

Тематическое моделирование

Воронцов Константин Вячеславович
ФИЦ ИУ РАН • МФТИ • МГУ •

Научный семинар • ФКН НИУ ВШЭ
30 сентября 2016

1 Разведочный информационный поиск

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

2 Тематическое моделирование

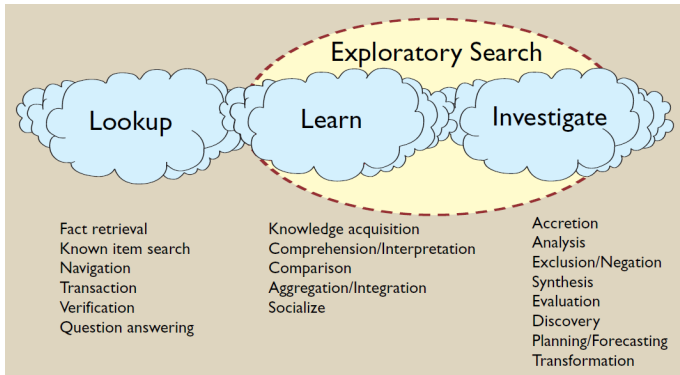
- Вероятностные тематические модели
- Теория аддитивной регуляризации (ARTM)
- Примеры тематических моделей

3 За пределами «мешка слов»

- Короткие тексты и тематическая сегментация
- Битермы, сети слов и когерентность
- Обсуждение и резюме

Концепция разведочного поиска (exploratory search)

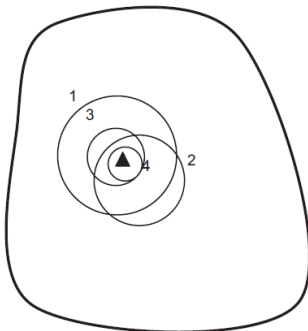
- пользователь может не знать ключевых терминов
- пользователя может интересовать множество ответов



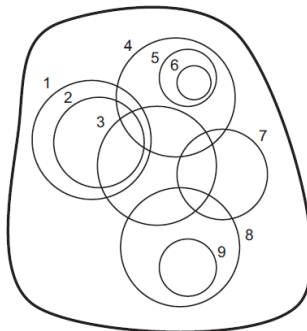
Gary Marchionini. Exploratory Search: from finding to understanding. 2006.

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)

Information Seeking Mantra [B.Shneiderman, 1996]

«Overview first, **zoom and filter, details on demand**»

Понятие *дальнего чтения* [Franco Moretti, 2005]

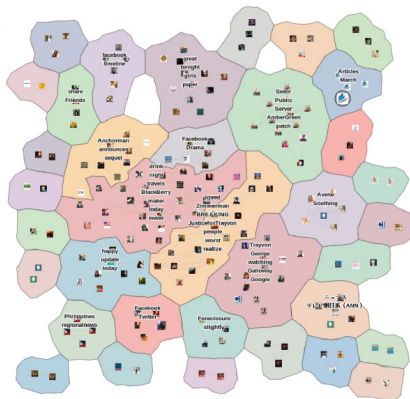
«*Distant reading* is not an obstacle but a specific form of knowledge: fewer elements, hence a sharper sense of their overall interconnection. Shapes, relations, structures. Forms. Models.»

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.

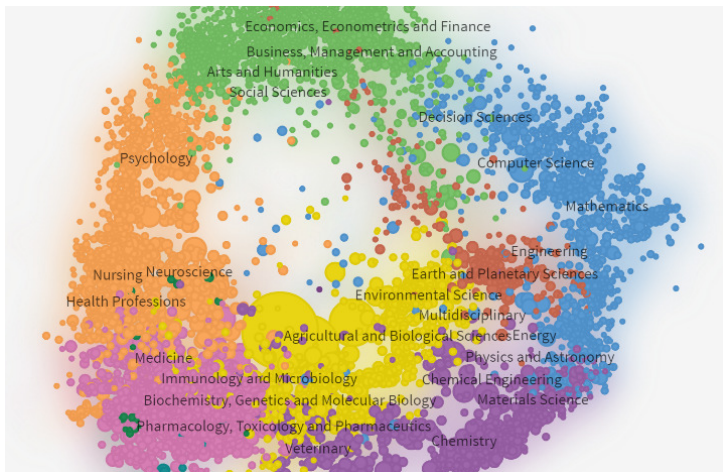
Географическая метафора: карта кластеризации документов



«A map metaphor visualization (left) seems more appealing than a plain graph layout (right), and clusters seem easier to identify.»

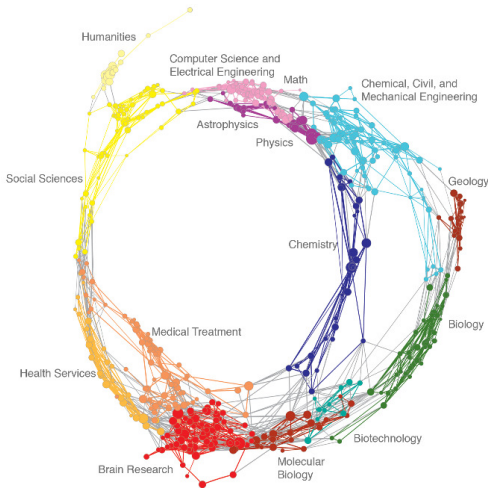
E.R.Gansner, Y.Hu, S.North. Visualizing Streaming Text Data with Dynamic Maps. 2012.

