

# Тематическое моделирование (семинар)

К. В. Воронцов  
vokov@forecsys.ru

Этот курс доступен на странице вики-ресурса  
<http://www.MachineLearning.ru/wiki>  
«Машинное обучение (курс лекций, К.В.Воронцов)»

ШАД Яндекс • 6 октября 2020

- 1 Библиотека тематического моделирования BigARTM**
  - Библиотека тематического моделирования BigARTM
  - Универсальность и гибкость
  - Скорость и масштабируемость
- 2 Эксперименты с тематическим поиском**
  - Методика измерения качества поиска
  - Тематическая модель для документного поиска
  - Оптимизация гиперпараметров
- 3 Другие эксперименты с тематическими моделями**
  - Поиск этно-релевантных тем в социальных сетях
  - Анализ банковских транзакций
  - Визуализация тематических моделей

# 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>

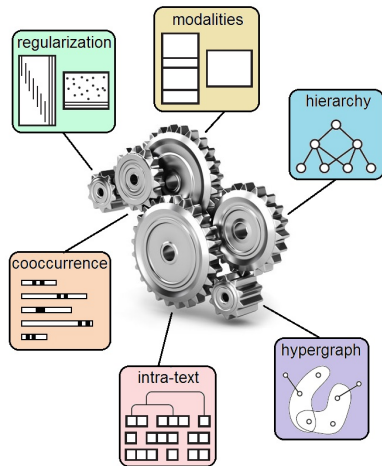


## Лицензия и среда разработки:

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

# Шесть ключевых механизмов BigARTM

- 1 регуляризация
- 2 модальности
- 3 иерархия тем
- 4 сочетаемость термов
- 5 внутритекстовые связи
- 6 гиперграфовые данные



# BigARTM упрощает разработку тематических моделей

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

## Этапы моделирования

### Bayesian TM

### ARTM

	Bayesian TM	ARTM
	Анализ требований	Анализ требований
<i>Формализация:</i>	Вероятностная порождающая модель данных	Стандартные критерии   Свои критерии
<i>Алгоритмизация:</i>	Байесовский вывод для данной порождающей модели (VI, GS, EP)	Общий регуляризованный EM-алгоритм для любых моделей
<i>Реализация:</i>	Исследовательский код (Matlab, Python, R)	Промышленный код BigARTM (C++, Python API)
<i>Оценивание:</i>	Исследовательские метрики, исследовательский код	Стандартные метрики   Свои метрики
	Внедрение	Внедрение

  -- нестандартизируемые этапы, уникальная разработка для каждой задачи

  -- стандартизуемые этапы

## 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)		<b>8m</b> (5648)	<b>5m</b> (6220)
8	100	107m (5380)		<b>10m</b> (4660)	<b>6m</b> (5119)
8	200	288m (4263)		<b>15m</b> (3929)	<b>10m</b> (4309)

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

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

## Две коллекции новостей про технологии

### Habrahr.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 ассессорам

- 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 (библиотека) написанная распределенно: выделены для больших объемов данных и разноразмерных шардов, представляющих собой набор Lucene-кластеров и исполняемых узлах для создания и обработки данных на параллельно обработке.

Основные компоненты Поиск MapReduce можно сгруппировать так:

- обработка вычислений больших объемов данных;
- масштабируемость;
- автоматическая распределенная нагрузка;
- работа на неоднородном оборудовании;
- автоматическая обработка отказов вычислений узлов.

Поиск – популярная программная платформа (облачные сервисы) построена распределенных приложений для высоко-параллельной обработки (разные работы, ресурсы, CPU) данных.

Поиск включает в себе следующие компоненты:

1. HDFS – распределенная файловая система;

2. Поиск MapReduce – программная модель (библиотека) написанная распределенно: выделены для больших объемов данных и разноразмерных шардов;

Ключевые особенности в архитектуре Поиск MapReduce и структуру HDFS, стали привычной реляционной базе данных, в том числе и основные точки отказа. Это, в конечном итоге, определило ограничение платформ Поиск и целом. К сожалению можно отметить:

Ограничение масштабируемости кластера Поиск – не вычислительных узлов, – не К масштабируемых узлов;

Сильная зависимость от распределенно вычислений и элементов вычисления распределенно распределенной нагрузки. Как следствие:

Отсутствие поддержки альтернативной программной модели написанной распределенно: вычислений в Поиск v1.0 поддерживается только модель написанной шардов;

Многие вычисления, точки отказа и как следствие, необходимость масштабирования в средстве вычисления требовались в масштабе;

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

Пример запроса для разведочного поиска



## Векторный поиск тематически близких документов

$\theta_{tq} = p(t|q)$  — тематический вектор запроса  $q$

$\theta_{td} = p(t|d)$  — тематические векторы документов  $d \in D$

Косинусная мера близости документа  $d$  и запроса  $q$ :

$$\text{sim}(q, d) = \frac{\sum_t \theta_{tq} \theta_{td}}{(\sum_t \theta_{tq}^2)^{1/2} (\sum_t \theta_{td}^2)^{1/2}}.$$

Ранжируем документы коллекции  $d \in D$  по убыванию  $\text{sim}(q, d)$

Выдача тематического поиска —  $k$  первых документов.

Реализация: *векторный индекс* для быстрого поиска документов  $d$  по каждой из тем  $t$  запроса

---

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

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

## 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 A4):

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

## Оценивание качества поиска

Precision — доля релевантных среди найденных

Recall — доля найденных среди релевантных

$$P = \frac{TP}{TP + FP} \text{ — точность (precision)}$$

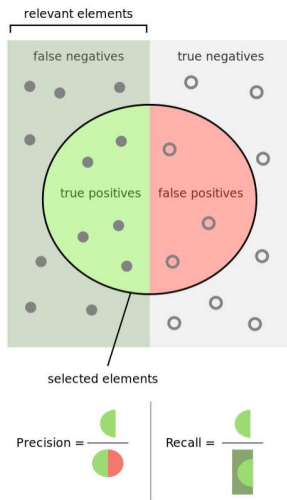
$$R = \frac{TP}{TP + FN} \text{ — полнота, (recall)}$$

$$F_1 = \frac{P + R}{2PR} \text{ — F1-мера}$$

TP (true positive) — найденные релевантные

FP (false positive) — найденные нерелевантные

FN (false negative) — не найденные релевантные



## Какие модели поиска сравнивались

- **assessors**: результаты поиска, выполненного ассессорами
- **TF-IDF, BM25**: сравнение документов по частотам слов
- **word2vec**: нетематические векторные представления слов
- **PLSA**: Probabilistic Latent Semantic Analysis (1999)
- **LDA**: Latent Dirichlet Allocation (2003)
- **ARTM**: тематическая модель с тремя регуляризаторами
- **hARTM**: двухуровневая иерархическая модель ARTM

