

# Вероятностные тематические модели

## Лекция 4. Тематические иерархии и разведочный информационный поиск

К. В. Воронцов  
vokov@forecsys.ru

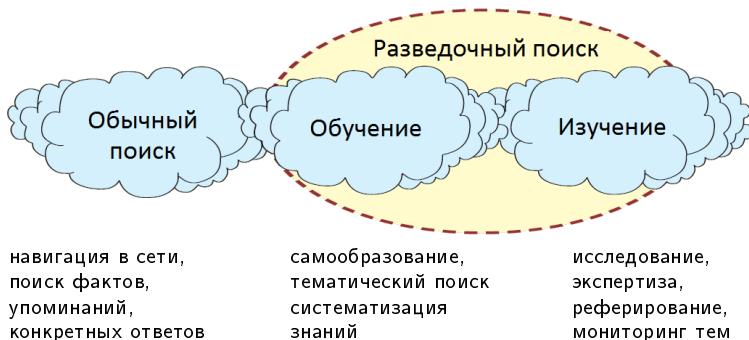
Этот курс доступен на странице вики-ресурса  
<http://www.MachineLearning.ru/wiki>  
«Вероятностные тематические модели (курс лекций, К.В.Воронцов)»

ВМК МГУ • 7 марта 2019

- 1 Разведочный информационный поиск**
  - Концепция разведочного поиска
  - Концепция distant reading и идеи визуализации
  - Сценарии использования разведочного поиска
- 2 Иерархические тематические модели**
  - Визуализация тематических иерархий
  - Метод нисходящего послойного построения иерархии
  - Спектр тем
- 3 Эксперименты с тематическим поиском**
  - Методика измерения качества поиска
  - Тематическая модель для документного поиска
  - Оптимизация гиперпараметров

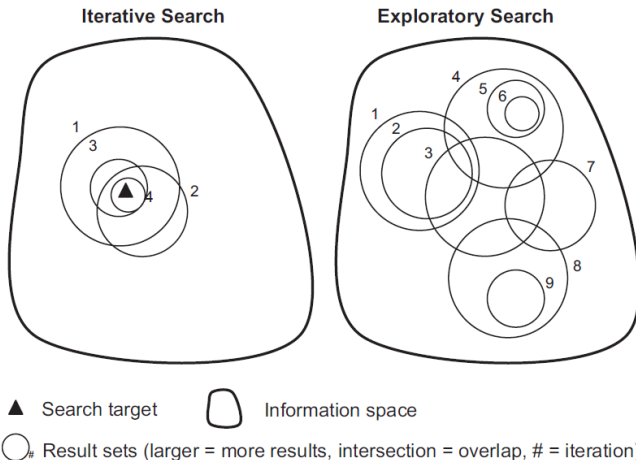
## Концепция разведочного поиска (exploratory search)

- 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»

