

О некоторых технологиях информационного поиска в текстовых массивах

Воронцов Константин Вячеславович

зав. лаб. Машинного интеллекта и семантического анализа
Института ИИ МГУ; проф., и.о. зав. каф. ММП ВМК МГУ;
проф., зав. каф. МОЦГ МФТИ; г.н.с. ФИЦ ИУ РАН

Круглый стол ВНИИДАД «Практические задачи внедрения
технологий искусственного интеллекта в деятельность архивов»

• 10 апреля 2023 •

1 Разведочный поиск

- Векторный документный поиск
- Полуавтоматическое реферирование
- Distant reading и визуализация

2 Тематический поиск

- Тематическое моделирование
- Поиск документов по документам или словарям
- Темпоральные тематические модели

3 Лингвистический поиск

- Конкурс ПРО//ЧТЕНИЕ: поиск смысловых ошибок
- Поиск обмана, фейков, противоречий
- Унификация разметки, моделирования и оценивания

Концепция разведочного поиска (exploratory search)

- пользователь может не знать ключевых терминов
- запросом может быть текст произвольной длины
- информационная потребность — систематизация знаний



Gary Marchionini. Exploratory Search: from finding to understanding. 2006.

Концепция проекта «Мастерская 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



От поиска информации к «Мастерской знаний»

Недостатки обычного поиска:

- как искать новые знания?
- что делать с найденным?



Мастерская знаний — инструментарий для автоматизации **последующих этапов** работы с профессиональными знаниями:

- ищу – чтобы накапливать
- накапливаю – чтобы анализировать
- анализирую – чтобы понимать
- понимаю – чтобы применять и передавать

Эти задачи связаны с *автоматической обработкой текстов* (только применение знаний остаётся за пределами системы)

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*

Прототип поисково-рекомендательной системы

Тематическая подборка пользователя:

https://arxiv.aitheta.com/collections/Q29sbGVjdGVjbjozUFTUEFxaHBH

FEEDS | SEARCH | **COLLECTIONS** | About | FAQ | Konstantin Vorontsov

MOOC (massive open online course)

PAPERS | RECOMMENDED

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, Osmar Zaiane

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...

Citations: 0

Прототип поисково-рекомендательной системы

Список статей, рекомендуемых для добавления в подборку:

← → ↻ 🔒 https://arxiv.aitha.com/collections/Q29sbGVjZGJlbnJozUFVTUEFxaHBH

FEEDS | SEARCH | COLLECTIONS | About | FAQ | Konstantin Vorontsov

MOOC (massive open online course)

PAPERS → **RECOMMENDED**

2 JUN 2019

A Survey of Natural Language Generation Techniques with a Focus on Dialogue Systems - Past, Present and Future Directions

Sashank Santhanam, Samira Shaikh

One of the hardest problems in the area of Natural Language Processing and Artificial Intelligence is automatically generating language that is coherent and understandable to humans. Teaching machines how to converse as humans do falls under the broad umbrella of Natural Language Generation. Recent years have seen unprecedented growth in the number of research articles published on this subject in conferences and journals both by academic and industry researchers. There have...

Citations: 6

🔖 👍 🔄

20 SEP 2014

Capturing "attrition intensifying" structural traits from didactic interaction sequences of MOOC learners

Tanmay Sinha, Nan Li, Patrick Jermann, Pierre Dillenbourg

This work is an attempt to discover hidden structural configurations in learning activity sequences of students in Massive Open Online Courses (MOOCs). Leveraging combined representations of video clickstream interactions and forum activities, we seek to fundamentally understand traits that are predictive of decreasing engagement over time. Grounded in the interdisciplinary field of network science, we follow a graph based approach to successfully extract indicators of active and...

