

Историческая информатика

Правильная ссылка на статью:

Галушко И.Н. — Применение тематического моделирования для оптимизации процесса поиска релевантных исторических документов (на примере биржевой прессы начала XX в.) // Историческая информатика. – 2023. – № 2. DOI: 10.7256/2585-7797.2023.2.43466 EDN: SKBPNS URL: https://nbpublish.com/library_read_article.php?id=43466

Применение тематического моделирования для оптимизации процесса поиска релевантных исторических документов (на примере биржевой прессы начала XX в.)

Галушко Илья Николаевич

магистр, кафедра исторической информатики, исторический факультет, Московский государственный университет имени М.В. Ломоносова (МГУ)

119234, Россия, г. Москва, ул. Ломоносовский Проспект, 27, корп.4

✉ i.galushko15@gmail.com



[Статья из рубрики "Новые методы и технологии обработки исторических источников"](#)

DOI:

10.7256/2585-7797.2023.2.43466

EDN:

SKBPNS

Дата направления статьи в редакцию:

30-06-2023

Аннотация: Ключевой задачей представленной статьи является апробация методики анализа информационного потенциала коллекции исторических источников с помощью тематического моделирования. Некоторые современные коллекции оцифрованных исторических материалов насчитывают десятки тысяч документов, и на уровне отдельного исследователя охват всего доступного наследия представляется затруднительным. Вслед за рядом исследователей мы предполагаем, что тематическое моделирование может стать удобным инструментом предварительной оценки содержания коллекции исторических документов; инструментом отбора только тех документов, в которых присутствует информация, релевантная поставленным исследовательским задачам. В нашем случае в качестве основной коллекции исторических документов была выбрана подборка газеты «Биржевые ведомости». На данном этапе мы можем подтвердить, что в рамках нашего исследования применение тематического моделирования оказалось продуктивным решением для оптимизации процесса поиска исторических документов в объемной коллекции оцифрованных исторических материалов. В то же время необходимо подчеркнуть, что в нашей работе тематическое

моделирование применялось исключительно как прикладной инструмент ускорения поиска и первичной оценки информационного потенциала коллекции документов через анализ выделенных топики. Наш опыт показал, что по крайней мере для «Биржевых ведомостей» тематическое моделирование с использованием LDA не позволяет делать выводы с позиции применяемой нами методологии содержательного анализа. Данные наших моделей слишком фрагментарны, их можно использовать только для первичной оценки тематик информации, содержащейся в источнике.

Ключевые слова:

тематическое моделирование, латентное размещение Дирихле, Биржевые ведомости, поведенческие финансы, обработка естественного языка, распознавание исторических документов, исторические газеты, поиск исторических документов, машинное обучение, фондовый рынок

Ключевой задачей представленной статьи является апробация методики анализа информационного потенциала коллекции исторических источников с помощью тематического моделирования. Некоторые современные коллекции оцифрованных исторических материалов насчитывают десятки тысяч документов (как, например, «Электронная библиотека исторических документов», созданная Российским историческим обществом (РИО), содержит 294 тысячи распознанных исторических документов [\[1\]](#) – и на уровне отдельного исследователя охват всего доступного наследия представляется затруднительным. Вслед за рядом исследователей [\[2\]](#) мы предполагаем, что тематическое моделирование может стать удобным инструментом предварительной оценки содержания коллекции исторических документов; инструментом отбора только тех документов, в которых присутствует информация, релевантная поставленным исследовательским задачам.

Наше исследование, для которого и была разработана описываемая в статье методика, посвящено изучению доходности ценных бумаг на Санкт-Петербургской фондовой бирже в начале XX в. с позиции поведенческих финансов. Нас интересовали принципы инвестиционной оценки публичных компаний – как определялись приемлемые или недостаточные уровни капитализации; как определялись ценные бумаги, представляющие хороший выбор для помещения капиталов, насколько широко данные методики (если они существовали) применялись в практике биржевой торговли. В качестве одного из основных источников была выбрана газета «Биржевые ведомости», в ежедневных выпусках которой велась биржевая колонка, где печатался комментарий хроникера, в котором описывался настрой участников торгов и нередко приводился подробный анализ текущей ситуации в экономике Российской империи. В колонках «Биржевых ведомостей» часто встречаются аналитические заметки о доходности ценных бумаг: под какой процент размещается очередная эмиссия государственных долговых бумаг; на каком уровне относительно номинала торгуются эти бумаги; соответствует ли предлагаемый процент актуальной статистике денежного рынка и как объяснить курсовую динамику последних дней. Для исследования было решено собрать коллекцию таких заметок, чтобы на ее основе выделить устойчивые аналитические паттерны, характерные для биржевой прессы в вопросах, касающихся доходности ценных бумаг. Мы воспользовались материалами оцифрованного комплекта «Биржевых ведомостей» с сайта Российской национальной библиотеки (447 номеров за 1905 и 1913 года, утренние и вечерние выпуски). В рамках общего исследования доходности ценных бумаг на

фондовом рынке Российской империи нас интересовал ограниченный набор проблем, связанных с поведенческими и институциональными аспектами функционирования фондового рынка Российской империи. Содержание «Биржевых ведомостей», напротив, разнообразно. И, если не считать колонку биржевого хроникера, то встречаются номера, полностью лишенные нужной нам информации. Содержание таких номеров заполнено военными новостями, театральными и литературными обзорами, экономическими рассуждениями небиржевого характера и другими подобными статьями широкого профиля (см. рис. 1 и 2). И в этом контексте логично обратиться к возможностям тематического моделирования в качестве прикладного инструмента для автоматического поиска тех номеров (страниц) из нашей коллекции оцифрованных газетных материалов, которые содержат информацию, касающуюся особенностей функционирования рынка ценных бумаг.



