

Вероятностные тематические модели

Лекция 4.

Тематический поиск. Эксперименты с тематическими моделями

К. В. Воронцов
vokov@forecsys.ru

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

ВМК МГУ • 15 марта 2018

- 1 Разведочный информационный поиск**
 - Концепция разведочного поиска
 - Визуализация больших текстовых коллекций
 - Сценарий разведочного поиска
- 2 Эксперименты с тематическим поиском**
 - Методика эксперимента
 - Построение тематической модели
 - Оптимизация гиперпараметров
- 3 Эксперименты с тематическими моделями**
 - Измерение качества тематической модели
 - Многокритериальное оценивание качества модели
 - Проблема определения числа тем

Напоминание. Задача тематического моделирования

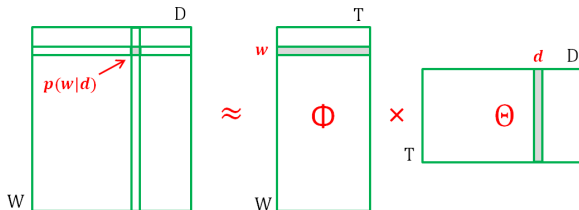
Дано: коллекция текстовых документов, $p(w|d) = \frac{n_{dw}}{n_d}$

Вероятностная тематическая модель:

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

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

Это задача стохастического матричного разложения:



Максимизация \log правдоподобия с регуляризатором R :

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

EM-алгоритм: метод простой итерации для системы уравнений

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

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

Максимизация \log правдоподобия с k регуляризаторами R_i :

$$\sum_{d,w} n_{dw} \ln \sum_t \phi_{wt} \theta_{td} + \sum_{i=1}^k \tau_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

где τ_i — коэффициенты регуляризации.

EM-алгоритм: метод простой итерации для системы уравнений

$$\begin{cases} \text{E-шаг:} & p_{tdw} = \operatorname{norm}_{t \in T}(\phi_{wt} \theta_{td}) \\ \text{M-шаг:} & \begin{cases} \phi_{wt} = \operatorname{norm}_{w \in W} \left(\sum_{d \in D} n_{dw} p_{tdw} + \phi_{wt} \sum_{i=1}^k \tau_i \frac{\partial R_i}{\partial \phi_{wt}} \right) \\ \theta_{td} = \operatorname{norm}_{t \in T} \left(\sum_{w \in W} n_{dw} p_{tdw} + \theta_{td} \sum_{i=1}^k \tau_i \frac{\partial R_i}{\partial \theta_{td}} \right) \end{cases} \end{cases}$$

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

Максимизация суммы \log правдоподобий с регуляризацией:

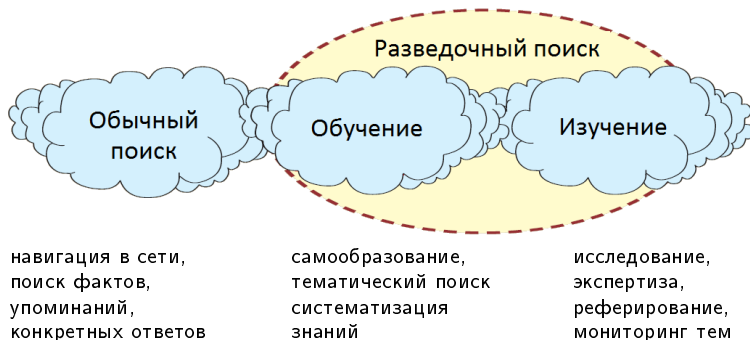
$$\sum_{m \in M} \tau_m \sum_{d \in D} \sum_{w \in W^m} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

EM-алгоритм: метод простой итерации для системы уравнений

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

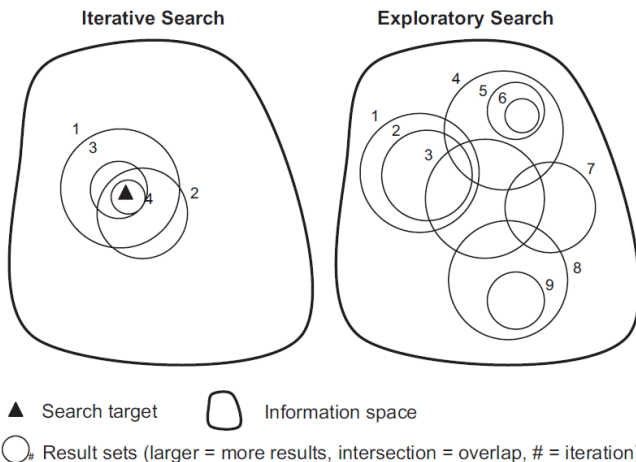
Концепция разведочного поиска (exploratory search)

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



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

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



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* — не ограничение, а способ представления знаний: меньше элементов, чётче понимание их взаимосвязей, акцент на формах, отношениях, структурах, моделях»

Мантра Шнейдермана

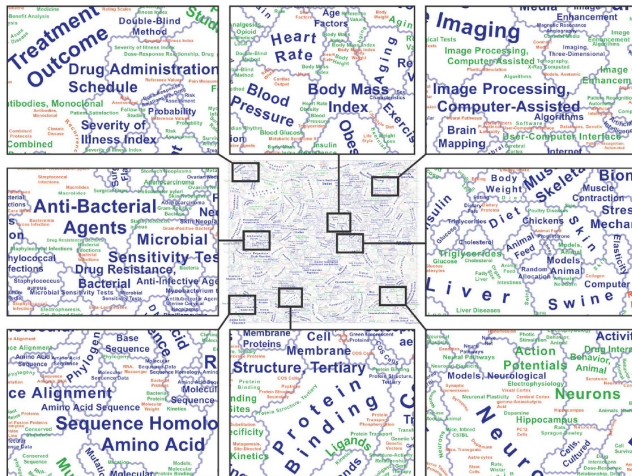
«Сначала крупный план, затем масштабирование и фильтрация, детали по требованию»

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.