Задачи регуляризаторов в ARTM и hARTM:

- сделать темы как можно более различными
- сделать векторы  $p(t|d)$  как можно более разреженными
- не допустить вырожденности распределений  $p(w|t)$

## Стратегия регуляризации

Последовательное применение трёх регуляризаторов

- 1 декоррелирование тем:

$$R(\Phi) = -\tau \sum_{s,t \in T} \sum_{w \in W} \phi_{wt} \phi_{ws}$$

- 2 разреживание распределений  $p(t|d)$ :

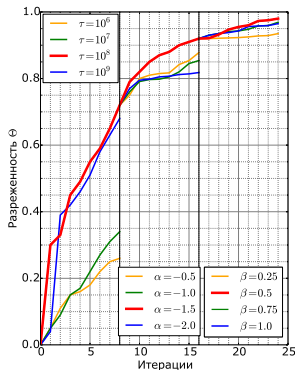
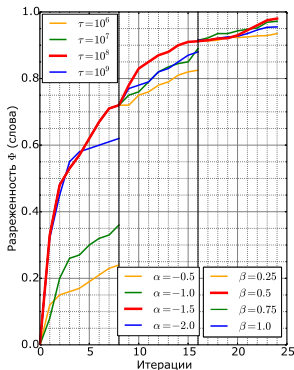
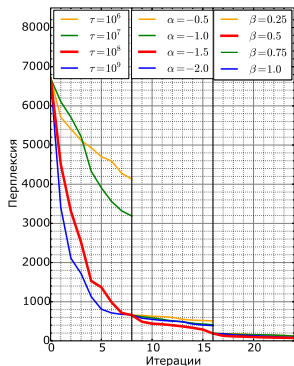
$$R(\Theta) = -\alpha \sum_{d,t} \ln \theta_{td}$$

- 3 сглаживание распределений  $p(w|t)$ :

$$R(\Phi) = \beta \sum_{t,w} \ln \phi_{wt}$$

## Последовательный подбор коэффициентов регуляризации

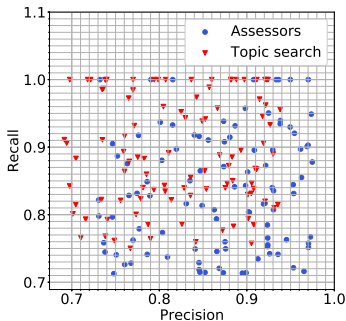
- декоррелирование распределений термов в темах ( $\tau$ ),
- разреживание распределений тем в документах ( $\alpha$ ),
- сглаживание распределений термов в темах ( $\beta$ ).



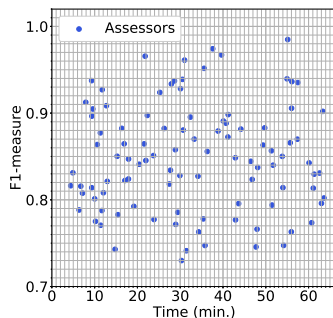
## Результаты измерения точности и полноты по запросам

100 запросов, 3 ассессора на запрос

точность и полнота поиска



время и  $F_1$ -мера (ассессоры)

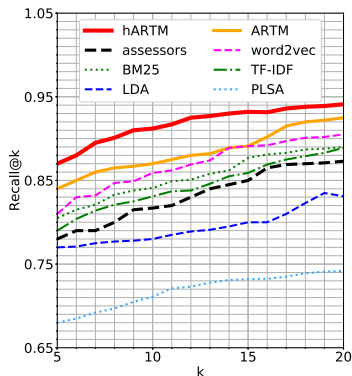
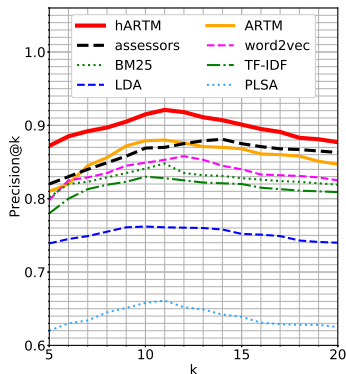


- среднее время обработки запроса ассессором — 30 минут
- точность выше у ассессоров, полнота — у поисковика



## Сравнение с ассессорами по качеству поиска

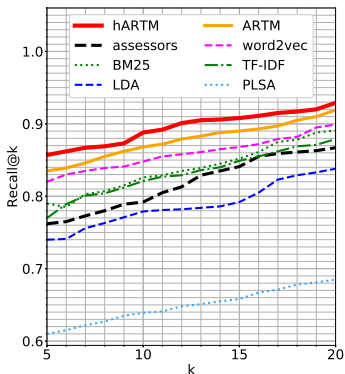
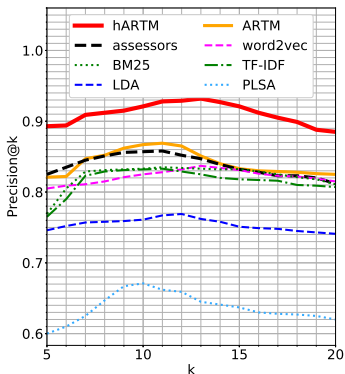
Точность и полнота по первым  $k$  позициям поисковой выдачи (коллекция Habrahabr.ru)



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

## Сравнение с ассессорами по качеству поиска

Точность и полнота по первым  $k$  позициям поисковой выдачи (коллекция TechCrunch.com)