Рисунок 1. Пример организации страницы номера «Биржевых ведомостей».

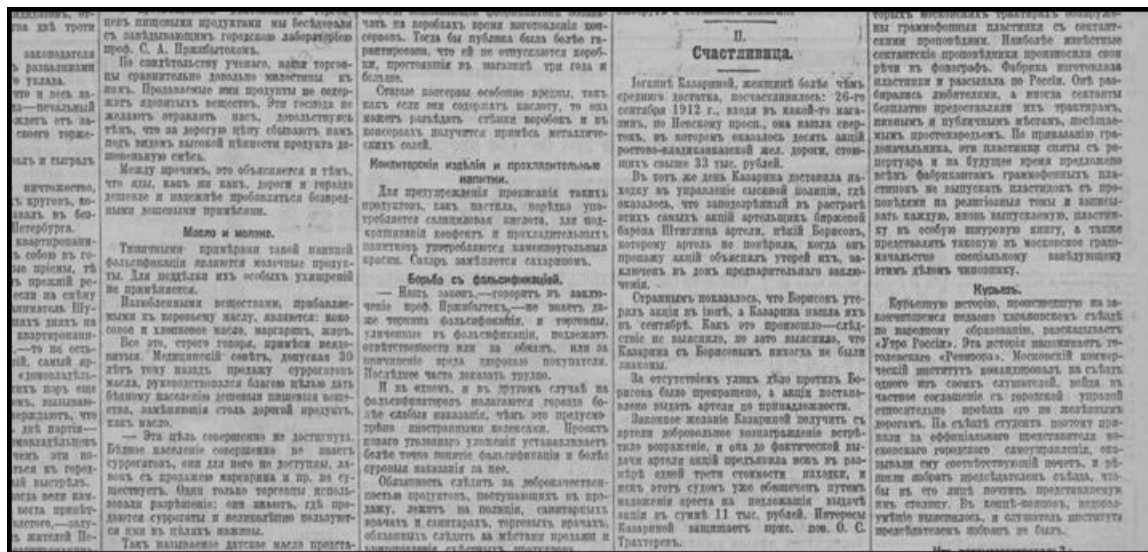


Рисунок 2. Пример искомого текста. Колонка «Счастливица». Рассказ о том, как некая женщина нашла в магазине свертков с акциями, которые потерял банкир Борисов (Выпуск № 13621 от 28 июня (11 июля)).

Тематическое моделирование — это метод машинного обучения без учителя, применяемый для определения основных тем коллекции документов (или тем предложений одного документа, который рассматривается в таком случае как совокупность предложений) на основе выделения топиков (про разницу понятий «тема» и «топик» см. Приложение 1). Как правило, топик представляет собой взвешенный по вероятности список слов, которые вместе выражают общее содержание предполагаемой темы [2]. Чем выше коэффициент слова, тем большее значение модель придает этому слову при формировании топика. Так, в Таблице 2 представлен пример двух топиков, определенных нашей моделью для третьей страницы выпуска «Биржевых ведомостей» от 29 апреля 1913 г. (№13521):

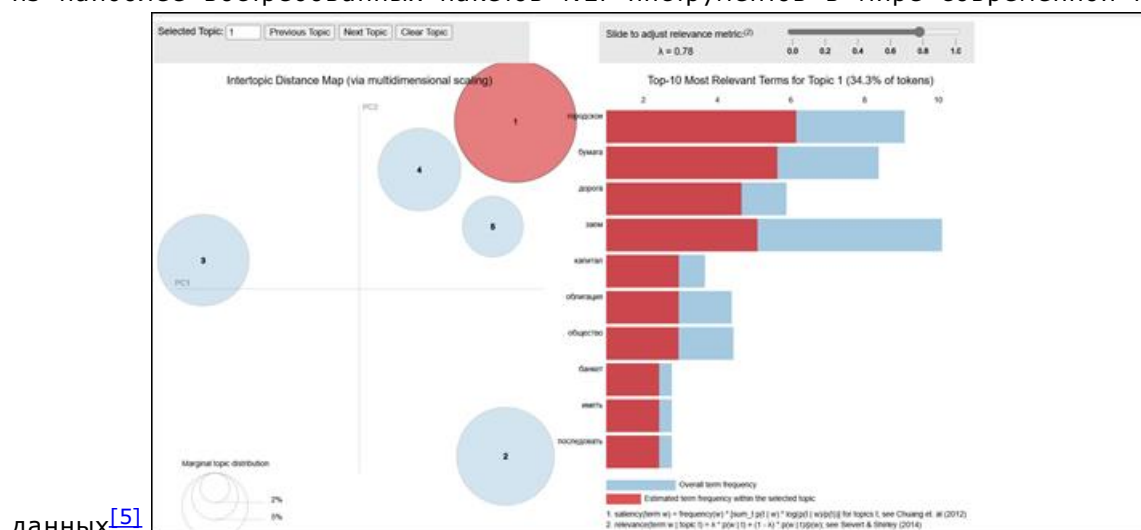
Таблица 1. Примеры двух определенных моделью топиков, представленных набором ключевых слов и их вероятностями (см. Приложение 2.1)