Пример карты науки



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

Ещё один пример карты науки



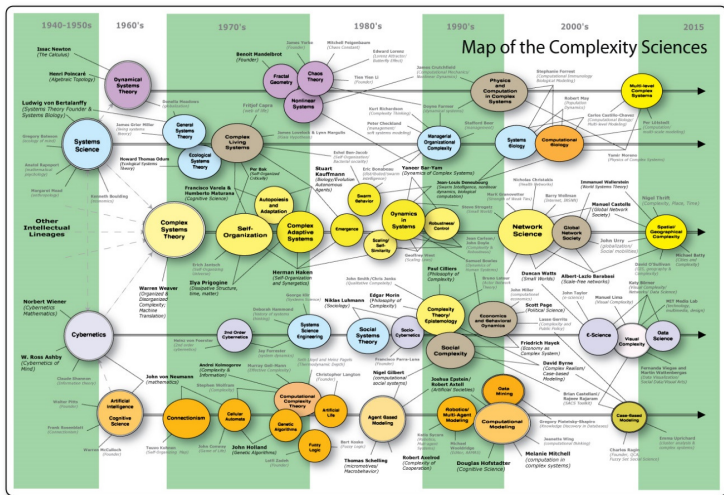
Важное наблюдение:
области знания
самопроизвольно
располагаются по кругу,
значит,
их можно располагать
и вдоль прямой линии.

Недостатки:

- оси не имеют интерпретации
- искажение сходства при двумерном проецировании

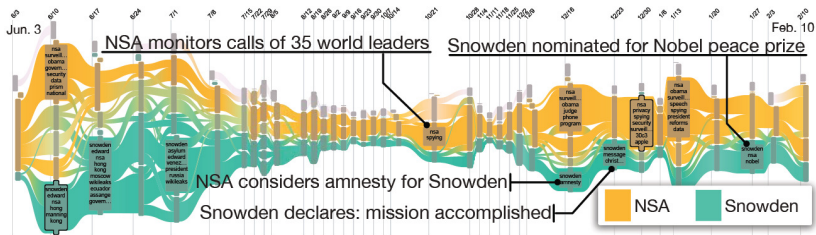
<http://scimaps.org>

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,
- генерирует отчёт.

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 тема–время
- Ось тем: гуманитарные → естественные → точные
- Темы делятся на подтемы иерархически
- Интерактивность: zoom / filter / details
- При любом масштабе на карте достаточно много текста



<http://textvis.lnu.se> — обзор 330 средств визуализации текстов



Айсина Р.М. Обзор средств визуализации тематических моделей коллекций текстовых документов. Машинное обучение и анализ данных. 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 — технология BigARTM
- 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.

Что такое «тема»?

- *тема* — семантически однородный кластер текстов
- *тема* — специальная терминология предметной области
- *тема* — набор терминов (слов или словосочетаний), совместно часто встречающихся в документах
- тем много меньше, чем терминов и чем документов

Более формально,

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

Тематическая модель выявляет латентные темы по наблюдаемым частотам $p(w|d)$ слов w в документах d .

Прямая задача — порождение коллекции по $p(w|t)$ и $p(t|d)$

Вероятностная тематическая модель коллекции документов D описывает появление терминов w в документах d темами t :

$$p(w|d) = \sum_t p(w|t)p(t|d), \quad d \in D$$



w_1, \dots, w_{n_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 дупликаций и мегасателлитные участки в геноме, районы синтении 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: W — словарь терминов

D — коллекция текстовых документов $d = \{w_1 \dots w_{n_d}\}$

n_{dw} = сколько раз термин w встречается в документе d

Найти: параметры модели $p(w|d) = \sum_t \phi_{wt} \theta_{td}$:

$\phi_{wt} = p(w|t)$ — вероятности терминов w в каждой теме t

$\theta_{td} = p(t|d)$ — вероятности тем t в каждом документе d

Задача стохастического матричного разложения является некорректно поставленной — решение не единственно:

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

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

Модель: разумные дополнительные ограничения на Φ, Θ .

Тематическая модель для разведочного поиска должна быть...

- 1 Темпоральная: отображение динамики развития тем
- 2 Иерархическая: систематизация областей знания
- 3 Интерпретируемая: каждая тема понятна для людей
- 4 Мультиграммная: выделение тематичных словосочетаний
- 5 Мультимодальная: авторы, связи, тэги, пользователи,...
- 6 Мультиязычная: кросс- и много-языковой поиск
- 7 Разреженная: для эффективности поискового индекса
- 8 Сегментирующая: выделение тем внутри документа
- 9 Обучаемая: учёт обратной связи с пользователями
- 10 Создающая и именующая темы автоматически
- 11 Онлайновая: обрабатывающая коллекцию за 1 проход
- 12 Параллельная, распределённая для больших коллекций

Некоторые тематические модели

- PLSA (1999) вероятностный латентный семантический анализ
- LDA (2003) латентное размещение Дирихле
- ATM (2004) авторы документов
- TOT (2006) метки времени документов
- HDP (2006) определение числа тем
- TNG (2007) группирование слов в мультиграммы
- CTM (2007) корреляции между темами
- NetPLSA (2008) граф связей между документами
- ML-LDA (2009) многоязычные параллельные тексты
- ssLDA (2012) частичное обучение
- Dependency-LDA (2012) классификация
- BitermTM (2013) битермы в коротких документах
- mLDA (2013) метаданные с тремя и более модальностями
- WNTM (2014) локальные контексты слов

Байесовское обучение — доминирующий подход в ВТМ

Основа подхода — байесовский вывод:

$$\text{Prior}(\Phi, \Theta) + \text{Data} \rightarrow \text{Posterior}(\Phi, \Theta).$$

Проблемы:

- Нам нужны лишь значения Φ, Θ , а не распределения
- Dir — удобный, но лингвистически не обоснованный Prior
- Байесовский вывод уникален для каждой модели
- Технически трудно комбинировать модели
- Невозможно реализовать тысячи моделей в одном коде
- Сложно для понимания, не спасает даже плоская нотация

Rob Zinkov. Stop using Plate Notation.