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

## Влияние числа тем на качество поиска

Все регуляризаторы и модальности, **плоская модель**

	Habrahabr						TechCrunch					
	асесс	100	150	<b>200</b>	250	400	асесс	350	400	450	<b>475</b>	500
Pr@5	0.821	0.662	0.721	<b>0.810</b>	0.761	0.693	0.822	0.653	0.725	0.752	<b>0.819</b>	0.777
Pr@10	0.869	0.761	0.812	<b>0.879</b>	0.825	0.673	0.851	0.663	0.732	0.762	<b>0.867</b>	0.811
Pr@15	0.875	0.733	0.795	<b>0.868</b>	0.791	0.651	0.835	0.682	0.743	0.787	<b>0.833</b>	0.793
Pr@20	0.863	0.724	0.795	<b>0.847</b>	0.792	0.642	0.813	0.650	0.743	0.773	<b>0.825</b>	0.793
R@5	0.780	0.732	0.807	<b>0.840</b>	0.821	0.721	0.762	0.731	0.762	0.793	<b>0.835</b>	0.817
R@10	0.817	0.771	0.843	<b>0.870</b>	0.851	0.751	0.792	0.763	0.793	0.812	<b>0.868</b>	0.855
R@15	0.850	0.824	<b>0.895</b>	0.891	0.871	0.773	0.835	0.782	0.807	0.855	<b>0.890</b>	0.882
R@20	0.873	0.857	0.905	<b>0.925</b>	0.892	0.771	0.867	0.792	0.823	0.862	<b>0.919</b>	0.903

- существует оптимальное число тем
- чем больше коллекция, тем больше оптимум числа тем

## Влияние числа тем на качество поиска

Nabrahabr. Все регуляризаторы и модальности, **два уровня**

$ T_1 $	20		25			30	
$ T_2 $	150	200	250	275	300	400	450
Pr@5	0.621	0.742	0.839 0.850	0.865 <b>0.869</b> <b>0.869</b>	0.803 0.769	0.701 0.670	
Pr@10	0.645	0.749	0.850 0.861	0.879 <b>0.911</b> 0.895	0.809 0.796	0.719 0.689	
Pr@15	0.635	0.751	0.848 0.869	0.873 <b>0.893</b> 0.887	0.807 0.781	0.721 0.701	
Pr@20	0.630	0.745	0.841 0.855	0.864 0.874 <b>0.875</b>	0.800 0.775	0.709 0.675	
R@5	0.628	0.773	0.843 0.865	0.881 <b>0.881</b> 0.868	0.849 0.839	0.715 0.691	
R@10	0.652	0.782	0.855 0.871	0.902 <b>0.918</b> 0.877	0.871 0.845	0.745 0.699	
R@15	0.671	0.801	0.870 0.889	0.929 <b>0.939</b> 0.901	0.883 0.861	0.781 0.722	
R@20	0.680	0.819	0.886 0.892	<b>0.955</b> <b>0.955</b> 0.907	0.901 0.872	0.801 0.729	

- существует оптимальное число тем на каждом уровне
- два уровня лучше, чем один
- увеличивается оптимальное число тем на нижнем уровне

## Влияние числа тем на качество поиска

Nabrahabr. Все регуляризаторы и модальности, три уровня

$ T_1 $	20		25			30					
$ T_2 $	150	200	250		275		300		400	450	
$ T_3 $	750	800	1200	1300	1300	<b>1400</b>	1500	1500	1600	3000	3500
Pr@5	0.625	0.743	0.840	0.852	0.869	<b>0.872</b>	0.870	0.805	0.771	0.705	0.672
Pr@10	0.648	0.754	0.851	0.867	0.882	<b>0.915</b>	0.901	0.811	0.799	0.722	0.694
Pr@15	0.632	0.752	0.850	0.872	0.878	<b>0.895</b>	0.889	0.809	0.785	0.729	0.703
Pr@20	0.629	0.745	0.845	0.861	0.871	0.877	<b>0.882</b>	0.803	0.778	0.710	0.681
R@5	0.632	0.780	0.845	0.869	0.883	<b>0.889</b>	0.872	0.851	0.841	0.721	0.695
R@10	0.654	0.792	0.859	0.873	0.905	<b>0.922</b>	0.881	0.873	0.850	0.749	0.703
R@15	0.675	0.805	0.874	0.892	0.932	<b>0.942</b>	0.905	0.889	0.863	0.787	0.725
R@20	0.684	0.824	0.889	0.901	0.958	<b>0.961</b>	0.912	0.904	0.878	0.805	0.734

- существует оптимальное число тем на каждом уровне
- три уровня лучше, чем один или два
- увеличивается оптимальное число тем на нижнем уровне

## Влияние числа тем на качество поиска

TechCrunch. Все регуляризаторы и модальности, **два уровня**

$ T_1 $	80		100			120					
$ T_2 $	300	350	500	550	600	700	750				
Pr@5	0.651	0.701	0.749	0.789	0.883	<b>0.889</b>	<b>0.889</b>	0.785	0.721	0.701	0.675
Pr@10	0.675	0.709	0.771	0.821	0.891	<b>0.918</b>	0.902	0.803	0.738	0.718	0.691
Pr@15	0.687	0.712	0.773	0.827	0.899	<b>0.919</b>	0.905	0.817	0.741	0.721	0.701
Pr@20	0.683	0.707	0.759	0.815	0.885	0.888	<b>0.895</b>	0.805	0.732	0.716	0.679
R@5	0.749	0.791	0.801	0.854	0.868	<b>0.875</b>	0.861	0.849	0.829	0.731	0.701
R@10	0.765	0.809	0.823	0.873	0.890	<b>0.904</b>	0.875	0.867	0.835	0.745	0.708
R@15	0.771	0.820	0.841	0.882	0.909	<b>0.921</b>	0.895	0.890	0.848	0.769	0.717
R@20	0.778	0.825	0.851	0.887	0.928	<b>0.942</b>	0.929	0.901	0.869	0.785	0.728

- существует оптимальное число тем на каждом уровне
- два уровня лучше, чем один
- увеличивается оптимальное число тем на нижнем уровне

## Влияние числа тем на качество поиска

TechCrunch. Все регуляризаторы и модальности, **три уровня**

$ T_1 $	80		100						120		
$ T_2 $	300	350	500		550		600		700	750	
$ T_3 $	1500	1700	2500	2600	2600	<b>2800</b>	3000	3000	3200	4500	4700
Pr@5	0.655	0.707	0.751	0.792	0.887	<b>0.893</b>	0.890	0.789	0.722	0.703	0.678
Pr@10	0.678	0.712	0.773	0.823	0.895	<b>0.922</b>	0.905	0.805	0.741	0.722	0.692
Pr@15	0.692	0.715	0.775	0.831	0.902	<b>0.921</b>	0.907	0.821	0.743	0.725	0.703
Pr@20	0.687	0.709	0.761	0.819	0.889	0.885	<b>0.898</b>	0.809	0.736	0.719	0.683
R@5	0.751	0.795	0.802	0.856	0.871	<b>0.877</b>	0.863	0.852	0.831	0.738	0.705
R@10	0.767	0.812	0.825	0.875	0.892	<b>0.908</b>	0.879	0.871	0.842	0.751	0.711
R@15	0.772	0.824	0.841	0.887	0.912	<b>0.927</b>	0.901	0.893	0.854	0.772	0.721
R@20	0.783	0.830	0.854	0.892	0.931	<b>0.949</b>	0.935	0.905	0.871	0.790	0.732

