым происходит неявная максимизация меры TC-PMI, поскольку похожие слова склонны встречаться в одних и тех же документах. Поэтому было принято решение модифицировать данную метрику для учета не всех топ-10 слов, а только топ-10 непохожих слов в темах (в дальнейшем в статье данная метрика будет обозначаться как TC - PMI - nSIM). В табл. 3 подытожены результаты добавления похожих слов в тематические модели с использованием описанных выше критериев и введённой новой метрики.

Таблица 3 Результаты наилучших способов добавления похожих слов

Алгоритм	Перплексия	TC-PMI-nSIM
Исходный	1694	78.3
PLSA		
PLSA-SIM	1376	87.8

Как видно, модифицированная версия алгоритма PLSA-SIM показывает результаты лучше оригинального алгоритма PLSA по обеим целевым метрикам.

5.2. Добавление схожих биграмм в тематические модели

Для применения подхода, представленного в разделе 5.1 к топ-1000 биграммам, упорядоченными в соответствии с различными ассоциативными мерами, описанными в разделе 4, было решено ввести дополнительный критерий схожести биграмм и униграмм. Биграмма (w_1, w_2) считается похожей на униграмму w_3 , если выполнен один из следующих критериев:

- слово w_3 похоже на w_1 или w_2 в соответствии с критериями, описанными в разделе 5.1;
- слово w_3 совпадает с w_1 или w_2 и длина w_3 больше трех букв.

Хотя эксперименты были проведены для всех ассоциативных мер, описанных в разделе 4, в табл. 4 представлены только наиболее характерные результаты интеграции биграмм и добавлению похожести униграмм и биграмм наряду с результатами алгоритмов PLSA и PLSA-SIM.

Как видно, добавление в тематическую модель похожих униграмм и топ-1000 биграмм, упорядоченных в соответствии с большинством ассоциативных мер, приводит к улучшению качества получающихся тем по сравнению с алгоритмом PLSA-SIM.

6. ИТЕРАТИВНЫЙ АЛГОРИТМ ДЛЯ ВЫБОРА НАИБОЛЕЕ ПОДХОДЯЩИХ БИГРАММ

На последнем этапе экспериментов было сделано предположение, что темы могут сами выбирать себе наиболее подходящие биграммы. Для проверки данной гипотезы был предложен новый итеративный алгоритм выбора биграмм исходя из вида верхушек тем.

При описании предлагаемого алгоритма будут использоваться следующие дополнительные обозначения:

B – множество всех биграмм в коллекции документов D;

 $B_{\scriptscriptstyle A}$ – множество биграмм, добавленных в тематическую модель;

 $S_{\scriptscriptstyle A}$ – множество потенциальных кандидатов на похожие слова;

 $(u_1^t,...,u_{10}^t)$ – топ-10 униграмм в теме t; $f(u_i^t,u_i^t)$ – частота биграммы (u_i^t,u_i^t) .

Псевдокод предлагаемого алгоритма представлен в Алгоритме 2. На каждой итерации алгоритм добавляет в множество кандидатов в похожие слова топ-10 униграмм из каждой темы. Также в это же множество и в саму тематическую модель добавляются все биграммы, которые могут быть образованы с помощью этих топ-10 униграмм. Было принято решение анализировать только первые топ-10 слов в темах, поскольку одной из целевой метрик является согласованность тем, использующая

Таблица 4

Результаты добавления похожих униграмм и биграмм

Алгоритм	Перплексия	TC-PMI-nSIM
PLSA	1694	78.3
PLSA-SIM	1376	87.8
PLSA-SIM + MI	1411	106.2
PLSA-SIM + Настоящая MI	1204	177.8
PLSA-SIM + Кубическая MI	1186	151.7
PLSA-SIM + DC	1288	99
PLSA-SIM + Модифицированный DC	1163	156.2
PLSA-SIM + T-Score	1222	171.5
PLSA-SIM + Лексическая Связность	1208	125.6
PLSA-SIM + Хи-квадрат	1346	122.9

Итеративный алгоритм

Вход: коллекция документов D, число тем |T|, множество биграмм B Выход: полученные темы T Запуск оригинального PLSA на коллекции текстов D для получения тем T $B_A = \emptyset$ While не выполнится критерий остановки \mathbf{do} $S_A = \emptyset$ For $t \in T$ do $S_A = S_A \cup \{u_1^t, \dots, u_{10}^t\}$ For $u_i^t, u_j^t \in (u_1^t, \dots, u_{10}^t)$ \mathbf{do} If $(u_i^t, u_i^t) \in B$ and $f(u_i^t, u_i^t) > f(u_i^t, u_i^t)$ \mathbf{do}

 $S_{\scriptscriptstyle A} = S_{\scriptscriptstyle A} \cup B_{\scriptscriptstyle A}$ Запуск PLSA-SIM с множествами $S_{\scriptscriptstyle A}$ и $B_{\scriptscriptstyle A}$ для получения тем T

именно это множество (см. определение метрики в разделе 3). В соответствии с данным алгоритмом темы могут выбирать себе только те биграммы, которые образуются с помощью топ-10 униграмм в темах, а такие биграммы с большей вероятностью могут оказаться наиболее подходящими.