Топики	
(1,	'0.023*"склад" + 0.021*"акция" + 0.016*"предприятие" + 0.012*"общество" + 0.012*"пароходных" + 0.012*"товарищество" + 0.012*"транспорт" + 0.012*"страх" + 0.012*"железный" + 0.007*"баланс")
(2,	'0.015*"городской" + 0.013*"бумага" + 0.012*"заем" + 0.011*"дорога" + 0.007*"капитал" + 0.007*"облигация" + 0.007*"общество" + 0.007*"акция" + 0.006*"специальный" + 0.006*"иметь")

Одним из наиболее популярных методов тематического моделирования, используемых в настоящее время, является латентное распределение Дирихле (LDA), которое представляет собой «генеративную вероятностную модель для коллекций дискретных данных, таких как текстовые корпуса» [3]. Этот метод используется в рамках Digital Humanities для извлечения тем из набора текстов [4]. В этой модели документ (в нашем случае – отдельная страница газетного номера «Биржевых ведомостей») представляет собой смесь топиков, а топик — распределение вероятностей по словарю. Под словарем понимается список всех слов изучаемой коллекции документов – именно словарь задает модели пространство слов, в котором нужно распределить документы таким образом, чтобы сформировать заданное исследователем количество топиков. Важной особенностью LDA является заранее определенное количество топиков, которое устанавливается до начала обучения. Удобной функцией моделей LDA является установление весов, равных 0.001 для слов, определяющих избыточные топик. То есть, если исследователь изначально установил для модели поиск 5 топиков, а по результатам тематического моделирования для трех тем модель установила разные веса слов, изменяющиеся в границах, например, от 0.005 до 0.0021, а для двух других тем – вес 0.001 для всех слов – значит, LDA при заданных настройках не может распределить текст более чем на 3 топика. В ходе наших экспериментов мы установили, что для решения наших исследовательских задач достаточным является порог в 5 топиков.

Создание словаря может дополняться применением различных методов предобработки текстов. В нашей работе использовался классический для компьютерной лингвистики метод лемматизации – приведения всех слов к основной форме. Другой широко применяемый метод предобработки подразумевает поиск и добавление в словарь модели биграмм – слов с двойной основой, выделяемых в тексте по частоте совместной встречаемости (например, «акционерный_капитал», «государственный_рента»). Данный подход мы тестировали изначально, когда планировали создавать LDA модель всей

имеющейся коллекции. При таком подходе в качестве одного документа рассматривается целый номер биржевой газеты. Однако результаты оказались крайне непрактичными – ввиду слишком широкого разброса тем на один выпуск газеты мы не смогли создать интерпретируемые модели. Добавление фильтров (TF-IDF (см. Приложение 1), изменение параметров обучения (например, размера батча – порции данных для итерации обучения), увеличение количества тем (разные варианты в диапазоне от 5 до 30 топиков) – не смогли улучшить ситуацию. Тогда мы решили создавать LDA модели для отдельных номеров газеты, рассматривая их в качестве коллекции предложений. Результаты улучшились, но все еще остались недостаточными для использования тематического моделирования в качестве стабильного поискового алгоритма. Последним и наиболее продуктивным вариантом оказалось построение LDA для каждой отдельной страницы в нашей коллекции «Биржевых ведомостей». Подобный подход позволил нам создать алгоритм отбора релевантных для исследования материалов на основе тематического моделирования, но ввиду ограниченного словаря нам пришлось отказаться от создания биграмм и фильтрации по TF-IDF, что в конечном итоге никак не повлияло на качество итогового моделирования, поскольку мы не ставили задачу исчерпывающего семантического анализа; выявление тем на основе одинарных слов оказалось достаточным для отбора всех номеров газеты, обладающих информацией, релевантной нашим исследовательским задачам. В качестве основной библиотеки для реализации программного кода был выбран Python-модуль Gensim, являющийся одним из наиболее востребованных пакетов NLP-инструментов в мире современной науки о



данных [5].

Рисунок 3. Визуализация LDA модели, созданной для 3-ей страницы выпуска №13521 «Биржевых ведомостей» от 29 апреля 1913 г.