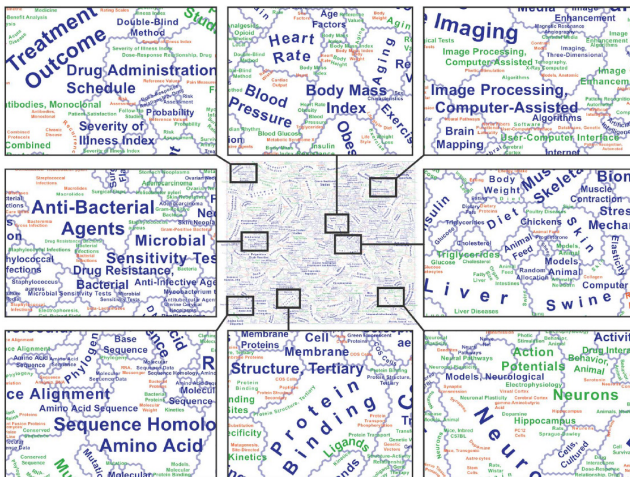
---

*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



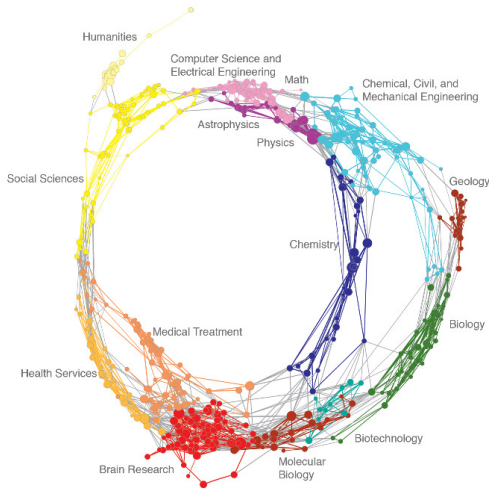
Skupin, Biberstine, Borner. Visualizing the Topical Structure of the Medical Sciences: A Self-Organizing Map Approach. PLoS ONE, 2013.







## Ещё один пример карты науки



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

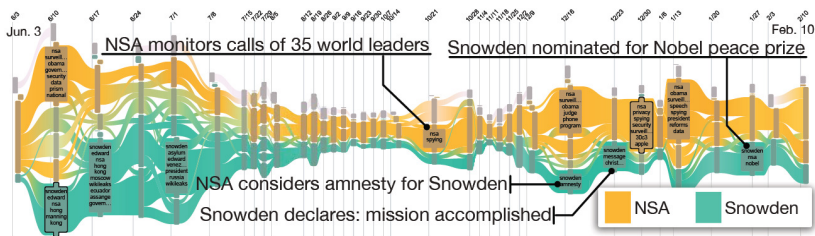
### Недостатки:

- оси не имеют интерпретации
- искажение сходства при двумерном проецировании

<http://scimaps.org>



## 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)

- Интерпретация осей: время–темы или сложность–темы
- Иерархичность: темы делятся на подтемы
- Спектр тем: гуманитарные → естественные → точные
- Интерактивность: реализация мантры Шнейдермана
- Суммаризация: на карте любого масштаба много текста



<http://textvis.lnu.se>

## Интерактивный обзор 430 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?
- какова тематическая структура предметной области?
- в каком порядке читать, чтобы лучше разобраться в теме?
- какие области являются смежными?
- какие фрагменты текста наиболее релевантны теме?

### Сценарий поиска:

- 1 имея текст под рукой,
- 2 получить карту содержащихся в нём тем-подтем
- 3 и карту предметной области в целом

# Документ-запрос и результат тематического поиска (концепт)

Тематическая сегментация: структура документа-запроса

Тематическая карта: кластеризация релевантных документов

The screenshot shows the BigARTM web interface. The main content area displays the title "Теоретическое введение" and a paragraph of text. Below the text, there is a mathematical formula for the joint probability distribution of topics and documents:

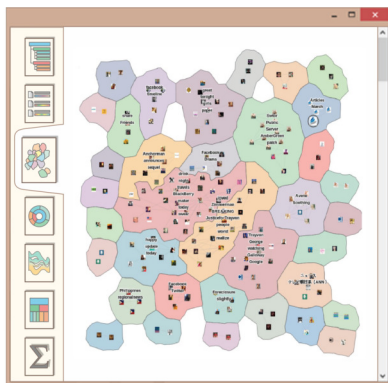
$$p(\theta, d) = \prod_{t \in T} p(\theta_t) p(d|t)$$

Below the formula, there is a definition of the unknowns and the parameters of the model:

где  $T$  — множество тем;  
 $\theta_t = p(\theta|t)$  — неизвестное распределение темов в теме  $t$ ;  
 $p(d|t) = p(d|\theta)$  — неизвестное распределение тем в документе  $d$ ;  
 Параметры тематической модели — матрица  $\Phi = (\phi_{wt}) \times \Theta = (\phi_{wt})$  — матрица нулей  
 решение задачи максимизации правдоподобия

$$\sum_{d \in D} \sum_{t \in T} \sum_{w \in V} \sum_{t \in T} \phi_{wt} \theta_t \rightarrow \max_{\Phi}$$