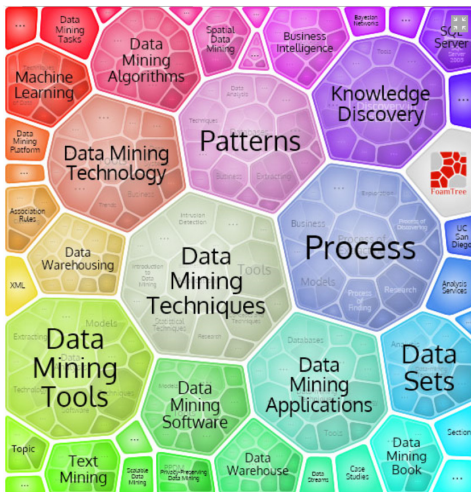
S.Janicke, G.Franzini, M.F.Cheema, G.Scheuermann. On Close and Distant Reading in Digital Humanities: A Survey and Future Challenges. EuroVis, 2015.

Пример карты медицинских знаний



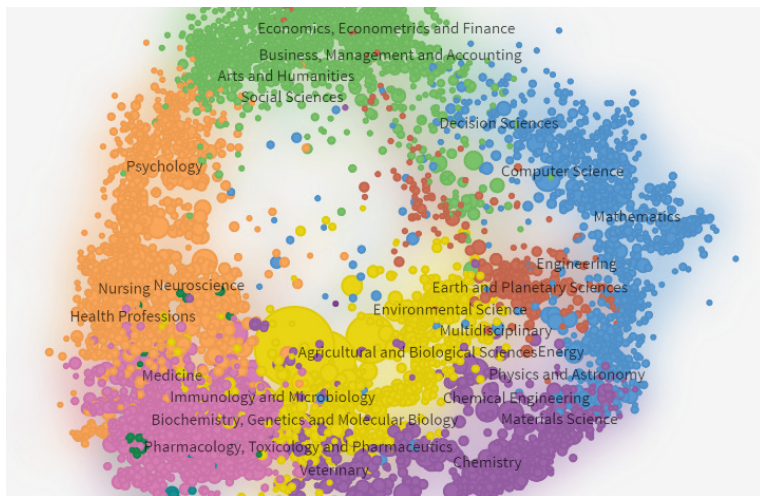
Skupin, Biberstine, Borner. Visualizing the Topical Structure of the Medical Sciences: A Self-Organizing Map Approach. PLoS ONE, 2013.

Пример иерархической карты области *Data Mining*



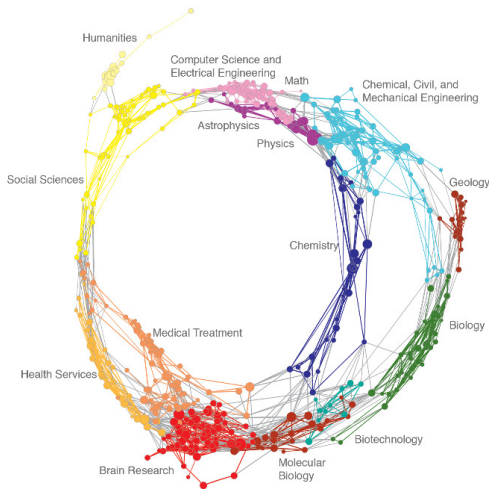
FoamTree: <https://carrotsearch.com/foamtree>

Пример карты науки



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

Ещё один пример карты науки



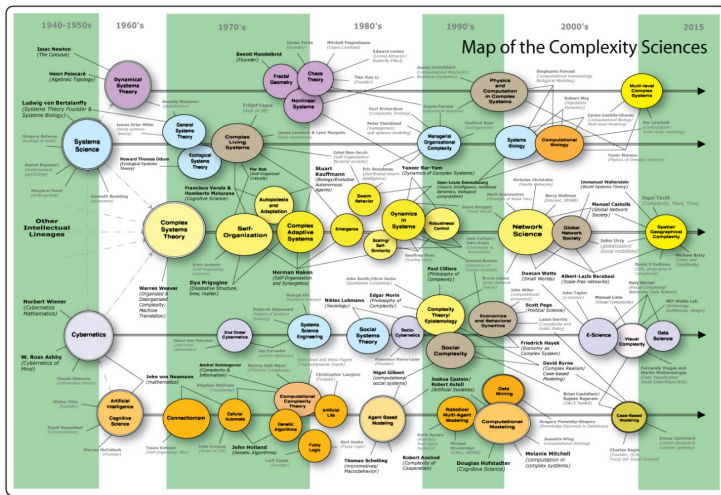
Важное наблюдение:
области знания
самопроизвольно
располагаются по кругу,
значит,
их можно располагать
и вдоль прямой линии.