<http://zinkov.com/posts/2013-07-28-stop-using-plates>

Байесовское обучение — доминирующий подход в ВТМ

$$p(\Theta|\alpha) = \prod_{d=1}^D p(\theta_{d,\cdot}|\alpha) = \prod_{d=1}^D \frac{1}{B(\alpha)} \prod_{k=1}^K \alpha_k^{\theta_{d,k}-1}$$

$$p(Z|\Theta) = \prod_{d=1}^D \theta_{d,\cdot} = \prod_{d=1}^D \prod_{k=1}^K \theta_{d,k}^{I(d,k)}$$

$$p(Z|\alpha) = \int p(Z|\Theta)p(\Theta|\alpha)d\Theta$$

$$= \prod_{d=1}^D \left(\int \frac{1}{B(\alpha)} \prod_{k=1}^K \alpha_k^{I(d,k)+n_{d,k}-1} d\theta_d \right)$$

$$= \prod_{d=1}^D \frac{B(I(d,\cdot)+\alpha)}{B(\alpha)}$$

$$B(d,k) = \sum_{i=1}^K I(d,i) \{d_i = m \wedge \lambda_i = z\}$$

$$p(Z,W|\alpha, \beta) = \prod_{d=1}^D \prod_{i=1}^V \theta_{d,i} \phi_{i,w}$$

$$p(z_i = k | Z, \alpha) = \frac{\theta_{d,k}}{\sum_{l=1}^K \theta_{d,l}}$$

$$p(w = t | z = k, W, Z, \beta) = \frac{\phi_{k,t}}{\sum_{s=1}^V \phi_{k,s}}$$

$$p(w = t, z = k | W, Z, \alpha, \beta) = p(z = k | Z, \alpha) p(w = t | z = k, W, Z, \beta)$$

$$= \frac{\theta_{d,k} \phi_{k,t}}{\sum_{l=1}^K \theta_{d,l} \sum_{s=1}^V \phi_{k,s}}$$

$$p(z_i = k | Z, \alpha)$$

$$p(w = t | z = k, W, Z, \beta)$$

$$p(w = t, z = k | W, Z, \alpha, \beta)$$

Graphical models showing hierarchical structures with nodes for topics, words, and documents, illustrating the probabilistic relationships between them.

ARTM — альтернатива байесовскому обучению

The central white box contains the following equations:

$$\left\{ \begin{array}{l} p_{tdw} = \text{norm}_t(\phi_{wt}\theta_{td}) \\ \phi_{wt} = \text{norm}_w\left(\sum_d n_{dw} p_{tdw} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}\right) \\ \theta_{td} = \text{norm}_t\left(\sum_w n_{dw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}\right) \end{array} \right.$$

Other visible formulas include:

- $p(\Theta|\alpha) = \prod_{d=1}^D p(\theta_{d,\cdot}|\alpha) = \prod_{d=1}^D \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_{d,k}^{\alpha_k - 1}$
- $p(\Theta) = \prod_{d=1}^D \theta_{d,\cdot}$
- $p(Z|\Theta) = \int p(Z|\Theta)p(\Theta|\alpha)\Theta$
- $p(Z, W|\alpha, \beta) = \prod_{d=1}^D \prod_{t=1}^T \prod_{w=1}^W p(z_{dt}, w_{dt}|\alpha, \beta)$
- $p(z_{dt} = k|Z_{-dt}, W_{-dt}) = \frac{\alpha_k + \sum_{d'=1}^D \sum_{t'=1}^T \mathbb{1}(z_{d't'} = k)}{\sum_{k=1}^K (\alpha_k + \sum_{d'=1}^D \sum_{t'=1}^T \mathbb{1}(z_{d't'} = k))}$
- $p(w_{dt} = l|Z_{-dt}, W_{-dt}) = \frac{\beta_l + \sum_{d'=1}^D \sum_{t'=1}^T \mathbb{1}(w_{d't'} = l)}{\sum_{l=1}^L (\beta_l + \sum_{d'=1}^D \sum_{t'=1}^T \mathbb{1}(w_{d't'} = l))}$

Diagram elements include:

- Plate notation for $p(z_{dt} = k|Z_{-dt}, W_{-dt})$ and $p(w_{dt} = l|Z_{-dt}, W_{-dt})$.
- A hierarchical tree diagram with nodes labeled $\alpha, \beta, \tau, \sigma, \tau, \sigma, \tau, \sigma$ and a note "These trees grouped into M documents".
- Other mathematical expressions like $p(z_{dt} = k|Z_{-dt}, W_{-dt})$ and $p(w_{dt} = l|Z_{-dt}, W_{-dt})$.