В рамках тематического моделирования эксперт должен сам определить название для топика таким образом, чтобы оно отражало семантику отдельных слов, формирующих тему. Для примера по рисунку 8 мы можем однозначно утверждать, что слова из топика №1 явно тяготеют к теме городских займов, на что указывают такие слова, как «городской», «заем», «облигация», «капитал». Архитектура существующих LDA моделей позволяет быстро анализировать содержание выделенных тем на основе поиска по ключевым словам с учетом их «веса». В качестве алгоритма поиска номеров газеты, содержащих нужную нам информацию, мы решили использовать поиск по словам, формирующим топик, с пороговым значением весов слова > 0.01 . В качестве смысловых единиц, маркирующих присутствие биржевых сведений в моделируемом документе, мы использовали слова, перечисленные в Таблице 1. Перечень паттернов Fuzzy-match и соответствующих им форм слов [6].

Для всей отобранной коллекции текстов «Биржевых ведомостей» было создано 2411 LDA моделей. Из них наш алгоритм поиска определил в группу «содержащих биржевую информацию» – 457. Разумеется, в значительной степени к таковым относятся страницы с колонкой биржевого хроникера. В качестве определенного доказательства применимости LDA-моделей в подобных источниковедческих задачах отметим, что анализ случайной выборки из 100 номеров «Биржевых ведомостей» показал, что ни одна колонка хроникера не была пропущена поисковым алгоритмом – каждая из них была помещена в сводную таблицу. В качестве иллюстрации мы приведем малую выборку из полей итоговой таблицы, включающую только те страницы и номера, в которых удалось обнаружить биржевые сведения вне колонки биржевого хроникера. Всего таких — 29 моделей. В Таблице 2 представлены примеры определенных моделью наборов топиков. Таблица 3 содержит краткое описание материалов, найденных на страницах, перечисленных в Таблице 2. В приложении мы поместили фотографии страниц, представленных в обеих таблицах.

Таблица 2. Примеры топиков и ключевых слов LDA моделей для выборки номеров из коллекции «Биржевых ведомостей» за 1913 год.

	Топик 1	Топик 2	Топик 3
Биржевые ведомости. 1913, №13521 (29 апр. (12 мая))3.	'0.021*"место" + 0.012*"завод" + 0.012*"съезд" + 0.009*"час" + 0.009*"станция" + 0.009*"октябрист" + 0.009*"академик" + 0.009*"совещание" + 0.009*"вновь" + 0.009*"вода"	+ '0.023*"склад" + 0.021*"акция" + 0.016*"предприятие" + 0.012*"общество" + 0.012*"пароходных" + 0.012*"товарищество" + 0.012*"транспорт" + 0.012*"страх" + 0.012*"железный" + 0.007*"баланс"	+ '0.015*"городской" + 0.013*"бумага" + 0.012*"заем" + 0.011*"дорога" + 0.007*"капитал" + 0.007*"облигация" + 0.007*"общество" + 0.007*"акция" + 0.006*"специальный" + 0.006*"иметь"
Биржевые ведомости. 1913, №13529 (3 (16) мая) 2.	'0.020*"акция" + 0.014*"заем" + 0.011*"рабочий" + 0.009*"май" + 0.007*"русский" + 0.007*"банк" + 0.007*"русско" + 0.007*"рынок" + 0.007*"предприятие" + 0.007*"новый"	+ '0.015*"остров" + 0.013*"венгерский" + 0.008*"новый" + 0.008*"представитель" + 0.008*"округ" + 0.006*"акция" + 0.006*"май" + 0.006*"арестовать" + 0.006*"вопрос" + 0.005*"запрос"	+ '0.010*"время" + 0.007*"май" + 0.007*"ленский" + 0.007*"товарищество" + 0.007*"рабочий" + 0.007*"договор" + 0.007*"держава" + 0.006*"последний" + 0.006*"подписать" + 0.006*"принять"
Биржевые ведомости. 1913, №13543 (11 (24) мая) 1.	'0.027*"коп" + 0.019*"опер" + 0.019*"лето" + 0.010*"мир" + 0.010*"новый" + 0.010*"торговый" + 0.010*"николаевский" + + 0.010*"эстрада" + 0.010*"шт"	+ '0.022*"год" + 0.017*"вклад" + 0.017*"вопрос" + 0.017*"прирост" + 0.012*"министр" + 0.012*"проект" + 0.012*"бумага" + 0.012*"правительство" + 0.012*"сотрудничать" + 0.012*"ответить" + +	+ '0.022*"театр" + 0.015*"касса" + 0.015*"невский" + 0.015*"май" + 0.015*"вход" + 0.015*"сад" + 0.013*"проспект" + 0.008*"выпуск" + 0.008*"день" + 0.008*"адрес"

