

# Тематика НИР: математические методы анализа текстов и информационного поиска

Константин Вячеславович Воронцов

- Лаборатория машинного обучения и семантического анализа  
Института искусственного интеллекта МГУ •
- Кафедра математических методов прогнозирования ВМК МГУ •
- Кафедра машинного обучения и цифровой гуманитаристики,  
кафедра интеллектуальных систем МФТИ •
  - ФИЦ ИУ РАН •

[voron@mlsa-iai.ru](mailto:voron@mlsa-iai.ru)

<http://www.MachineLearning.ru/wiki?title=User:Vokov>

AI Masters, ИИИ МГУ • 18 ноября 2023

## 1 Проект «Мастерская знаний»

- Концепция «мастерской знаний»
- Полуавтоматическая суммаризация
- Тематизация и визуализация

## 2 Проект «Новостной коллаيدر»

- Детекция постправды
- Технологический конкурс ПРО//ЧТЕНИЕ
- Задачи автоматизации разметки текста

## 3 Вероятностное тематическое моделирование

- Векторизация текста и тематическое моделирование
- Аддитивная регуляризация (ARTM)
- Приложения ВТМ и проект «Тематизатор»

## Концепция «мастерской знаний»

«Огромное и все возрастающее богатство знаний разбросано сегодня по всему миру. Этих знаний, вероятно, было бы достаточно для решения всего громадного количества трудностей наших дней, но они рассеяны и неорганизованы. Нам необходима очистка мышления в *своеобразной мастерской*, где можно **получать, сортировать, суммировать, усваивать, разъяснять и сравнивать** знания и идеи»  
— Герберт Уэллс, 1940

“An immense and ever-increasing wealth of knowledge is scattered about the world today; knowledge that would probably suffice to solve all the mighty difficulties of our age, but it is dispersed and unorganized. We need a sort of mental clearing house for the mind: a depot where knowledge and ideas are **received, sorted, summarized, digested, clarified and compared**”  
— Herbert Wells, 1940



## Концепция сервисов «Мастерской знаний»

*Подборка* — долгосрочный поисковый интерес пользователя

### Поисково-рекомендательные функции:

- поиск тематически близких документов по *подборке*
- мониторинг новых документов для *подборки*
- контекстные рекомендации по документу из *подборки*

### Аналитические функции:

- автоматизация реферирования *подборки*
- кластеризация трендов, аспектов, отношений в *подборке*
- рекомендация порядка чтения внутри *подборки*
- выделение «важных мест» в документе из *подборки*

### Коммуникативные функции:

- совместное составление и использование *подборок*
- интерактивная визуализация и инфографика по *подборке*

## Поисково-рекомендательная система SciSearch.ai

Тематическая подборка пользователя:

The screenshot shows a web browser window with the URL <https://arxiv.aithea.com/collections/Q29sbGVjdGVhbjozUFVTUEFxaIBH>. The navigation menu includes FEEDS, SEARCH, COLLECTIONS (circled in red), About, and FAQ. The user's name, Konstantin Vorontsov, is displayed in the top right corner (circled in red). The main content area is titled "MOOC (massive open online course)" (highlighted in green) and features a "PAPERS" section (circled in red) with a "RECOMMENDED" filter. Two papers are listed:

- 19 JUL 2014**  
Towards Feature Engineering at Scale for Data from Massive Open Online Courses  
Kalyan Veeramachaneni, Una-May O'Reilly, Colin Taylor  
We examine the process of engineering features for developing models that improve our understanding of learners' online behavior in MOOCs. Because feature engineering relies so heavily on human insight, we argue that extra effort should be made to engage the crowd for feature proposals and even their operationalization. We show two approaches where we have started to engage the crowd. We also show how features can be evaluated for their relevance in predictive accuracy. When we...  
Citations: 6
- 2 JUL 2017**  
Reciprocal Recommender System for Learners in Massive Open Online Courses (MOOCs)  
Sankalp Prabhakar, Gerasimos Spanakis, Ozmar Zalane  
Massive open online courses (MOOC) describe platforms where users with completely different backgrounds subscribe to various courses on offer. MOOC forums and discussion boards offer learners a medium to communicate with each other and maximize their learning outcomes. However, oftentimes learners are hesitant to approach each other for different reasons (being shy, don't know the right match, etc.). In this paper, we propose a reciprocal recommender system which matches...

Разработка: <http://aithea.com>, <http://ddecisions.ai>, <http://machine-intelligence.ru>

# Поисково-рекомендательная система SciSearch.ai

Список статей, рекомендуемых для добавления в подборку:

The screenshot shows the SciSearch.ai interface. At the top, there are navigation tabs: FEEDS, SEARCH, and COLLECTIONS. The current page is titled 'MOOC (massive open online course)'. Below the title, there are two tabs: 'PAPERS' and 'RECOMMENDED'. A red arrow points from 'PAPERS' to 'RECOMMENDED', which is circled in red. The 'RECOMMENDED' tab displays a list of articles. The first article is titled 'A Survey of Natural Language Generation Techniques with a Focus on Dialogue Systems - Past, Present and Future Directions' by Sashank Santhanam and Samira Shalikh, dated 2 JUN 2019. The second article is titled 'Capturing "attrition intensifying" structural traits from didactic interaction sequences of MOOC learners' by Tanmay Sinha, Nan Li, Patrick Jermann, and Pierre Dillenbourg, dated 20 SEP 2014. Each article includes a brief description and a 'Citations' count.

Разработка: <http://aithea.com>, <http://ddecisions.ai>, <http://machine-intelligence.ru>

## Поисково-рекомендательная система SciSearch.ai

Добавление статьи из списка рекомендаций в подборку:

The screenshot shows a web browser window displaying a paper titled "A Survey of Natural Language Generation Tasks" by Sashank Santhanam and Samira Shaikh. A modal dialog titled "Add to collections" is open, showing a list of collection options: "Exploratory Search", "MOOC (massive open online course)", "Opinion Mining and Sentiment Analysis with Topic Modeling", "Textual Complexity and Readability", and "Topic modeling of genomic data". The "MOOC" option is selected. A "SAVE CHANGES" button is visible at the bottom of the dialog. A "RECOMMENDED" badge is circled in red on the right side of the page. A red arrow points from the "MOOC" option in the dialog to the "SAVE CHANGES" button.

Разработка: <http://aithea.com>, <http://ddecisions.ai>, <http://machine-intelligence.ru>

# Полуавтоматическое реферирование тематических подборок

Рекомендации фраз для реферата с помощью сифлэров:

The screenshot displays a web application interface for paper summarization and phrase recommendation. It is organized into three main columns: **PAPERS**, **RECOMMENDED**, and **SUMMARIZATION**.