Citations: 0

🔖 👍 🔄

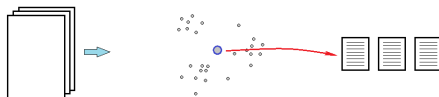
Прототип поисково-рекомендательной системы

Добавление статьи из списка рекомендаций в подборку:

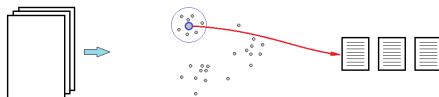
The screenshot shows a web browser window with the URL <https://arxiv.aitha.com/collections/Q29sbGVjZjJGJVbjozUjVVTUEFxaHBH>. The page title is "MOOC (massive open online course)". The main content area displays a list of papers under the heading "PAPERS". The first paper is "A Survey of Natural Language Generation T...", by Sashank Santhanam and Samira Shaikh, dated 2 JUN 2019. A red circle highlights the bookmark icon next to this paper. A red arrow points from this icon to the "Add to collections" dialog box. The dialog box is open, showing a list of collections: "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 (massive open online course)" option is selected with a radio button. A red circle highlights this option. A red arrow points from the selected option to the "SAVE CHANGES" button at the bottom of the dialog box. The "SAVE CHANGES" button is also highlighted with a red circle. Below the dialog box, there is a "NEW COLLECTION" link. In the background, a "RECOMMENDED" section is visible, with the word "RECOMMENDED" circled in red. The page also shows a navigation bar with "FEEDS", "SEARCH", and "COLLECTIONS" tabs, and a user profile "Konstantin Vorontsov".

Стратегии поиска документов по тематическим векторам

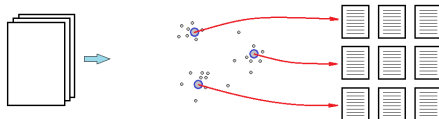
Поиск по среднему вектору **подборки** (неудачная стратегия):



Поиск по части **подборки**, близкой к выбранному документу:

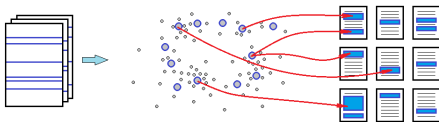


Поиск по тематике кластеров, на которые делится **подборка**:

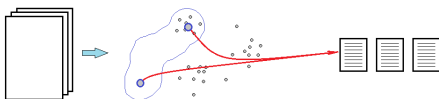


Стратегии поиска документов по тематическим векторам

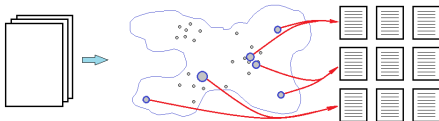
Поиск по тематике сегментов документов подборки:



Поиск по тематике, смежной для части подборки:



Поиск по тематике, смежной для всей подборки:



Полуавтоматическое реферирование тематических подборок

Рекомендации фраз для реферата с помощью сугфлёров:

The screenshot shows a web interface with three main sections: PAPERS, RECOMMENDED, and SUMMARIZATION.

- PAPERS:** A list of papers with titles and authors. The paper "SummaRuNNer: A Recurrent Neural Network based..." is highlighted with a red arrow pointing to the "Promoters" section.
- RECOMMENDED:** A summary of the selected paper, including an abstract and a goal statement. A red arrow points from the "Promoters" section to the "Recommended phrases" section.
- Promoters:** A set of buttons labeled "Annotate", "Idea", "Theory", "Method", "Citation", "Dataset", "Experiment", "Result", and "Conclusion". The "Theory" button is highlighted with a red arrow pointing to the "Recommended phrases" section.
- Recommended phrases:** A list of phrases extracted from the paper, such as "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."

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

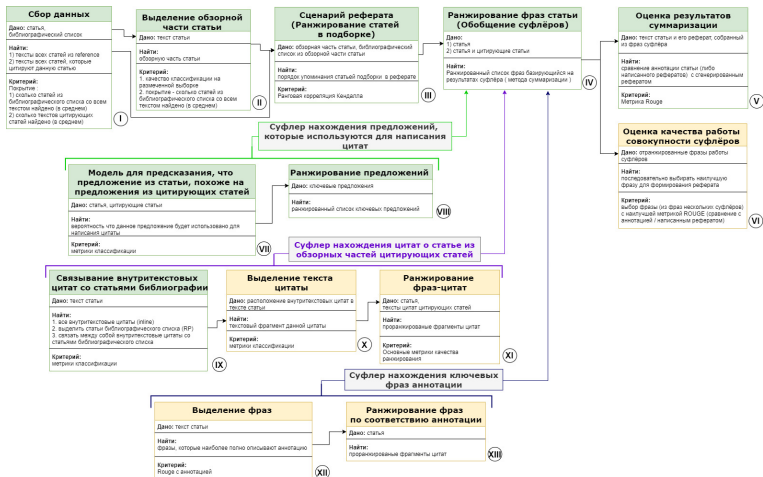
Задачи машинного обучения для МАНС:

- 1 Формирование обучающей выборки: paper \rightarrow (refs, survey)
- 2 Ранжирование статей подборки для сценария реферата
- 3 Выбор релевантных фраз из текста статьи для сфлёрра
- 4 Ранжирование выбранных фраз для каждого сфлёрра
- 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.