Недостатки:

- оси не имеют интерпретации
- искажение сходства при двумерном проецировании

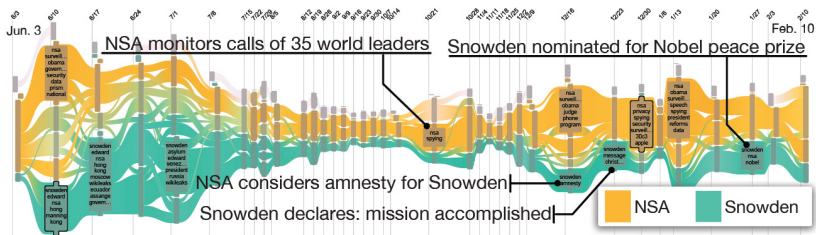
<http://scimaps.org>

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)

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.

Визуализация тематического разведочного поиска (концепт)

- Интерпретируемые оси: время–темы или сложность–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 400 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

Документ-запрос и результат тематического поиска (концепт)

Тематическая сегментация: структура документа-запроса

Дорожная карта: кластеризация релевантных документов

BigARTM

BigARTM — открытая библиотека для тематического моделирования больших коллекций текстовых документов. Параллельная распределенная реализация метода вероятностного тематического моделирования на основе аддитивной регуляризации.

ARTM (англ.) — Additive Regularization for Topic Modeling.

ARTM (рус.) — Аддитивная Регуляризация Тематических Модели.

Теоретическое введение

Основная идея: Тематическое моделирование

Вероятностное тематическое моделирование — это обобщенный инструмент статистического анализа текстов, предназначенный для выявления тематик в коллекциях документов. Тематическая модель описывает каждую тему дисперсным распределением на множестве термов, каждый документ — дисперсным распределением на множестве тем. Тематические модели используются для информационного поиска, классификации, категоризации, аннотирования, сегментации текстов.

Тематическая модель — это предположение наблюдений условного распределения $p(w|d)$ термов (слов или словосочетаний) w в документе d коллекции D :

$$p(w|d) = \sum_{t \in T} p(w|t) p(t|d),$$

где T — множество тем;

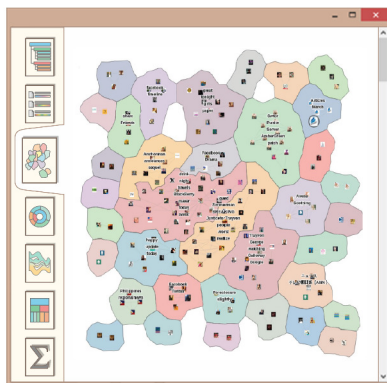
$\phi_{wt} = p(w|t)$ — неизвестное распределение термов в теме t ;

$\theta_{dt} = p(t|d)$ — неизвестное распределение тем в документе d .

Параметры тематической модели — матрица $\Phi = (\phi_{wt})_{w \in W, t \in T}$ задает пути решения задачи максимизации правдоподобия

$$\sum_{d \in D} \sum_{w \in W} n_{dw} \log \sum_{t \in T} \phi_{wt} \theta_{dt} \rightarrow \max_{\Phi, \Theta},$$

при ограничениях нормировки и неотрицательности



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

Тематическая модель — ключевое звено разведочного поиска.

Теория ARTM позволяет комбинировать тематические модели и строить композитные модели с требуемыми свойствами.

Тематическая модель для разведочного поиска должна быть...

- 1 Интерпретируемая: каждая тема понятна людям
- 2 Мультиграммная: термины-словосочетания неразрывны
- 3 Мультимодальная: авторы, связи, теги, пользователи, ...
- 4 Мультиязычная: для кросс- и много-языкового поиска
- 5 Иерархическая: выявление иерархических связей тем
- 6 Динамическая: прослеживание истории развития 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

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

Поиск тематически близких документов

$q = (w_1, \dots, w_{n_q})$ — 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

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

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

Поисковик MapReduce

Поисковик **MapReduce** – программа поиска (**Байду**) вычислено распределенно: вычисления для больших объемов данных в рамках параллельных шардов, представляющих собой набор Java-классов и исполняемых узлов для создания и обработки данных на параллельной обработке.

Основные компоненты MapReduce можно сформулировать как:

- обработка вычисленными большими объемами данных;
- масштабируемость;
- автоматическое распределение заданий;
- работа по минимальным объемам данных;
- автоматическая обработка отказов вычислений заданий.

MapReduce – популярная программная платформа (**Байду**, **Амазон**) построена распределенными приложениями для массово-параллельной обработки (**раздел**, **раздел**, **раздел**, **МР**) данных.

MapReduce включает в себе следующие компоненты:

1. **HDFS** – распределенная файловая система;

2. **MapReduce** – программная платформа (**Байду**) вычислено распределенно: вычисления для больших объемов данных в рамках параллельных шардов.