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

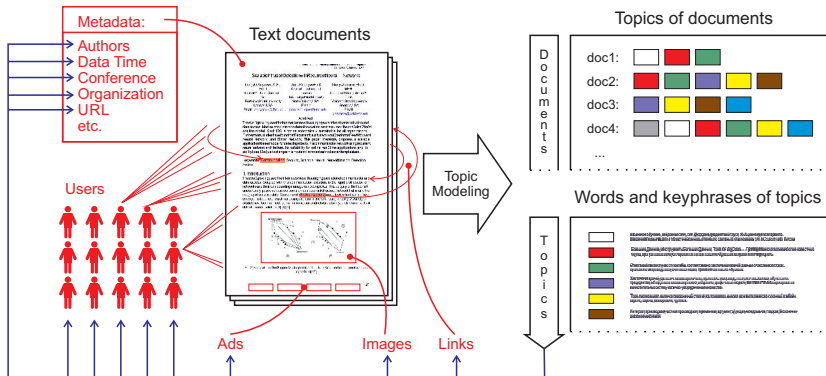
$$\begin{cases} \text{E-шаг:} & p_{tdw} = \operatorname{norm}_{t \in T}(\phi_{wt} \theta_{td}) \\ \text{M-шаг:} & \begin{cases} \phi_{wt} = \operatorname{norm}_{w \in W} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right), & n_{wt} = \sum_{d \in D} n_{dw} p_{tdw} \\ \theta_{td} = \operatorname{norm}_{t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

Модель PLSA: $R(\Phi, \Theta) = 0$

Модель LDA: $R(\Phi, \Theta) = \sum_{t,w} \beta_w \ln \phi_{wt} + \sum_{d,t} \alpha_t \ln \theta_{td}$

ARTM легко обобщается на мультимодальные задачи

Выявление тематики документов $p(t|d)$ и терминов $p(t|w)$,
 а также модальностей: $p(t|автор)$, $p(t|время)$, $p(t|ссылка)$,
 $p(t|баннер)$, $p(t|элемент изображения)$, $p(t|пользователь)$,...



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

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

Библиотека тематического моделирования BigARTM

Ключевые возможности:

- Комбинирование требований, моделей, модальностей
- (благодаря теории аддитивной регуляризации, ARTM)
- 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#

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

 -- стандартизуемые этапы

Тесты производительности

- 3.7M статей английской Вики, 100K 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 100K тестовых документов
- *perplexity* вычислена на тестовой выборке документов

№1. 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.

№1. 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.

№2. Биграммы радикально улучшают интерпретируемость тем

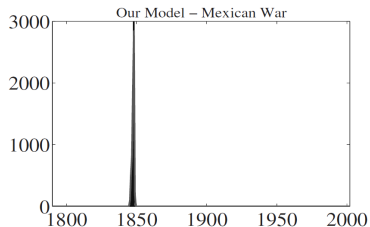
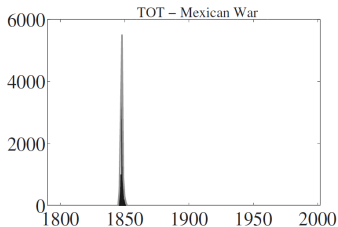
Коллекция 850 статей конференций ММРО, ИОИ на русском

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	тупиковое множество	комитет	задача MASC

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

№3. Совмещение динамической и n -граммной модели

По коллекции выступлений президентов США



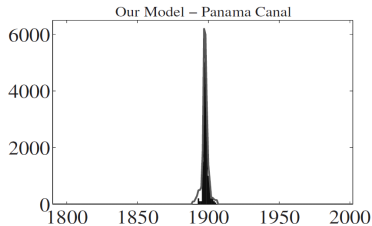
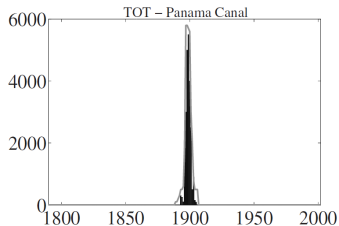
1. mexico	8. territory
2. texas	9. army
3. war	10. peace
4. mexican	11. act
5. united	12. policy
6. country	13. foreign
7. government	14. citizens

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.

№3. Совмещение динамической 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

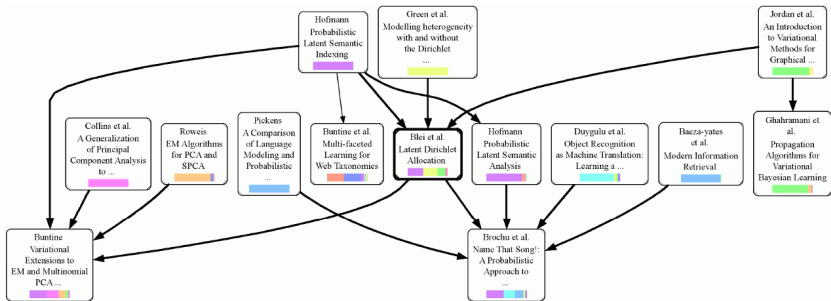
1. panama canal	8. united states senate
2. isthmian canal	9. french canal company
3. isthmus panama	10. caribbean sea
4. republic panama	11. panama canal bonds
5. united states government	12. panama
6. united states	13. american control
7. state panama	14. canal

Shoaib Jameel, Wai Lam. An N -Gram Topic Model for Time-Stamped Documents. ECIR 2013.

№4. Модели, учитывающие цитирования или гиперссылки

Учёт ссылок уточняет тематическую модель

Тематическая модель выявляет самые влиятельные ссылки



Laura Dietz, Steffen Bickel, Tobias Scheffer. Unsupervised prediction of citation influences. ICML-2007.