 $B_A = B_A \cup \{(u_i^t, u_i^t)\}$

В табл. 5 представлены первые несколько итераций предложенного итеративного алгоритма наряду с результатами оригинального алгоритма PLSA (в таблице обозначен как нулевая итерация).

Таблица 5 Результаты итеративного алгоритма

Итерация	Перплексия	TC-PMI-nSIM
0 (PLSA)	1694	78.3
1	936	180.5
2	934	210.2
3	933	230
4	940	235.8
5	931	193.5

Как видно, после первой итерации наблюдается существенное улучшение качества получаемых тем по обеим целевым метрикам. Однако на следующих итерациях результаты начинают колебаться вокруг примерно тех же самых уровней перплексии и согласованности тем (с незначительным улучшением последней). Поэтому мы считаем, что согласно

результатам первой итерации выбор необходимых биграмм и кандидатов в похожие слова самими темами приводит к наилучшим значениям перплексии и согласованности тем.

8. ЗАКЛЮЧЕНИЕ

В работе представлены эксперименты по добавлению биграмм в тематические модели. Эксперименты, проведённые на русскоязычных статьях из электронных банковских журналов, показывают, что большинство ассоциативных мер упорядочивает биграммы таким образом, что при добавлении верхушки этих списков в тематические модели ухудшается перплексия и улучшается согласованность тем. Затем в статье предлагается новый алгоритм PLSA-SIM, добавляющий схожесть униграмм и биграмм в тематические модели. Проведённые эксперименты показывают значительное улучшение перплексии и согласованности тем для этого алгоритма. В конце статьи предлагается еще один новый итеративный алгоритм, основанный на идее, что темы сами могут выбирать себе наиболее подходящие биграммы и похожие слова. Эксперименты показывают дальнейшее улучшение качества по обеим целевым метрикам.

СПИСОК ЛИТЕРАТУРЫ

- 1. *Blei D.*, *Ng A. and Jordan M.* Latent Dirichlet Allocation // Journal of Machine Learning Research. № 3, 2003, P. 993–1002.
- 2. *Hofmann T.* Probabilistic Latent Semantic Indexing // In the Proceedings of the 22nd Annual International SIGIR Conference on Research and Development in Information Retrieval. 1999, P. 50–57.
- 3. Wallach H. Topic Modeling: beyond bag-of-words // In the Proceedings of the 23rd International Conference on Machine Learning. 2006, P. 977–984.
- 4. *Griffiths T., Steyvers M. and Tenenbaum J.* Topics in semantic representation // Psychological Review. 144, 2, 2007, P. 211–244.
- 5. Wang X., McCallum A. and Wei X. Topical n-grams: Phrase and topic discovery, with an application to information retrieval // In the Proceedings of the 2007 Seventh IEEE

Нокель М. А. – аспирант факультета ВМК Московского государственного университета им. М. В. Ломоносова.

Тел.: +7-926-121-61-57 E-mail: mnokel@yandex.ru

- International Conference on Data Mining. 2007, P. 697 702.
- 6. Newman D., Lau J. H., Grieser K. and Baldwin T. Automatic evaluation of topic coherence // In the Proceedings of Human Language Technologies: The 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics. P. 100-108.
- 7. Lau J. H., Baldwin T. and Newman D. On Collocations and Topic Models // In ACM Transactions on Speech and Language Processing. 10 (3), 2013, P. 1–14.
- 8. *Mimno D.*, *Wallach H.*, *Talley E.*, *Leenders M. and McCallum A.* Optimizing semantic coherence in topic models // In the Proceedings of EMNLP'2011. 2011, P. 262–272.
- 9. *Vorontsov K. and Potapenko A.* EM-like algorithms for probabilistic topic modeling // Machine Learning and Data Analysis. 2013, vol. 1 (6), P. 657–686.

Nokel M. A. – PhD student, Faculty of Computational Mathematics and Cybernetics, Moscow State University.

Tel.: +7-926-121-61-57 E-mail: mnokel@yandex.ru