

Тематическое моделирование и разведочный информационный поиск

Воронцов Константин Вячеславович

ФИЦ ИУ РАН • МФТИ • МГУ • ШАД Яндекс • НИУ ВШЭ

Научный семинар • ФКН НИУ ВШЭ

31 января 2017

1 Философия

- Разведочный информационный поиск
- Дальнее чтение и визуализация
- Сценарий разведочного поиска

2 Теория

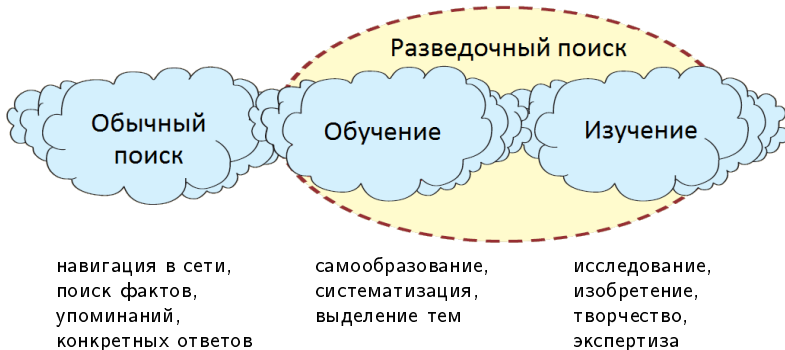
- Вероятностные тематические модели
- Теория аддитивной регуляризации (ARTM)
- Алхимия тематического моделирования

3 Практика

- Библиотека с открытым кодом BigARTM
- Поиск этно-релевантных тем в социальных сетях
- Разведочный поиск в коллективном блоге

Концепция разведочного поиска (exploratory search)

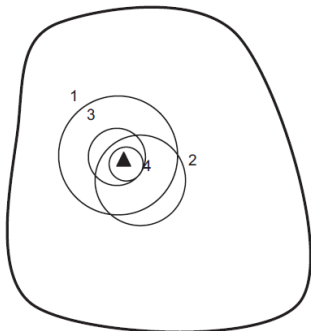
- пользователь может не знать ключевых терминов
- пользователя может интересовать множество ответов



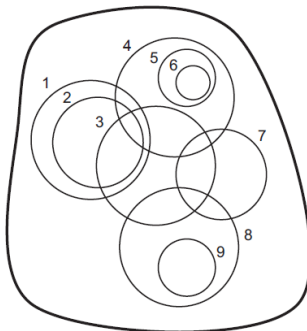
Gary Marchionini. Exploratory Search: from finding to understanding. 2006.

От итераций «query-browse-refine» к разведочному поиску

Iterative Search



Exploratory Search



- ▲ Search target ◊ Information space
○# Result sets (larger = more results, intersection = overlap, # = iteration)

R.W.White, R.A.Roth. Exploratory Search: beyond the Query-Response paradigm. San Rafael, CA: Morgan and Claypool, 2009.

От ближнего чтения (close reading) к дальнему (distant reading)

Концепция дальнего чтения Франко Моретти

Дальнее чтение — это специальная форма представления знаний: меньше элементов, грубее смысл их взаимосвязей, важны лишь общие очертания и структуры.

Мантра Шнейдермана

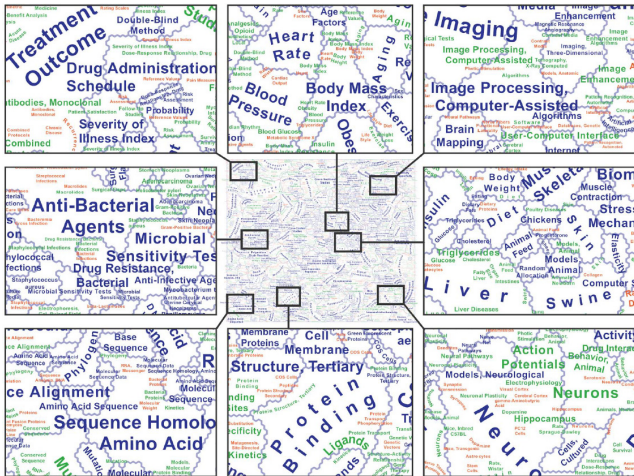
«Сначала крупный план, затем масштабирование и фильтрация, детали по требованию»

B.Shneiderman. The Eyes Have It: A Task by Data Type Taxonomy for Information Visualizations. Visual Languages, 1996.