Ключевые, **значение** и архитектура **MapReduce** и структура **MapReduce** стали популярной реляционной моделью вычислений, в том числе и в качестве точки отказа. Число, в конечном итоге, определенное ограниченными платформами **MapReduce** и числом K последующих можно отметить:

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

Сильная связность **MapReduce** распределенно вычислений и элементов вычисления, реализованных распределенно алгоритмом. Как следствие:

Отсутствие поддержки альтернативной программной модели вычислений распределенно вычислений в **MapReduce** поддерживается только модель вычислений шардов.

Модель вычислений, точки отказа и как следствие, масштабируемость вычислений в среде с высокими требованиями к надежности.

Проблема **MapReduce** совместности: требование по единственному объекту вычисления всех вычислительных узлов кластера при обилии платформ **MapReduce** (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 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

Оценивание качества поиска

Precision — доля релевантных среди найденных

Recall — доля найденных среди релевантных

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

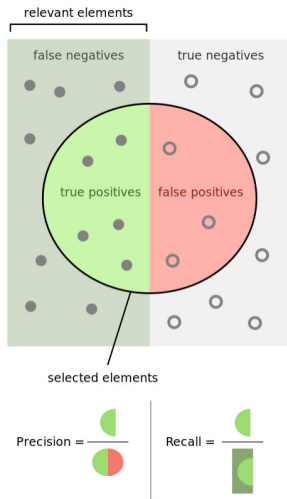
$$R = \frac{TP}{TP + FN} \text{ — полнота, (recall)}$$

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

TP (true positive) — найденные релевантные

FP (false positive) — найденные нерелевантные

FN (false negative) — ненайденные релевантные



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 $p(t|d)$:

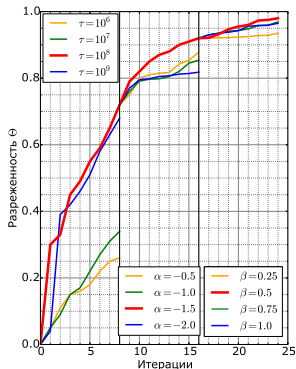
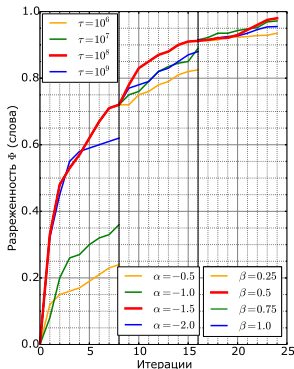
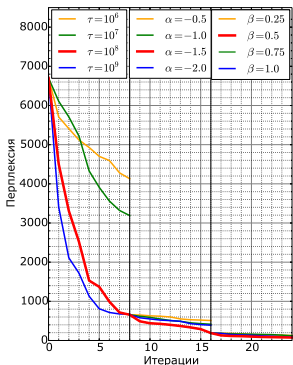
$$R(\Theta) = -\alpha \sum_{d,t} \ln \theta_{td}$$

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

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

Последовательный подбор коэффициентов регуляризации

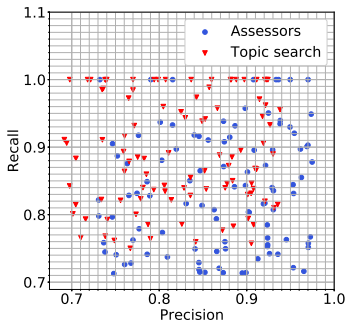
- декоррелирование распределений терминов в темах (τ),
- разреживание распределений тем в документах (α),
- сглаживание распределений терминов в темах (β).



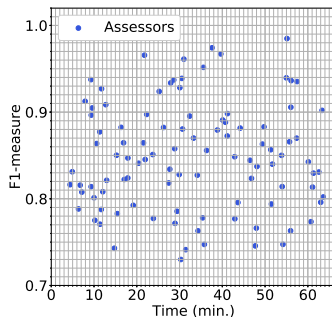
Результаты измерения точности и полноты по запросам

100 запросов, 3 ассессора на запрос

точность и полнота поиска



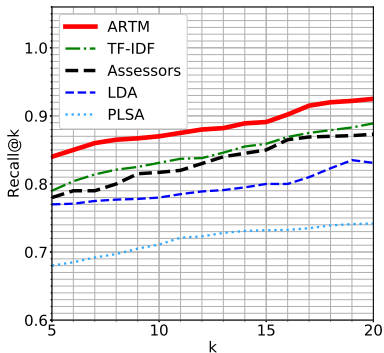
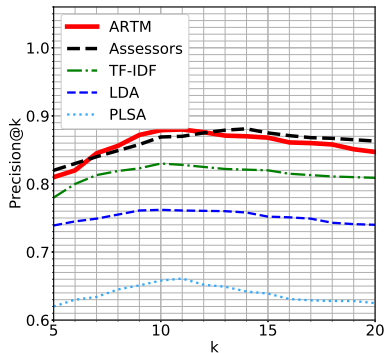
время и F_1 -мера (ассессоры)