- PAPERS:** A list of papers is shown. The paper "SummaRuNNer: A Recurrent Neural Network based..." is highlighted in green. Other papers include "BanditSum: Extractive Summarization as a Contextual Bandit Problem", "A Survey on Neural Network-Based Summarization...", "A Deep Reinforced Model for Abstractive Summarization...", "Neural Extractive Summarization with Side Information", and "Get To The Point: Summarization with Pointer-Generator Networks".
- RECOMMENDED:** A summary of the selected paper is displayed. The text includes: "A novel method for training neural networks to perform single-document extractive summarization without heuristically-generated extractive labels. We call our approach BANDITSUM as it treats extractive summarization as a contextual bandit (CB) problem, where the model receives a document to summarize (the context), and chooses a sequence of sentences to include in the summary (the action). A policy gradient reinforcement learning algorithm is used to train the model to select sequences of sentences that maximize ROUGE score. The aim of this literature review is to survey the recent work on neural-based models in automatic text summarization. We examine in detail ten state-of-the-art neural-based..."
- SUMMARIZATION:** A list of recommended phrases is shown. The phrases include: "SummaRuNNer, a Recurrent Neural Network (RNN) based sequence model for extractive summarization of documents and show that it achieves performance better than or comparable to state-of-the-art.", "Our model has the additional advantage of being very interpretable, since it allows visualization of its predictors broken up by abstract features, such as information content, salience and novelty.", and "Another novel contribution of our work is abstractive training of our extractive model that can train on human generated reference summaries alone, eliminating the need for sentence-level extractive labels."
- Prompters:** A row of buttons is located at the bottom, labeled "Annotate", "Idea", "Theory", "Method", "Citation", "Dataset", "Experiment", "Result", and "Conclusion". The "Theory" button is highlighted in green.

Red arrows indicate the flow of information from the paper title in the PAPERS section to the summary in the RECOMMENDED section, and then to the recommended phrases in the SUMMARIZATION section.

*А.Власов.* Методы полуавтоматической суммаризации подборок научных статей. 2020. ФУПМ МФТИ.

*С.Крыжановская.* Технология полуавтоматической суммаризации тематических подборок научных статей. 2022. ВМК МГУ.



## Концепция MAHS (Machine Aided Human Summarization)

- 1 Система рекомендует *сценарий реферата* — список статей **подборки**, ранжированный в порядке упоминания
- 2 **Пользователь** может скорректировать сценарий в соответствии со своими целями и творческим замыслом
- 3 В цикле по ранжированному списку статей **подборки**:
  - **пользователь** запрашивает аспекты статьи у суфлёров: «как другие авторы ссылаются на эту статью», «цель», «идея», «подход», «достижение», «недостаток», «результат», «вывод» и т.д.
  - **суфлёр** выдаёт ранжированный список найденных фраз
  - **пользователь** добавляет фразу из поисковой выдачи и корректирует её в соответствии с целями и замыслом

---

А.Власов. Методы полуавтоматической суммаризации подборок научных статей. 2020. ФУПМ МФТИ.

С.Крыжановская. Технология полуавтоматической суммаризации тематических подборок научных статей. 2022. ВМК МГУ.

# Полуавтоматическое реферирование тематических подборок

## Задачи машинного обучения для МАНС:

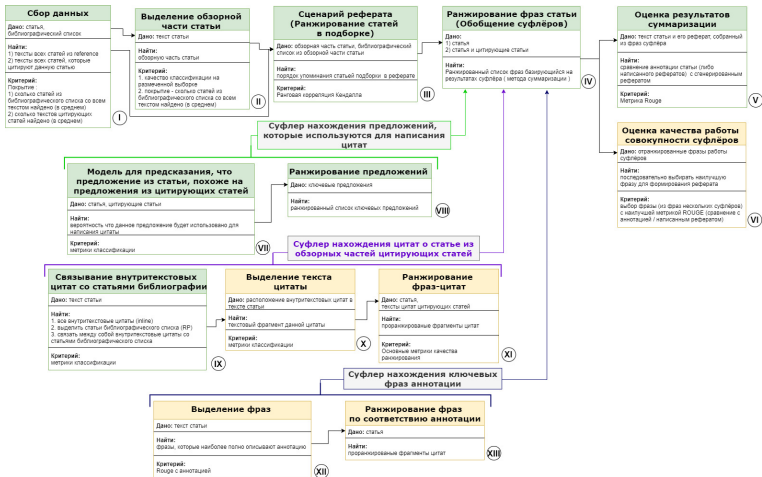
- 1 Формирование обучающей выборки: paper  $\rightarrow$  (refs, survey)
- 2 Ранжирование статей подборки для сценария реферата
- 3 Выбор релевантных фраз из текста статьи для сфлэра
- 4 Ранжирование выбранных фраз для каждого сфлэра
- 5 Выбор начала и конца контекста фразы, в частности, выбор релевантного контекста вокруг ссылки:

Few contextual citation graphs are publicly available. The ACL Anthology Network (AAN) (Radev et al., 2009) is one such contextual citation graph built from the ACL Anthology corpus (Bird et al., 2008), consisting of 24.6K papers manually augmented with citation information. CiteSeer (Giles et al., 1998) provides a large corpus consisting of 1.0M papers with full text and bibliography entries parsed from PDFs. Saier and Farber (2019) introduces a contextual citation graph of approximately 1.0M arXiv papers with full text LaTeX parses where citations are linked to papers in the Microsoft Academic Graph.

---

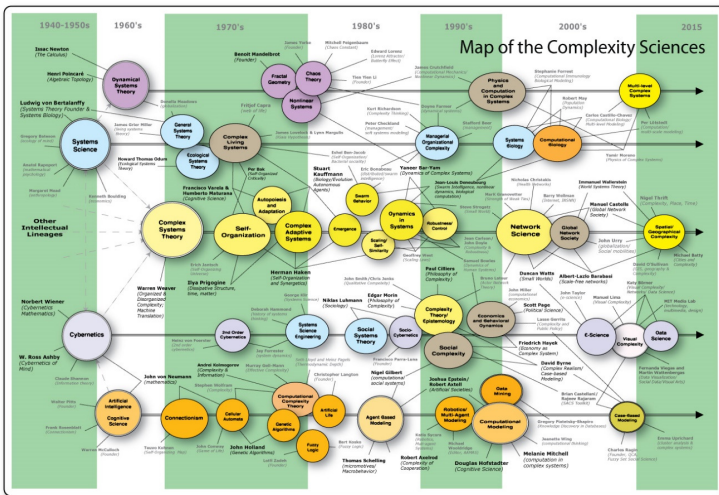
*M. Yasunaga et al.* ScisummNet: A large annotated corpus and content-impact models for scientific paper summarization with citation networks. 2019.