№6. 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

№6. 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

№6. 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,

№6. 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,

№6. 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,

№6. Результаты: 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**

№7. 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)

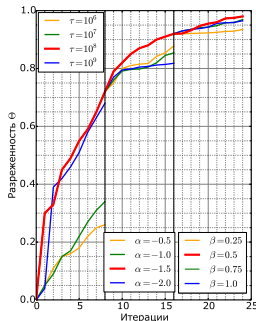
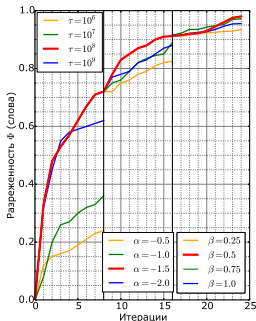
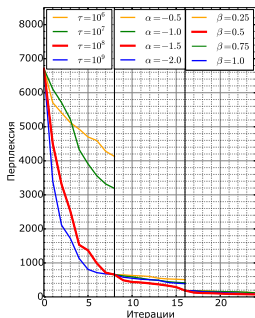
Предобработка текстов

- отброшены 5% наиболее частотных слов (общая лексика)
- удаление пунктуации
- нижний регистр, ё→е
- лемматизация `rumorphy2`

№7. Подбор коэффициентов регуляризации

Последовательное добавление регуляризаторов:

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



№7. 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

№7. 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) вычислений распределенных вычислений для больших объемов данных в рамках параллельных вычислений, представляющая собой набор Java-классов и исполняемых утилит для создания и обработки данных на параллельной обработке.

Основные возможности Поиск MapReduce можно сформулировать как:

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

Поиск – популярная программная платформа (язык Java, библиотека) построения распределенных приложений для массово-параллельной обработки (разное разбиение, репликация, МР) данных.

Поиск включает в себе следующие компоненты:

1. HDFS – распределенная файловая система;

2. **Поиск MapReduce** – программная модель (библиотека) вычислений распределенных вычислений для больших объемов данных в рамках параллельных вычислений.

Ключевые, **значение** и архитектура **Поиск MapReduce** и структура HDFS, стали примером того, как можно работать с данными, в том числе и с большими объемами данных. Это, в конечном итоге, определило направление платформ **Поиск** в целом. К последним можно отнести:

Ограничение масштабируемости кластера **Поиск** –4K вычислительных узлов, –4K параллельных заданий.

Сильная зависимость **Поиск** от распределенных вычислений и клиентских вычислений, реализованных распределенной архитектурой. Как следствие:

Отсутствие поддержки альтернативной программной модели вычислений распределенных вычислений в **Поиск v1.0** подразумевает только модель вычислений **MapReduce**.

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

Проблема **высокой** совместности требований по единственному объектно-ориентированному языку вычислений при обращении платформ **Поиск** (установка новой версии или пакета обновлений).

Пример запроса для разведочного поиска

№7. 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.

№7. Тематика запросов разведочного поиска

Примеры заголовков разведочных запросов к Хабру
(объём каждого запроса — около одной страницы A4):

Алгоритмы раскраски графов	Система IBM Watson
Рекомендательная система Netflix	3D-принтеры
Методики быстрого набора текста	CERN-кластер
Космические проекты Илона Маска	АВ-тестирование
Технологии Hadoop MapReduce	Облачные сервисы
Беспилотный автомобиль Google car	Контекстная реклама
Криптосистемы с открытым ключом	Марсоход Curiosity
Обзор платформ онлайн-курсов	Видеокарты NVIDIA
Data Science Meetups в Москве	Распознавание образов
Образовательные проекты mail.ru	Сервисы Google scholar
Межпланетная станция New horizons	MIT MediaLab Research
Языковая модель word2vec	Платформа Microsoft Azure

№7. Оценки качества поиска

Precision — доля релевантных среди найденных

Recall — доля найденных среди релевантных

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

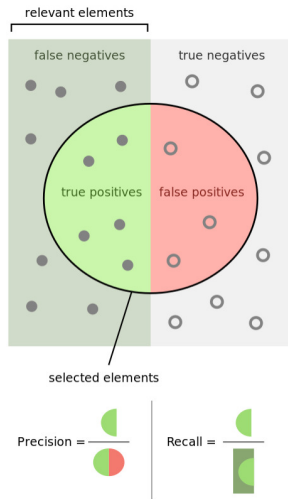
$$R = \frac{TP}{TP + FN} \text{ — полнота, (recall)}$$

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

TP (true positive) — найденные релевантные

FP (false positive) — найденные нерелевантные

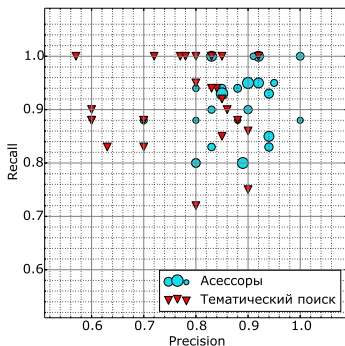
FN (false negative) — не найденные релевантные



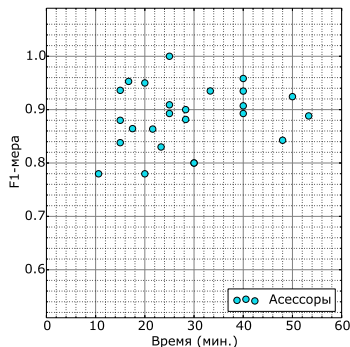
№7. Результаты измерения точности и полноты по запросам

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

№7. 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

№7. 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 асессоров по полноте

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

Проблема коротких текстов

Короткие тексты (short text): Twitter и другие микроблоги, социальные медиа, заголовки статей и новостных сообщений.