- существует оптимальное число тем на каждом уровне
- три уровня лучше, чем один или два
- увеличивается оптимальное число тем на нижнем уровне

## Влияние модальностей на качество поиска

Все регуляризаторы и модальности, 3 уровня, оптимальное  $|T|$

Модальности: Words, Bigrams, Authors, Comments, Tags, Hubs, Categories

	Habrahabr						TechCrunch					
	асесс	W	Com	WB	WBTH	All	асесс	W	C	WB	WBC	All
Pr@5	0.821	0.621	0.558	0.673	0.871	<b>0.872</b>	0.822	0.718	0.569	0.795	0.891	<b>0.893</b>
Pr@10	0.869	0.645	0.567	0.712	0.911	<b>0.915</b>	0.851	0.729	0.592	0.807	0.919	<b>0.922</b>
Pr@15	0.875	0.631	0.532	0.693	0.894	<b>0.895</b>	0.835	0.737	0.603	0.803	0.920	<b>0.921</b>
Pr@20	0.863	0.628	0.531	0.688	0.877	<b>0.877</b>	0.813	0.729	0.594	0.792	0.883	<b>0.885</b>
R@5	0.780	0.725	0.645	0.797	0.888	<b>0.889</b>	0.762	0.754	0.659	0.775	0.874	<b>0.877</b>
R@10	0.817	0.748	0.652	0.812	0.921	<b>0.922</b>	0.792	0.778	0.671	0.808	0.908	<b>0.908</b>
R@15	0.850	0.782	0.679	0.842	0.941	<b>0.942</b>	0.835	0.783	0.679	0.825	0.927	<b>0.927</b>
R@20	0.873	0.789	0.672	0.852	0.960	<b>0.961</b>	0.867	0.785	0.711	0.837	0.949	<b>0.949</b>

- лучше использовать все модальности
- биграммы и категории выигрывают у ассессоров
- авторы и комментаторы наименее важны



## Влияние регуляризаторов на качество поиска

Все регуляризаторы и модальности, 3 уровня, оптимальное |T|

Регуляризаторы: Decorrelation, Θ-sparsing, Φ-smoothing, Hierarchy

	Habrahabr					TechCrunch				
	нет	D	DΘ	DΘΦ	DΘΦH	нет	D	DΘ	DΘΦ	DΘΦH
Pr@5	0.628	0.772	0.771	0.865	<b>0.872</b>	0.652	0.777	0.779	0.879	<b>0.893</b>
Pr@10	0.653	0.781	0.812	0.883	<b>0.915</b>	0.679	0.788	0.819	0.895	<b>0.922</b>
Pr@15	0.642	0.785	0.792	0.891	<b>0.895</b>	0.669	0.791	0.798	0.901	<b>0.921</b>
Pr@20	0.643	0.771	0.783	0.875	<b>0.877</b>	0.673	0.775	0.792	<b>0.892</b>	0.885
R@5	0.692	0.820	0.805	0.875	<b>0.889</b>	0.673	0.825	0.812	0.869	<b>0.877</b>
R@10	0.714	0.831	0.834	0.905	<b>0.922</b>	0.685	0.856	0.845	0.881	<b>0.908</b>
R@15	0.725	0.847	0.867	0.921	<b>0.942</b>	0.712	0.877	0.869	0.912	<b>0.927</b>
R@20	0.735	0.873	0.891	0.943	<b>0.961</b>	0.723	0.892	0.895	0.934	<b>0.949</b>

- Лучше использовать все регуляризаторы
- Модели со слабой регуляризацией (PLSA, LDA) слабы

## Влияние функции близости на качество поиска

Все регуляризаторы и модальности, 3 уровня, оптимальное |T|

Функции близости: Euclidean, Cosine, Manhattan, Hellinger, KL-div

	Habrahabr					TechCrunch				
	Eu	cos	Ma	He	KL	Eu	cos	Ma	He	KL
Pr@5	0.652	<b>0.872</b>	0.772	0.725	0.741	0.647	<b>0.893</b>	0.752	0.742	0.735
Pr@10	0.693	<b>0.915</b>	0.798	0.749	0.772	0.658	<b>0.922</b>	0.794	0.758	0.751
Pr@15	0.695	<b>0.895</b>	0.803	0.737	0.751	0.672	<b>0.921</b>	0.801	0.745	0.742
Pr@20	0.671	<b>0.877</b>	0.789	0.731	0.738	0.652	<b>0.885</b>	0.793	0.739	0.738
R@5	0.693	<b>0.889</b>	0.721	0.742	0.833	0.688	<b>0.877</b>	0.708	0.733	0.858
R@10	0.715	<b>0.922</b>	0.732	0.775	0.868	0.692	<b>0.908</b>	0.715	0.753	0.872
R@15	0.732	<b>0.942</b>	0.739	0.791	0.892	0.724	<b>0.927</b>	0.719	0.785	0.895
R@20	0.741	<b>0.961</b>	0.721	0.812	0.902	0.732	<b>0.949</b>	0.711	0.808	0.901

- косинусная функция близости уверенно лидирует

## Выводы по результатам экспериментов

- Ассессорские данные относятся не к темам, а к коллекции; поэтому с их помощью можно оценивать новые модели
- Небольших ассессорских данных хватает для оценивания тематических моделей, т. к. они обучаются *без учителя*
- Регуляризаторы, улучшающие интерпретируемость модели, повышают также и качество поиска
- Иерархия улучшает качество поиска (в основном точность) благодаря постепенному сужению области поиска
- Подбор траектории регуляризации и оптимизация коэффициентов регуляризации влияет на качество поиска
- При тщательной оптимизации тематический поиск превосходит как ассессоров, так и конкурирующие модели

---

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

## Поиск этно-релевантных тем в социальных сетях

### Основные задачи проекта:

- Разведочный поиск этнических тем в социальных медиа (сколько различных тем, и что это за темы)
- Мониторинг этих тем во времени и по регионам
- Сентимент-анализ и оценивание конфликтности