# Систематизация задач машинного обучения для MAHS



А. Власов. Методы полуавтоматической суммаризации подборок научных статей. 2020. ФУПМ МФТИ.

# Пример карты предметной области, построенной вручную



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

## Тематическая карта научных знаний (концепт)

- **Интерпретация осей:** темы/время/важность/сложность
- **Иерархичность:** темы делятся на подтемы
- **Спектр тем:** гуманитарные → естественные → точные
- **Интерактивность:** реализация мантры Шнейдермана
- **Суммаризация:** масштаб карты определяет объём текста



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.*

## Направления исследований

- Участие в разработке «Мастерской знаний»
- Эффективный векторный поиск в коллекциях 100M+
- Мультиязычные нейросетевые модели языка (ruSciBERT)
- Сервис полуавтоматической суммаризации
  - задача выделения обзорных разделов из публикации
  - задача ранжирования статей (сценарий реферата)
  - задачи ранжирования фраз для суфлёров
  - задача выделения цитирующего фрагмента
- Сервис выявления новых тем и трендов
- Сервис визуализации «карты знаний»
- Сервис контекстных подсказок

## Концепция «новостного коллайдера»

Цель создания адронного коллайдера — сталкивая потоки частиц, узнать больше о строении материи

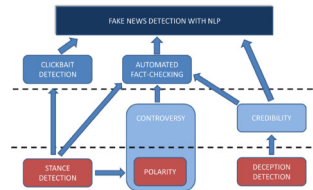


Цель создания новостного коллайдера — сталкивая потоки новостей, защитить общество от угроз эпохи постправды и информационных войн



## Область исследований «Fake News Detection»

- 1 **Deception Detection**  
выявление обмана в тексте
- 2 **Automated Fact-Checking**  
автоматическая проверка фактов
- 3 **Stance Detection**  
выявление позиции за или против
- 4 **Controversy Detection**  
выявление и кластеризация разногласий
- 5 **Polarization Detection**  
выявление полярных позиций
- 6 **Clickbait Detection**  
противоречия заголовка и текста
- 7 **Credibility Scores**  
оценка достоверности источников



*E.Saquete, D.Tomas, P.Moreda,  
P.Martinez-Barco, M.Palomar.*

**Fighting post-truth using  
natural language processing:  
a review and open challenges.**

Expert Systems With  
Applications, Elsevier, 2020.

# Задачи Propaganda/Manipulation/Persuasion Detection

## Базовая разметка: «фрагмент, метка класса»



Gallia est omnis divisa in partes tres, quarum unam incolunt Belgae, aliam Aquitani, tertiam qui ipsorum lingua Celtae, nostra Galli appellantur. Hi omnes lingua, institutis, legibus inter se differunt. Gallos ab Aquitania Garumna flumen, a Belgis Matrona et Sequana dividit. Horum **omnium fortissimi** sunt Belgae, propterea quod a cultu atque humanitate provinciae longissime absunt, minimeque ad eos mercatores saepe comeant atque ea quae ad effeminandos **animos pertinent important**, proximique sunt Germanis, qui trans Rhenum incolunt, quibuscum continenter bellum gerunt. Qua de causa **Helvetii quoque reliquos Gallos virtute praecedunt, quod fere cotidianis proelis cum Germanis contendunt**, cum aut suis finibus eos prohibent aut ipsi in eorum finibus bellum gerunt. Eorum una pars, quam Gallos obtinere dictum est, initium capit a flumine Rhodano, continetur Garumna flumine, Oceano, finibus Belgarum, attingit etiam ab Sequanis et Helvetiis flumen Rhenum, vergit ad septentriones. Belgae ab extremis Galliae finibus oriuntur, pertinent

Manipulative Wording: Loaded Language

Attack on Reputation: Smears

Manipulative Wording: Exaggeration

Justification: Appeal to Values



Commissio  
PopulusQue  
Europaea

## Упрощённая разметка: «предложение, метка класса»

## Продвинутая разметка: «фрагмент, мишень, метка класса»













*SemEval-2023 task 3*. Detecting the genre, the framing, and the persuasion techniques in online news in a multi-lingual setup.

<https://propaganda.math.unipd.it/semEval2023task3>

G.Martino, P.Nakov et al. A survey on computational propaganda detection. 2020.

## Типология угроз и задачи их автоматической детекции

воздействия → фейки → пропаганда → инф.война

1.  детекция приёмов манипулирования
2.  детекция замалчивания
3.  детекция обмана (deception detection), слухов (rumors d.), мистификаций (hoaxes d.)
4.  детекция кликбэйта (clickbait detection)
5.  автоматическая проверка фактов (auto fact-checking)
6.  детекция позиции (stance d.), противоречий (controversy d.), поляризации (polarization d.)
7.  выявление конструктов картины мира: идеологем, мифологем
8.  оценивание возможных психо-эмоциональных реакций
9.  выявление целевых аудиторий воздействия
10.  оценивание и предсказание скорости распространения (virality prediction)
11.  оценивание достоверности источников (credibility scores)
12.  детекция прямой агрессии (угрозы, призывы, провокации, вербовка, экстремизм)

*E.Saquete, D.Tomas, P.Moreda, P.Martinez-Barco, M.Palomar. Fighting post-truth using natural language processing: A review and open challenges // Expert Systems With Applications, Elsevier, 2020.*

## Типы задач ML/NLU для мониторинга медиа-пространства

- 1. Классификация текста (сообщения/предложения) целиком**
  - deception detection, fact-checking, text credibility
- 2. Классификация пары текстов**
  - stance, controversy, polarization, clickbait detection
  - выявление противоречий, разногласий, замалчивания
- 3. Разметка текста (выделение и классификация фрагментов)**
  - поиск лингвистических маркеров (linguistic-based cues) в тексте
  - детекция приёмов манипулирования
  - выявление конструктов картины мира: мифологем, идеологем
  - выявление психо-эмоциональных реакций и целевых аудиторий
- 4. Кластеризация или тематическое моделирование**
  - кластеризация мнений по заданной теме (controversy detection)
  - выявление поляризованных мнений (polarization detection)
  - выявление мнений как сочетаний слов, семантических ролей и тональностей
  - выявление «картин мира» – устойчивых сочетаний суждений и идеологем

# ПРО//ЧТЕНИЕ — технологический конкурс Up Great

**Задача:** поиск смысловых ошибок в сочинениях ЕГЭ по русскому, литературе, истории, обществознанию, английскому

**Период:** декабрь 2019 — декабрь 2022

**Призовой фонд:**

— 100М руб. русский язык

— 100М руб. английский язык

**Типов ошибок:** 152

(р:70 л:16 о:23 и:20 а:23)

**Подтипов ошибок:** 236

(р:112 л:19 о:29 и:26 а:50)

Алгоритм должен выделять ошибки и давать их объяснения.