Основные проблемы коротких сообщений:

- огромный объём ($\sim 10^9$ твитов в день)
- опечатки и намеренное искажение слов языка
- концентрация распределения $p(t|d)$ в одной теме
- выделение редких тем на фоне основных тем микроблогов (личная переписка, life style, репосты новостей)
- раннее обнаружение новых тем

Тривиальные подходы и их недостатки

- Считать каждое сообщение отдельным документом
 - для коротких сообщений $p(t|d)$ оценивается не надёжно
- Разреживать $p(t|d)$ вплоть до единственной темы; добавить модальности авторов, времени, регионов и т.п.
 - решение в духе АРТМ, пока не попробовали...
- Объединить сообщения по автору (времени, региону и т.п.)
 - появится дисбаланс документов по длине
 - появятся тематически неоднородные документы
- Объединить посты с комментариями
 - комментарии могут отсутствовать у большинства постов
- Дополнить коллекцию длинными текстами (Википедия и др.)
 - часть тем может не покрываться внешней коллекций
 - лексикон социальной сети может существенно отличаться

Модель Twitter-LDA

Предположения:

1. Каждый автор $a \in A$ написал множество сообщений $d \in D_a$.
2. Каждое сообщение d относится к одной теме $p(t|d) \in \{0, 1\}$.
3. Есть фоновая тема $b \in T$ с распределением $p(w|b)$.
4. Вероятность фона одинакова для документов, $p(b|d) = \pi$.

Порождающий процесс:

Input: распределения $p(w|t)$, $p(t|a)$

- 1 **for all** авторов $a \in A$
- 2 **for all** сообщений $d \in D_a$ автора a
- 3 выбрать тему t из $p(t|a)$, кроме фоновой, $t \neq b$;
- 4 **for all** позиций слов $i = 1, \dots, n_d$ в сообщении d
- 5 выбрать слово w_i из $(1 - \pi)p(w|t) + \pi p(w|b)$;

Wayne Xin Zhao, Jing Jiang, Jianshu Weng, Jing He, Ee Peng Lim et al.
Comparing Twitter and traditional media using topic models // ECIR 2011.

Тематическая модель предложений

Аналог Twitter-LDA с точностью до переобозначений:

1. Каждый документ d состоит из предложений $s \in S_d$.
2. Каждое предложение относится к одной теме $p(t|s) \in \{0, 1\}$.
3. Наблюдаемая выборка образуется тройками $(d_i, s_i, w_i)_{i=1}^n$.
4. Гипотеза условной независимости: $p(s, w|t) = p(s|t)p(w|t)$.

Тематическая модель сегментированного текста:

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

Критерий максимума регуляризованного правдоподобия:

$$\sum_{d \in D} \sum_{s \in S_d} \sum_{w \in d} n_{dsw} \ln \sum_{t \in T} \phi_{wt} \psi_{st} \theta_{td} + R(\Phi, \Psi, \Theta) \rightarrow \max_{\Phi, \Psi, \Theta}$$

где n_{dsw} — частота термина w в предложении $s \in S_d$.

EM-алгоритм для модели сегментированного текста

Максимизация \log правдоподобия с регуляризатором R :

$$\sum_{d \in D} \sum_{s \in S_d} \sum_{w \in d} n_{dsw} \ln \sum_{t \in T} \phi_{wt} \psi_{st} \theta_{td} + R(\Phi, \Psi, \Theta) \rightarrow \max_{\Phi, \Psi, \Theta}$$

EM-алгоритм: метод простой итерации для системы уравнений

$$\begin{cases} \text{E-шаг:} & p_{tdsw} \equiv p(t|d, s, w) = \operatorname{norm}_{t \in T}(\phi_{wt} \psi_{st} \theta_{td}); \\ \text{M-шаг:} & \left\{ \begin{array}{l} \phi_{wt} = \operatorname{norm}_{w \in W} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right); \quad n_{wt} = \sum_{d,s} n_{dsw} p_{tdsw} \\ \psi_{st} = \operatorname{norm}_{s \in S_d} \left(n_{st} + \psi_{st} \frac{\partial R}{\partial \psi_{st}} \right); \quad n_{st} = \sum_w n_{dsw} p_{tdsw} \\ \theta_{td} = \operatorname{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad n_{td} = \sum_{s,w} n_{dsw} p_{tdsw} \end{array} \right. \end{cases}$$

Битермы: модель совстречаемости слов в коротких текстах

Битерм — пара слов, встречающихся рядом:
в одном коротком сообщении / предложении / окне $\pm h$ слов.

Тематическая модель битермов (Biterm topic model):

$$p(u, w) = \sum_{t \in T} p(w|t)p(u|t)p(t) = \sum_{t \in T} \phi_{wt}\phi_{ut}\pi_t,$$

где $\phi_{wt} = p(w|t)$, $\pi_t = p(t)$ — параметры модели.

Критерий максимума логарифма правдоподобия:

$$\sum_{u,w} n_{uw} \ln \sum_t \phi_{wt}\phi_{ut}\pi_t \rightarrow \max_{\Phi, \pi},$$
$$\phi_{wt} \geq 0; \quad \sum_w \phi_{wt} = 1; \quad \pi_t \geq 0; \quad \sum_t \pi_t = 1$$

Xiaohui Yan, Jiafeng Guo, Yanyan Lan, Xueqi Cheng. A Biterm Topic Model for Short Texts // WWW 2013.