### Вспомогательные задачи:

- Фильтрация (обогащение) потока данных
- Обеспечение полноты поиска этнических тем
- Выявление тематических сообществ
- Выделение событийных и региональных тем
- Решение проблемы коротких сообщений

## Примеры этнонимов

османский	руси ч
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

## Примеры этнических 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,

## Примеры этнических 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,

## Примеры этнических тем

**(японцы)**: японский, япония, корея, китайский, жилища, авария, фукусиму, цунами, общаться, 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,



## Результаты: модель ARTM находит намного больше этно-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	38	42	30	104

Регуляризаторы ARTM-1:

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

**фоновые темы:** сглаживание, разреживание этнонимов

Регуляризаторы ARTM-2:

ARTM-1 + **модальность этнонимов**

*Apishev M., Koltcov S., Koltsova O., Nikolenko S., Vorontsov K. Additive Regularization for Topic Modeling in Sociological Studies of User-Generated Texts. 2016.*

## Анализ транзакций розничных клиентов банка

**Дано** (Sberbank Data Science Contest):

$D$  — множество клиентов (15 000)

$W$  — категории = MCC-коды (Merchant Category Code) (328)

$n_{dw}$  — сумма транзакций клиента  $d$  по категории  $w$

**Найти:** темы — типы потребительского поведения клиентов

$\phi_{wt} = p(w|t)$  — структура потребления для темы  $t$

$\theta_{td} = p(t|d)$  — типы потребления клиента  $d$

**Регуляризаторы:**

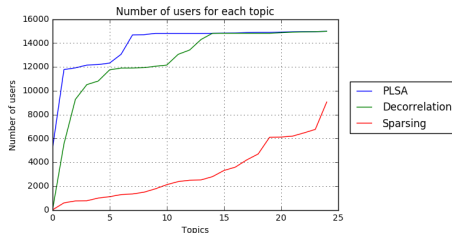
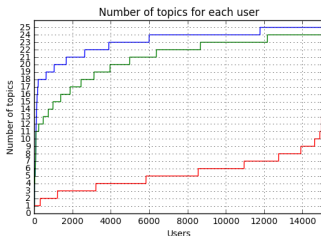
- повышение различности тем
- разреживание  $p(t|d)$
- учёт модальностей времени, типа транзакции, терминала

---

*Egorov E., Nikitin F., Goncharov A., Alekseev V., Vorontsov K. Topic Modelling for Extracting Behavioral Patterns from Transactions Data // IC-AIAI 2019.*

## Построение модели ARTM, 25 тем

- 30 итераций PLSA — без регуляризаторов
- 10 итераций — декоррелирование тем
- 10 итераций — разреживание  $p(t|d)$



Декоррелирование  $\Phi$  и разреживание  $\Theta$  определяют минимальное число типов экономического поведения каждого клиента, достаточное для описания его расходов.

## Пользуюсь картой только чтобы снять наличные

$\phi_{wt},\%$  МСС-код (категория расходов)

72 Финансовые институты — снятие наличности вручную

27 Финансовые институты — снятие наличности автоматически

0.23 Денежные переводы MasterCard MoneySend

0.1 Денежные переводы

0.012 Финансовые институты — снятие наличности вручную

0.0055 Легковой и грузовой транспорт: продажа, сервис, ремонт, лизинг

0.0027 Магазины игрушек

## Наличные + авто, спорт, компьютеры

$\phi_{wt},\%$  МСС-код (категория расходов)

- 55 Финансовые институты — снятие наличности автоматически
- 44 Денежные переводы
- 0.111 Станции техобслуживания
- 0.105 Автозапчасти и аксессуары
- 0.094 Компьютерная сеть/информационные услуги
- 0.043 Спортивная одежда, одежда для верховой езды и езды на мотоцикле
- 0.024 Финансовые институты — снятие наличности вручную
- 0.020 СТО общего назначения
- 0.018 Горючее топливо — уголь, нефть, разжиженный бензин, дрова
- 0.015 Магазины мужской и женской одежды
- 0.015 Финансовые институты — снятие наличности вручную
- 0.013 Магазины спорттоваров
- 0.012 Садовые принадлежности (в том числе для ухода за газонами) в розницу
- 0.011 Паркинги и гаражи
- 0.011 Бакалейные магазины, супермаркеты
- 0.010 Различные магазины одежды и аксессуаров

## Цивилизованный потребитель: разные магазины, связь, авто

- $\phi_{wt},\%$  МСС-код (категория расходов)
- 27 Станции техобслуживания
  - 20 Различные продовольственные магазины, рынки, полуфабрикаты
  - 15 Звонки с использованием телефонов, считывающих магнитную ленту
  - 12 Финансовые институты — снятие наличности автоматически
  - 4.7 Горючее топливо — уголь, нефть, разжиженный бензин, дрова
  - 4.1 Универсальные магазины
  - 3.4 Автозапчасти и аксессуары
  - 1.4 Аптеки
  - 1.2 Магазины с продажей спиртных напитков на вынос
  - 1.1 Бакалейные магазины, супермаркеты
  - 0.57 Автошины
  - 0.37 Прямой маркетинг — торговля через каталог
  - 0.35 Товары для дома
  - 0.33 Универмаги
  - 0.32 Плавательные бассейны — распродажа
  - 0.21 Места общественного питания, рестораны

Всего 24 категории с  $\phi_{wt} > 0.1\%$ ; 61 категория с  $\phi_{wt} > 0.01\%$

## Продвинутые мамки

$\phi_{wt}, \%$  МСС-код (категория расходов)

- 56 Бакалейные магазины, супермаркеты
- 8.6 Финансовые институты — снятие наличности автоматически
- 5.4 Аптеки
- 4.0 Звонки с использованием телефонов, считывающих магнитную ленту
- 2.2 Рестораны, закусочные
- 1.8 Обувные магазины
- 1.5 Различные продовольственные магазины — рынки, полуфабрикаты
- 1.4 Магазины спорттоваров
- 1.4 Детская одежда, включая одежду для самых маленьких
- 1.3 Магазины игрушек
- 1.3 Места общественного питания, рестораны
- 1.1 Магазины мужской и женской одежды
- 1.1 Магазины с продажей спиртных напитков на вынос
- 1.1 Магазины косметики
- 1.0 Садовые принадлежности в розницу
- 0.73 Одежда для всей семьи