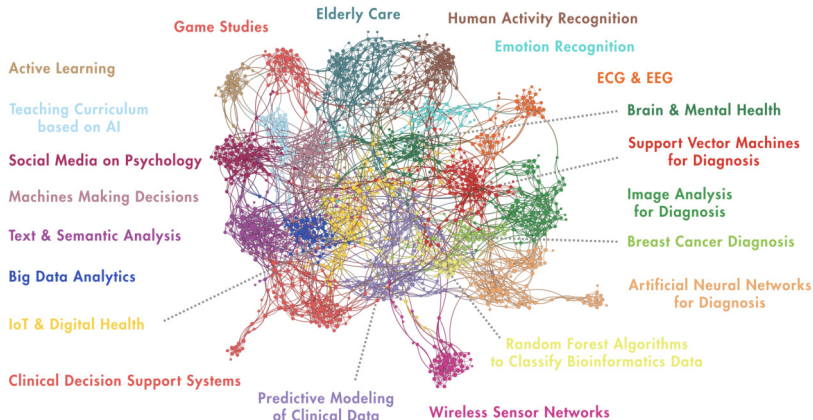
Полная систематизация задачи машинного обучения для MAHS



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

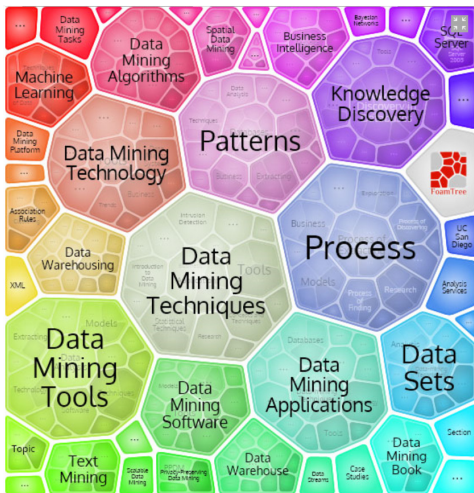
Пример тематической карты «ИИ в биомедицине»

Academic papers on AI in Healthcare published in 2016



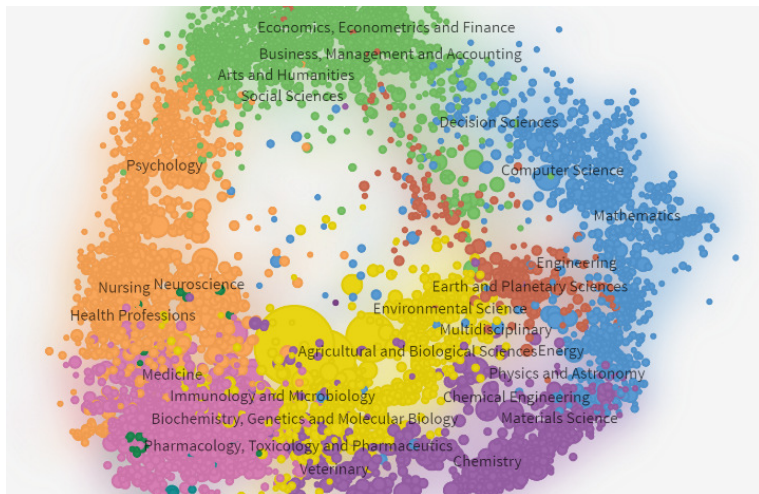
C.Folgar, J.McCuan. The 3 most-cited studies in healthcare and AI. Quid, 2017.