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

Сравнение с ассессорами по качеству поиска

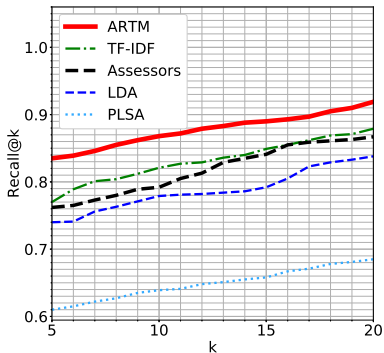
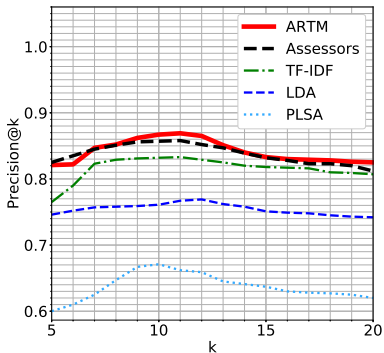
Точность и полнота по первым k позициям поисковой выдачи (коллекция Nabrahabr.ru)



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

Сравнение с ассессорами по качеству поиска

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



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

Влияние меры близости документа и запроса на качество поиска

Меры близости распределений:

Euclidean, Cosine, Manhattan, Hellinger, Kullback–Leibler

	Коллекция Habrahabr.ru					Коллекция TechCrunch.com				
	Eu	cos	Ma	He	KL	Eu	cos	Ma	He	KL
Prec@5	0.612	0.810	0.682	0.709	0.721	0.635	0.819	0.673	0.732	0.715
Prec@10	0.657	0.879	0.697	0.735	0.749	0.665	0.867	0.683	0.752	0.732
Prec@15	0.627	0.868	0.635	0.727	0.711	0.643	0.833	0.642	0.742	0.724
Prec@20	0.619	0.847	0.627	0.728	0.707	0.638	0.825	0.638	0.729	0.708
Recall@5	0.672	0.840	0.692	0.721	0.803	0.658	0.835	0.669	0.733	0.775
Recall@10	0.682	0.870	0.707	0.775	0.856	0.671	0.868	0.682	0.753	0.787
Recall@15	0.705	0.891	0.725	0.791	0.878	0.715	0.890	0.708	0.785	0.809
Recall@20	0.703	0.925	0.732	0.812	0.888	0.712	0.919	0.715	0.808	0.812

- Наилучшее качество поиска — при косинусной мере
- Одни и те же ассессорские оценки можно использовать для оценивания новых моделей и поисковых движков

Влияние комбинаций регуляризаторов на качество поиска

Декоррелирование, Θ-разреживание, Φ-сглаживание

	Коллекция Habrahabr.ru				Коллекция TechCrunch.com			
	$R = 0$	Д	ДΘ	ДΘΦ	$R = 0$	Д	ДΘ	ДΘΦ
Prec@5	0.628	0.748	0.771	0.810	0.652	0.775	0.779	0.819
Prec@10	0.653	0.776	0.812	0.879	0.679	0.787	0.819	0.867
Prec@15	0.642	0.765	0.792	0.868	0.669	0.773	0.798	0.833
Prec@20	0.643	0.759	0.783	0.847	0.673	0.777	0.792	0.825
Recall@5	0.692	0.784	0.805	0.840	0.673	0.812	0.812	0.835
Recall@10	0.714	0.814	0.834	0.870	0.685	0.821	0.845	0.868
Recall@15	0.725	0.835	0.867	0.891	0.712	0.859	0.869	0.890
Recall@20	0.735	0.862	0.891	0.925	0.723	0.882	0.895	0.919

- Комбинирование регуляризаторов улучшает качество поиска,
- хотя исходно все регуляризаторы нацелены на улучшение интерпретируемости тем и не оптимизируют поиск явно

Влияние сочетания модальностей на качество поиска

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

	ассессоры	С	К	СБ	СБТХ	все
Prec@5	0.821	0.612	0.549	0.654	0.737	0.810
Prec@10	0.869	0.635	0.568	0.701	0.752	0.879
Prec@15	0.875	0.625	0.532	0.685	0.682	0.868
Prec@20	0.863	0.616	0.533	0.682	0.687	0.847
Recall@5	0.780	0.722	0.636	0.797	0.827	0.840
Recall@10	0.817	0.744	0.648	0.812	0.875	0.870
Recall@15	0.850	0.778	0.677	0.842	0.893	0.891
Recall@20	0.873	0.803	0.685	0.852	0.898	0.925

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

Влияние сочетания модальностей на качество поиска

Коллекция TechCrunch.com. Число тем $|T| = 450$.

Модальности: Слова, Категории, Биграмммы, Авторы.

	асессоры	С	К	СБ	СБК	все
Prec@5	0.822	0.711	0.557	0.767	0.808	0.819
Prec@10	0.851	0.721	0.581	0.783	0.818	0.867
Prec@15	0.835	0.733	0.594	0.793	0.833	0.833
Prec@20	0.813	0.727	0.566	0.772	0.822	0.825
Recall@5	0.762	0.752	0.657	0.775	0.825	0.835
Recall@10	0.792	0.776	0.669	0.808	0.855	0.868
Recall@15	0.835	0.782	0.684	0.825	0.877	0.890
Recall@20	0.867	0.825	0.702	0.837	0.901	0.919

- Наилучшее качество поиска — по всем модальностям
- Наиболее полезные модальности — слова и категории