F.Moretti. Graphs, Maps, Trees: Abstract Models for a Literary History. 2005.

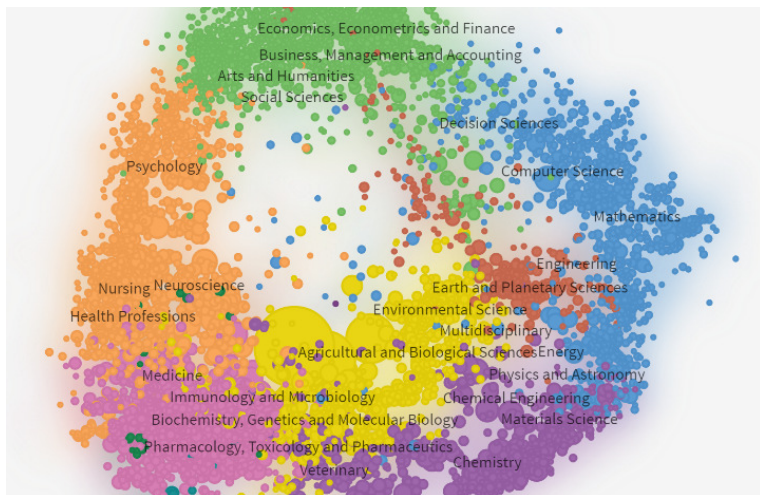
S.Janicke, G.Franzini, M.F.Cheema, G.Scheuermann. On Close and Distant Reading in Digital Humanities: A Survey and Future Challenges. EuroVis, 2015.

Пример карты медицинских знаний



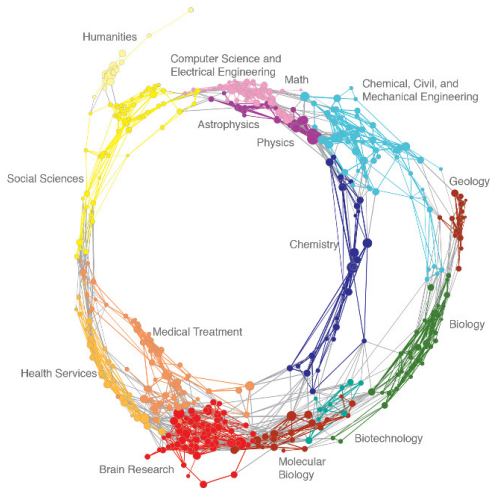
Skupin, Biberstine, Borner. Visualizing the Topical Structure of the Medical Sciences: A Self-Organizing Map Approach. PLoS ONE, 2013.

Пример карты науки



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

Ещё один пример карты науки



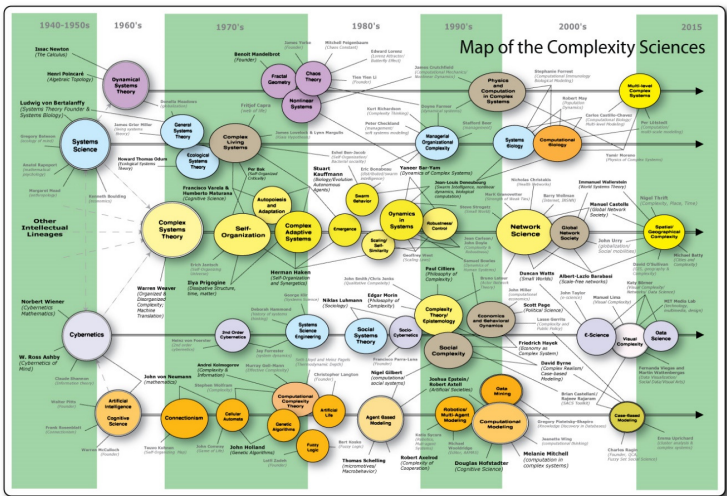
Важное наблюдение:
области знания
самопроизвольно
располагаются по кругу,
значит,
их можно располагать
и вдоль прямой линии.