**ФАКТИЧЕСКАЯ ОШИБКА**  
 автор высказывания А.Франц

В своем высказывании «Если человек зависит от природы, то и она от него зависит» Д. Мережковский **говорит** о необходимости защиты природы.

**ЛОГИЧЕСКАЯ ОШИБКА**  
 тезис не обоснован

Официальный сайт конкурса: <http://ai.upgreat.one>



## Модели разметки для формализации гуманитарных знаний

**Цель** — автоматизация обработки текстовых источников (контент-анализа и др.) в социогуманитарных исследованиях

**Гипотеза:** достаточно четырёх базовых операций разметки:

- 1 выделить фрагмент
- 2 классифицировать (тегировать) фрагмент по рубрикатору
- 3 связать несколько фрагментов
- 4 дать комментарий (затекст) к фрагменту или связи

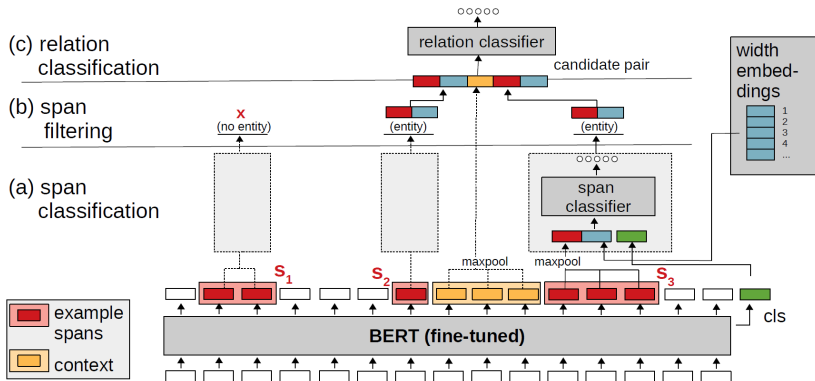
**Задачи** универсализации обучаемой модели разметки:

- 1 унификация правил разметки и инструментария разметки
- 2 унификация нейросетевой архитектуры модели разметки
- 3 унификация методики оценивания моделей разметки





# Унификация нейросетевых архитектур моделей разметки

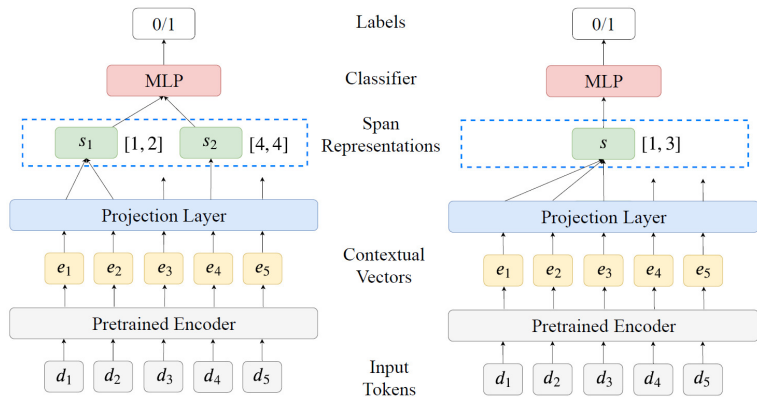


*M.Eberts, A.Ulges.* Span-based joint entity and relation extraction with transformer pre-training. 2020.

*L.Anisiutin, T.Batura, N.Shvarts.* Information extraction from news texts using a joint deep learning model. 2021.

*Wayne Xin Zhao et al.* A Survey of Large Language Models. ArXiv, 29 Jun 2023.

## Сравнение методов формирования эмбедингов фрагментов



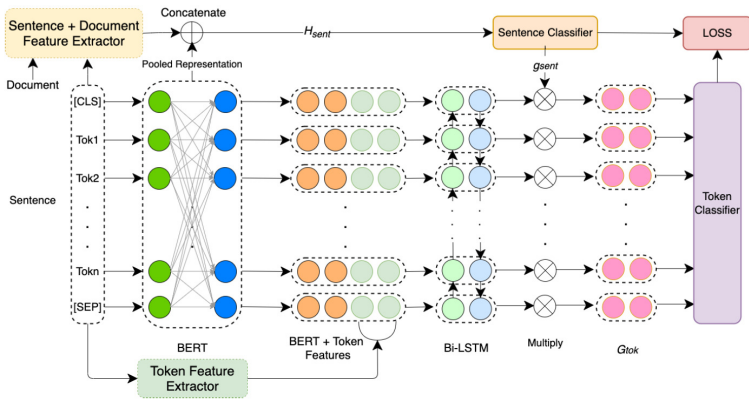
Полнота (recall) до 90% на задачах NER, SRL, Mention Detection

Xiaoya Li et al. A Unified MRC Framework for Named Entity Recognition. 2022.

S. Toshniwal et al. A Cross-Task Analysis of Text Span Representations. 2020.

## Извлечение признаков предложений и документов

## Задача детекции фрагментов с приёмами пропаганды



Sopan Khosla et al. LTIatCMU at SemEval-2020 Task 11: Incorporating Multi-Level Features for Multi-Granular Propaganda Span Identification. 2020.

