

# Технологии и приложения тематического моделирования в цифровых гуманитарных исследованиях

Воронцов Константин Вячеславович

д.ф.-м.н., профессор РАН

Институт ИИ МГУ • ВМК МГУ • МФТИ • ФИЦ ИУ РАН

Методологический семинар МГИМО — ИСП РАН по анализу данных  
10 марта 2023

- 1 Вероятностное тематическое моделирование**
  - Математическая технология
  - Инструментарий
  - Способы и средства визуализации
- 2 Примеры приложений**
  - Тематический поиск
  - «Классификация иголок в стоге сена»
  - Темпоральные модели
- 3 Тематизатор для гуманитарных исследований (проект)**
  - Тематические модели в политологии
  - Тематические модели в исторических исследованиях
  - Проект «Тематизатор»

## Задача тематического моделирования

### Дано:

- коллекция текстовых документов

### Найти:

- $T$  — множество тем, составляющих эту коллекцию
- $\phi_{wt} = p(w|t)$  — вероятности слов  $w$  в каждой теме  $t$
- $\theta_{td} = p(t|d)$  — вероятности тем  $t$  в каждом документе  $d$

### Критерий:

- вероятностная тематическая модель  $p(w|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$  должна предсказывать появление слов  $w$  в документах  $d$ ,
- заодно максимизируя сумму регуляризаторов  $R_i(\Phi, \Theta)$ :

$$\sum_{d \in D} \sum_{w \in d} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + \sum_i \tau_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$









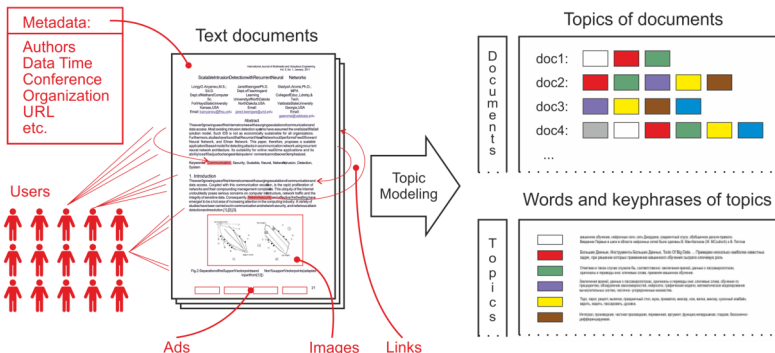




# Мультимодальные тематические модели

Тема  $t$  может порождать термины различных *модальностей*:

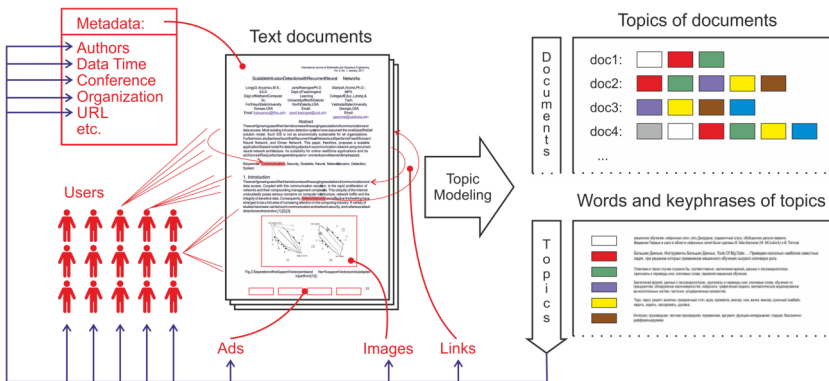
$p(\text{слово}|t)$ ,  $p(n\text{-грамма}|t)$ ,  $p(\text{автор}|t)$ ,  $p(\text{время}|t)$ ,  $p(\text{источник}|t)$ ,  
 $p(\text{объект}|t)$ ,  $p(\text{ссылка}|t)$ ,  $p(\text{баннер}|t)$ ,  $p(\text{пользователь}|t)$



# Мультимодальные тематические модели

Тема  $t$  может порождать термины различных *модальностей*:

$p(\text{слово} | t)$ ,  $p(n\text{-грамма} | t)$ ,  $p(\text{автор} | t)$ ,  $p(\text{время} | t)$ ,  $p(\text{источник} | t)$ ,  
 $p(\text{объект} | t)$ ,  $p(\text{ссылка} | t)$ ,  $p(\text{баннер} | t)$ ,  $p(\text{пользователь} | t)$



## Пример 1. Мультиязычная модель Википедии

216 175 русско-английских пар статей. Языки — модальности.  
Первые 10 слов и их частоты  $p(w|t)$  в %:

Тема №68				Тема №79			
research	4.56	институт	6.03	goals	4.48	матч	6.02
technology	3.14	университет	3.35	league	3.99	игрок	5.56
engineering	2.63	программа	3.17	club	3.76	сборная	4.51
institute	2.37	учебный	2.75	season	3.49	фк	3.25
science	1.97	технический	2.70	scored	2.72	против	3.20
program	1.60	технология	2.30	cup	2.57	клуб	3.14
education	1.44	научный	1.76	goal	2.48	футболист	2.67
campus	1.43	исследование	1.67	apps	1.74	гол	2.65
management	1.38	наука	1.64	debut	1.69	забивать	2.53
programs	1.36	образование	1.47	match	1.67	команда	2.14

Ассессор оценил 396 тем из 400 как хорошо интерпретируемые.

*Vorontsov, Frei, Apishev, Romov, Suvorova.* BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections. AIST-2015.

## Пример 1. Мультиязычная модель Википедии

216 175 русско-английских пар статей. Языки — модальности.  
Первые 10 слов и их частоты  $p(w|t)$  в %:

Тема №88				Тема №251			
opera	7.36	опера	7.82	windows	8.00	windows	6.05
conductor	1.69	оперный	3.13	microsoft	4.03	microsoft	3.76
orchestra	1.14	дирижер	2.82	server	2.93	версия	1.86
wagner	0.97	певец	1.65	software	1.38	приложение	1.86
soprano	0.78	певица	1.51	user	1.03	сервер	1.63
performance	0.78	театр	1.14	security	0.92	server	1.54
mozart	0.74	партия	1.05	mitchell	0.82	программный	1.08
sang	0.70	сопрано	0.97	oracle	0.82	пользователь	1.04
singing	0.69	вагнер	0.90	enterprise	0.78	обеспечение	1.02
operas	0.68	оркестр	0.82	users	0.78	система	0.96

Ассессор оценил 396 тем из 400 как хорошо интерпретируемые.

*Vorontsov, Frei, Apishev, Romov, Suvorova.* BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections. AIST-2015.

## Пример 2. Биграммная модель научных конференций

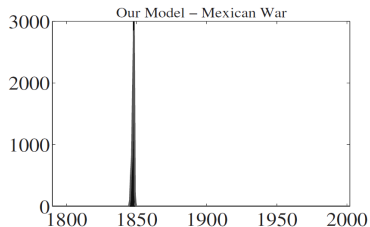
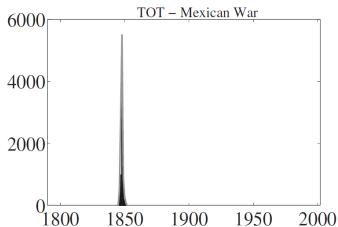
Коллекция 1000 статей конференций ММРО, ИОИ на русском

распознавание образов в биоинформатике		теория вычислительной сложности	
униграммы	биграммы	униграммы	биграммы
объект	задача распознавания	задача	разделять множества
задача	множество мотивов	множество	конечное множество
множество	система масок	подмножество	условие задачи
мотив	вторичная структура	условие	задача о покрытии
разрешимость	структура белка	класс	покрытие множества
выборка	распознавание вторичной	решение	сильный смысл
маска	состояние объекта	конечный	разделяющий комитет
распознавание	обучающая выборка	число	минимальный аффинный
информативность	оценка информативности	аффинный	аффинный комитет
состояние	множество объектов	случай	аффинный разделяющий
закономерность	разрешимость задачи	покрытие	общее положение
система	критерий разрешимости	общий	множество точек
структура	информативность мотива	пространство	случай задачи
значение	первичная структура	схема	общий случай
регулярность	тупиковое множество	комитет	задача MASC

*Сергей Стенин.* Мультиграммные аддитивно регуляризованные тематические модели // Магистерская диссертация, МФТИ, 2015.