Недостатки:

- оси не имеют интерпретации
- искажение сходства при двумерном проецировании

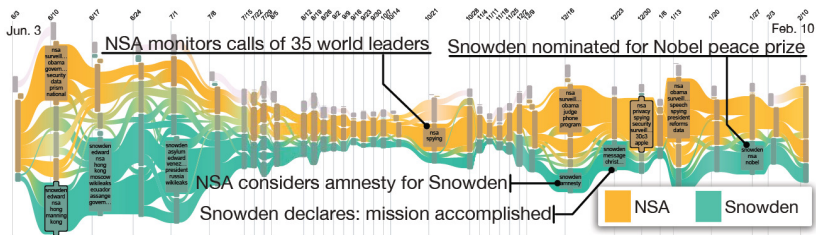
<http://scimaps.org>

Пример карты предметной области, построенной вручную



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

Динамика тем: эволюция предметной области



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

- 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.

Визуализация тематического разведочного поиска (концепт)

- Двумерная карта в интерпретируемых осях время–темы
- Спектр тем: гуманитарные → естественные → точные
- Темы делятся на подтемы иерархически
- Интерактивность: реализация мантры Шнейдермана
- При любом масштабе на карте достаточно много текста



<http://textvis.lnu.se>

Интерактивный обзор 365 средств визуализации текстов



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

Возможный сценарий разведочного поиска

Поисковый запрос:

- документ любой длины или даже коллекция документов

Цели поиска:

- к каким темам относится мой запрос?
- что ещё известно по этим темам?
- какова тематическая структура этой предметной области?
- какие области являются смежными?
- что ещё есть понятного, обзорного, важного, свежего?

Сценарий поиска:

- 1 имея любой текст под рукой, в любом приложении,
- 2 получаем картину содержащихся в нём тем-подтем
- 3 и «дорожную карту» предметной области в целом

Технологические элементы разведочного поиска

По всем технологиям имеются готовые решения:

- 1 интернет-краулинг
- 2 фильтрация контента
- 3 тематическое моделирование
- 4 инвертированный индекс
- 5 ранжирование
- 6 визуализация
- 7 персонализация

Наша научная группа развивает теорию и технологии тематического моделирования как ключевой и наиболее наукоёмкий элемент разведочного поиска.

Тематическая модель для разведочного поиска должна быть...

- 1 Интерпретируемая: каждая тема понятна людям
- 2 Мультиграммная: термины-словосочетания неразрывны
- 3 Мультимодальная: авторы, связи, тэги, пользователи, ...
- 4 Мультиязычная: для кросс- и много-языкового поиска
- 5 Динамическая: выявление истории развития тем
- 6 Иерархическая: выявление иерархических связей тем
- 7 Сегментирующая: выделение тем внутри документа
- 8 Обучаемая по оценкам ассессоров и логам пользователей
- 9 Определяющая число тем автоматически
- 10 Создающая и именующая новые темы автоматически
- 11 Онлайновая: обрабатывающая коллекцию за 1 проход
- 12 Параллельная, распределённая для больших коллекций

Что такое «тема» в коллекции текстовых документов?

Неформально,

- *тема* — семантически однородный кластер текстов
- *тема* — специальная терминология предметной области
- *тема* — набор часто совместно встречающихся терминов

Более формально,

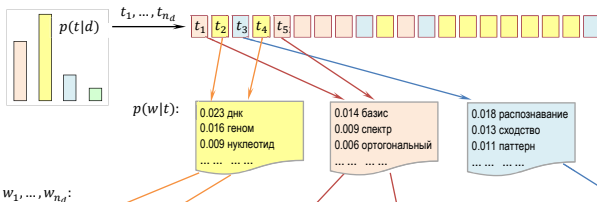
- *тема* — условное распределение на множестве терминов,
 $p(w|t)$ — вероятность (частота) термина w в теме t ;
- *тематика* документа — условное распределение
 $p(t|d)$ — вероятность (частота) темы t в документе d .

Когда автор писал термин w в документ d , он думал о теме t , и мы хотели бы выявить, о какой именно.

Тематическая модель оценивает вероятности $p(w|t)$ и $p(t|d)$ по наблюдаемым частотам $p(w|d)$ слов w в документах d .