## Унификация методики оценивания моделей разметки

- В основе методики — сравнение пар разметок текста: «алгоритм – эксперт», «эксперт-1 – эксперт-2», путём оптимального сопоставления их элементов
- Вводятся меры согласованности пары разметок  $Con_k(A, B)$
- Вводится их средневзвешенная согласованность  $Con(A, B)$
- СТАР (Средняя Точность Алгоритмической Разметки) — средняя по размеченной выборке согласованность  $Con(A, E)$  разметки модели  $A$  и разметки эксперта  $E$
- СТЭР (Средняя Точность Экспертной Разметки) — средняя по размеченной выборке согласованность  $Con(E_1, E_2)$  разметок двух экспертов,  $E_1$  и  $E_2$
- ОТАР = СТАР / СТЭР (Относительная Точность Алгоритмической Разметки) — если выше 100%, то это означает, что алгоритм работает не хуже экспертов

## Направления исследований

- Технология формализации гуманитарных знаний (разметка + языковая модель + оценивание)
- Универсальная нейросетевая модель для автоматической разметки текстов (на основе архитектуры трансформера)
- Выявление ценностей социокультурного кода
- Выявление манипулятивных приёмов
- Выявление поляризации общественного мнения
- Выявление психо-эмоциональных воздействий
- Обнаружение новых тем и событий (topic detection and tracking)

## Эволюция подходов машинного обучения в анализе текстов

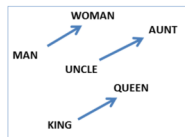
### Декомпозиция задач по уровням пирамиды NLP

- морфологический анализ, лемматизация, опечатки
- синтаксический анализ, выделение терминов, NER
- семантический анализ, выделение фактов, тем



### Модели векторных представлений (эмбедингов) слов на основе матричных разложений

- модели дистрибутивной семантики: word2vec [Mikolov, 2013], FastText [Bojanowski, 2016]
- тематические модели LDA [Blei, 2003], ARTM [2014]



### Нейросетевые модели локальных контекстов

- рекуррентные нейронные сети
- модели внимания и трансформеры: BERT [2018], GPT-3 [2020] и др.

$$\text{softmax} \left( \frac{\begin{matrix} Q \\ \text{grid} \end{matrix} \times \begin{matrix} K^T \\ \text{grid} \end{matrix}}{\sqrt{d}} \right) \begin{matrix} V \\ \text{grid} \end{matrix}$$

## Постановка задачи тематического моделирования

**Дано:**

- $W$  — конечное множество (словарь) термов (слов, токенов)
- $D$  — конечное множество (коллекция) документов
- $n_{dw}$  — частота термина  $w \in W$  в документе  $d \in D$

**Найти:** вероятностную тематическую языковую модель

$$p(w|d) = \sum_{t \in T} p(w | \cancel{d}, t) p(t|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$$

где  $\phi_{wt} = p(w|t)$ ,  $\theta_{td} = p(t|d)$  — параметры модели

**Критерий:** максимум логарифма правдоподобия

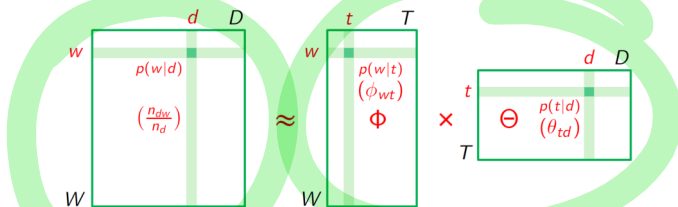
$$L(\Phi, \Theta) = \ln \prod_{d,w} p(w|d)^{n_{dw}} = \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} \rightarrow \max_{\Phi, \Theta}$$

при ограничениях  $\phi_{wt} \geq 0$ ,  $\sum_w \phi_{wt} = 1$ ,  $\theta_{td} \geq 0$ ,  $\sum_t \theta_{td} = 1$

*Hofmann T.* Probabilistic Latent Semantic Indexing. ACM SIGIR, 1999.

## Некорректно поставленная задача матричного разложения

Низкоранговое стохастическое матричное разложение:



Если  $\Phi, \Theta$  — решение, то стохастические  $\Phi', \Theta'$  — тоже решения

- $\Phi' \Theta' = (\Phi S)(S^{-1} \Theta)$ ,  $\text{rank} S = |T|$
- $L(\Phi', \Theta') = L(\Phi, \Theta)$  — линейно не зависимые решения
- $L(\Phi', \Theta') \geq L(\Phi, \Theta) - \varepsilon$  — приближённые решения

**Регуляризация** — стандартный приём доопределения решения с помощью добавления дополнительных критериев.



## ARTM: аддитивная регуляризация тематических моделей

Максимизация логарифма правдоподобия с регуляризатором:

$$\sum_{d,w} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}; \quad R(\Phi, \Theta) = \sum_i \tau_i R_i(\Phi, \Theta)$$

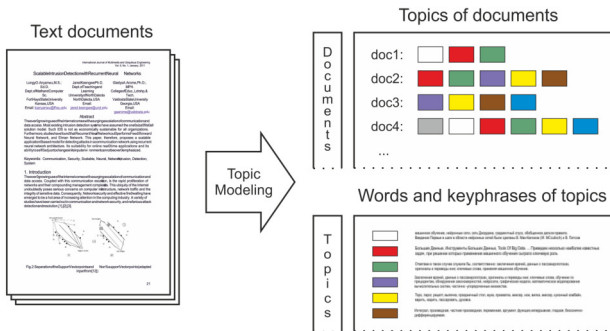
EM-алгоритм: метод простой итерации для системы уравнений

$$\begin{cases} \text{E-шаг:} & p_{tdw} \equiv p(t|d, w) = \operatorname{norm}_{t \in T}(\phi_{wt} \theta_{td}) \\ \text{M-шаг:} & \begin{cases} \phi_{wt} = \operatorname{norm}_{w \in W} \left( n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right), & n_{wt} = \sum_{d \in D} n_{dw} p_{tdw} \\ \theta_{td} = \operatorname{norm}_{t \in T} \left( n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right), & n_{td} = \sum_{w \in D} n_{dw} p_{tdw} \end{cases} \end{cases}$$

Воронцов К. В. Аддитивная регуляризация тематических моделей коллекций текстовых документов. Доклады РАН, 2014.