## Пример 3. Совмещение темпоральной и $n$ -граммной модели

### Коллекция еженедельных выступлений президентов США



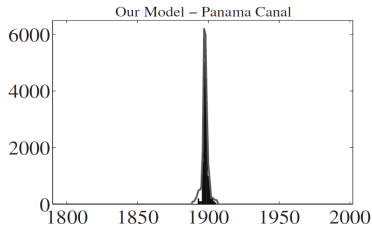
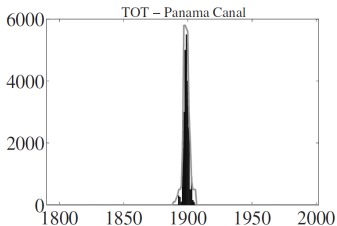
1. mexico	8. territory
2. texas	9. army
3. war	10. peace
4. mexican	11. act
5. united	12. policy
6. country	13. foreign
7. government	14. citizens

1. east bank	8. military
2. american coins	9. general herrera
3. mexican flag	10. foreign coin
4. separate independent	11. military usurper
5. american commonwealth	12. mexican treasury
6. mexican population	13. invaded texas
7. texan troops	14. veteran troops

Shoaib Jameel, Wai Lam. An  $N$ -gram topic model for time-stamped documents. 2013.

## Пример 3. Совмещение темпоральной и $n$ -граммной модели

### Коллекция еженедельных выступлений президентов США



1. government	8. spanish
2. cuba	9. island
3. islands	10. act
4. international	11. commission
5. powers	12. officers
6. gold	13. spain
7. action	14. rico

1. panama canal	8. united states senate
2. isthmian canal	9. french canal company
3. isthmus panama	10. caribbean sea
4. republic panama	11. panama canal bonds
5. united states government	12. panama
6. united states	13. american control
7. state panama	14. canal

Shoaib Jameel, Wai Lam. An  $N$ -gram topic model for time-stamped documents. 2013.

## Цели и не-цели тематического моделирования

### Цели:

- Выяснить тематическую кластерную структуру текстовой коллекции, сколько в ней тем и какие они
- Получать интерпретируемые тематические векторные представления (эмбединги) документов  $p(t|d)$ , слов  $p(t|w)$ , фрагментов и прочих объектов  $p(t|x)$
- Решать задачи поиска, категоризации, сегментации, суммаризации с помощью тематических эмбедингов

### Не-цели:

- Угадывать следующие слова (ТМ — слабые модели языка)
- Генерировать связный текст
- Понимать смысл текста



## Некоторые приложения тематического моделирования

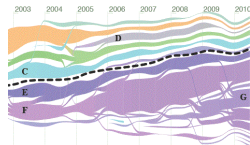
разведочный поиск в  
электронных библиотеках



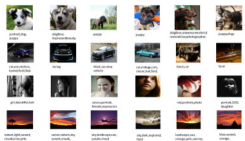
поиск тематического  
контента в соцсетях



выявление и отслеживание  
цепочек новостей



мультимодальный поиск  
текстов и изображений



анализ банковских  
транзакционных данных



управлением диалогом в  
разговорном интеллекте



*J.Boyd-Graber, Yuening Hu, D.Mimno. Applications of Topic Models. 2017.*  
*H.Jelodar et al. Latent Dirichlet allocation (LDA) and topic modeling: models, applications, a survey. 2019.*

## BigARTM: библиотека тематического моделирования

### Ключевые возможности:

- Большие данные: коллекция не хранится в памяти
- Онлайн-овый параллельный мультимодальный ARTM
- Встроенная библиотека регуляризаторов и мер качества

### Сообщество:

- Открытый код <https://github.com/bigartm>  
(discussion group, issue tracker, pull requests)
- Документация <http://bigartm.org>



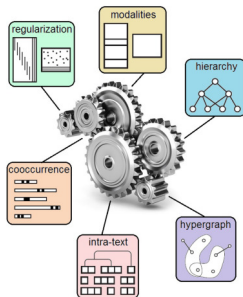
### Лицензия и среда разработки:

- Свободная коммерческая лицензия (BSD 3-Clause)
- Кросс-платформенность: Windows, Linux, MacOS (32/64 bit)
- Интерфейсы API: command-line, C++, and Python

## Ключевые возможности библиотек BigARTM и TopicNet

### BigARTM

- библиотека регуляризаторов
- мультимодальные модели
- иерархические модели
- гиперграфовые модели
- модели связности текста



### TopicNet

- Перебор сценариев регуляризации для выбора моделей
- Автоматическое протоколирование экспериментов
- Построение «банка тем» из множества моделей
- Визуализация результатов тематического моделирования

*V. Bulatov, E. Egorov, E. Veselova, D. Polyudova, V. Alekseev, A. Goncharov, K. Vorontsov.*  
TopicNet: making additive regularisation for topic modelling accessible. LREC-2020

## Качество и скорость: BigARTM vs Gensim и Vowpal Wabbit

3.7М статей Википедии, 100К слов:

время min (перплексия)

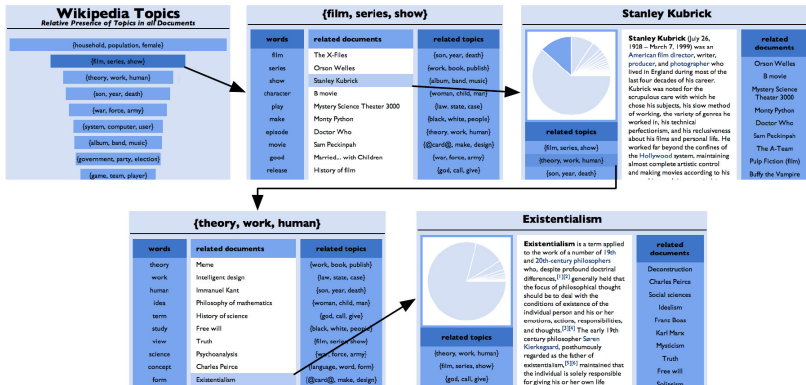
проц.	$ T $	Gensim	Vowpal Wabbit	BigARTM	BigARTM асинхрон
1	50	142m (4945)	50m (5413)	42m (5117)	25m (5131)
1	100	287m (3969)	91m (4592)	52m (4093)	32m (4133)
1	200	637m (3241)	154m (3960)	83m (3347)	53m (3362)
2	50	89m (5056)		22m (5092)	13m (5160)
2	100	143m (4012)		29m (4107)	19m (4144)
2	200	325m (3297)		47m (3347)	28m (3380)
4	50	88m (5311)		12m (5216)	7m (5353)
4	100	104m (4338)		16m (4233)	10m (4357)
4	200	315m (3583)		26m (3520)	16m (3634)
8	50	88m (6344)		8m (5648)	5m (6220)
8	100	107m (5380)		10m (4660)	6m (5119)
8	200	288m (4263)		15m (3929)	10m (4309)

*D.Kochedykov, M.Apishev, L.Golitsyn, K.Vorontsov.*

Fast and Modular Regularized Topic Modelling. FRUCT ISMW, 2017.

# Система TMVE — Topic Model Visualization Engine

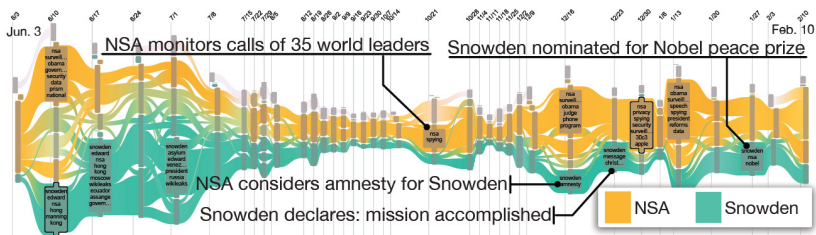
Тематический навигатор с веб-интерфейсом:



<https://github.com/ajbc/tmv>

Chaney A., Blei D. Visualizing Topic Models, 2012.

## Динамика тем: эволюция предметной области



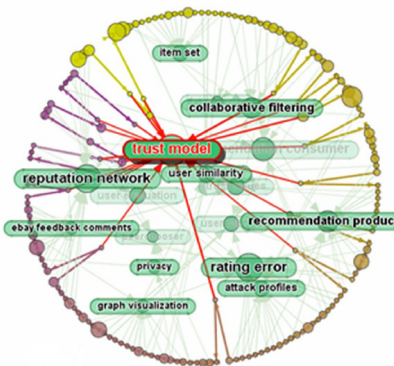
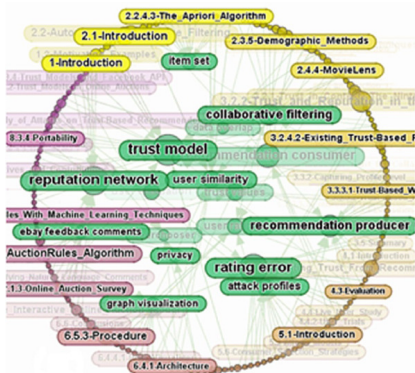
Эволюция выбранных тем иерархии. Данные Prism (2013/06/03–2014/02/09)

Стратегия визуализации в системах TextFlow и RoseRiver:

- эксперт задаёт сечение иерархии (дерева) тем,
- интерактивно выбирает подмножество тем и событий,
- получает сгенерированный отчёт с инфографикой.