Влияние числа тем на качество поиска

Коллекция Nabrabr.ru

Используем все 5 модальностей, меняем $|T|$

	ассессоры	100	150	200	250	400
Prec@5	0.821	0.662	0.721	0.810	0.761	0.693
Prec@10	0.869	0.761	0.812	0.879	0.825	0.673
Prec@15	0.875	0.733	0.795	0.868	0.791	0.651
Prec@20	0.863	0.724	0.795	0.847	0.792	0.642
Recall@5	0.780	0.732	0.807	0.840	0.821	0.721
Recall@10	0.817	0.771	0.843	0.870	0.851	0.751
Recall@15	0.850	0.824	0.895	0.891	0.871	0.773
Recall@20	0.873	0.857	0.905	0.925	0.892	0.771

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

Влияние числа тем на качество поиска

Коллекция TechCrunch.com

Используем все 4 модальности, меняем $|T|$

	ассессоры	350	400	450	475	500
Prec@5	0.822	0.653	0.725	0.752	0.819	0.777
Prec@10	0.851	0.663	0.732	0.762	0.867	0.811
Prec@15	0.835	0.682	0.743	0.787	0.833	0.793
Prec@20	0.813	0.650	0.743	0.773	0.825	0.793
Recall@5	0.762	0.731	0.762	0.793	0.835	0.817
Recall@10	0.792	0.763	0.793	0.812	0.868	0.855
Recall@15	0.835	0.782	0.807	0.855	0.890	0.882
Recall@20	0.867	0.792	0.823	0.862	0.919	0.903

- Наилучшее качество поиска — при 475 темах
- Тематический поиск превосходит ассессоров по полноте

Выводы

- Регуляризаторы, улучшающие интерпретируемость модели, повышают также и качество поиска
- Тщательный подбор траектории регуляризации важен для повышения качества поиска
- Только при тщательной оптимизации тематический поиск превосходит как ассессоров, так и конкурирующие модели
- Ассессорские данные относятся не к темам, а к коллекции; поэтому их можно собрать один раз и использовать для оценивания многих моделей
- Ассессорских данных хватает для оценивания моделей, поскольку тематическая модель обучается *без учителя*

Стандартная методика оценивания моделей языка

Перplexия тестовой коллекции D' (hold-out perplexity):

$$\mathcal{P}(D') = \exp \left(- \frac{\sum_{d \in D'} \sum_{w \in d''} n_{dw} \ln p(w|d)}{\sum_{d \in D'} \sum_{w \in d''} n_{dw}} \right),$$

$d = d' \sqcup d''$ — случайное разбиение тестового документа на две половины равной длины;

параметры ϕ_{wt} оцениваются по обучающей коллекции D ;

параметры θ_{td} оцениваются по первой половине d' ;

перplexия вычисляется по второй половине d'' .

Интерпретации перplexии:

- 1) $\mathcal{P}(D') \rightarrow |d''|$ при $n \rightarrow \infty$, если слова равновероятны;
- 2) насколько хорошо мы предсказываем слова в документах (чем меньше перplexия, тем лучше).

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.

Оценки разреженности темы

- Разреженность:
 - доля нулевых элементов в Φ
 - доля нулевых элементов в Θ
- Характеристики различности тем:
 - размер ядра темы: $|W_t|$, ядро $W_t = \{w : p(t|w) > 0.25\}$
 - чистота темы: $\sum_{w \in W_t} p(w|t)$
 - контрастность темы: $\frac{1}{|W_t|} \sum_{w \in W_t} p(t|w)$
- Вырожденность тематической модели:
 - доля фоновых слов: $\frac{1}{n} \sum_{d \in D} \sum_{w \in d} \sum_{t \in B} p(t|d, w)$
 - доля нетематичных документов: $\frac{1}{|D|} \sum_{d \in D} \left[\sum_{t \in B} p(t|d) > 0.95 \right]$
 - доля нетематичных терминов: $\frac{1}{|W|} \sum_{w \in W} \left[\sum_{t \in B} p(t|w) > 0.95 \right]$

Разреживание, сглаживание, декоррелирование, отбор тем

M-шаг при комбинировании 6 регуляризаторов:

$$\phi_{wt} = \text{norm}_w \left(n_{wt} + \tau_1 \underbrace{\beta_w[t \in B]}_{\substack{\text{сглаживание} \\ \text{фоновых} \\ \text{тем}}} - \tau_2 \underbrace{\beta_w[t \in S]}_{\substack{\text{разреживание} \\ \text{предметных} \\ \text{тем}}} - \tau_3 \underbrace{\phi_{wt} \sum_{s \in S \setminus t} \phi_{ws}}_{\text{декоррелирование}} \right)$$