Пример иерархической карты области *Data Mining*



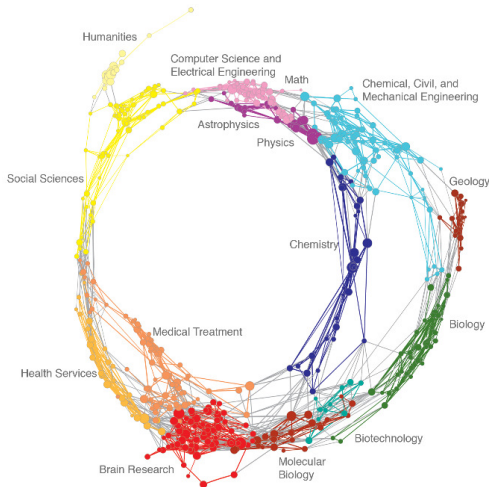
FoamTree: <https://carrotsearch.com/foamtree>

Пример карты науки



<http://onlinelibrary.wiley.com/browse/subjects>

Ещё один пример карты науки

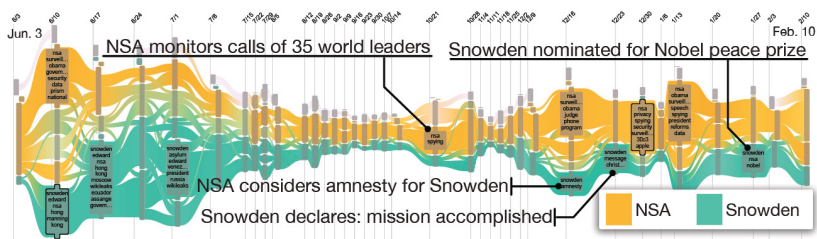


Важное наблюдение:
области знания
самопроизвольно
располагаются по кругу,
значит,
их можно располагать
и вдоль прямой линии.

Недостатки:

- оси не имеют интерпретации
- искажение сходства при двумерном проецировании

Динамика тем: эволюция предметной области



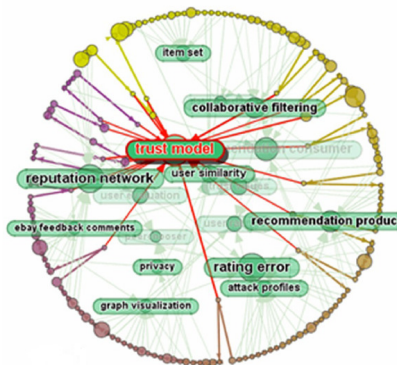
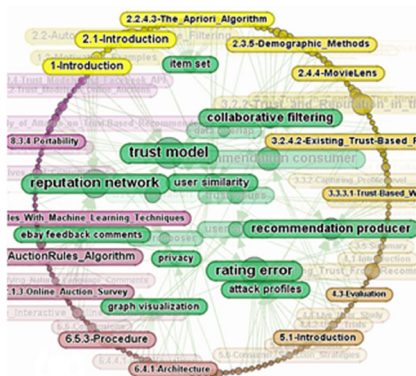
Эволюция выбранных тем иерархии. Данные Prism (2013/06/03–2014/02/09)

Стратегия визуализации в системах TextFlow и RoseRiver:

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

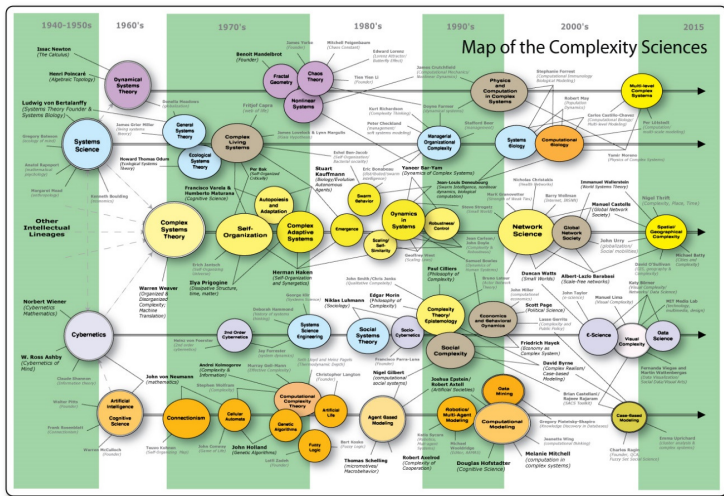
Weiwei Cui, Shixia Liu, Zhuofeng Wu, Hao Wei. How hierarchical topics evolve in large text corpora. 2014.

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.

Пример карты предметной области (построено 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.

Тематическое моделирование: «о чём все эти тексты?»