Вероятностная модель порождения текстовой коллекции объясняет появление терминов w в документах d темами t :

$$p(w|d) = \sum_t p(w|t)p(t|d)$$



Разработан спектрально-аналитический подход к выявлению размытых протяженных повторов в **геномных** последовательностях. Метод основан на разномасштабном оценивании сходства **нуклеотидных** последовательностей в пространстве коэффициентов разложения фрагментов кривых GC- и GA-содержания по классическим **ортогональным базисам**. Найдены условия оптимальной аппроксимации, обеспечивающие автоматическое **распознавание** повторов различных видов (прямых и инвертированных, а также **тандемных**) на спектральной матрице сходства. Метод одинаково хорошо работает на разных масштабах данных. Он позволяет выявлять следы **сегментных дубликаций** и **мегасателлитные** участки в **геноме**, районы **синтезии** при сравнении пары **геномов**. Его можно использовать для детального изучения фрагментов **хромосом** (поиска размытых участков с умеренной длиной повторяющегося **паттерна**).

Обратная задача — восстановление $p(w|t)$ и $p(t|d)$ по коллекции

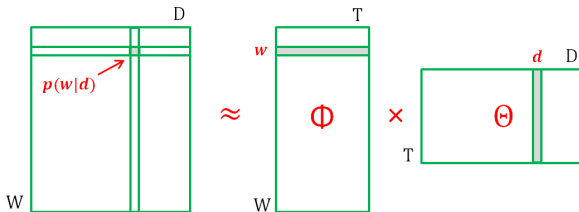
Дано: коллекция текстовых документов

- n_{dw} — частоты терминов в документах, $p(w|d) = \frac{n_{dw}}{n_d}$

Найти: параметры тематической модели $p(w|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$

- $\phi_{wt} = p(w|t)$ — вероятности терминов w в каждой теме t
- $\theta_{td} = p(t|d)$ — вероятности тем t в каждом документе d

Это задача стохастического матричного разложения:



Задачи, некорректно поставленные по Адамару

Задача *корректно поставлена*, если её решение

- существует,
- единственно,
- устойчиво.



Жак Саломон Адамар
(1865–1963)

Задача стохастического матричного разложения *некорректно поставлена*, так как имеет бесконечно много решений:

$$\Phi\Theta = (\Phi S)(S^{-1}\Theta) = \Phi'\Theta'$$

для невырожденных $S_{T \times T}$ таких, что Φ', Θ' — стохастические.

Регуляризация — стандартный приём, введение новых ограничений или критериев, доопределяющих решение.

ARTM — Аддитивная Регуляризация Тематических Моделей

Максимизация \log правдоподобия с регуляризатором R :

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

EM-алгоритм: метод простой итерации для системы уравнений

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

где $\operatorname{norm}_{t \in T} x_t = \frac{\max\{x_t, 0\}}{\sum_{s \in T} \max\{x_s, 0\}}$ — операция нормирования вектора.

Мультимодальная ARTM [Vorontsov et al, 2015]

W^m — словарь токенов m -й модальности, $m \in M$

$W = W^1 \sqcup \dots \sqcup W^M$ — объединённый словарь всех модальностей

Максимизация суммы \log правдоподобий с регуляризацией:

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

EM-алгоритм: метод простой итерации для системы уравнений

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

25 полезных свойств тематических моделей

interpretable	sparse	robust	decorrelated	multigram
multimodal	multilingual	hierarchical	temporal	spacio-temporal
short-text	sentence	segmentation	relational	sentiment
supervised	classification	semi-supervised	auto-labeled	summarization
fast	online	extendable	parallel	distributed

- Для каждого свойства X имеется масса литературы по байесовским моделям « X Topic Model» (погуглите!)
- Про комбинации двух свойств « X Y Topic Model» литературы намного меньше.
- Комбинирование трёх и более моделей крайне редко.