*Weiwei Cui, Shixia Liu, Zhuofeng Wu, Hao Wei.* How hierarchical topics evolve in large text corpora. 2014.

## Динамика тем внутри документа: тематическая сегментация



Gretarsson B., O'Donovan J., Bostandjiev S., Hollerer T., Asuncion A., Newman D., Smyth P. TopicNets: visual analysis of large text corpora with topic modeling. ACM Trans. on Intelligent Systems and Technology. 2012.

## Визуализация иерархической тематической модели



Chirkova N. A., Vorontsov K. V. Additive regularization for hierarchical multimodal topic modeling. JMLDA, 2016.



## Визуализация иерархической тематической модели

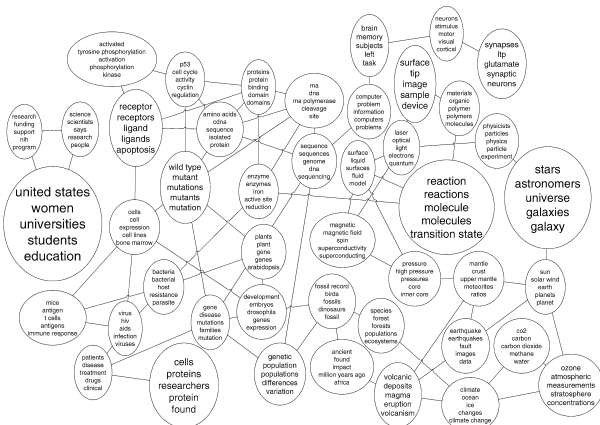


Belyy A.V., Seleznova M.S., Sholokhov A.K., Vorontsov K.V. Quality Evaluation and Improvement for Hierarchical Topic Modeling. Dialogue 2018.

# Модель коррелированных тем СТМ

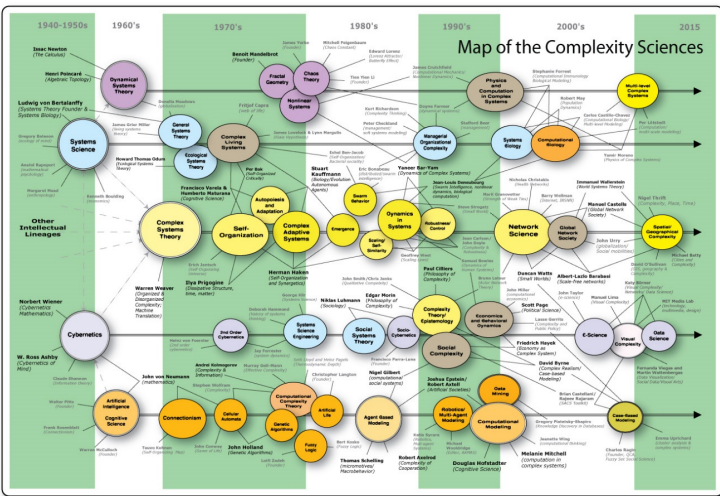
**Мотивация:** выявлять междисциплинарные связи.

Статьи по археологии чаще связаны с историей и геологией, чем с генетикой



David Blei, John Lafferty. A Correlated Topic Model of SCIENCE, 2007.

# Пример карты предметной области, построенной вручную



<http://www.theoryculturesociety.org/brian-castellani-on-the-complexity-sciences>

## Источники вдохновения: <http://textvis.lnu.se>

### Интерактивный обзор 440 средств визуализации текстов



*Shixia Liu, Weiwei Cui, Yingcai Wu, Mengchen Liu. A survey on information visualization: recent advances and challenges. 2014.*

*Айсина Р. М. Обзор средств визуализации тематических моделей коллекций текстовых документов // JMLDA, 2015.*

## Декоррелирование, сглаживание, разреживание

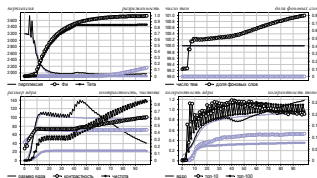
**Цель:** найти комбинацию регуляризаторов, улучшающую интерпретируемость тем по совокупности критериев.

**Регуляризаторы:**

$$\mathcal{L}\left(\begin{array}{|c|} \hline \text{PLSA} \\ \hline \Phi \quad \Theta \\ \hline \end{array}\right) + R\left(\begin{array}{|c|} \hline \text{decorrelated} \\ \hline \begin{array}{|c|} \hline \diagdown \\ \hline \end{array} \quad \square \\ \hline \end{array}\right) + R\left(\begin{array}{|c|} \hline \text{sparse} \\ \hline \begin{array}{|c|} \hline \cdot \\ \hline \end{array} \quad \begin{array}{|c|} \hline \cdot \\ \hline \end{array} \\ \hline \end{array}\right) + R\left(\begin{array}{|c|} \hline \text{background} \\ \hline \begin{array}{|c|} \hline \text{|||||} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{|||||} \\ \hline \end{array} \\ \hline \end{array}\right) \rightarrow \max$$

**Результаты:**

- разреженность  $0 \rightarrow 95\%$ , когерентность  $0.25 \rightarrow 0.96$ , чистота  $0.14 \rightarrow 0.89$ , контрастность  $0.43 \rightarrow 0.52$ ,
- без заметного ущерба для перплексии:  $1920 \rightarrow 2020$
- выработаны рекомендации по стратегии регуляризации



## Разведочный поиск в технологических блогах — 1

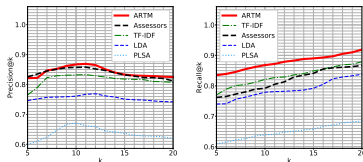
**Цель:** поиск документов по длинным текстовым запросам  
 — Habr.ru (175К документов),  
 — TechCrunch.com (760К док.).

### Регуляризаторы:

$$\mathcal{L} \left( \begin{matrix} \text{PLSA} \\ \Phi \quad \Theta \end{matrix} \right) + R \left( \begin{matrix} \text{interpretable} \\ \text{[diagram]} \end{matrix} \right) + R \left( \begin{matrix} \text{multimodal} \\ \text{[diagram]} \end{matrix} \right) + R \left( \begin{matrix} \text{n-gram} \\ \text{[diagram]} \end{matrix} \right) \rightarrow \max$$

### Результаты:

- Точность и полнота 88%, превосходит ассессоров и другие методы (tf-idf, word2vec, PLSA, LDA).
- Векторный поиск мгновенный, ассессоры тратили 5–65 мин.



*A. Ianina, L. Golitsyn, K. Vorontsov. Multi-objective topic modeling for exploratory search in tech news. AINL, 2017.*

## Разведочный поиск в технологических блогах — 2

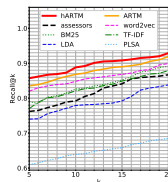
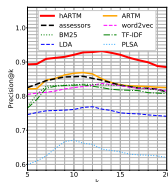
**Цель:** улучшение качества поиска с помощью иерархической тематической модели hARTM и отсекация нерелевантных тем.

**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{|c|} \hline \text{PLSA} \\ \hline \Phi \quad \Theta \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{hierarchy} \\ \hline \text{graph} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{interpretable} \\ \hline \text{matrix} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{multimodal} \\ \hline \text{img} \quad \text{text} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{n-gram} \\ \hline \text{grid} \\ \hline \end{array} \right) \rightarrow \max$$

**Результаты:**

- Точность и полнота **93%**, превосходит ассессоров и другие методы (tf-idf, BM25, word2vec, PLSA, LDA, ARTM).
- Увеличилась оптимальная размерность векторов:  
200 → 1400 (Habr.ru), 475 → 2800 (TechCrunch.com).



A. Ianina, K. Vorontsov. Regularized multimodal hierarchical topic model for document-by-document exploratory search. FRUCT-ISMW, 2019.

## Две коллекции новостей про технологии

### Habr.ru

175 143 статей на русском  
10 552 слов (униграмм)  
742 000 биграмм  
524 авторов статей  
10 000 авторов комментариев  
2546 тегов  
123 хаба (категории)

### TechCrunch.com

759 324 статей на английском  
11 523 слов (униграмм)  
1.2 млн. биграмм  
605 авторов  
184 категорий

### Предобработка текстов

- отброшены 5% наиболее частотных слов (общая лексика)
- удаление пунктуации
- нижний регистр, ё→е
- лемматизация r morphology2



# Методика оценивания качества разведочного поиска

## Поисковый запрос

набор ключевых слов или фрагментов текста, около одной страницы A4

## Поисковая выдача

документы  $d$  с распределением  $p(t|d)$ , близким к распределению  $p(t|q)$  запроса

## Два задания асессорам

- 1 найти как можно больше статей, пользуясь любыми средствами поиска (и засечь время)
- 2 оценить релевантность поисковой выдачи на том же запросе

### Поиск MapReduce

**Поиск MapReduce** – программа поиска (анализов) вычислений распределенных вычислений для больших объемов данных в рамках параллельных вычислений, представляющих собой набор Java-классов и исполняемых заданий для создания и обработки данных на параллельной обработке.