Необходимые условия точки максимума правдоподобия

Максимизация \log правдоподобия с регуляризатором R :

$$\sum_{u,w} n_{uw} \ln \sum_t \phi_{wt} \phi_{ut} \pi_t + R(\Phi, \pi) \rightarrow \max_{\Phi, \pi}$$

n_{uw} — частота битерма (u, w) в документах коллекции.

EM-алгоритм: метод простой итерации для системы уравнений

$$\begin{cases} \text{E-шаг:} & p_{tuw} \equiv p(t|u, w) = \operatorname{norm}_{t \in T}(\phi_{wt} \phi_{ut} \pi_t) \\ \text{M-шаг:} & \begin{cases} \phi_{wt} = \operatorname{norm}_{w \in W} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right), & n_{wt} = \sum_{u \in W} n_{uw} p_{tuw} \\ \pi_t = \operatorname{norm}_{t \in T} \left(n_t + \pi_t \frac{\partial R}{\partial \pi_t} \right), & n_t = \sum_{u, w \in W} n_{uw} p_{tuw} \end{cases} \end{cases}$$

Битермы как регуляризатор для обычной $\Phi\Theta$ -модели

1. Регуляризатор битермов для матрицы Φ :

$$R(\Phi) = \tau \sum_{u,w \in W} n_{uw} \ln \sum_{t \in T} n_t \phi_{ut} \phi_{wt} \rightarrow \max$$

Подставляем в формулу M-шага, получаем сглаживание:

$$\phi_{wt} = \text{norm}_w \left(n_{wt} + \tau \sum_{u \in W} n_{uw} p_{tuw} \right);$$

$$p_{tuw} \equiv p(t|u, w) = \text{norm}_{t \in T} (n_t \phi_{wt} \phi_{ut}).$$

2. Регуляризатор разреживания для матрицы Θ :

$$R(\Theta) = -\tau' \sum_{d \in D} \sum_{t \in T} \alpha_t \ln \theta_{td} \rightarrow \max.$$

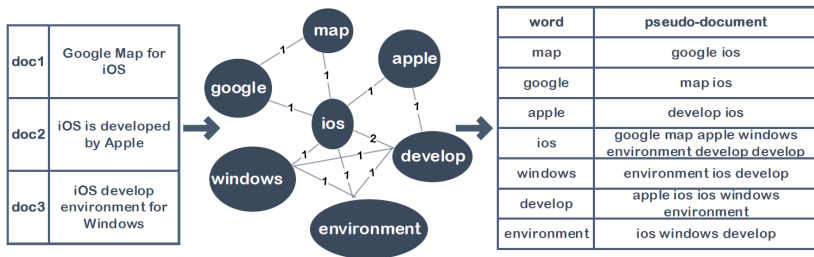
Модель сети слов WNTM для коротких текстов

Идея: моделировать не документы, а связи между словами.

d_w — псевдо-документ, объединение всех контекстов слова w .

n_{wu} — число вхождений слова u в псевдо-документ d_w .

Контекст — короткое сообщение / предложение / окно $\pm h$ слов.



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

Модели WNTM и WTM (Word Topic Model)

Тематическая модель контекстов, разложение $W \times W$ -матрицы:

$$p(u|d_w) = \sum_{t \in T} p(u|t)p(t|d_w) = \sum_{t \in T} \phi_{ut}\theta_{tw},$$

где d_w — псевдо-документ слова w .

Максимизация логарифма правдоподобия:

$$\sum_{u, w \in W} n_{wu} \log \sum_{t \in T} \phi_{ut}\theta_{tw} \rightarrow \max_{\Phi, \Theta},$$

где n_{wu} — встречаемость слов w, u .

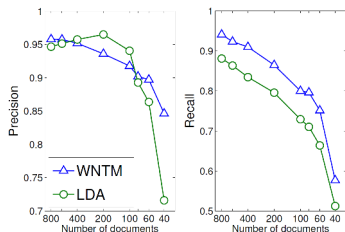
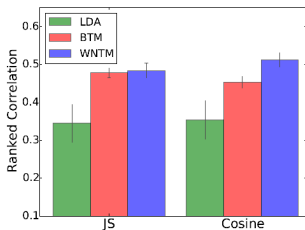
Отличие от модели битермов: там $\Theta = \text{diag}(p_1, \dots, p_t)\Phi^T$.

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

Berlin Chen. Word Topic Models for spoken document retrieval and transcription // ACM Trans., 2009.

Результаты оценивания модели WNTM

- Когерентность на коротких текстах лучше, чем у LDA и Bitern TM; на длинных текстах преимуществ нет.
- Слева: оценивание семантической близости слов по $p(t|w)$, корреляция с 10-балльными экспертными оценками.
- Справа: полнота и точность распознавания новой темы в зависимости от числа документов.



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

Интерпретируемости и когерентность

Тема интерпретируемая, если по топовым словам темы эксперт может определить, о чём эта тема, и дать ей название.

- Экспертные оценки:
 - интерпретируемость темы по балльной шкале;
 - каждую тему оценивают несколько экспертов.
- Метод интрузий (intrusion):
 - в список топовых слов внедряется лишнее слово;
 - измеряется доля ошибок экспертов его при определении

Нужна автоматически вычисляемая мера интерпретируемости, коррелирующая с экспертными оценками.

Ею оказалась *когерентность* (согласованность, coherence).

Newman D., Lau J.H., Grieser K., Baldwin T. Automatic evaluation of topic coherence // Human Language Technologies, HLT-2010, Pp. 100–108.