	0.010*"дебюты"		
Биржевые ведомости. 1913, №13553 (17 (30) мая) 1.	0.011*"день"	+ 0.026*"акция"	+ 0.024*"коп"
	0.011*"новый"	+ 0.016*"собор"	+ 0.019*"сегодня"
	0.011*"друг"	+ 0.014*"час"	+ 0.016*"май"
	0.011*"специально"	+ 0.010*"день"	+ 0.013*"акция"
	0.011*"представление"	0.010*"новый"	+ 0.013*"маг"
	+ 0.011*"выпуск"	+ 0.010*"город"	+ 0.013*"невский"
	0.011*"театр"	+ 0.008*"императорский"	+ 0.013*"проспект"
	0.011*"славянский"	+ 0.008*"августейший"	+ 0.013*"театр"
	0.011*"вход"	+ 0.008*"местный"	+ 0.013*"начать"
	0.011*"сербский"	0.008*"лицо"	0.013*"сад"
Биржевые ведомости. 1913, №13565 (25 мая (7 июня)) 7.	0.014*"личный"	+ 0.013*"год"	+ 0.023*"книга"
	0.014*"мыло"	+ 0.012*"акция"	+ 0.016*"корова"
	0.014*"поэтому"	+ 0.012*"часть"	+ 0.012*"местность"
	0.008*"иметь"	+ 0.009*"иметь"	+ 0.012*"иметь"
	0.008*"убийство"	+ 0.009*"собрание"	+ 0.012*"продажа"
	0.008*"время"	+ 0.009*"общий"	+ 0.012*"право"
	0.008*"преступление"	0.009*"акционер"	+ 0.008*"лес"
	+ 0.008*"деньга"	+ 0.009*"склад" +	0.008*"невский"
	0.007*"станция"	+ 0.009*"час"	0.008*"мочь"
	0.007*"хороший"	0.009*"правление"	+ 0.008*"назначить"
Биржевые ведомости. 1913, №13587 (8 (21) июня) 2.	0.020*"пират"	+ 0.011*"комитет"	+ 0.013*"торговля"
	0.014*"заем"	+ 0.011*"неделя"	+ 0.011*"министерство"
	0.012*"акция"	+ 0.011*"коп"	+ 0.011*"капитал"
	0.010*"дорога"	+ 0.006*"иметь"	+ 0.007*"предприятие"
	0.010*"лодка"	+ 0.006*"боль"	+ 0.007*"промышленнос"
	0.008*"нефтяной"	+ 0.006*"принять"	+ 0.007*"правительство"
	0.008*"казна"	+ 0.006*"заседание"	+ 0.004*"страна"
	0.008*"железный"	+ 0.006*"денежный"	+ 0.001*"пират"
	0.006*"синдикат"	+ 0.006*"поклон"	+ 0.001*"акционерный"
	0.006*"особенно"	0.006*"разрешать"	0.001*"политика"
Биржевые ведомости. 1913, №13597 (14 (27) июня) 2.	0.016*"президент"	+ 0.022*"акция"	+ 0.018*"депутат"
	0.013*"присутствовать"	0.016*"русский"	+ 0.014*"дорога"
	0.011*"июнь"	0.013*"предприятие"	+ 0.007*"июнь"
	0.011*"духовенство"	+ 0.013*"правление"	+ 0.007*"кредитор"
	0.011*"право"	+ 0.013*"банк"	+ 0.007*"роспись"
	0.008*"согласие"	+ 0.010*"общество"	+ 0.006*"дело"
	0.008*"запрос"	+ 0.010*"дивиденд"	+ 0.006*"принять"
	0.008*"мир"	+ 0.010*"доклад"	+ 0.006*"также"
	0.008*"выбор"	+ 0.010*"зав"	+ 0.006*"член"
	0.006*"дело"	+ 0.010*"русско")	0.006*"вчера")
Биржевые ведомости. 1913, №13621 (28 июня (11 июля)) 3.	0.010*"пристав"	+ 0.016*"акция"	+ 0.012*"июнь"
	0.009*"дело"	+ 0.011*"рубль"	+ 0.008*"год"
	0.009*"московский"	+ 0.007*"день"	+ 0.008*"общественный"
	0.008*"справка"	+ 0.007*"дело"	+ 0.005*"городской"
	0.007*"кислота"	+ 0.007*"чиновник"	+ 0.005*"закон" +
	0.006*"фальсификация"	0.007*"артели"	+ 0.005*"выпуск"
	0.006*"проц"	+ 0.005*"пристав"	+ 0.005*"пройти"
		0.005*"оказаться"	

	0.006*"салициловый" + 0.005*"пластинка" + 0.005*"чиновник"	0.005*"помощник" + 0.005*"служба"	0.005*"управление" 0.005*"жизнь" 0.005*"издать"
Биржевые ведомости. 1913, № 13696 (12 (25) авг.) 1.	(0.012*"акционер" 0.012*"день" 0.012*"помещение" 0.012*"акция" 0.012*"кризис" 0.012*"класс" 0.007*"сведение" 0.007*"третий" 0.007*"выпуск" 0.007*"подлинный")	+ 0.009*"дипломатический" + 0.008*"август" + 0.008*"театр" + 0.008*"дело" + 0.008*"великий" + 0.008*"июль" + 0.008*"иметь" + 0.008*"переговоры" + 0.008*"вопрос" 0.008*"близкий")	0.002*"дипломатическ 0.002*"переговоры" + 0.002*"близкий" + 0.002*"вопрос" + 0.002*"посол" + 0.002*"союз" + 0.002*"круг" + 0.002*"великий" + 0.002*"балканский" + 0.002*"сад"

Таблица 3. Содержание отобранных примеров с указателем на приложение 2 (соответствующие сканы страниц «Биржевых ведомостей»).