Всего 41 категория с  $\phi_{wt} > 0.1\%$ ; 95 категорий с  $\phi_{wt} > 0.01\%$

## Бизнес-леди: забыла про наличку — всё по карте

$\phi_{wt}, \%$  МСС-код (категория расходов)

- 12 Магазины мужской и женской одежды
- 7.3 Оборудование, мебель и бытовые принадлежности
- 7.0 Места общественного питания, рестораны
- 5.6 Магазины по продаже часов, ювелирных изделий и изделий из серебра
- 5.3 Обувные магазины
- 4.7 Магазины косметики
- 4.6 Одежда для всей семьи
- 3.8 Универмаги
- 3.2 Готовая женская одежда
- 2.8 Практикующие врачи, медицинские услуги
- 1.8 Прямой маркетинг — торговля через каталог
- 1.5 Салоны красоты и парикмахерские
- 1.3 Детская одежда, включая одежду для самых маленьких
- 1.3 Аптеки
- 1.0 Изготовление и продажа меховых изделий
- 1.0 Центры здоровья

Всего 70 категорий с  $\phi_{wt} > 0.1\%$ ; 134 категории с  $\phi_{wt} > 0.01\%$



## Бизнес-класс: авиа, отели, казино, рестораны, ценные бумаги

- $\phi_{wt}, \%$  МСС-код (категория расходов)
- 28 Авиалинии, авиакомпании
  - 19 Финансовые институты — торговля и услуги
  - 9.5 Отели, мотели, базы отдыха, сервисы бронирования
  - 8.6 Транзакции по азартным играм (плюс)
  - 5.2 Финансовые институты — торговля и услуги
  - 3.2 Места общественного питания, рестораны
  - 3.1 Не-финансовые институты: ин.валюта, переводы, дорожн.чеки, квази-кэш
  - 2.2 Пассажирские железнодорожные перевозки
  - 1.7 Бизнес-сервис
  - 1.4 Жилье — отели, мотели, курорты
  - 1.3 Галереи/учреждения видеоигр
  - 1.3 Транзакции по азартным играм (минус)
  - 0.6 Ценные бумаги: брокеры/дилеры
  - 0.5 Туристические агентства и организаторы экскурсий
  - 0.3 Лимузины и такси
  - 0.3 Беспшлинные магазины Duty Free

Всего 50 категорий с  $\phi_{wt} > 0.1\%$ ; 103 категории с  $\phi_{wt} > 0.01\%$

## Провинциальный малый бизнес

$\phi_{wt}, \%$  МСС-код (категория расходов)

- 27 Финансовые институты — снятие наличности автоматически
- 8.5 Лесо- и строительный материал
- 8.4 Бытовое оборудование
- 6.6 Плавательные бассейны — распродажа
- 5.5 Продажа электронного оборудования
- 4.1 Бакалейные магазины, супермаркеты
- 3.3 Универсальные магазины
- 3.0 Садовые принадлежности в розницу
- 2.6 Телекоммуникационное оборудование, включая продажу телефонов
- 2.4 Легковой и грузовой транспорт: продажа, сервис, ремонт, лизинг
- 2.2 Товары для дома
- 2.1 Пассажирские железнодорожные перевозки
- 1.5 Оборудование, мебель и бытовые принадлежности
- 1.3 Скобяные товары в розницу
- 1.2 Магазины спорттоваров
- 1.1 Аптеки

Всего 54 категории с  $\phi_{wt} > 0.1\%$ ; 104 категории с  $\phi_{wt} > 0.01\%$

## Анализ транзакций корпоративных клиентов банка

### Данные:

лесная отрасль, 2016 г., 10.7М транзакций, 1М компаний.

Транзакция — это тройка ⟨покупатель, продавец, текст⟩.

Некоторые *тексты* платёжных поручений (далеко не все!) содержат названия товаров и услуг.

Документ — это история транзакций одной компании

### Семь модальностей:

- компании: поставщики / покупатели
- слова в платёжных поручениях: поставщики / покупатели
- ОКВЭДы данной компании
- ОКВЭДы контрагентов: поставщики / покупатели

## Примеры тем — видов деятельности компаний

покупка	продажа
0.11: услуга	0.12: лдсп
0.07: классик	0.08: дсп
0.05: дрова	0.03: мдф
0.05: пиловочник	0.03: поставка
0.05: материал	0.02: услуга
0.03: порода	0.02: охранный
0.03: лесоматериал	0.02: ламинировать
0.03: сертум	0.02: хдф
0.02: хвойный	0.02: материал
0.01: дерево	0.01: накл
0.01: транспортный	0.01: товар

покупка	продажа
0.19: право	0.16: арендный
0.17: сбис	0.10: часть
0.16: использование	0.08: плата
0.03: аккаунт	0.04: минимальный
0.02: электронный	0.04: участок
0.02: лицевой	0.04: использование
0.02: устный	0.02: земля
0.01: устройство	0.02: лесничество
0.01: генерация	0.02: земельный
0.01: хранение	0.01: фонд
0.01: ключевой	0.01: федеральный

## Примеры тем — видов деятельности компаний

покупка	продажа
0.09: ткань	0.16: мебель
0.09: поставка	0.05: плёнка
0.02: мебельный	0.04: стул
0.02: деревянный	0.03: кресло
0.02: транспортный	0.03: изделие
0.02: фанера	0.02: краска
0.02: поролон	0.02: фанера
0.01: механизм	0.01: лкм
0.01: плата	0.01: лакокрасочный
0.01: частичный	0.01: лак
	0.01: материал
	0.01: клеить

покупка	продажа
0.06: лдсп	0.37: товар
0.05: фурнитура	0.15: мебель
0.02: плёнка	0.04: поставка
0.02: материал	0.04: накладный
0.02: мебельный	0.03: накл
0.02: стекло	0.03: рубль
0.02: мдф	
0.02: кромка	
0.01: транспортный	
0.01: клеить	
0.01: профиль	
0.01: пвх	

## Примеры тем — видов деятельности компаний

покупка	продажа
0.52: гсм	0.14: вывоз
0.43: далее	0.09: тбо
	0.04: мусор
	0.03: отход
	0.02: утилизация
	0.01: тко