ARTM:

- Переформулировка: байесовская модель \rightarrow регуляризатор
- Складывая регуляризаторы, можем комбинировать модели

Разреживание, сглаживание и декоррелирование тем

interpretable	sparse	robust	decorrelated	multigram
multimodal	multilingual	hierarchical	temporal	spacio-temporal
short-text	sentence	segmentation	relational	sentiment
supervised	classification	semi-supervised	auto-labeled	summarization
fast	online	extendable	parallel	distributed

Сглаживание фоновых тем $t \in B \subset T$ делает модель робастной:

$$R(\Phi, \Theta) = \sum_{t \in B} \sum_{w \in W} \beta_{wt} \ln \phi_{wt} + \sum_{d \in D} \sum_{t \in B} \alpha_{td} \ln \theta_{td} \rightarrow \max.$$

Разреживание тем $t \in S = T \setminus B$ улучшает интерпретируемость:

$$R(\Phi, \Theta) = - \sum_{t \in S} \sum_{w \in W} \beta_{wt} \ln \phi_{wt} - \sum_{d \in D} \sum_{t \in S} \alpha_{td} \ln \theta_{td} \rightarrow \max.$$

Декоррелирование увеличивает различность тем:

$$R(\Phi) = - \frac{\tau}{2} \sum_{t, s \in S} \sum_{w \in W} \phi_{wt} \phi_{ws} \rightarrow \max.$$

Частичное обучение для коррекции тем

interpretable	sparse	robust	decorrelated	multigram
multimodal	multilingual	hierarchical	temporal	spacio-temporal
short-text	sentence	segmentation	relational	sentiment
supervised	classification	semi-supervised	auto-labeled	summarization
fast	online	extendable	parallel	distributed

Сглаживание по «чёрным» и «белым» спискам документов и терминов, указанных ассессорами для каждой темы:

$$R(\Phi, \Theta) = \sum_{t \in T} \sum_{w \in W} \beta_{wt} \ln \phi_{wt} + \sum_{d \in D} \sum_{t \in T} \alpha_{td} \ln \theta_{td} \rightarrow \max.$$

$\beta_{wt} = [w \in W_t^+]$, W_t^+ — *белый список терминов* в теме t

$\alpha_{td} = [d \in D_t^+]$, D_t^+ — *белый список документов* в теме t

$\beta_{wt} = -[w \in W_t^-]$, W_t^- — *чёрный список терминов* в теме t

$\alpha_{td} = -[d \in D_t^-]$, D_t^- — *чёрный список документов* в теме t

Частичное обучение для поиска релевантных тем

interpretable	sparse	robust	decorrelated	multigram
multimodal	multilingual	hierarchical	temporal	spacio-temporal
short-text	sentence	segmentation	relational	sentiment
supervised	classification	semi-supervised	auto-labeled	summarization
fast	online	extendable	parallel	distributed

Применение: выявление в социальных медиа всех тем про болезни / катастрофы / межэтнические отношения / терроризм / страну / компанию / продукт / персону и т.д.

Сглаживание тем из $T_0 \subset T$ по семантическому ядру W_0 :

$$R(\Phi) = \tau \sum_{t \in T_0} \sum_{w \in W_0} \ln \phi_{wt} \rightarrow \max.$$

Paul, M.J., Dredze, M. Discovering health topics in social media using topic models. 2014.

Тематическая модель сети слов (WNTM) для коротких текстов

interpretable	sparse	robust	decorrelated	multigram
multimodal	multilingual	hierarchical	temporal	spacio-temporal
short-text	sentence	segmentation	relational	sentiment
supervised	classification	semi-supervised	auto-labeled	summarization
fast	online	extendable	parallel	distributed

Короткие тексты: сообщения твиттера, заголовки статей.

Тематическая модель сети слов (word network topic model):

$$R(\Phi, \Psi) = \tau \sum_{u, w \in W} n_{uw} \log \sum_{t \in T} \phi_{ut} \psi_{tw} \rightarrow \max_{\Phi, \Psi}$$

где n_{uw} — число со-встречаемостей пары терминов (u, w) в коротком контексте (предложении или окне в 10 слов).

Это модель дистрибутивной семантики, аналог word2vec.

Yuan Zuo, Jichang Zhao, Ke Xu. Word Network Topic Model: a simple but general solution for short and imbalanced texts. 2014.

Мультимодальные и иерархические тематические модели

interpretable	sparse	robust	decorrelated	multigram
multimodal	multilingual	hierarchical	temporal	spacio-temporal
short-text	sentence	segmentation	relational	sentiment
supervised	classification	semi-supervised	auto-labeled	summarization
fast	online	extendable	parallel	distributed

Все эти свойства являются частными случаями модальностей.