## Мультимодальная тематическая модель

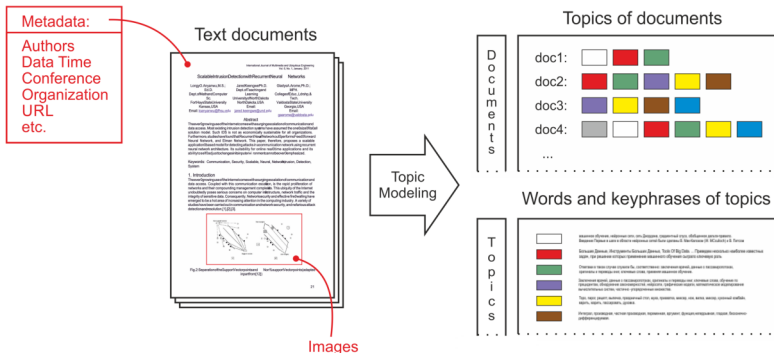
Тема может порождать термины различных *модальностей*:  
 $p(\text{слово} | t)$ ,  $p(n\text{-грамма} | t)$ ,





## Мультимодальная тематическая модель

Тема может порождать термины различных *модальностей*:  
 $p(\text{слово} | t)$ ,  $p(n\text{-грамма} | t)$ ,  $p(\text{автор} | t)$ ,  $p(\text{время} | t)$ ,  $p(\text{источник} | t)$ ,  
 $p(\text{объект} | t)$ ,

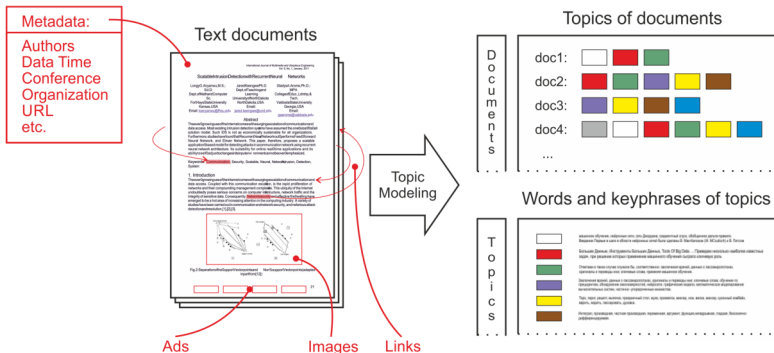




## Мультимодальная тематическая модель

Тема может порождать термины различных *модальностей*:

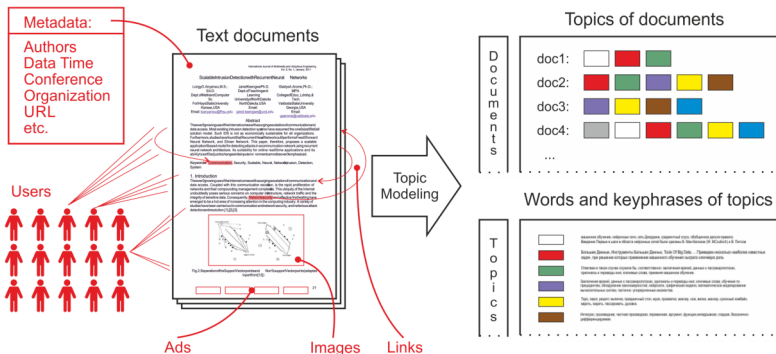
$p(\text{слово} | t)$ ,  $p(n\text{-грамма} | t)$ ,  $p(\text{автор} | t)$ ,  $p(\text{время} | t)$ ,  $p(\text{источник} | t)$ ,  
 $p(\text{объект} | t)$ ,  $p(\text{ссылка} | t)$ ,  $p(\text{баннер} | t)$ ,



## Мультимодальная тематическая модель

Тема может порождать термины различных *модальностей*:

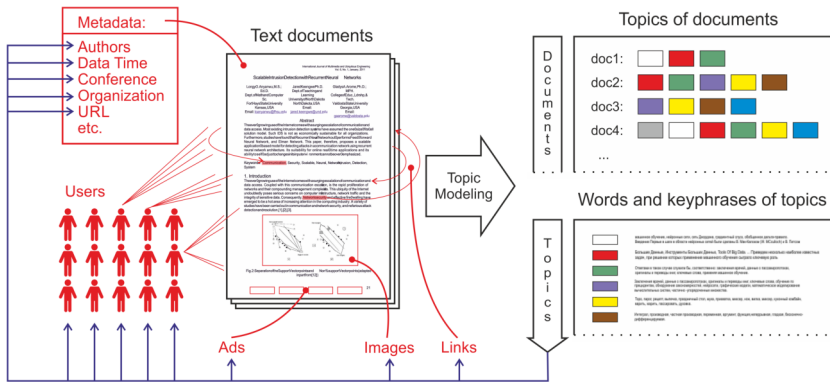
$p(\text{слово} | t)$ ,  $p(n\text{-грамма} | t)$ ,  $p(\text{автор} | t)$ ,  $p(\text{время} | t)$ ,  $p(\text{источник} | t)$ ,  
 $p(\text{объект} | t)$ ,  $p(\text{ссылка} | t)$ ,  $p(\text{баннер} | t)$ ,  $p(\text{пользователь} | t)$



# Мультимодальная тематическая модель

Тема может порождать термины различных *модальностей*:

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





## Мультимодальная ARTM

$W_m$  — словарь термов  $m$ -й модальности,  $m \in M$

Максимизация суммы log-правдоподобий с регуляризацией:

$$\sum_{m \in M} \tau_m \sum_{d \in D} \sum_{w \in W^m} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

EM-алгоритм: метод простой итерации для системы уравнений

$$\begin{cases} \text{E-шаг:} & p_{tdw} = \operatorname{norm}_{t \in T}(\phi_{wt} \theta_{td}) \\ \text{M-шаг:} & \begin{cases} \phi_{wt} = \operatorname{norm}_{w \in W^m} \left( \sum_{d \in D} \tau_m(w) n_{dw} p_{tdw} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right) \\ \theta_{td} = \operatorname{norm}_{t \in T} \left( \sum_{w \in W^m} \tau_m(w) n_{dw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right) \end{cases} \end{cases}$$

*K. Vorontsov, O. Freij, M. Apishev et al.* Non-Bayesian additive regularization for multimodal topic modeling of large collections. CIKM TM workshop, 2015.