**Основные компоненты MapReduce** можно сформулировать как:

- обработка вычислений больших объемов данных;
- масштабируемость;
- автоматическое распределение заданий;
- работа по итерационным обработкам;
- автоматическая обработка отказов вычислений заданий.

**MapReduce** – популярная программная платформа (язык Java, библиотека) построения распределенных приложений для массово-параллельной обработки (задачи, задачи, ресурсы, МР) данных.

**MapReduce** включает в себе следующие компоненты:

1. HDFS – распределенная файловая система;
2. **MapReduce** – программная модель (библиотека) вычислений распределенных вычислений для больших объемов данных в рамках параллельных вычислений;

**Ключевые особенности** архитектуры **MapReduce** и структуры HDFS, такие как принцип разделения задач на кластеры, а также точки отказа. Это, в конечном итоге, определяет ограничения платформ **MapReduce** в целом. К недостаткам можно отнести:

Ограничение масштабируемости кластера **MapReduce** – не масштабируемый упор «МР» параллельных заданий.

Сильная связность **MapReduce** распределенных вычислений и элементов вычисления, реализованных распределенной программой. Как следствие:

Отсутствие поддержки алгоритмической программы вычисления распределенных вычислений в **MapReduce** поддерживается только модель вычислений **MapReduce**.

Модель вычислений, точки отказа и как следствие, необходимость использования в среде с высокими требованиями к надежности;

Проблема совместности требований по единообразному обслуживанию всех вычислительных узлов кластера при обслуживании платформ **MapReduce** (установка новых версий или пакета обновлений).

Пример запроса для разведочного поиска

## Пример: фрагмент запроса «Система IBM Watson»

IBM Watson — суперкомпьютер фирмы IBM, оснащённый вопросно-ответной системой искусственного интеллекта, созданный группой исследователей под руководством Дэвида Феруччи. Его создание — часть проекта DeepQA. Основная задача Уотсона — понимать вопросы, сформулированные на естественном языке, и находить на них ответы в базе данных. Назван в честь основателя IBM Томаса Уотсона.

IBM Watson представляет собой когнитивную систему, которая способна понимать, делать выводы и обучаться. Она также позволяет преобразовывать целые отрасли, различные направления науки и техники. Например, предсказывать появление эпидемий или возникновения очагов природных катастроф в различных регионах, вести мониторинг состояния атмосферы больших городов, оптимизировать бизнес-процессы, узнавать, какие товары будут в тренде в ближайшее время.

... ..

**Релевантные тексты:** примеры сервисов и приложений, основа которых — когнитивная платформа IBM Watson, используемые в IBM Watson технологии, вопрос-ответные системы, сопоставление IBM Watson с Wolfram-Alpha.

**Нерелевантные тексты:** общие вопросы искусственного интеллекта, другие коммерческие решения на рынке бизнес-аналитики.

## Тематика запросов разведочного поиска

Примеры заголовков разведочных запросов к Хабру  
(объём каждого запроса — около одной страницы А4):

Алгоритмы раскраски графов  
Рекомендательная система Netflix  
Методики быстрого набора текста  
Космические проекты Илона Маска  
Технологии Hadoop MapReduce  
Беспилотный автомобиль Google car  
Криптосистемы с открытым ключом  
Обзор платформ онлайн-курсов  
Data Science Meetups в Москве  
Образовательные проекты mail.ru  
Межпланетная станция New horizons  
Языковая модель word2vec

Система IBM Watson  
3D-принтеры  
CERN-кластер  
АВ-тестирование  
Облачные сервисы  
Контекстная реклама  
Марсоход Curiosity  
Видеокарты NVIDIA  
Распознавание образов  
Сервисы Google scholar  
MIT MediaLab Research  
Платформа Microsoft Azure

## Поиск и рубрикация научных публикаций на 100 языках

**Цель:** мультязыковой поиск и классификация научных публикаций по рубрикам УДК, ГРНТИ, ОЭСР, ВАК

модель	ср.ч. УДК	ср.% УДК	ср.ч. ГРНТИ	ср.% ГРНТИ
Базовая ТМ	0.558	0.165	0.536	0.220
XLM-RoBERTa	0.835	0.179	0.832	0.288
ARTM	<b>0.995</b>	<b>0.225</b>	<b>0.852</b>	<b>0.366</b>

**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{c} \text{PLSA} \\ \left[ \begin{array}{|c|} \hline \Phi \\ \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \Theta \\ \hline \end{array} \right] \end{array} \right) + R \left( \begin{array}{c} \text{interpretable} \\ \left[ \begin{array}{|c|} \hline \text{bar chart} \\ \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \text{matrix} \\ \hline \end{array} \right] \end{array} \right) + R \left( \begin{array}{c} \text{multimodal} \\ \left[ \begin{array}{|c|} \hline \text{img} \\ \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \text{text} \\ \hline \end{array} \right] \end{array} \right) + R \left( \begin{array}{c} \text{multilanguage} \\ \left[ \begin{array}{|c|} \hline \text{img} \\ \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \text{text} \\ \hline \end{array} \right] \end{array} \right) + R \left( \begin{array}{c} \text{supervised} \\ \left[ \begin{array}{|c|} \hline \text{scatter plot} \\ \hline \end{array} \right] \left[ \begin{array}{|c|} \hline \text{decision tree} \\ \hline \end{array} \right] \end{array} \right) \rightarrow \max$$

**Результаты:**

- точность мультязычного поиска 94%
- сокращение модели 128 Гб → 4.8 Гб при редукции словарей (BPE-токенизация) до 11К токенов на каждый язык.

П.Потапова, А.Грабовой, О.Бахтеев, Е.Егоров, Н.Зиновкин, Ю.Чехович, К.Воронцов и др. Мультязыковая автоматическая рубрикация научных документов. 2023.



## Задача 1. Поиск этно-релевантных тем в социальных сетях

- **Дано:**

- 1) данные социальных медиа (ВК и др.)
- 2) словарь этнонимов (около 300)

- **Найти:**

- 1) как можно больше тем про этничности
- 2) темы с сочетанием этничностей (возможные конфликты)

- **Критерий:**

- 1) интерпретируемость всех тем
- 2) точность и полнота поиска этно-релевантных тем

### Используемые регуляризаторы:

- сглаживание этно-релевантных тем по словарю этнонимов
- декоррелирование этно-релевантных тем
- модальность этнонимов

## Задача 1. Примеры этнонимов (всего около 300)

османский	русич
восточноевропейский	сингапурец
эвенк	перуанский
швейцарская	словенский
аланский	вепсский
саамский	ниггер
латыш	адыги
литовец	сомалиец
цыганка	абхаз
ханты-мансийский	темнокожий
карачаевский	нигериец
кубинка	лягушатник
гагаузский	камбоджиец

## Задача 1. Примеры этно-релевантных тем

**(русские)**: русский, князь, россия, татарин, великий, царить, царь, иван, император, империя, грозить, государь, век, московская, екатерина, москва,

**(русские)**: акция, организация, митинг, движение, активный, мероприятие, совет, русский, участник, москва, оппозиция, россия, пикет, протест, проведение, националист, поддержка, общественный, проводить, участие,

**(славяне, византийцы)**: славянский, святослав, жрец, древние, письменность, рюрик, летопись, византия, мефодий, хазарский, русский, азбука,

**(сирийцы)**: сирийский, асад, боевик, район, террорист, уничтожить, группировка, дамаск, оружие, алесию, оппозиция, операция, селение, сша, нусра, турция,

**(турки)**: турция, турецкий, курдский, эрдоган, стамбул, страна, кавказ, горин, полиция, премьер-министр, регион, курдистан, ататюрк, партия,

**(иранцы)**: иран, иранский, сша, россия, ядерный, президент, тегеран, сирия, оон, израиль, переговоры, обама, санкция, исламский,

**(палестинцы)**: террорист, израиль, терять, палестинский, палестинец, террористический, палестина, взрыв, территория, страна, государство, безопасность, арабский, организация, иерусалим, военный, полиция, газ,

**(ливанцы)**: ливанский, боевик, район, ливан, армия, террорист, али, военный, хизбалла, раненый, уничтожить, сирия, подразделение, квартал, армейский,

**(ливийцы)**: ливан, демократия, страна, ливийский, каддафи, государство, алжир, война, правительство, сша, арабский, али, муаммар, сирия,



## Задача 1. Примеры этно-релевантных тем

**(евреи)**: израиль, израильский, страна, война, нетаньяху, тель-авив, время, сша, сирия, египет, случай, самолет, еврейский, военный, ближний,