Эксперимент. Связь когерентности и интерпретируемости

Измерялась ранговая
корреляция Спирмена
между 15 метрикам
и экспертными оценками
интерпретируемости.

PMI — лучшая метрика.

Gold-standard — средняя
корреляция Спирмена
между оценками
разных экспертов.

Resource	Method	Median	Mean
WordNet	HSO	0.15	0.59
	JCN	-0.20	0.19
	LCH	-0.31	-0.15
	LESK	0.53	0.53
	LIN	0.09	0.28
	PATH	0.29	0.12
	RES	0.57	0.66
	VECTOR	-0.08	0.27
	WuP	0.41	0.26
	RACO	0.62	0.69
Wikipedia	MiW	0.68	0.70
	DOCsim	0.59	0.60
	PMI	0.74	0.77
Google	TITLES	0.51	
	LOGHITS	-0.19	
Gold-standard	IAA	0.82	0.78

Вывод: когерентность близка к «золотому стандарту».

Newman D., Lau J.H., Grieser K., Baldwin T. Automatic evaluation of topic coherence // Human Language Technologies, HLT-2010, Pp. 100–108.

Когерентность как внутренняя мера интерпретируемости

Когерентность (согласованность) темы t по k топовым словам:

$$\text{PMI}_t = \frac{2}{k(k-1)} \sum_{i=1}^{k-1} \sum_{j=i}^k \text{PMI}(w_i, w_j)$$

где w_i — i -й термин в порядке убывания ϕ_{wt} .

$\text{PMI}(u, v) = \ln \frac{|D|N_{uv}}{N_u N_v}$ — поточечная взаимная информация (pointwise mutual information),

N_{uv} — число документов, в которых термины u, v хотя бы один раз встречаются рядом (в окне 10 слов),

N_u — число документов, в которых u встретился хотя бы 1 раз.

Newman D., Lau J.H., Grieser K., Baldwin T. Automatic evaluation of topic coherence // Human Language Technologies, HLT-2010, Pp. 100–108.

Регуляризатор для максимизации когерентности тем

Гипотеза: тема лучше интерпретируется, если она содержит *когерентные* (часто встречающиеся рядом) слова $u, w \in W$.

Пусть C_{uw} — оценка когерентности, например $\hat{p}(w|u) = \frac{N_{uw}}{N_u}$.
Согласуем ϕ_{wt} с оценками $\hat{p}(w|t)$ по когерентным словам,

$$\hat{p}(w|t) = \sum_u p(w|u)p(u|t) = \frac{1}{n_t} \sum_u C_{uw} n_{ut};$$

$$R(\Phi) = \tau \sum_{t \in T} n_t \sum_{w \in W} \hat{p}(w|t) \ln \phi_{wt} \rightarrow \max.$$

Подставляем в формулу M-шага, получаем сглаживание:

$$\phi_{wt} = \operatorname{norm}_w \left(n_{wt} + \tau \sum_{u \in W} C_{uw} n_{ut} \right).$$

Mimno D., Wallach H. M., Talley E., Leenders M., McCallum A. Optimizing semantic coherence in topic models // EMNLP-2011. — Pp. 262–272.

Альтернативный регуляризатор когерентности

Квадратичный регуляризатор Quad-Reg:

$$R(\Phi) = \tau \sum_{t \in T} \ln \sum_{u, v \in W} C_{uv} \phi_{ut} \phi_{vt} \rightarrow \max,$$

$$C_{uv} = N_{uv} [\text{PMI}(u, v) > 0],$$

N_{uv} — число документов, в которых u, v хотя бы раз встречаются на расстоянии не более 10 слов,

N_u — число документов, в которых u встречается хотя бы раз,

$\text{PMI}(u, v) = \ln \frac{|D| N_{uv}}{N_u N_v}$ — поточечная взаимная информация.

В литературе пока не выработан окончательный вариант регуляризатора когерентности.

Newman D., Bonilla E. V., Buntine W. L. Improving topic coherence with regularized topic models. 2011.

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

Элементы лингвистики в тематическом моделировании

Предобработка:

- Морфология, лемматизация
- Выделение терминов: POST, подчинительные связи, NER
- Выделение общей лексики (контрастные корпуса)

Лингвистическая регуляризация:

- Выделение общей лексики
(сглаживание, разреживание, декоррелирование)
- Выделение терминов (оценивание тематичности)
- Короткие тексты, битермы, когерентность, WTM, WNTM
- Тематическая сегментация
- Тематические модели предложений
- Синтаксические тематические модели (не взлетело)

Открытые проблемы

- Что такое «тема»? Тема у слова, предложения или абзаца?
- Полнота и устойчивость
- Отсев зависимых/дублирующих/слитых/расщеплённых тем
- Автоматическое обнаружение новых тем (в новостях)
- Автоматическое выделение сюжетов (в новостях)
- Именованние новых тем
- Тематическая суммаризация
- Выделение тематических фактов, понятий, определений
- Тематическая сегментация
- Внутритекстовая когерентность
- Интерпретируемые векторные представления (W2V+TM)
- Адаптивный выбор коэффициентов регуляризации
- Алгоритмы инициализации

-  *K. Vorontsov*. Additive regularization for topic models of text collections. 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. (в печати)
-  *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. (в печати)
-  *А.О.Янина, К.В.Воронцов*. Мультимодальные тематические модели для разведочного поиска в коллективном блоге. ИОИ 2016. (на рецензии)