Пример: построение следующего уровня иерархии тем:

$$R(\Phi, \Psi) = \tau \sum_{a,w} n_{aw} \ln \sum_t \phi_{wt} \psi_{ta} \rightarrow \max_{\Phi, \Psi}$$

$\psi_{ta} = p(t|a)$ связывают подтемы t с родительскими темами a . Родительские темы a представимы как «псевдодокументы».

N. A. Chirkova, K. V. Vorontsov. Additively Regularized Multimodal Topic Hierarchies. JMLDA. 2016.

Библиотека тематического моделирования BigARTM

Ключевые возможности:

- Комбинирование моделей, регуляризация, модальности
- Онлайн-параллельный EM-алгоритм
- Большие данные: коллекция не хранится в памяти
- Встроенная библиотека регуляризаторов и мер качества

Сообщество:

- Открытый код <https://github.com/bigartm>
- Документация <http://bigartm.org>



Лицензия и среда разработки:

- Freely available for commercial usage (BSD 3-Clause license)
- Cross-platform — Windows, Linux, Mac OS X (32 bit, 64 bit)
- Programming APIs: command-line, Python, C++, C#

Тесты производительности

- 3.7М статей английской Вики, 100К уникальных слов

	procs	train	inference	perplexity
BigARTM	1	35 min	72 sec	4000
Gensim.LdaModel	1	369 min	395 sec	4161
VowpalWabbit.LDA	1	73 min	120 sec	4108
BigARTM	4	9 min	20 sec	4061
Gensim.LdaMulticore	4	60 min	222 sec	4111
BigARTM	8	4.5 min	14 sec	4304
Gensim.LdaMulticore	8	57 min	224 sec	4455

- *procs* = число параллельных потоков
- *inference* = время тематизации 100К тестовых документов
- *perplexity* вычислена на тестовой выборке документов

BigARTM: унификация разработки тематических моделей

На практике чаще всего используют устаревшую модель LDA. Причина — байесовские модели приходится строить «с нуля».

Этапы моделирования

Bayesian TM

ARTM

	Bayesian TM	ARTM
	Анализ требований	Анализ требований
Формализация:	Вероятностная порождающая модель данных	Стандартные критерии Свои критерии
Алгоритмизация:	Байесовский вывод для данной порождающей модели (VI, GS, EP)	Общий регуляризованный EM-алгоритм для любых моделей
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

-- стандартизируемые этапы

Мультиязычная модель Википедии

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

Дударенко М. А. Регуляризация многоязычных тематических моделей.
Вычислительные методы и программирование. 2015. Т. 16. С. 26–36.

Мультиязычная модель Википедии

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

Биграммы радикально улучшают интерпретируемость тем

Коллекция 850 статей конференций ММРО, ИОИ на русском

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

Стенин С. С. Мультиграммные аддитивно регуляризованные тематические модели. Магистерская диссертация, МФТИ, 2015.

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

Примеры этно-релевантных тем

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

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

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

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

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

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

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

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

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

Примеры этно-релевантных тем

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

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

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

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

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

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

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

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

Примеры этно-релевантных 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: ARTM находит больше этно-релевантных тем

Число этно-релевантных тем, найденных моделью:

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

Данные коллективного блога Хабрахабр.ру

Данные

- 132 157 статей
- Модальности:
 - 52 354 терминов (слов)
 - 524 авторов статей
 - 10 000 комментаторов (авторов комментариев к статьям)
 - 2546 тегов
 - 123 хаба (категории)

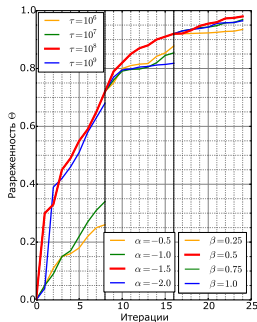
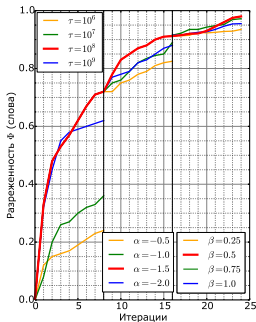
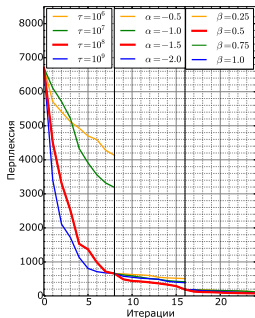
Предобработка текстов