**(американцы)**: американский, американка, война, россия, военный, страна, вашингтон, америка, армия, конгресс, сирия, союзный, российский, обама, войска, русский, оружие, операция,

**(немцы)**: армия, война, войска, советский, военный, дивизия, немец, фронт, немецкий, генерал, борт, операция, оборона, русский, бог, победа,

**(немцы)**: германий, немец, германский, ссср, немецкий, война, старое, советский, россия, береза, русский, правительство, территория, полный, документ, вопрос, сорт, договор, отношение, франция,

**(евреи, немцы)**: еврей, еврейский, холодный, германий, антисемитизм, гетра, немец, синагога, сша, израиль, малиновского, комиссия, нацбол, документ, война, еврейка, миллион, украина,

**(украинцы, немцы)**: украинский, унс, оун, немец, немецкий, ковальков, хохол, волынский, бандера, организация, россиянин, советский, русский, польский, армия, шухевича, ровенский,

**(таджики, узбеки)**: мигрант, страна, россия, миграция, азия, нелегальный, миграционный, таджикистан, гастарбайтер, гражданка, трудовой, рабочий, фмс, коренево, среднее, узбекистан, таджик, проблема, русский, население,

**(канадцы)**: команда, игра, игрок, канадский, сезон, хоккей, сборная, играть, болельщик, победа, кубок, счет, забирать, хоккейный, выигрывать, хоккеист, чемпионат, шайба,

## Задача 1. Примеры этно-релевантных тем

**(японцы)**: японский, япония, корея, китайский, жилища, авария, фукусиму, цунами, сообщать, океан, станция, хатико, район, правительство, атомный,

**(норвежцы)**: дитя, ребенок, родиться, детский, семья, воспитанный, право, возраст, отец, воспитание, норвежский, родительский, родить, мальчик, взрослый, опека, сын,

**(венесуэльцы)**: куба, кастро, венесуэла, чавес, президент, уго, мадура, боливия, фидель, глава, латинский, венесуэльский, лидер, боливарианской, президентский, альенде, гевару,

**(китайцы)**: китайский, россия, производство, китаи, продукция, страна, предприятие, компания, технология, военный, регион, производить, производственный, промышленность, российский, экономический, кнр,

**(азербайджанцы)**: русский, азербайджан, азербайджанец, россия, азербайджанский, таксист, диаспора, анапа, народ, москва, страна, армянин, слово, рынок,

**(грузины)**: грузинский, спецназ, военный, август, баташева, российский, спецназовец, миротворец, операция, румын, бригада, миротворческий, абхазия, группа, войска, русский, цхинвале,

**(осетины)**: конституция, осетия, аминат, русский, осетинский, южный, северный, россия, война, республика, вопрос, алахай, российский, население, конфликт,

**(цыгане)**: наркотик, цыган, цыганка, хороший, место, страна, деньга, время, работать, жизнь, жить, рука, дом, цыганский, наркоманка,

## Задача 1. Результат: ARTM находит больше этно-релевантных тем

Число этно-релевантных тем, найденных моделью:

модель	этно-тем	фон.тем	++	+-	-+	всего
PLSA	300		9	11	18	38
PLSA	400		12	15	17	44
ARTM-1	200	100	18	33	20	71
ARTM-1	250	150	21	27	20	68
ARTM-2	200	100	28	23	23	74
ARTM-2	250	150	<b>38</b>	<b>42</b>	<b>30</b>	<b>104</b>

Регуляризаторы ARTM-1:

**этно темы:** разреживание, декоррелирование, сглаживание этнонимов

**фоновые темы:** сглаживание, разреживание этнонимов

Регуляризаторы ARTM-2:

ARTM-1 + **модальность этнонимов**

---

*M. Apishev, S. Koltcov, O. Koltsova, S. Nikolenko, K. Vorontsov. Additive regularization for topic modeling in sociological studies of user-generated text content. MICAI 2016.*

## Аналогичные по структуре исследования

### Задачи «поиска и классификации иголок в стоге сена»

- поиск и кластеризация новостей [1]
- поиск в социальных медиа информации, связанной с болезнями, симптомами и методами лечения [2]
- поиск чатов, связанных с преступностью и экстремизмом [3, 4]
- поиск выступлений о правах человека в ООН [5]

- 
1. *J.Jagarlamudi, H.Daumé III, R.Udupa*. Incorporating lexical priors into topic models. 2012.
  2. *M.Paul, M.Dredze*. Discovering health topics in social media using topic models. 2014.
  3. *M.A.Basher, A.Rahman, B.C.M.Fung*. Analyzing topics and authors in chat logs for crime investigation. 2014.
  4. *A.Sharma, M.Pawar*. Survey paper on topic modeling techniques to gain useful forecasting information on violant extremist activities over cyber space. 2015.
  5. *Kohei Watanabe, Yuan Zhou*. Theory-driven analysis of large corpora: semisupervised topic classification of the UN speeches. 2022.

## Задача 2. Анализ программ развития российских вузов

**Цель** — выявить закономерности в стратегиях развития вузов, не читая всех этих документов (Distant Reading)

- **Дано:**

программам развития ВУЗов: 396 файлов, 284 вуза

- **Найти:**

полный тематический спектр направлений развития

- **Критерий:**

интерпретируемость тем;

чёткого количественного критерия нет :(

## Задача 2. Пример интерпретации темы

**(слова)**: инновационный исследование результат региональный предприятие проведение основа среда внедрение уровень рамка сфера исследовательский научно научно-исследовательский участие приоритетный специалист цель выполнение международный прикладной ведущий взаимодействие

**(биграммы)**: научный \_ исследование инновационный \_ деятельность приоритетный \_ направление научно \_ исследовательский исследование \_ разработка развитие \_ инновационный фундаментальный \_ прикладной разработка \_ внедрение направление \_ развитие мировой \_ уровень научно \_ образовательный исследовательский \_ деятельность инновационный \_ развитие малое \_ инновационный инновационный \_ предприятие научный \_ инновационный модернизация \_ научно-исследовательский прикладной \_ исследование инновационный \_ проект развитие \_ научный инновационный \_ инфраструктура проведение \_ научный

**(ИНТЕРПРЕТАЦИЯ)**: научные исследования и инновационное развитие

## Задача 2. Пример интерпретации темы

**(слова)**: международный число количество участие конференция  
зарубежный увеличение учёный академический мобильность конкурс  
сотрудничество грант иностранный аспирант совместный молодая  
ведущий специалист привлечение преподаватель исследование школа  
сотрудник семинар

**(биграммы)**: увеличение \_ количество академический \_ мобильность  
увеличение \_ число международный \_ деятельность  
международный \_ сотрудничество международный \_ научный  
развитие \_ международный принять \_ участие российский \_ международный  
научный \_ мероприятие международный \_ образовательный  
участие \_ международный иностранный \_ студент количество \_ студент  
научный \_ проект университет \_ международный международный \_ уровень  
международный \_ академический количество \_ участник  
научный \_ конференция программа \_ академический участие \_ студент

**(ИНТЕРПРЕТАЦИЯ)**: академическая мобильность и международное  
сотрудничество

## Задача 2. Пример интерпретации темы

**(слова)**: общежитие корпус здание ремонт площадь инфраструктура комплекс помещение строительство объект капитальный кампус имущественный спортивный реконструкция безопасность территория сооружение место оборудование современный замена учебно-лабораторный комфортный

**(биграммы)**: учебный \_ корпус капитальный \_ ремонт имущественный \_ комплекс общий \_ площадь здание \_ сооружение студенческий \_ общежитие корпус \_ общежитие развитие \_ имущественный инфраструктура \_ университет создание \_ комфортный развитие \_ инфраструктура университетский \_ кампус комплекс \_ университет спортивный \_ комплекс студент \_ сотрудник объект \_ университет земельный \_ участок условие \_ проживание территория \_ университет объект \_ инфраструктура социальный \_ инфраструктура использование \_ имущественный строительство \_ новый ремонтный \_ работа общежитие \_ университет

**(ИНТЕРПРЕТАЦИЯ)**: инфраструктура, кампус, строительство



## Задача 2. Интерпретация всех 50 тем

- Для интерпретируемости тем важны биграммы
- Модель построили примерно с 10-й попытки (подбирали число тем, регуляризацию, добивались различности тем)
- Интерпретация 50 тем заняла примерно 20 минут работы
- Иногда выделялись темы исследований и разработок, но для этого нужна более гранулированная модель
- Темы были сгруппированы вручную по 5 категориям:
  - 1 16 тем про науку, инновации и сотрудничество
  - 2 14 тем про образование и кадровый потенциал
  - 3 11 тем про административное управление и хозяйство вуза
  - 4 3 темы «юридические», о самой стратегии развития
  - 5 6 тем «малые и мусорные», вместе не более 5% контента