## Необходимые условия экстремума и метод простых итераций

Операция нормировки вектора:  $p_i = \text{norm}_{i \in I}(x_i) = \frac{\max(x_i, 0)}{\sum_k \max(x_k, 0)}$

**Лемма.** Пусть  $f(\Omega)$  непрерывно дифференцируема по  $\Omega$ . Если  $\omega_j$  — вектор локального экстремума задачи  $f(\Omega) \rightarrow \max$  и  $\exists i: \omega_{ij} \frac{\partial f}{\partial \omega_{ij}} > 0$ , то  $\omega_j$  удовлетворяет системе уравнений

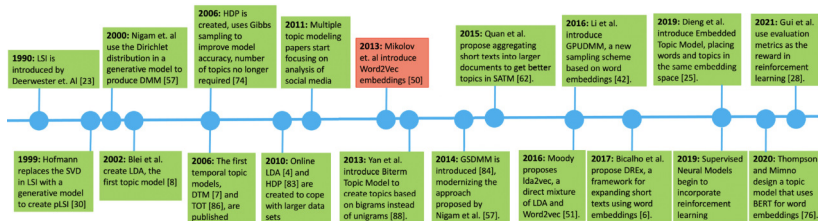
$$\omega_{ij} = \text{norm}_{i \in I_j} \left( \omega_{ij} \frac{\partial f}{\partial \omega_{ij}} \right).$$

- Численное решение системы — методом простых итераций
- Векторы  $\omega_j = 0$  отбрасываются как вырожденные решения
- Итерации похожи на градиентную оптимизацию:

$$\omega_{ij} := \omega_{ij} + \eta \frac{\partial f}{\partial \omega_{ij}},$$

но учитывают ограничения и не требуют подбора шага  $\eta$

## Эволюция тематического моделирования



Neural Topic Models — поток публикаций начиная с 2016

Как «объединить лучшее от двух миров»?

- **Neural:** качество, универсальность, генеративность
- **Topic:** скорость, интерпретируемость, простота

**Что объединяет:** векторизация, оптимизация, регуляризация, гомогенизация, локализация (контекст и внимание)

*Rob Churchill, Lisa Singh. The Evolution of Topic Modeling. November, 2022.*

## Направления исследований

- Строить любые тематические модели в `pyTorch`, где градиентный шаг на симплекс уже реализован
- Neural Topic Modeling (NTM) — как совмещать интерпретируемость тематических моделей с преимуществами глубоких нейросетевых моделей языка
- В том числе, создание тематических моделей внимания
- Открытая проблема — несбалансированность тем
- Автоматическое именование и аннотирование тем (каждая тема должна уметь «рассказать о себе»)
- Как приблизить долю интерпретируемых тем к 100%
- Динамическое обнаружение новых тем и трендов (Topic Detection & Tracking, First Story Detection)
- Автоматический подбор гиперпараметров, AutoML

# 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

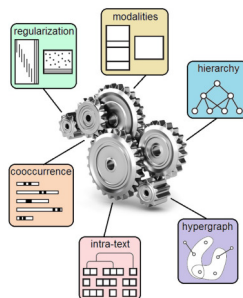
---

*K. Vorontsov, O. Frej, M. Apishev, P. Romov, M. Suvorova.* BigARTM: Open source library for regularized multimodal topic modeling of large collections. AIST 2015.

## 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

---

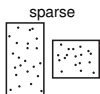
*V. Bulatov, E. Egorov, E. Veselova, D. Polyudova, V. Alekseev, A. Goncharov, K. Vorontsov.*  
TopicNet: making additive regularisation for topic modelling accessible. LREC-2020

## Регуляризаторы для улучшения интерпретируемости тем



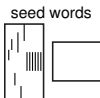
Сглаживание фоновых тем  $B \subset T$ :

$$R(\Phi, \Theta) = \beta_0 \sum_{t \in B} \sum_w \beta_w \ln \phi_{wt} + \alpha_0 \sum_d \sum_{t \in B} \alpha_t \ln \theta_{td}$$



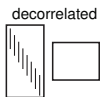
Разреживание предметных тем  $S = T \setminus B$ :

$$R(\Phi, \Theta) = -\beta_0 \sum_{t \in S} \sum_w \beta_w \ln \phi_{wt} - \alpha_0 \sum_d \sum_{t \in S} \alpha_t \ln \theta_{td}$$



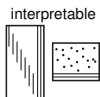
Сглаживание для выделения релевантных тем

с помощью словаря «затравочных» ключевых слов



Декоррелирование для повышения различности тем:

$$R(\Phi) = -\frac{\tau}{2} \sum_{t,s} \sum_w \phi_{wt} \phi_{ws}$$



Сглаживание + разреживание + декоррелирование  
для улучшения интерпретируемости тем



## Регуляризаторы для учёта дополнительной информации

regression



Линейная модель регрессии  $\hat{y}_d = \langle v, \theta_d \rangle$  документов:

$$R(\Theta, v) = -\tau \sum_{d \in D} \left( y_d - \sum_{t \in T} v_t \theta_{td} \right)^2$$

biterm



Связи сочетаемости слов ( $n_{uv}$  — частота битерма):

$$R(\Phi) = \tau \sum_{u \in W} \sum_{v \in W} n_{uv} \ln \sum_{t \in T} n_t \phi_{ut} \phi_{vt}$$

relational



Связи или ссылки между документами:

$$R(\Theta) = \tau \sum_{d, c \in D} n_{dc} \sum_{t \in T} \theta_{td} \theta_{tc}$$

hierarchy



Связи родительских тем  $t$  с дочерними подтемами  $s$ :

$$R(\Phi, \Psi) = \tau \sum_{t \in T} \sum_{w \in W} n_{wt} \ln \sum_{s \in S} \phi_{ws} \psi_{st}$$

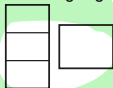
# Регуляризаторы для мультимодальных тематических моделей

supervised



Модальности меток классов или категорий для задач классификации и категоризации текстов.

multilanguage

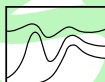


Модальность языков и регуляризация со словарём

$\pi_{uwt} = p(u|w, t)$  переводов с языка  $k$  на  $\ell$ :

$$R(\Phi, \Pi) = \tau \sum_{u \in W^k} \sum_{t \in T} n_{ut} \ln \sum_{w \in W^\ell} \pi_{uwt} \phi_{wt}$$

temporal



Темпоральные модели с модальностью времени  $i$ :

$$R(\Phi) = -\tau \sum_{i \in I} \sum_{t \in T} |\phi_{it} - \phi_{i-1,t}|.$$

geospatial



Модальность геолокаций  $g$  с близостью  $S_{gg'}$ :

$$R(\Phi) = -\frac{\tau}{2} \sum_{g, g' \in G} S_{gg'} \sum_{t \in T} n_t^2 \left( \frac{\phi_{gt}}{n_g} - \frac{\phi_{g't}}{n_{g'}} \right)^2$$

## Регуляризаторы для моделирования последовательного текста

sentence



Тематические модели, учитывающие границы предложений, абзацев и секций документов

n-gram



Модели с модальностями  $n$ -грамм, коллокаций, именованных сущностей (используем TopMine)

syntax



Модели, учитывающие результаты автоматического синтаксического разбора (используем UDPipe)

sentiment



Модели выделения мнений на основе тональностей, фактов, семантических ролей именованных сущностей

segmentation



Тематические модели сегментации с автоматическим определением границ сегментов

## ARTM — модульный подход к синтезу требуемых моделей

Для построения композитных моделей в BigARTM не нужны ни математические выкладки, ни программирование «с нуля».