- отброшены 5% наиболее частотных слов (общая лексика)
- удаление пунктуации
- нижний регистр, ё→е
- лемматизация руморphy2

Подбор коэффициентов регуляризации

Последовательное добавление регуляризаторов:

- декоррелирование распределений терминов в темах (τ),
- разреживание распределений тем в документах (α),
- сглаживание распределений терминов в темах (β).



Разведочный поиск

$q = (w_1, \dots, w_{n_q})$ — текст запроса произвольной длины n_q

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

$\theta_{td} = p(t|d)$ — тематические профили 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: *инвертированный индекс* для быстрого поиска документов d по каждой из тем t запроса

Методика оценивания качества разведочного поиска

Поисковый запрос

набор ключевых слов или фрагментов текста, около одной страницы A4

Поисковая выдача

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

Два задания асессорам

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

Набор MapReduce

Набор MapReduce – программа поиска (библиотека) вычислений распределенных вычислений для больших объемов данных и реляционных шардов, представляющих собой набор Java-классов и исполняемых утилит для создания и обработки данных на параллельной обработке.

Основные возможности Набор MapReduce можно сформулировать как:

- обработка вычислений больших объемов данных;
- масштабируемость;
- автоматическое распределение заданий;
- работа на неидеальных оборудовании;
- автоматическая обработка отказов вычислений заданий.

Набор – популярная программная платформа (язык, библиотека) построения распределенных приложений для высоко-параллельной обработки (раздел работы, решение, МЭУ) данных.

Набор включает в себе следующие компоненты:

1. HDFS – распределенная файловая система;

2. Набор MapReduce – программная модель (библиотека) вычислений распределенных вычислений для больших объемов данных и реляционных шардов.

Ключевые особенности в архитектуре Набор MapReduce и структуре HDFS, такие как: принцип разделения задач на кластеры, а также число и название точек отказа. Это, в конечном итоге, определяет ограничение платформ Набор в целом. К сожалению можно отметить:

Ограничение масштабируемости кластера Набор – не вычислительных узлов, – не К параллельных заданий.

Сильная зависимость от распределенных вычислений и клиентских вычислений, реализованных распределенной программой. Как следствие:

Отсутствие поддержки альтернативной программной модели вычислений распределенных вычислений в Набор v1.0 поддерживается только модель вычислений шардов.

Модель вычислений точек отказа и как следствие, неопределенность масштабов и средств с высоким требованиями к надежности.

Проблема совместности требований по единственному объектно-ориентированному интерфейсу всех вычислительных узлов кластера при обращении платформ Набор (установка новой версии или пакета обновлений).

Пример запроса для разведочного поиска

Пример: фрагмент запроса «Система IBM Watson»

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

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

... ..

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

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

Тематика запросов разведочного поиска

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

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

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

Оценки качества поиска

Precision — доля релевантных среди найденных

Recall — доля найденных среди релевантных

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

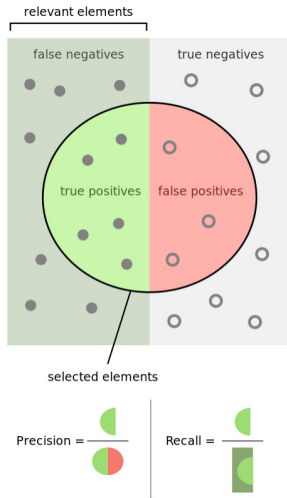
$$R = \frac{TP}{TP + FN} \text{ — полнота, (recall)}$$

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

TP (true positive) — найденные релевантные

FP (false positive) — найденные нерелевантные

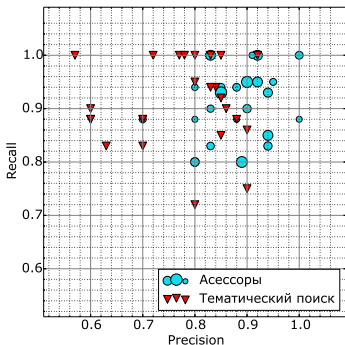
FN (false negative) — ненайденные релевантные