## Задача 2. Интерпретация всех 50 тем

доля контента	доля вузов		название темы
	более 2%	более 5%	
7	95	67	научные исследования и инновационное развитие
12	92	39	стратегия развития
15	84	23	академическая мобильность и международное сотрудничество
19	82	17	кадровой потенциал и кадровая политика
22	80	14	иностранные студенты
27	75	30	образовательные программы
30	75	13	повышение квалификации и переподготовка кадров
33	70	10	система управления вузом
36	68	16	учебный процесс
39	62	15	финансы и бюджет
43	62	21	бюрократия
45	56	3	подготовка высококвалифицированных кадров
48	47	9	инфраструктура, кампус, строительство
50	44	4	меры повышения качества образования
52	42	4	влияние на экономику региона
54	41	8	молодежная политика
56	41	6	центры компетенций и технологического превосходства
58	40	6	отсылки к стратегическим документам и НПА
60	36	1	работа со школьниками и талантливой молодежью
62	34	7	ректорат и органы управления вузом
64	30	5	материально-техническая база вуза
65	29	2	связь с общественностью, имидж вуза
67	29	8	исследования с/х, лес, химия, ит
69	29	1	публикационная активность и защиты диссертаций
71	29	2	взаимодействие с региональной властью

доля контента	доля вузов		название темы
	более 2%	более 5%	
72	27	1	образовательные программы, аккредитация, профстандарты
74	25	3	спортивная и культурная жизнь вуза
75	21	5	стратегия развития и региональная среда
77	20	1	образовательный процесс и образовательные технологии
78	19	1	международное сотрудничество и договорные отношения
79	19	2	цифровизация и цифровые технологии
81	18	2	медицинское обеспечение, обучение инвалидов
82	18	5	блоки мероприятий и показатели результативности
84	18	5	работа структурных подразделений вуза
85	17	2	выход в мировые рейтинги университетов
86	14	1	технологии транспорта и искусственного интеллекта
87	13	1	публикационная и издательская деятельность
88	12	1	финансовое и ресурсное обеспечение программы развития
89	11	1	мониторинг показателей эффективности
90	11	0	сетевые образовательные программы, ворлдскиллс
92	11	1	региональные особенности приёма и рынка труда
93	10	1	приём абитуриентов
93	10	0	исследования в экологии и медицине
94	9	1	образовательные программы (частные вопросы)
95	8	1	частные и региональные проблемы
96	8	2	авиационные технологии
97	8	0	смесь тем
98	7	0	образовательные программы & урбанистика и туризм (смесь тем)
99	7	1	смесь тем
100	7	1	частные юридические вопросы

- 16 тем — наука и инновации
- 14 тем — образование и кадры
- 11 тем — управление и хозяйство
- 3 темы — о стратегии развития
- 6 тем — мелкие мусорные



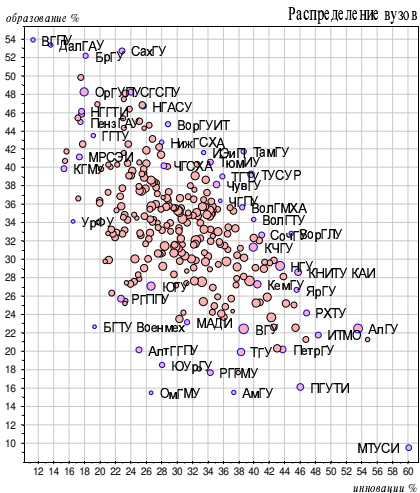
## Задача 2. Тематическая карта вузов

### По осям:

- объёмная доля тем
- про инновации
- про образование

### Вывод:

объёмные доли тем, возможно, показывают баланс приоритетов развития ...хотя... это похоже на оценивание научного отчёта толщиной в сантиметрах :)



## Тематические модели коротких текстов

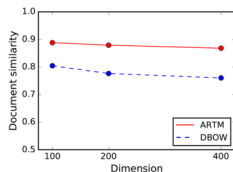
**Цель:** интерпретируемые разреженные тематические эмбединги на основе дистрибутивной семантики, аналоги word2vec и WNTM.

**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{|c|} \hline \text{PLSA} \\ \hline \Phi \quad \Theta \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{co-occurrence} \\ \hline \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{interpretable} \\ \hline \text{---} \\ \text{---} \\ \text{---} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{multimodal} \\ \hline \text{---} \\ \text{---} \\ \text{---} \\ \hline \end{array} \right) \rightarrow \max$$

**Результаты:**

- Точность поиска схожих документов:  $0.8 \rightarrow 0.9$
- Когерентность тем:  $0.08 \rightarrow 0.33$
- Семантическая близость слов:  $0.53 \rightarrow 0.58$ ,  $0.38 \rightarrow 0.61$



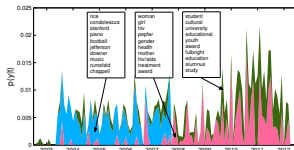
*A.Potapenko, A.Popov, K.Vorontsov. Interpretable probabilistic embeddings: bridging the gap between topic models and neural networks. AINL, 2017.*

## Выявление динамики тем в новостных потоках

**Цель:** выделение тем в коллекции пресс-релизов МИДов 4х стран, с привязкой ко времени.

**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{c} \text{PLSA} \\ \Phi \quad \Theta \end{array} \right) + R \left( \begin{array}{c} \text{interpretable} \\ \text{[bar chart]} \quad \text{[stacked bar chart]} \end{array} \right) + R \left( \begin{array}{c} \text{temporal} \\ \text{[line graph]} \end{array} \right) + R \left( \begin{array}{c} \text{multimodal} \\ \text{[stacked bar chart]} \quad \text{[square]} \end{array} \right) \\ + R \left( \begin{array}{c} \text{n-gram} \\ \text{[grid]} \end{array} \right) + R \left( \begin{array}{c} \text{multilanguage} \\ \text{[stacked bar chart]} \quad \text{[square]} \end{array} \right) \rightarrow \max$$



**Результаты:**

- разделение тем на событийные и перманентные
- когерентность тем: 5.5  $\rightarrow$  6.5

*Н. Дойков.* Адаптивная регуляризация вероятностных тематических моделей.  
ВКР бакалавра, ВМК МГУ, 2015.

## Выделение поляризованных мнений в политических новостях

**Цель:** найти признаки, по которым событийная тема разделяется на кластеры-мнения

Modalities	<i>Pr</i>	<i>Rec</i>	<i>F1</i>
TF-IDF	0.51	0.95	0.67
SPO	0.59	0.7	0.64
FR	0.86	0.49	0.65
Sent	0.69	0.57	0.66
SPO+FR	0.86	0.68	0.76
SPO+Sent	0.83	0.78	0.81
FR+Sent	0.9	0.52	0.67
<b>All</b>	<b>0.77</b>	<b>0.97</b>	<b>0.86</b>

**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{c} \text{PLSA} \\ \left( \begin{array}{|c|} \hline \Phi \\ \hline \end{array} \begin{array}{|c|} \hline \Theta \\ \hline \end{array} \right) \end{array} \right) + R \left( \begin{array}{c} \text{interpretable} \\ \left( \begin{array}{|c|} \hline \text{matrix} \\ \hline \end{array} \begin{array}{|c|} \hline \text{matrix} \\ \hline \end{array} \right) \end{array} \right) + R \left( \begin{array}{c} \text{multimodal} \\ \left( \begin{array}{|c|} \hline \text{matrix} \\ \hline \end{array} \begin{array}{|c|} \hline \text{matrix} \\ \hline \end{array} \right) \end{array} \right) + R \left( \begin{array}{c} \text{n-gram} \\ \left( \begin{array}{|c|} \hline \text{matrix} \\ \hline \end{array} \begin{array}{|c|} \hline \text{matrix} \\ \hline \end{array} \right) \end{array} \right) + R \left( \begin{array}{c} \text{syntax} \\ \left( \begin{array}{|c|} \hline \text{tree} \\ \hline \end{array} \right) \end{array} \right) \rightarrow \max$$

**Результаты:**

- выделение мнений внутри тем: F1-мера = 0.86%
- совместное использование трёх модальностей: факты «субъект–предикат–объект», семантические роли слов по Филлмору, тональности именованных существей

*D.Feldman, T.Sadekova, K.Vorontsov. Combining facts, semantic roles and sentiment lexicon in a generative model for opinion mining. Dialogue 2020.*

## Выделение поляризованных мнений в политических новостях

... Президент Петр Порошенко заявил, что Россия де-факто конфисковала украинские предприятия, которые находятся на неподконтрольной Киеву территории. Сегодня ДНР и ЛНР "национализировали" украинские предприятия ... При этом Кремль защитил конфискацию предприятий в ЛДНР ... Украина потребует расширить санкции ... За все эти действия обязательно наступит наказание. Украина потребует расширения санкций на тех, кто украл украинские предприятия ... (*Kiev opinion*)

... По словам Захарченко, Киев встретит свой "ужасный конец" ... Киев возьмется за ум, и в целях спасения собственной промышленности снимет блокаду ... Обстановка, которую искусственно создала Украина с блокадой Донбасса, вынудила ... кошмарить свой народ ... если в Киеве были приняты какое-либо постановление ... положительные результаты, как в республиках, так и в России ... Если им удастся сместить Порошенко и при этом не развалить Украину, то все вернется на свои места ... (*Moscow opinion*)

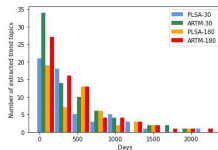


Слова «Порошенко», «Россия», «Украина» встречаются в тексте-1 и тексте-2 одинаково часто, однако:

- «Порошенко» — субъект в тексте-1 и объект в тексте-2;
- «Россия» — агент в тексте-1 и локация в тексте-2;
- негативная тональность: «Россия», «Кремль» в тексте-1, «Киев», «Украина» в тексте-2.