Дано:

- коллекция текстовых документов

Найти:

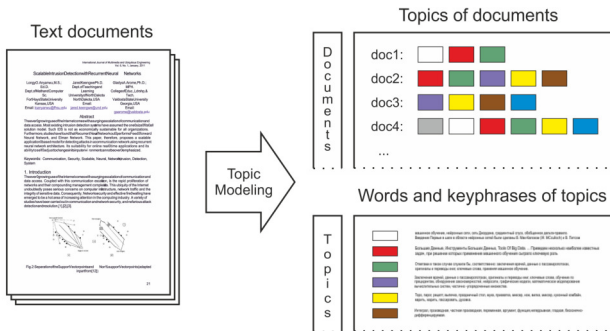
- T — множество тем, составляющих эту коллекцию
- $p(w|t) = \phi_{wt}$ — вероятности слов w в каждой теме t
- $p(t|d) = \theta_{td}$ — вероятности тем t в каждом документе d
- $p(w|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$ — вероятностная тематическая модель

Критерий: правдоподобие предсказания слов w в документах d с дополнительными критериями-регуляризаторами $R_i(\Phi, \Theta)$:

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

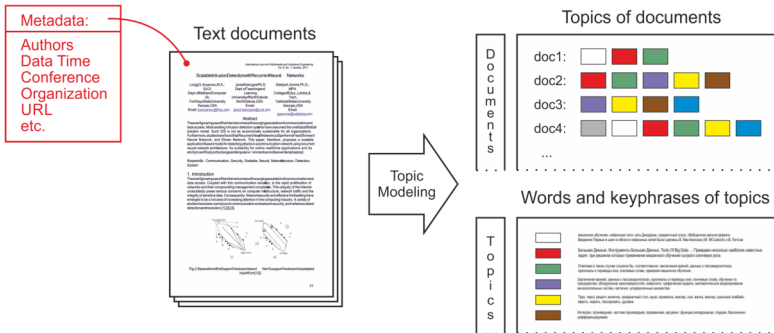
Мультимодальная тематическая модель

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



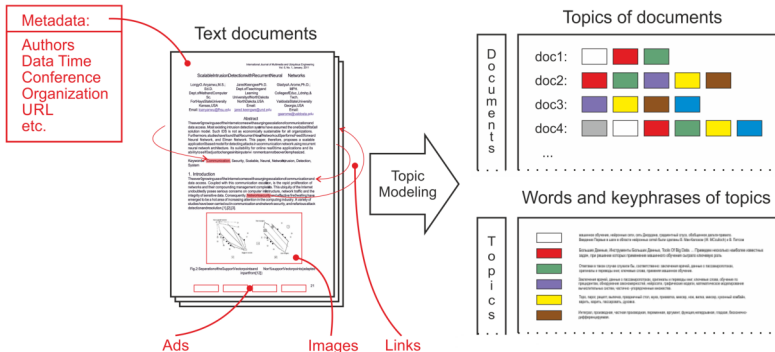
Мультимодальная тематическая модель

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



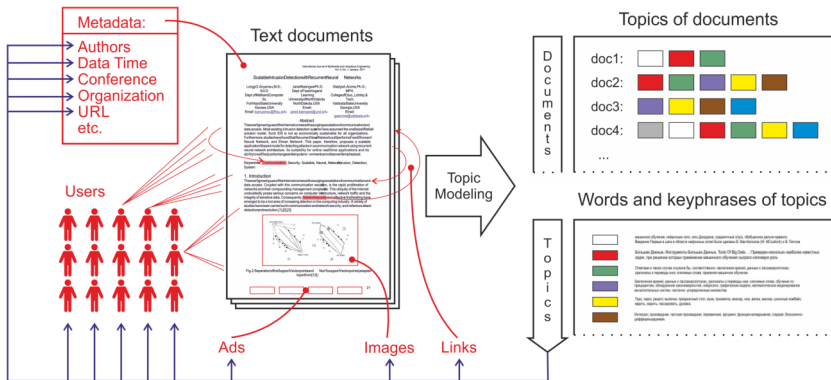
Мультимодальная тематическая модель

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



Мультимодальная тематическая модель

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



Пример 1. Мультиязычная тематическая 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.

K. Vorontsov, O. Freij, M. Apishev, P. Romov, M. Suvorova. BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections. AIST-2015.