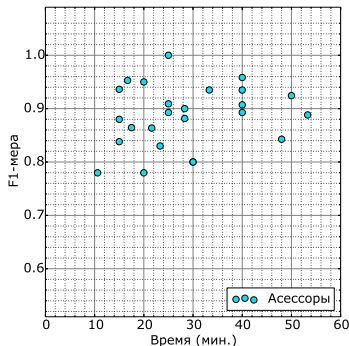
Результаты измерения точности и полноты по запросам

25 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

Выбор модальностей по критериям точности и полноты

Модальности: Слова, Авторы, Комментаторы, Теги, Хабы.
Число тем $|T| = 200$.

	ассессоры	С	К	ТХ	СТ	СХ	СТХ	все
Precision@5	0.82	0.63	0.54	0.59	0.74	0.73	0.73	0.74
Precision@10	0.87	0.67	0.56	0.58	0.77	0.74	0.75	0.77
Precision@15	0.86	0.65	0.53	0.55	0.67	0.67	0.68	0.68
Precision@20	0.85	0.64	0.53	0.54	0.66	0.67	0.68	0.68
Recall@5	0.78	0.77	0.63	0.69	0.82	0.81	0.82	0.82
Recall@10	0.84	0.79	0.64	0.71	0.88	0.82	0.87	0.88
Recall@15	0.88	0.82	0.67	0.73	0.90	0.84	0.89	0.90
Recall@20	0.88	0.85	0.68	0.74	0.91	0.85	0.89	0.91

- Наилучшее качество поиска — по всем модальностям
- Наиболее полезные модальности — термины и теги

Выбор числа тем по критериям точности и полноты

Теперь используем все 5 модальностей, меняем число тем $|T|$

	асессоры	100	200	300	400	500
Precision@5	0.82	0.61	0.74	0.71	0.69	0.59
Precision@10	0.87	0.65	0.77	0.72	0.67	0.61
Precision@15	0.86	0.67	0.68	0.67	0.65	0.62
Precision@20	0.85	0.64	0.68	0.67	0.64	0.60
Recall@5	0.78	0.62	0.82	0.80	0.72	0.63
Recall@10	0.84	0.63	0.88	0.81	0.75	0.64
Recall@15	0.88	0.67	0.90	0.82	0.77	0.67
Recall@20	0.88	0.69	0.91	0.85	0.77	0.68

- Наилучшее качество поиска — при 200 темах
- Тематический поиск превосходит ассессоров по полноте

Янина А. О., Воронцов К. В. Мультимодальные тематические модели для разведочного поиска в коллективном блоге. JMLDA, 2016.









Обсуждение. Цели тематического моделирования

Цели:

- тематический разведочный поиск
- кросс-язычный и мультязычный поиск
- рубрикация, визуализация, систематизация контента
- определение фронтов исследований
- поиск экспертов

Не цели:

- понимание смысла текста
- синтаксический разбор
- машинный перевод

-  *К. В. Воронцов.* Аддитивная регуляризация тематических моделей коллекций текстовых документов. Доклады РАН, Т. 455., №3. 2014.
-  *K. Vorontsov, A. Potapenko.* Tutorial on probabilistic topic modeling: Additive regularization for stochastic matrix factorization. AIST 2014.
-  *K. Vorontsov, A. Potapenko.* Additive regularization of topic models. Machine Learning, 2015.
-  *K. Vorontsov, O. Frei, M. Apishev, P. Romov, M. Suvorova, A. Yanina.* Non-bayesian additive regularization for multimodal topic modeling of large collections. 2015.
-  *K. Vorontsov, A. Potapenko, A. Plavin.* Additive regularization of topic models for topic selection and sparse factorization. SLDS 2015.
-  *O. Frei, M. Apishev.* Parallel non-blocking deterministic algorithm for online topic modeling. AIST 2016.
-  *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.
-  *А. О. Янина, К. В. Воронцов.* Мультимодальные тематические модели для разведочного поиска в коллективном блоге. JMLDA, 2016.