## Выявление трендов в коллекции научных публикаций

**Цель:** ранее обнаружение трендовых тем с начальным экспоненциальным ростом в области AI/ML 2009–2021 гг.



**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{c} \text{PLSA} \\ \Phi \quad \Theta \end{array} \right) + R \left( \begin{array}{c} \text{interpretable} \\ \text{[Bar Chart Icon]} \quad \text{[Scatter Plot Icon]} \end{array} \right) + R \left( \begin{array}{c} \text{dynamic} \\ \text{[Line Graph Icon]} \end{array} \right) + R \left( \begin{array}{c} \text{multimodal} \\ \text{[Stacked Boxes Icon]} \quad \text{[Square Icon]} \end{array} \right) + R \left( \begin{array}{c} \text{n-gram} \\ \text{[Grid Icon]} \end{array} \right) \rightarrow \max$$

**Результаты:**

- выделение 90 из 91 тренда в области машинного обучения
- 63% тем выделяется за год, 79% за два года

*Н.Герасименко, А.Чернявский, М.Никифорова, М.Никитин, К.Воронцов.*

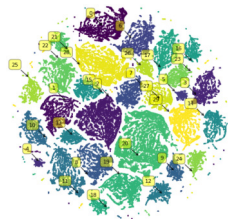
Инкрементальное обучение тематических моделей для поиска трендовых тем в научных публикациях. Доклады РАН, 2022.



## Тематическая модель банковских транзакционных данных

**Цель:** Выявление паттернов потребительского поведения клиентов банка, причём

- документы = клиенты,
- слова = МСС-коды продавцов.



**Регуляризаторы:**

$$\mathcal{L}\left(\begin{array}{|c|} \hline \text{PLSA} \\ \hline \Phi \quad \Theta \\ \hline \end{array}\right) + R\left(\begin{array}{|c|} \hline \text{interpretable} \\ \hline \text{[Bar chart icon]} \quad \text{[Scatter plot icon]} \\ \hline \end{array}\right) + R\left(\begin{array}{|c|} \hline \text{multimodal} \\ \hline \text{[Stacked bar chart icon]} \quad \text{[Box icon]} \\ \hline \end{array}\right) + R\left(\begin{array}{|c|} \hline \text{supervised} \\ \hline \text{[Decision tree icon]} \\ \hline \end{array}\right) \rightarrow \max$$

**Результаты:**

- темы — паттерны потребительского поведения
- предсказание пола, возраста, достатка клиентов

---

*E.Egorov, F.Nikitin, A.Goncharov, V.Alekseev, K.Vorontsov. Topic modelling for extracting behavioral patterns from transactions data. 2019.*

## Анализ публичных выступлений

- [1] Выступления (210К) в Европарламенте, 1999–2014:
  - выявление событийных тем и эволюции перманентных тем;
  - как члены и комитеты ЕП влияют на формирование тем
- [2] Модель контрастных мнений (Contrastive Opinion Modeling)
  - выступления в Сенате США ([www.votesmart.org](http://www.votesmart.org));
  - СМИ: New York Times, Xinhua News, The Hindu, 2009–2010
- [3] Выступления в Совбезе США по Афганистану, 2001–2017:
  - динамика отношения разных стран к проблеме Афганистана

---

[1] *D. Greene, J.P. Cross*. Unveiling the political agenda of the European Parliament plenary: a topical analysis. 2015.

[2] *Fang, Y., et al*. Mining contrastive opinions on political texts using cross-perspective topic model. 2012.

[3] *M. Schönfeld*. Discursive landscapes and unsupervised topic modeling in IR: a validation of text-as-data approaches through a new corpus of UN Security Council speeches on Afghanistan. 2018.

## Анализ СМИ и социальных медиа

- [1] Тематика изменения климата в СМИ Пакистана, 2010–2021  
— выявление, группирование и динамика тем
- [2] Выявление поляризации новостей (AYLIEN COVID-19)  
— 1,5М новостей, 440 источников СМИ, 11.2019–07.2020
- [3] Выявление политических взглядов пользователей Twitter
- [4] Что пишет NYT о ядерных технологиях с 1945 по н/в

---

[1] *W.Ejaz et al.* Politics triumphs: A topic modeling approach for analyzing news media coverage of climate change in Pakistan. 2023

[2] *Zihao He.* Detecting polarized topics using partisanship-aware contextualized topic embeddings. 2021

[3] *R.Cohen, D.Ruths.* Classifying Political Orientation on Twitter: It's Not Easy! 2013.

[4] *C.Jacobi.* Quantitative analysis of large amounts of journalistic texts using topic modelling. 2015.

*H.Jelodar et al.* Latent Dirichlet allocation (LDA) and topic modeling: models, applications, a survey. 2019.

## Исследования газетных архивов

[1] Корпус *Pennsylvania Gazette* 1728–1800, 25М слов:

- выделение последовательности событийных тем;
- изучение синхронности событий;
- комбинирование автоматического анализа и ручного.

[2] *Газеты Техаса* от гражданской войны до наших дней:

- выделение всех тем, связанных с хлопком;
- построение серии моделей в скользящих окнах;
- важность качественной предобработки текстов.

[3] Газеты и периодика Финляндии (1854–1917):

- выделение тем о церкви, религии, образовании;
- тренды модернизации и секуляризации финского общества.

---

1. *D.Newman, S.Block*. Probabilistic topic decomposition of an eighteenth-century American newspaper. 2006.

2. *Tze-I Yang, A.J.Torget, R.Mihalcea*. Topic modeling on historical newspapers. 2011.

3. *J.Marjanen et al*. Topic modelling discourse dynamics in historical newspapers. 2021.

## Исследования документальной литературы и дневников

- [1] Двухязычный корпус книг на английском и немецком:  
— все темы, связанные с эпистемологией
- [2] Корпус текстов на китайском языке (1644–1912):  
— все темы, связанные с бандитизмом, преступлениями;  
— необходим контекст для установления типа преступления;  
— важность правильной токенизации для китайского языка.
- [3] Дневник Martha Ballard (1735–1812), охватывает 27 лет:  
— выделение событийных и перманентных тем;  
— выделение персональных и исторических тем;  
— специфичный английский XVIII века.

- 
1. *M. Erlin*. Topic modeling, epistemology, and the English and German novel. 2017.
  2. *Ian Matthew Miller*. Rebellion, crime and violence in Qing China, 2013.
  3. *Cameron Blevins*.  
<http://www.cameronblevins.org/posts/topic-modeling-martha-ballards-diary>.

## Исследования научной и литературно-художественной периодики

Статьи коллекции JSTOR доступны в виде «мешков слов».

[1] Научные журналы XX века:

- различия тематики на английском и немецком языках;
- особенно исследовались различия, связанные со 2МВ;
- для объединения тем использовались интервики Википедии.

[2] Более 100 лет литературно-художественной периодики:

- как менялись темы;
- как менялись значения слов внутри каждой темы;
- как менялась тема насилия (violence, power, fear, blood, death, murder, act, guilt).

---

1. *D.Mimno*. Computational historiography: Data mining in a century of classics journals. 2012.

2. *A.Goldstone, T.Underwood*. The quiet transformations of literary studies: What thirteen thousand scholars could tell us. 2014.