Пример 1. Мультиязычная тематическая модель Википедии

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

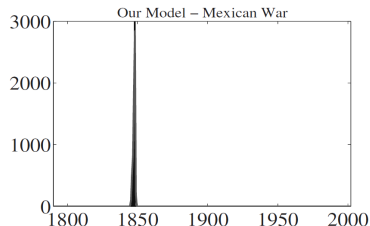
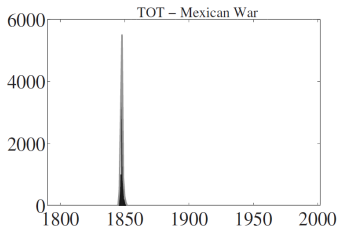
Тема №88				Тема №251			
opera	7.36	опера	7.82	windows	8.00	windows	6.05
conductor	1.69	оперный	3.13	microsoft	4.03	microsoft	3.76
orchestra	1.14	дирижер	2.82	server	2.93	версия	1.86
wagner	0.97	певец	1.65	software	1.38	приложение	1.86
soprano	0.78	певица	1.51	user	1.03	сервер	1.63
performance	0.78	театр	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.

K. Vorontsov, O. Freij, M. Apishev, P. Romov, M. Suvorova. BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections. AIST-2015.

Пример 2. Совмещение темпоральной и n -граммной модели

По коллекции выступлений президентов США



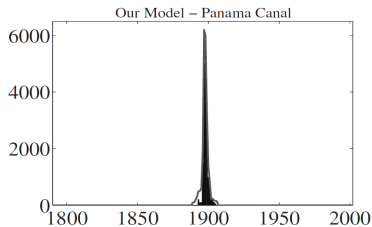
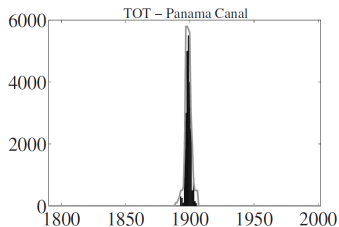
1. mexico	8. territory
2. texas	9. army
3. war	10. peace
4. mexican	11. act
5. united	12. policy
6. country	13. foreign
7. government	14. citizens

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. ECIR 2013.

Пример 2. Совмещение темпоральной и n -граммной модели

По коллекции 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. ECIR 2013.

BigARTM: библиотека тематического моделирования

Ключевые возможности:

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

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.

Разведочный поиск в технологических блогах

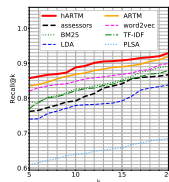
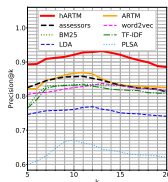
Цель: поиск документов по длинным текстовым запросам
 — Habr.ru (175К документов),
 — TechCrunch.com (760К док.).

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:

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

(японцы) японский, япония, япон, китайский, жилища, азия, фукусима, цунами, сакура, слики, слики-слики, озон, район, нана, гласиско, дзю-дзю, **(норвежцы)** дитя, ребенок, родился, детский, семья, воспитаный, повар, возраст, отец, воспитание, норвежский, родителский, родители, мальчик, взрослый, отец, сын, **(американцы)** айба, колдун, искусство, танец, предание, угл, издурю, божина, фидель, глаза, латинский, виртуальный, лидер, болгарская, призраческий, зелье, лидер, **(китайцы)** китайский, россия, производство, китай, продукция, страна, предприятие, компания, технология, военный, регион, производство, производственный, организованность, российский, экономика, кр, **(азербайджанцы)** русский, азербайджан, азербайджанец, россия, азербайджанский, тикет, дислока, анала, жарод, москва, страна, землянич, слово, рынок, **(германы)** германский, спецназ, военный, август, батальон, российский, специальность, мультимедиа, операция, ручны, братство, мультимедийский, абстракт, группа, война, русский, цинвале, **(осетины)** конституция, осетия, азиат, русский, осетинский, цинвал, осетинский, регион, майя, республика, мирот, азиат, российский, кр-ж-ж-ж, конфликт, **(бразильцы)** наркотики, азиат, шателю, ларинский, 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

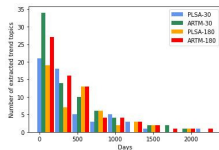
Задачи «поиска и классификации иголок в стоге сена»

- 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

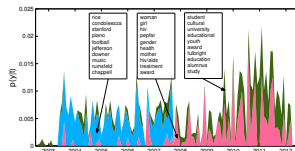
Н.Герасименко, А.Чернявский, М.Никифорова, М.Никитин, К.Воронцов.
 Инкрементальное обучение тематических моделей для поиска трендовых тем
 в научных публикациях. Доклады РАН, 2022.