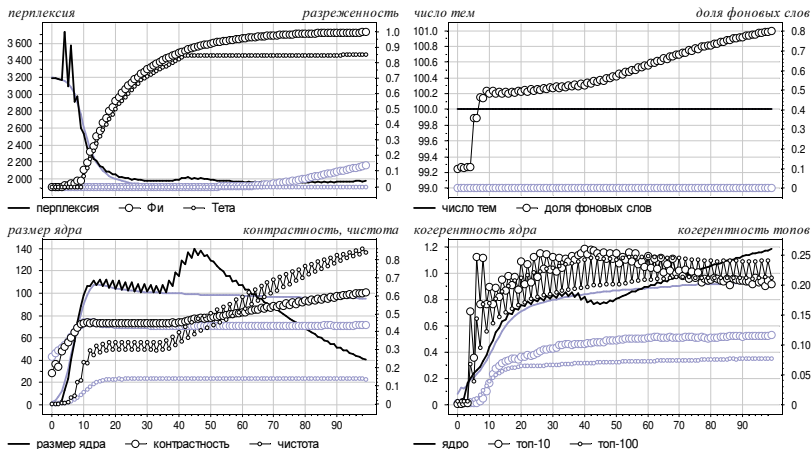
$$\theta_{td} = \text{norm}_t \left(n_{td} + \tau_4 \underbrace{\alpha_t[t \in B]}_{\substack{\text{сглаживание} \\ \text{фоновых} \\ \text{тем}}} - \tau_5 \underbrace{\alpha_t[t \in S]}_{\substack{\text{разреживание} \\ \text{предметных} \\ \text{тем}}} - \tau_6 \underbrace{\frac{n_d}{n_t} \theta_{td}}_{\text{удаление} \\ \text{малых тем}} \right)$$

Данные: статьи NIPS (Neural Information Processing System)
 $|D| = 1566$ статей, $n = 2.3$ М, $|W| = 13$ К,
 контрольная коллекция: $|D'| = 174$.

Vorontsov K. V., Potapenko A. A. Tutorial on Probabilistic Topic Modeling: Additive Regularization for Stochastic Matrix Factorization. AIST'2014.

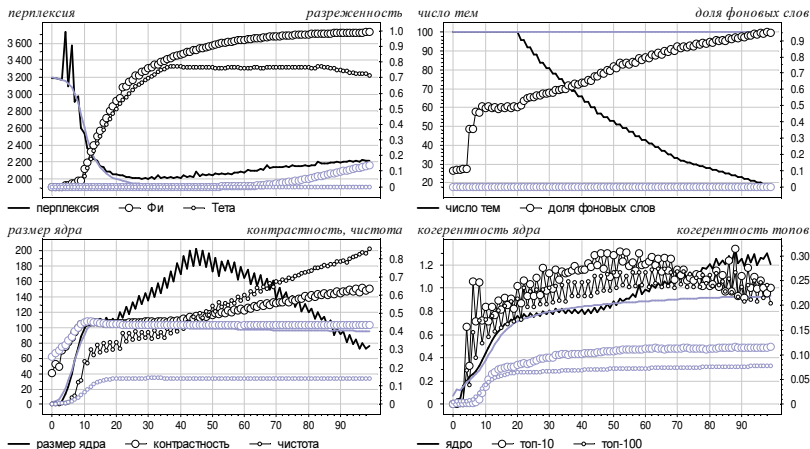
Разреживание, сглаживание, декоррелирование

Зависимости критериев качества от итераций EM-алгоритма
(серый — PLSA, чёрный — ARTM)



Те же регуляризаторы, плюс отбор тем

Зависимости критериев качества от итераций EM-алгоритма
(серый — PLSA, чёрный — ARTM)



Выводы

Одновременное улучшение многих критериев качества:

- разреженность выросла от 0 до 95%–98%
- когерентность тем выросла от 0.1 до 0.3
- чистота тем выросла от 0.15 до 0.8
- контрастность тем выросла от 0.4 до 0.6
- почти без потери перплексии (правдоподобия) модели

Подобраны траектории регуляризации:

- разреживание включать постепенно после 10-20 итераций
- сглаживание включать сразу
- декоррелирование включать сразу и как можно сильнее
- сокращение числа тем включать постепенно,
- не совмещая с декоррелированием на одной итерации

Регуляризатор для сокращения числа тем

Цель: избавиться от «мелких» незначимых тем.

Разрезаем распределение $p(t) = \sum_d p(d)\theta_{td}$, максимизируя кросс-энтропию между $p(t)$ и равномерным распределением:

$$R(\Theta) = -\tau \sum_{t \in S} \ln \sum_{d \in D} p(d)\theta_{td} \rightarrow \max.$$

Подставляем, получаем:

$$\theta_{td} = \operatorname{norm}_{t \in T} \left(n_{td} - \tau \frac{n_d}{n_t} \theta_{td} \right), \text{ вариант: } \theta_{td} = \operatorname{norm}_{t \in T} \left(n_{td} \left(1 - \frac{\tau}{n_t} \right) \right).$$

Эффект: обнуляются строки матрицы Θ с малыми n_t , заодно удаляются зависимые и расщеплённые темы.

Vorontsov K. V., Potapenko A. A., Plavin A. V. Additive Regularization of Topic Models for Topic Selection and Sparse Factorization. SLDS 2015.