### Этапы моделирования

#### Bayesian TM

#### ARTM

Формализация:

Анализ требований

Анализ требований

Вероятностная модель  
порождения данных

Стандартные  
критерии

Свои  
критерии

Алгоритмизация:

Байесовский вывод для  
данной порождающей модели  
(VI, GS, EP)

Единый регуляризованный  
EM-алгоритм для любых  
моделей и их композиций

Реализация:

Исследовательский код  
(Matlab, Python, R)

Промышленный код BigARTM  
(C++, Python API)

Оценивание:

Исследовательские метрики,  
исследовательский код

Стандартные  
метрики

Свои метрики

Внедрение

Внедрение

-- нестандартизируемые этапы, уникальная разработка для каждой задачи

-- стандартизуемые этапы

## Разведочный поиск в технологических блогах

**Цель:** поиск документов

по длинным текстовым запросам

— Habr.ru (175К документов),

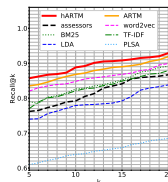
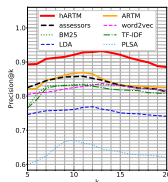
— TechCrunch.com (760К док.).

**Регуляризаторы:**

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

**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.*

## Поиск и классификация этно-релевантных тем в соцсетях

**Цель:** выявление как можно большего числа тем о национальностях и межнациональных отношениях (затравка — словарь 300 этнонимов).

**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{|c|} \hline \text{PLSA} \\ \hline \Phi \quad \Theta \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{seed words} \\ \hline \text{[Bar Chart]} \quad \square \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{interpretable} \\ \hline \text{[Bar Chart]} \quad \text{[Scatter Plot]} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{multimodal} \\ \hline \text{[Image]} \quad \square \\ \hline \end{array} \right) \\ + R \left( \begin{array}{|c|} \hline \text{temporal} \\ \hline \text{[Waveform]} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{geospatial} \\ \hline \text{[Map]} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|} \hline \text{sentiment} \\ \hline \text{[Sentiment Scale]} \\ \hline \end{array} \right) \rightarrow \max$$

**[японцы]** японский, япония, японя, корей, китайский, жилища, азияя, фукусима, цунами, собораи, якия, сланин, каино, рабон, цина, гласико, диланый,  
**[норвежцы]** дитя, ребенок, родился, детский, семья, воспитаный, поаво, возраст, отец, воспитаный, норвежский, родительский, родители, мальчик, взрослый, отец, сын,  
**[американцы]** айба, кастро, империализм, чикко, президент, итг, издурно, бонжани, фидель, гласа, катанский, коммунистический, лидер, болгарская, ирландский, зальвира, лидер,  
**[китайцы]** китайский, россия, производство, китай, продукция, страна, производство, компания, технология, азиатский, регион, производство, производственный, ориентация, российская, экономика, кит,  
**[азербайджанцы]** русский, азербайджан, азербайджанец, россия, азербайджанский, тиксист, диспоза, аналз, жарод, москва, страна, азербайджан, слово, рынок,  
**[германцы]** германский, спецназ, военный, август, батальон, российский, спецназ, министр, операция, румын, братство, микрофинансовый, абскал, группа, война, русский, цинвале,  
**[осетины]** конституция, осетия, азиат, русский, осетинский, цинвал, северный, регион, жабай, республика, история, азиат, республика, азия, азия, конфликт,  
**[американцы]** наркотики, америк, шатланд, лардшей, место, страна, деньги, время, работа, жизнь, жить, дуно, дин, цинвалский, наркотизма,

**Результаты:** число релевантных тем: 45 (LDA)  $\rightarrow$  83 (ARTM).

*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.*

*Mining ethnic content online with additively regularized topic models. 2016.*

## Аналогичные исследования по выделению узкой тематики

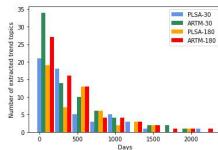
### Задачи «поиска и классификации иголок в стоге сена»

- 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.

## Выявление трендов в коллекции научных публикаций

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



**Регуляризаторы:**

$$\mathcal{L} \left( \begin{array}{c} \text{PLSA} \\ \Phi \quad \Theta \end{array} \right) + R \left( \begin{array}{c} \text{interpretable} \\ \text{[Bar chart icon]} \quad \text{[Scatter plot icon]} \end{array} \right) + R \left( \begin{array}{c} \text{dynamic} \\ \text{[Line graph icon]} \end{array} \right) + R \left( \begin{array}{c} \text{multimodal} \\ \text{[Stacked bar icon]} \quad \text{[Box icon]} \end{array} \right) + R \left( \begin{array}{c} \text{n-gram} \\ \text{[Grid icon]} \end{array} \right) \rightarrow \max$$

**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.*

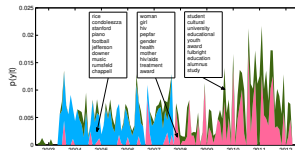
Инкрементальное обучение тематических моделей для поиска трендовых тем в научных публикациях. Доклады РАН, 2022.



## Выявление динамики тем в новостных потоках

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

**Регуляризаторы:**



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

**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	Pr	Rec	F1
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:**

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

**Результаты:**

- выделение мнений внутри тем: F1-мера = 0.86%
- совместное использование трёх модальностей:
  - SPO — факты как триплеты «субъект–предикат–объект»
  - FR — семантические роли слов по Филлмору
  - Sent — 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 ... кошмарит 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.

## 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 (2022)
- 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

## 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

- отбор моделей, оценивание и отбор релевантных тем
- рекомендации документов для тематических подборок

## Направления исследований

- Участие в разработке «Тематизатора»
- Встраивание модуля BigARTM в Orange
- Встраивание модуля BigARTM в PolyAnalyst
- Разработка сценариев разведочного анализа текстовых данных в социо-гуманитарных исследованиях
- Разработка или адаптация средств визуализации тематических моделей в модуле «Тематизатора»
- Разработка адаптивных стратегий AutoML для оптимизации гиперпараметров тематических моделей

---

*К.Воронцов.* Вероятностное тематическое моделирование: теория регуляризации ARTM и библиотека с открытым кодом BigARTM. 2023. (для изд-ва URSS)

<http://www.MachineLearning.ru/wiki/images/d/d5/Voron17survey-artm.pdf>