покупка	продажа
0.19: налог	0.11: бумага
0.06: услуга	0.08: гофроящик
0.04: макулатура	0.04: гофрокартон
0.03: поставка	0.03: гофрокороб
0.03: транспортный	0.03: поставка
0.02: лесопродукция	0.03: фактура
0.02: автоуслуга	0.02: гофропродукция
0.01: перевозка	0.02: гофротару
0.01: плата	0.02: гофрирование
	0.02: гофролоток
	0.02: товар
	0.01: лоток

## Примеры тем — видов деятельности компаний

покупка

0.15: программа

0.11: право

0.09: сертификат

0.07: эвм

0.07: использование

0.07: лицензия

0.04: криптопро

0.03: абонентский

0.02: обслужа

0.02: пользование

0.02: контур

0.01: проверка

продажа

0.13: фурнитура

0.09: материал

0.08: лдсп

0.04: кромка

0.04: мебельный

0.04: фрз

0.04: мдф

0.03: клеить

0.03: пвх

0.02: тмц

0.02: комплект

0.02: профиль

0.02: столешница

продажа

0.14: рекламный

0.13: размещение

0.09: материал

0.05: проект

0.05: яндекс

0.04: директ

0.04: реклама

0.02: рубль

0.01: стек

продажа

0.21: тмц

0.06: накл

0.04: инструмент

0.03: пила

0.02: заточка

0.02: нож

0.02: материал

0.02: фреза

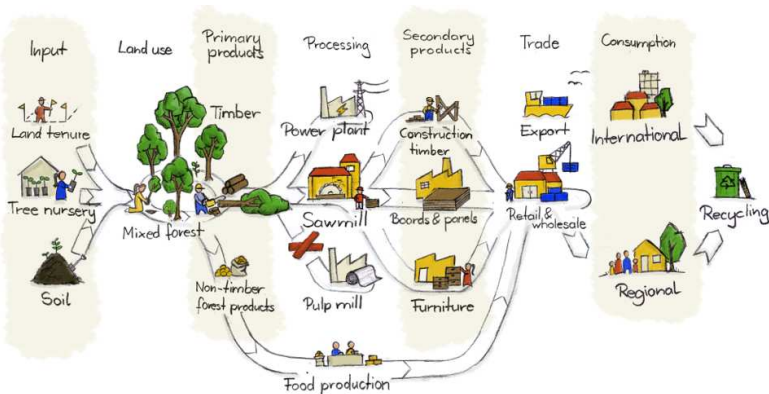
0.02: клеить

0.01: товар

0.01: перчатка

## Цели тематического моделирования банковских данных

- Получение векторных представлений компаний
- Поиск схожих и конкурирующих компаний
- Восстановление структуры товарных потоков отрасли





## Что можно визуализировать

Одна из целей тематического моделирования — систематизация результатов информационного поиска.

- текстовое представление темы: название, топ-слова, топ-термы, топ-документы, аннотация, близкие темы
- масштабируемая тематическая карта коллекции
- иерархия тем
- граф связей между темами
- текст документа: темы слов или термов, сегментация
- графическая тематическая сегментация документа
- динамика тем во времени: временные ряды, реки тем
- иерархия + динамика

# Система TMVE — Topic Model Visualization Engine

Тематический навигатор с веб-интерфейсом:

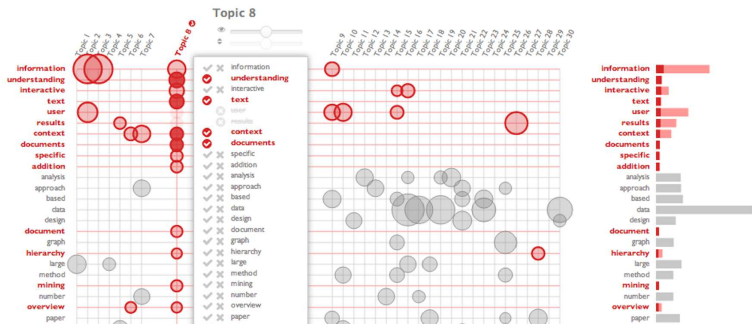


<https://github.com/ajbc/tmv>

Chaney A., Blei D. Visualizing Topic Models // Frontiers of computer science in China, 2012. — 55(4), pp. 77–84.

# Система Termite

Интерактивная визуализация матрицы  $\Phi$  и сравнение тем:

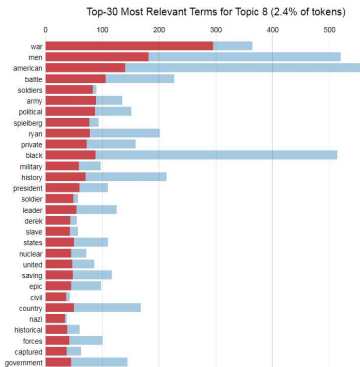
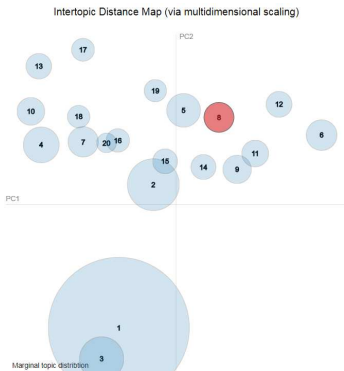


<https://github.com/uwdata/termite-visualizations>

Chuang J., Manning C., Heer J. Termite: Visualization Techniques for Assessing Textual Topic Models. IWCAVI 2012.

## Система LDAvis

Карта сходства тем и сравнение  $p(w|t)$  с  $p(w)$ :



<https://github.com/cpsievert/LDAvis>

C.Sievert, K.Shirley. LDAvis: A method for visualizing and interpreting topics. 2014.

## Система Serendip

Визуализация матриц  $\Phi$ ,  $\Theta$  и тематики слов в текстах:

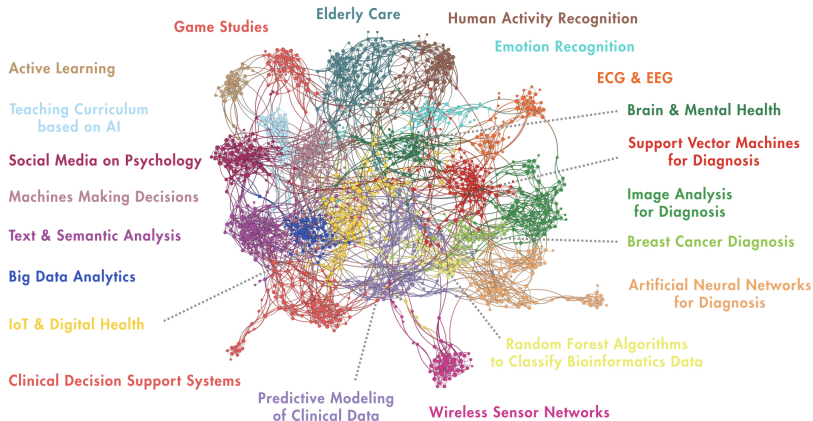


<http://vep.cs.wisc.edu/serendip>

E.Alexander, J.Kohlmann, R.Valenza, M.Witmore, M.Gleicher. Serendip: Topic Model-Driven Visual Exploration of Text Corpora. IEEE VAST 2014.