Эксперименты с регуляризатором отбора тем

Коллекция статей NIPS (Neural Information Processing System)

- $|D| = 1566$ обучающих документов; $|D'| = 174$ тестовых
- $|W| = 13\text{ K}$ — мощность словаря

Синтетическая коллекция:

- строим PLSA за 500 итераций, $|T_0| = 50$ тем на NIPS
- генерируем (n_{dw}^0) из полученных Φ и Θ :

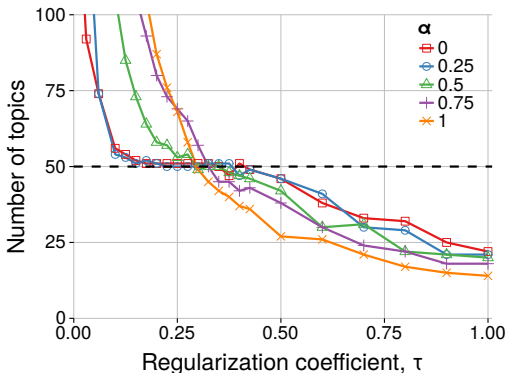
$$n_{dw}^0 = n_d \sum_{t \in T} \phi_{wt} \theta_{td}$$

Параметрическое семейство полусинтетических данных:

- n_{dw}^α — смесь синтетических данных n_{dw}^0 и реальных n_{dw} :

$$n_{dw}^\alpha = \alpha n_{dw} + (1 - \alpha) n_{dw}^0$$

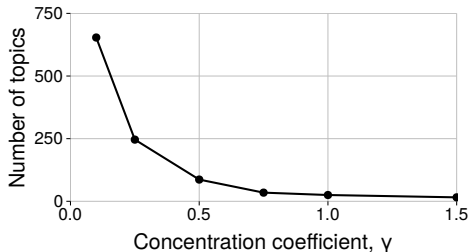
Попытка определения числа тем



- На синтетических данных надёжно находим $|T| = 50$,
- в широком интервале значений коэффициента τ ;
- однако на реальных данных нет столь чёткого интервала.

Сравнение с байесовской тематической моделью HDP

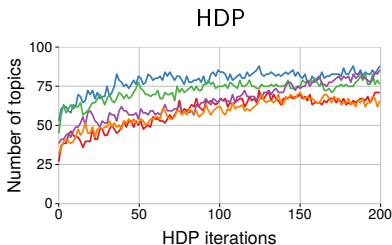
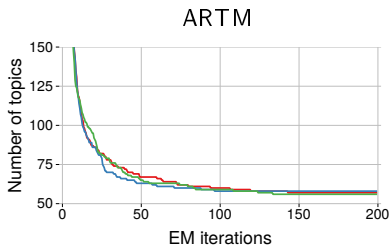
HDP, Hierarchical Dirichlet Process [Teh et.al, 2006] —
«state-of-the-art» байесовский подход к определению числа тем



- Коэффициент концентрации γ в HDP влияет на $|T|$ так же сильно, как выбор коэффициента τ в ARTM.

Сравнение ARTM и HDP по устойчивости

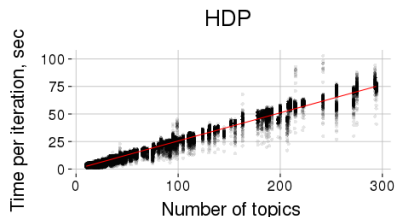
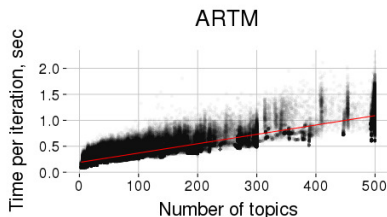
Запуск ARTM и HDP много раз из случайных инициализаций:



- HDP менее устойчив, причём в двух смыслах:
 - число тем сильнее флуктуирует от итерации к итерации;
 - результаты нескольких запусков различаются сильнее.
- «Рекомендуемые» значения параметров γ в HDP и τ в ARTM дают примерно равное число тем $|T| \approx 60$

Сравнение ARTM и HDP по времени вычислений

Сравнение времени одного прохода коллекции (sec)

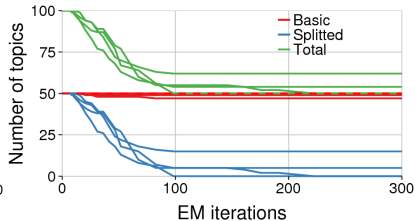
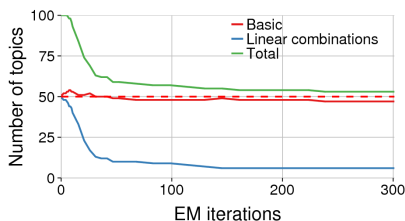


- ARTM в 100 раз быстрее!

Vorontsov K. V., Potapenko A. A., Plavin A. V. Additive regularization of topic models for topic selection and sparse factorization. SLDS 2015.

Удаление линейно зависимых и расщеплённых тем

Добавили 50 линейных комбинаций тем в модельную Φ .
Расщепили 50 тем, каждую на две подтемы в модельной Φ .



- Удаляются линейно зависимые и расщеплённые темы
- Остаются более различные темы исходной модели.

Vorontsov K. V., Potapenko A. A., Plavin A. V. Additive regularization of topic models for topic selection and sparse factorization. SLDS 2015.

Выводы

- Регуляризатор отбора тем — для удаления незначимых, зависимых, расщеплённых тем.
- Оптимального числа тем вообще не существует! Оно задаётся исходя из целей моделирования.
- Есть простой метод для удаления лишних тем, но пока в ARTM нет простых критериев добавления тем.
- **Открытая проблема:** почему этот регуляризатор удаляет линейно зависимые и расщеплённые темы?