Биржевые ведомости. 1913, № 13521 (29 апр. (12 мая)). 3 страница Приложение 2.1	Колонка «Может ли город обойтись без ссуды у банков» в разделе «Государственная Дума». На странице мы находим подробный доклад гласного петербургской городской думы по вопросу об организации очередного займа для покрытия муниципальных расходов Санкт-Петербурга. Автор высказывает сомнение в необходимости нового выпуска и считает, что продажа ценных бумаг, имеющих на балансе городской думы, сможет обеспечить финансовые потребности столицы. К докладу прилагаются подробные статистические выкладки со ссылкой на «отчет с.-петербургского городского общественного управления за 1911 г.»
Биржевые ведомости. 1913, № 13529 (3 (16) мая). 2 страница Приложение 2.2	Колонка «Разоблачение Керенского», посвященная оправданию забастовки Ленских рабочих. В тексте члены правления ленского товарищества (а также неназванные акционеры) обвиняются в использовании административных связей для сокрытия тяжелых условий труда рабочих.
Биржевые ведомости. 1913, № 13543 (11 (24) мая). 1 страница Приложение 2.3	Колонка «Бюджетные прения» в разделе «Государственная Дума». Приводятся основные тезисы речи А.И. Шингарева о медлительности железнодорожного строительства. Далее следует оценка недавних слов министра финансов В.Н. Коковцова о текущем состоянии фондового рынка, удержавшегося на своих уровнях в контексте общего падения цен на мировых рынках.
Биржевые ведомости. 1913, № 13553 (17 (30) мая). 1 страница	Объявление правления Волжского акционерного общества маслособойных и химических заводов «Салолин» об открытии подписки на акции дополнительного выпуска с приведением условий участия.

Приложение 2.4	
Биржевые ведомости. 1913, № 13565 (25 мая (7 июня)). 7 страница	Объявление правления Акционерного общества С.-Петербургских товарных складов о проведении чрезвычайного общего собрания акционеров с повесткой мероприятия.
Приложение 2.5	
Биржевые ведомости. 1913, № 13587 (8 (21) июня). 2 страница	Колонка «Речь Коновалова» в разделе «Государственная Дума». А.И. Коновалов поднял проблему несовершенства акционерного законодательства. Депутат говорил о длительности процедуры разрешения акционерных предприятий и высказывался против ограничений, устанавливаемых при разрешении акционерных обществ для поляков и евреев.
Приложение 2.6	
Биржевые ведомости. 1913, № 13597 (14 (27) июня). 2 страница	Обращение к владельцам облигаций Сестрорецкой железной дороги председателя конкурсного управления, юрисконсультанта министерства финансов К.К. Дыновского: сообщение о ликвидации текущего предприятия и организации продажи дороги в собственность другому обществу, готовому взять на себя обязательства по реорганизации дороги. Приводится сообщение представителя группы кредиторов дороги Л.Л. Балинского о том, что купоны не оплачивались дорогой с конца 1906 г.
Приложение 2.7	
Биржевые ведомости. 1913, № 13621 (28 июня (11 июля)). 3 страница	Колонка «Счастливица». Рассказ о том, как некая женщина нашла в магазине сверток с акциями, которые потерял банкир Борисов. Этот банкир находился в доме предварительного заключения по подозрению в растрате этих акций. Редакция высказывает подозрение, поскольку «утерянные» в июне акции были внезапно найдены в магазине в сентябре.
Приложение 2.8	
Биржевые ведомости. 1913, № 13696 (12 (25) авг.) 1 страница	Объявление от правления акционерного общества «Северный ломбард» об обмене временных свидетельств на подлинные акции.
Приложение 2.9	

Анализ материалов, найденных с помощью LDA моделей, был включен в основное содержание нашего исследования, посвященного изучению поведенческих практик, связанных с анализом доходности ценных бумаг. Таким образом, был реализован подход, при котором содержательные задачи исследования решались традиционными историческими методами, предполагающими детальное изучение текста источника и встраивание отдельного документа в единую систему с другим имеющимся материалом – архивными источниками и сочинениями биржевых практиков [\[7\]](#).

На данном этапе мы можем подтвердить, что в рамках нашего исследования применение тематического моделирования оказалось продуктивным решением для оптимизации процесса поиска исторических документов в объемной коллекции оцифрованных исторических материалов. В то же время необходимо подчеркнуть, что в нашей работе тематическое моделирование применялось исключительно как прикладной инструмент ускорения поиска и первичной оценки информационного потенциала коллекции документов через анализ выделенных топики. Наш опыт показал, что по крайней мере для «Биржевых ведомостей» тематическое моделирование с использованием LDA из библиотеки Gensim не позволяет делать выводы с позиции применяемой нами методологии содержательного анализа, предполагающей работу с внутренним содержанием источника на уровне аналитических практик. Данные наших моделей слишком фрагментарны, их можно использовать только для первичной оценки тематик информации, содержащейся в источнике. Безусловно, мы должны учитывать возможность дальнейшего усложнения моделей через применение аддитивной регуляризации [\[8\]](#) и совершенствование OCR-распознавания в сторону разбиения страницы на документы по каждой отдельной колонке, – что может существенно повысить способность модели улавливать более глубокую семантику газетного текста. Подводя итоги, мы бы хотели отметить прикладную значимость продолжения исследований практической применимости тематического моделирования в решении задач источниковедения.