Выявление динамики тем в новостных потоках

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

Регуляризаторы:

$$\mathcal{L} \left(\begin{array}{|c|} \hline \text{PLSA} \\ \hline \Phi \quad \Theta \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{interpretable} \\ \hline \text{[diagram: vertical bars and dots]} \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{temporal} \\ \hline \text{[diagram: wavy lines]} \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{multimodal} \\ \hline \text{[diagram: stacked boxes]} \\ \hline \end{array} \right) \\ + R \left(\begin{array}{|c|} \hline \text{n-gram} \\ \hline \text{[diagram: grid of boxes]} \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{multilanguage} \\ \hline \text{[diagram: stacked boxes]} \\ \hline \end{array} \right) \rightarrow \max$$



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
All	0.77	0.97	0.86

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:
 - факты как триплеты «субъект–предикат–объект»
 - семантические роли слов по Филлмору
 - тональности именованных существностей

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

ТМ в исторических исследованиях: газетные архивы

- [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.

ТМ в исторических исследованиях: летописи и дневники

- [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>.

ТМ в исторических исследованиях: журнальная периодика

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.

ТМ в политологии: 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);
 - 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.

ТМ в политологии: 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.

Эволюция подходов машинного обучения в анализе текстов

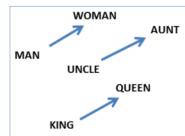
Декомпозиция задач по уровням пирамиды NLP

- 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], GPT-4 [2023]

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

Пример 1. Конкурс ПРО//ЧТЕНИЕ

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.



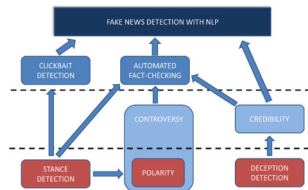
ФАКТИЧЕСКАЯ ОШИБКА
автор высказывания А.Франц

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

Пример 2. 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

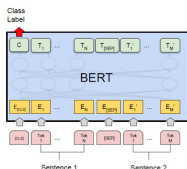


E.Saquete et al.
Fighting post-truth using
natural language processing:
a review and open
challenges // Expert Systems
With Applications, Elsevier,
2020.

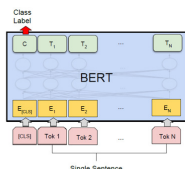
Унификация моделей разметки

Большие пред-обученные модели языка (трансформеры)

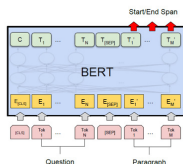
- обучены по терабайтам текстов, «они видели в языке всё»
- способны выделять и классифицировать фрагменты текста
- способны генерировать связный текст
- *мультиязычны*: обучаются на десятках языков
- *мультизадачны*: для каждой новой задачи NLP/NLU достаточно пред-обученной модели + дообучения на относительно небольшой размеченной выборке



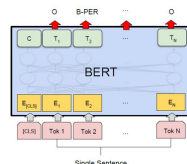
(a) Sentence Pair Classification Tasks:
MNL1, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA



(c) Question Answering Tasks:
SQuAD v1.1



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Унификация оценивания моделей разметки

- В основе методики — сравнение пар разметок текста: «алгоритм – эксперт», «эксперт-1 – эксперт-2», путём оптимального сопоставления их элементов
- Вводится мера согласованности пары разметок (A, B) ,
- если она измеряется несколькими критериями, то берётся 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 $Con(E1, E2)$ разметок двух экспертов, $E1$ и $E2$
- $OTAP = СТАР / СТЭР$,
если больше 100%, то алгоритм не хуже экспертов

Резюме

- *Разведочный поиск* (Exploratory Search)
— поиск по смыслу, а не по ключевым словам, реализуется через модели векторизации текста
- *Тематический поиск* — тоже векторный поиск, векторы интерпретируются как распределения вероятностей тем, каждая тема умеет рассказать о себе
- *Лингвистический поиск* — поиск в текстах фрагментов, имеющих определённый смысл, задаваемый выборкой текстов, размеченных экспертами
- Профессиональная экспертная разметка текстов
— магистральный путь формализации гуманитарных знаний для автоматизации аналитических задач