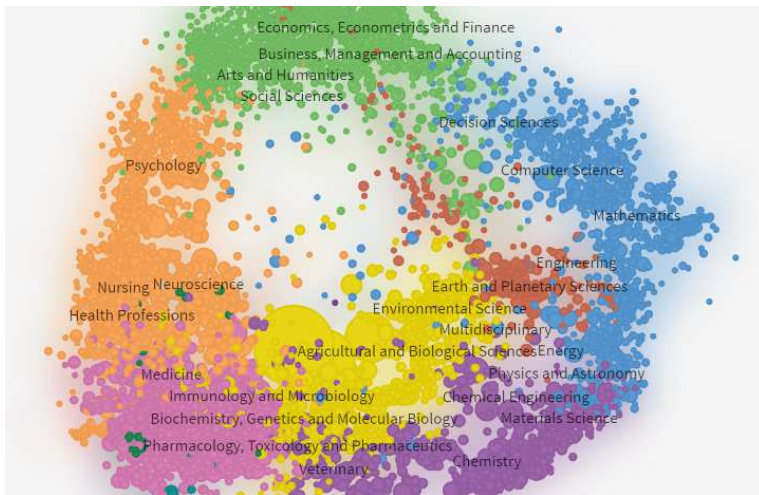
## Пример карты тематической структуры коллекции

Academic papers on AI in Healthcare published in 2016



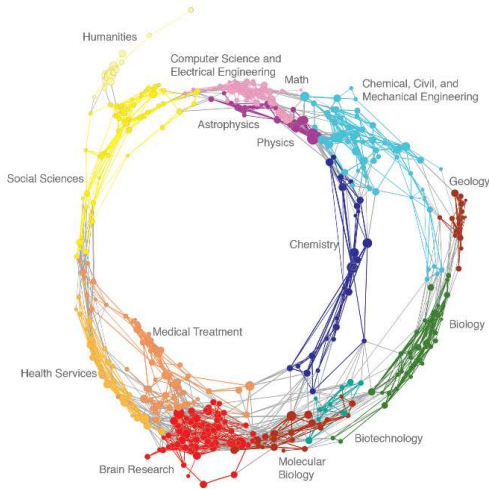
C.Folgar, J.McCuan. The 3 most-cited studies in healthcare and AI. Quid, 2017.

## Пример 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.

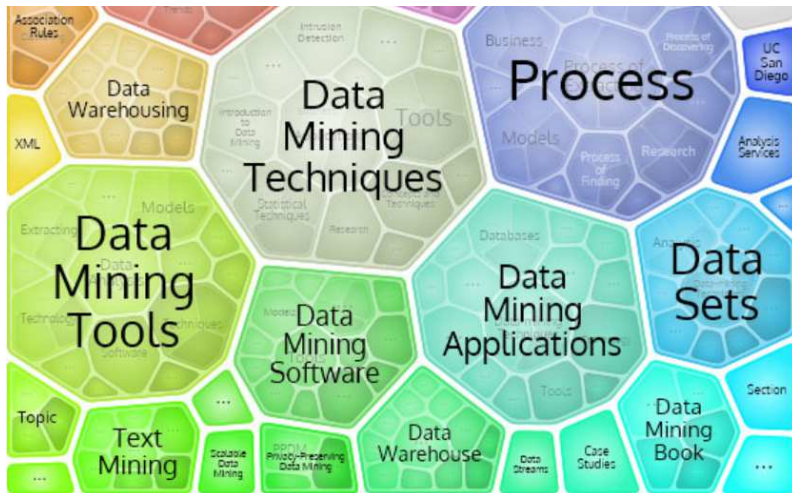
### Недостатки:

- оси не имеют интерпретации
- искажение сходства при двумерном проецировании

<http://scimaps.org>



## Пример иерархической карты области *Data Mining*



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

## Динамика тем: эволюция предметной области



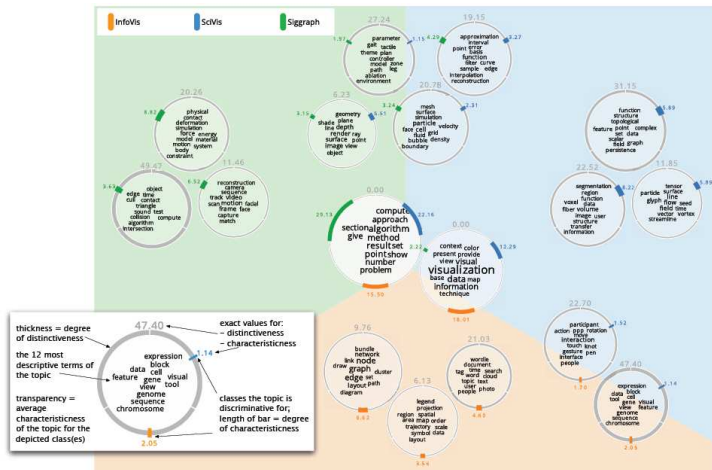
Эволюция выбранных тем иерархии. Данные Prism (2013/06/03–2014/02/09)

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.

## Тематический анализ источников



Oelke D., Stobelt H., Rohrdantz C., Gurevych I., Deussen O. Comparative exploration of document collections: a visual analytics approach. EuroVis. 2014.



<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.*

## Применения тематического моделирования:

- Разведочный информационный поиск
- Анализ тематической структуры коллекции
- Тематизация документов — формирование интерпретируемых тематических векторных представлений документов
- Классификация, суммаризация, сегментация текстов
- Обнаружение и отслеживание новостных цепочек
- Маршрутизация обращений в контактные центры
- Поиск экспертов и тематических сообществ
- ...