Для интересующихся читателей мы оставляем ссылку на доступ к предложенному автором программному коду и набору исходных данных [\[9\]](#).

Приложение 1: Словарь определений

Тематическое моделирование – это технология статистического анализа текстов для автоматического выявления тематики в больших коллекциях документов. Тематическая модель определяет, к каким темам относится каждый документ, и какими словами описывается каждая тема. Для этого не требуется ручная разметка текстов, обучение модели происходит без учителя. Этот процесс можно сравнить с кластеризацией, но тематическая кластеризация является «мягкой» и допускает, чтобы документ относился к нескольким кластерам-темам. Тематическое моделирование не претендует на понимание смысла текста, однако оно способно отвечать на вопросы «о чём этот текст» или «какие общие темы есть у этих текстов» [\[10\]](#).

Применяя тематическое моделирование, важно проводить разграничение между топики и, собственно, темами. Топики являются результатом статистической обработки коллекции документов и состоят из слов, которым в зависимости от выбранной статистической модели (LDA, LSA, BertTopic и т.д.) была присвоена определенная значимость, на основе которой моделью было сделано предположение, что эти топики (кластер значимых слов) формируют тему документа на семантическом уровне. Однако важно подчеркнуть, что процедура признания, что определенный топик действительно отражает присутствующую в тексте документа тему, является сугубо экспертной. Исследователь определяет, насколько семантически релевантным текстовой коллекции оказался сформированный моделью набор топики.

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 частоте употребления этого слова в документе и обратно

пропорционален частоте употребления слова во всех документах коллекции.

ПРИЛОЖЕНИЕ 2

Страницы газеты «Биржевые ведомости», выявленные методом автоматизированного поиска по содержанию комплекта LDA моделей



Рисунок 2. 1. Биржевые ведомости. 1913, № 13521 (29 апр. (12 мая)). Третья страница



Рисунок 2. 2. Биржевые ведомости. 1913, № 13529 (3 (16) мая). Вторая страница



Рисунок 2. 3. Биржевые ведомости. 1913, № 13543 (11 (24) мая). Первая страница

Рисунок 2. 4. Биржевые ведомости. 1913, № 13553 (17 (30) мая). Вторая страница

Рисунок 2. 5. Биржевые ведомости. 1913, № 13565 (25 мая (7 июня)). Седьмая страница



Рисунок 2. 6. Биржевые ведомости. 1913, № 13587 (8 (21) июня). Вторая страница



Рисунок 2. 7. Биржевые ведомости. 1913, № 13597 (14 (27) июня). Вторая страница



Рисунок 2. 8. Биржевые ведомости. 1913, № 13621 (28 июня (11 июля)). Третья страница