## Проект «Тематизатор»: общие требования

Переход от библиотек (BigARTM, VisARTM, TopicNet) к приложению «Тематизатор» для конечного пользователя — аналитика в области цифровых гуманитарных исследований

- 1 Цели пользователя — разведочный анализ, понимание тематической структуры данных, «о чём эта коллекция»
- 2 Пользователь не обязан знать
  - форматы исходных данных и способы их предобработки
  - теорию TM и ARTM, виды регуляризаторов
  - методики подбора гиперпараметров
  - критерии качества моделей
  - библиотеку BigARTM
- 3 Интуитивная визуальная среда, веб-интерфейс
- 4 Пользователю должны быть доступны настройки
- 5 Дефолтные настройки должны работать на любых данных

## Приложения и исследования, взятые для анализа требований

- 1 Поиск этно-релевантных тем в социальных медиа
- 2 Анализ программ развития российских вузов
- 3 Проекты Школы Прикладного Анализа Данных
- 4 Тематический поиск по длинному текстовому запросу
- 5 Составление тематических подборок
- 6 Поиск и рубрикация научных статей на 100 языках
- 7 Выявление трендов в коллекции научных публикаций
- 8 Тематизация научно-просветительского онлайн-журнала
- 9 Поиск похожих дел в актах арбитражных судов
- 10 Тематизация пресс-релизов внешнеполитических ведомств
- 11 Тематизация twitter о российско-украинских отношениях
- 12 Выявление событийных тем в новостных потоках



## Проекты Школы Прикладного Анализа Данных (ноябрь, 2022)

### Исходные данные ВКонтакте (через сервисы Крибрум)

- Анализ социального влияния на формирование образа правильного питания у студентов г. Томска
- Анализ научной и публикационной активности сотрудников университета или научной организации
- Анализ практик участия в читательских сообществах, формирующихся вокруг авторов или жанров
- Анализ социального и политического взаимодействия сетевых сообществ в регионах ресурсного типа
- Анализ туристической активности и оценка портрета потенциального туриста — путешественника по Камчатке
- Анализ корпуса текстов образовательных дисциплин: программа курса + материалы курса + отчёты студентов
- Анализ научной педагогической литературы для построения карт компетенций

## Основной пользовательский сценарий (без детализации)

### 1 Загрузка

- данные в различных «сырых» форматах
- возможна дозагрузка данных порциями

### 2 Предобработка

- автоматический выбор обработчиков на основании данных
- выделение модальностей: языков, времени, терминов и т.д.

### 3 Моделирование

- визуализация метрик качества в процессе обучения модели
- возможность перехода к анализу, не прерывая обучения

### 4 Визуализация

- каждая тема должна уметь «рассказать о себе»
- много разных графиков (distant reading)

### 5 Коррекция

- перебор моделей и накопление «банка тем»
- пользовательские темы как подборки с рекомендациями

## Функциональные требования (по приоритетности)

- 1 Визуализация множества всех тем и их характеристик
- 2 Визуализация каждой темы с её «рассказом о себе»
- 3 Возможность задавать словари затравок для (групп) тем
- 4 Определение динамики тем во времени
- 5 Выявление коротких тем-событий и долгих тем-трендов
- 6 Разбиение тем на подтемы иерархически
- 7 Возможность группировки тем вручную
- 8 Выявление связей тем по сочетаемости в документах
- 9 Возможность отбора и накопления «банка тем»
- 10 Тематическая фильтрация коллекции
- 11 Тематический поиск по документу или фрагменту
- 12 Рекомендательный поиск и построение подборок

## Требования к интерпретируемости (по приоритетности)

- 1 Доля интерпретируемых тем близка к 100%
- 2 Темы строятся более на терминах, чем на словах
- 3 Общая лексика выводится в отдельные фоновые темы
- 4 Нет мусорных тем, нет тем-дубликатов (декорреляция)
- 5 Решена проблема несбалансированности тем
- 6 Темы способны рассказать о себе словами и фразами
- 7 Нетекстовые термины способны рассказать о себе словами
- 8 Темы именуется автоматически
- 9 В иерархии имена дочерних тем уточняют родительские
- 10 Тематика слов согласуется с их локальными контекстами
- 11 Короткие тексты объяснимо наследуют тематику их слов
- 12 Длинные тексты разбиваются на тематические сегменты

## Требования к функциям Загрузки

- 1 Загрузка коллекций из различных сырых форматов
- 2 — txt, json, docx, odt, pdf и др.
- 3 — СМИ, соцмедиа, Википедия, статьи, патенты и др.
- 4 Представление метаданных и модальностей
- 5 Возможность загрузки как локально, так и из облака
- 6 Возможность дозагрузки данных из источника порциями
- 7 Текст как последовательность или как «мешок слов»
- 8 В одном файле один документ или много документов

## Требования к функциям Предобработки

- 1 Автоматическая токенизация и лемматизация
- 2 Автоматическое исправление опечаток (соцсети)
- 3 Автоматическое выделение терминов  $n$ -грамм
- 4 Метаданные: авторы, время, категории, заголовки и др.
- 5 Модальности: онимы, теги, ссылки, пользователи и др.
- 6 Настройка шаблонов для выделения модальностей
- 7 Сортировка по времени и нарезка по пакетам
- 8 Автоматическое определение коротких текстов
- 9 Автоматическая редукция словарей (по необходимости)
- 10 Автоматическое определение языков
- 11 Машинный перевод для получения параллельных текстов
- 12 Предобработка не должна идти дольше тематизации

## Требования к функциям Моделирования

- 1 Визуализация процесса обучения модели
- 2 Вывод метрик на графиках от #итерации, #пакета
- 3 Метрики перплексии, разреженности, вырожденности и др.
- 4 Автоматическая подстройка под короткие тексты
- 5 Автоматическая подстройка под длинные тексты
- 6 Темпоральная модель, если есть модальность времени
- 7 Подбор числа тем или построение иерархии тем
- 8 Автоматический подбор гиперпараметров, AutoML
- 9 Логирование информации о найденных аномалиях
- 10 Логирование данных о моделях, журнал экспериментов
- 11 Возможность перехода к анализу, не прерывая обучения
- 12 Возможность замены BigARTM на альтернативы

## Требования к функциям Визуализации

- 1 Визуальная навигация по темам, документам, терминам
- 2 XY-график тем в осях свойств тем
- 3 XY-график документов/объектов в осях объёмов тем/групп
- 4 Построение спектра тем по семантической близости
- 5 XY-график документов в осях «время–спектр тем»
- 6 Визуализация связей между словами и понятиями темы
- 7 Визуализация динамики тем в осях «время–объём темы»
- 8 Визуализация иерархии тем
- 9 Визуализация связей тем по их сочетаемости в документах
- 10 Визуализация тематической структуры документа
- 11 Выбор характеристик тем для осей XY-графиков
- 12 Выбор характеристик объектов и документов для осей



## Требования к функциям Коррекции

- 1 Разметка тем на релевантные, нерелевантные, мусорные
- 2 Разметка релевантных термов, документов в темах
- 3 Термы-затравки для «классификации иголок в стоге сена»
- 4 Обнаружение и расщепление неоднородных тем
- 5 Автоматический переход к тематической иерархии
- 6 Детекция новых событийных тем в темпоральных моделях
- 7 Накопление «банка тем» по множеству моделей
- 8 Многокритериальное оценивание качества моделей
- 9 Планирование экспериментов по улучшению моделей
- 10 Тематическая фильтрация коллекции и потока
- 11 Создание пользовательских тем — подборок документов
- 12 Ранжирование рекомендаций для пользовательских тем

## Требования к рабочему пространству проекта пользователя

- 1 Настройки входных данных — коллекций и потоков
- 2 Настройки модулей предобработки
- 3 Структура и гиперпараметры сравниваемых моделей
- 4 Структура и гиперпараметры финальной модели
- 5 Визуализации процесса обучения модели
- 6 Визуализации количественных результатов моделирования
- 7 Визуализации качественных результатов (аннотации тем)
- 8 Банк тем — множество тем, отобранных из моделей
- 9 Пользовательские темы — подборки документов
- 10 Настройка подробности отчёта по проекту
- 11 Настройка комментариев к пунктам отчёта по проекту
- 12 Сгенерированный отчёт по проекту

## Что точно войдёт в MVP

### 1 Загрузка

- несколько коллекций для тестирования
- «мешок слов» в формате Vowpal Wabbit (BigARTM)

### 2 Моделирование

- отображение статуса обработки пакетов и текущих метрик
- возможность прервать обучение и перейти к анализу
- регуляризатор декоррелирования — спрятан

### 3 Визуализация

- навигация по темам в духе TMVE
- спектр тем по семантической близости и релевантности

### 4 Коррекция

- разметка тем на релевантные, нерелевантные, мусорные
- перестроение модели с сохранением релевантных тем

- Тематическое моделирование — инструмент для поиска и систематизации больших текстовых коллекций
- Теория ARTM и библиотека BigARTM позволяют строить модели с нужным набором свойств
- Есть наработанные приёмы для
  - улучшения интерпретируемости тем
  - улучшения качества тематического поиска
  - исследования динамики тем во времени
  - выделения тем по большому списку слов-затравок
  - иерархического дробления тем на более мелкие подтемы
  - навигации по темам и их визуального анализа
  - учёта обратной связи с экспертом
- Эти приёмы активно используются для обработки больших текстовых массивов в гуманитарных исследованиях
- Фактор успеха — качественная предобработка текстов