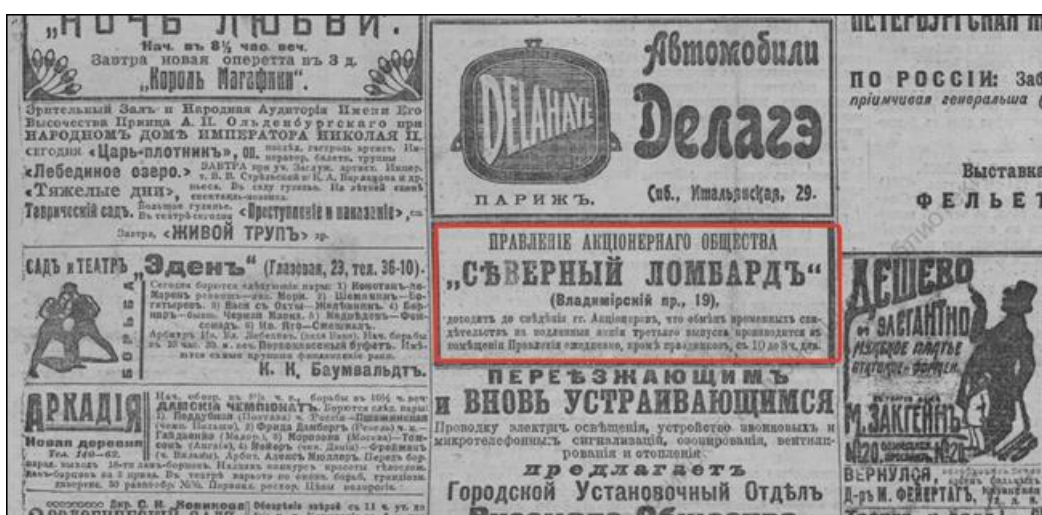


Рисунок 2. 9. Биржевые ведомости. 1913, № 13696 (12 (25) авг.). 1 страница

Библиография

1. URL: <http://docs.historyrussia.org/ru/nodes/1-glavnaya>
2. Tze-I Yang, A.J.Torget, R.Mihalcea (2011). Topic modeling in historical newspapers.
3. Marjanen, J., Zosa, E., Hengchen, S., Pivovarov, L., & Tolonen, M. (2020). Topic Modelling Discourse Dynamics in Historical Newspapers. DHN Post-Proceedings.
4. Koentges, Thomas (2020). Measuring Philosophy in the First Thousand Years of Greek Literature.
5. Egger, Roman (2020). A Topic Modeling Comparison Between LDA, NMF, Top2Vec, and BERTopic to Demystify Twitter Posts.
6. Галушко И.Н. Корректировка результатов OCR-распознавания текста исторического источника с помощью нечетких множеств (на примере газеты начала XX века) // Историческая информатика. – 2023. – № 1. – С. 102-113.
7. Представленная статья является частью моей магистерской диссертации по теме:

«Поведенческие аспекты анализа доходности ценных бумаг на фондовом рынке Российской империи в начале XX века: контент-анализ биржевых нарративов». Найденные LDA-алгоритмом выпуски «Биржевых ведомостей» в данной работе рассматривались в сочетании с материалами фонда №143 ЦГАМ (Московский биржевой комитет) и трудами биржевых практиков начала XX в. (Васильев А.А. Биржевая спекуляция, теория и практика. СПб., 1912.).

8. Воронцов К. В. Вероятностное тематическое моделирование: теория, модели, алгоритмы и проект BigARTM. 2020.
9. GitHub. URL: <https://github.com/iodinesky/Topic-modeling-in-historical-newspapers>
10. Воронцов К. В. Вероятностное тематическое моделирование: теория регуляризации ARTM и библиотека с открытым кодом BigARTM. 2023.

Результаты процедуры рецензирования статьи

В связи с политикой двойного слепого рецензирования личность рецензента не раскрывается.

Со списком рецензентов издательства можно ознакомиться [здесь](#).

Рецензируемая статья посвящена анализу информационного потенциала коллекции исторических источников с помощью тематического моделирования. В данном случае тематическое моделирование, пригодное для различных направлений анализа текстов, используется как инструмент предварительной оценки содержания коллекции исторических документов для отбора текстов, релевантных исследовательским запросам. В качестве массива текстов, на котором проводится апробация предлагаемой методики, используется пресса (газета «Биржевые ведомости» за 1905 и 1913 гг.).

Для исследования были взяты оцифрованные комплекты газет с сайта Российской национальной библиотеки (447 номеров). Автором была поставлена задача использования тематического моделирования для автоматического поиска тех газетных страниц, где есть информация о функционировании рынка ценных бумаг. В процессе поиска использовался хорошо апробированный в компьютерной лингвистике метод лемматизации. Использовались вероятностные LDA-модели, часть из которых была автоматически отобрана для дальнейшего анализа как содержащая биржевую информацию. Далее проводился содержательный анализ отобранного материала.

Актуальность работы не вызывает сомнений, поскольку поиск и отбор необходимой для дальнейшего исследования информации является на сегодняшний день серьезной проблемой, на решение которой уходит огромное количество времени. Любой способ адекватного решения подобных вопросов – это большой вклад в научно-исследовательскую практику.

Научная новизна работа также очевидна. Тематическое моделирование почти не освоено в исторической науке, появляются только первые опыты его использования, а вопрос о степени его полезности в рамках творческой лаборатории профессионального историка до сих пор остается открытым.

Статья построена не вполне традиционным образом, поскольку носит во многом экспериментальный характер. После постановки проблемы практически сразу начинается логический раздел, посвященный созданию методики использования тематического моделирования для поиска информации. Далее анализируются результаты поиска. Показательно, что, с одной стороны, автор констатирует продуктивность созданной методики для оптимизации поиска, с другой, – подчеркивает, что полученные модели слишком фрагментарны и их можно использовать только для первичной оценки тематики информации источника. Статья дополнена примерами

топиков и ключевых слов, фрагментами отобранных текстов, словарем определений и фотографиями газетных страниц. Следует отметить, что безусловно интересная и новаторская статья рассчитана на подготовленного читателя, знакомого с основами анализа текстов. В то же время продуманный стиль статьи облегчает понимание довольно сложных и не всегда привычных для традиционного взгляда историка вещей, о которых идет речь.

Библиография статьи содержит достаточный для подобных исследований список, хотя, думается, что не лишним было бы добавить работ на русском языке.

Статья фактически является приглашением к обмену мнениями и дискуссиям по рассматриваемой проблематике, скорее всего, ей обеспечена хорошая цитируемость в силу актуальности рассмотренных вопросов.

Публикация исследования, связанного с методами машинного обучения, рассчитана на определенный круг читателей, который, безусловно, будет достаточно широким, поскольку любая работа, связанная с интеллектуальной обработкой данных, вызывает сегодня большой интерес научной общественности. Статья рекомендуется к публикации.