при ограничениях: нормировка и неотрицательность



## 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

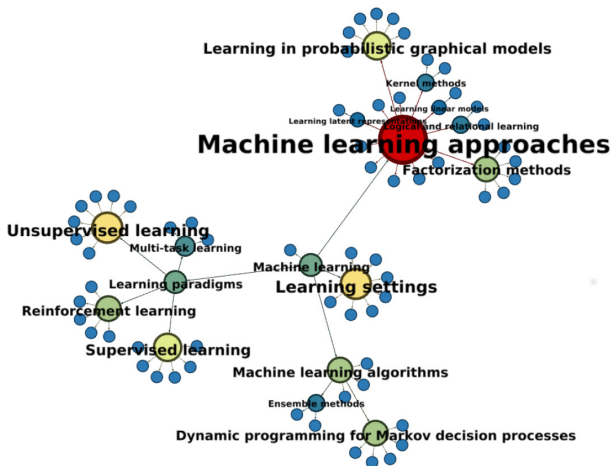
ARTM позволяет строить тематические модели с заданным набором требуемых свойств



## Тематическая модель для разведочного поиска должна быть...

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

## Пример древовидной тематической иерархии



Georgeta Bordea. Domain adaptive extraction of topical hierarchies for Expertise Mining. 2013.

## Пример тематической иерархии

Тексты научно-просветительского ресурса Postnauka.ru:  
2976 документов, 43196 слов, 1799 тэгов



Chirkova N.A., Vorontsov K.V. Additive regularization for hierarchical multimodal topic modeling. JMLDA, 2016.

Belyy A.V., Seleznova M.S., Sholokhov A.K., Vorontsov K.V. Quality Evaluation and Improvement for Hierarchical Topic Modeling. Dialogue 2018.

## Иерархические тематические модели

- структура иерархии: дерево / **многодольный граф**
- направление: снизу вверх / **сверху вниз** / одновременно
- наращивание: попершинное / **послойное**

### Открытые проблемы:

- “Despite recent activity in the field of HPTMs, determining the hierarchical model that best fits a given data set, in terms of the structure and size of the learned hierarchy, still remains a challenging task and an open issue.”
- “The evaluation of hierarchical PTMs is also an open issue.”

---

*Zavitsanos E., Paliouras G., Vouros G. A. Non-Parametric Estimation of Topic Hierarchies from Texts with Hierarchical Dirichlet Processes. 2011.*

## Регуляризатор $\Phi$ : родительские темы как псевдо-документы

Шаг 1. Строим модель с небольшим числом тем.

Шаг  $k$ . Пусть модель с множеством тем  $T$  уже построена.  
Строим множество дочерних тем  $S$  (subtopics),  $|S| > |T|$ .

Родительские темы приближаются смесями дочерних тем:

$$\sum_{t \in T} n_t \text{KL}_w \left( p(w|t) \parallel \sum_{s \in S} p(w|s)p(s|t) \right) \rightarrow \min_{\Phi, \tilde{\Psi}}$$

где  $\tilde{\Psi} = (\tilde{\psi}_{st})_{S \times T}$  — матрица связей,  $\tilde{\psi}_{st} = p(s|t)$ .

Родительская  $\Phi^P \approx \Phi \tilde{\Psi}$ , отсюда регуляризатор матрицы  $\Phi$ :

$$R(\Phi, \tilde{\Psi}) = \tau \sum_{t \in T} \sum_{w \in W} n_{wt} \ln \sum_{s \in S} \phi_{ws} \tilde{\psi}_{st} \rightarrow \max.$$

Родительские темы  $t$  — «документы» с частотами термов  $n_{wt}$ .

## Регуляризатор $\Theta$ : родительские темы как модальность

**Шаг 1.** Строим модель с небольшим числом тем.

**Шаг  $k$ .** Пусть модель с множеством тем  $T$  уже построена.  
Строим множество дочерних тем  $S$  (subtopics),  $|S| > |T|$ .

Родительские темы приближаются смесями дочерних тем:

$$\sum_{d \in D} n_d \text{KL}_t \left( p(t|d) \parallel \sum_{s \in S} p(t|s)p(s|d) \right) \rightarrow \min_{\Theta, \Psi}$$

где  $\Psi = (\psi_{ts})_{T \times S}$  — матрица связей,  $\psi_{ts} = p(t|s)$ .

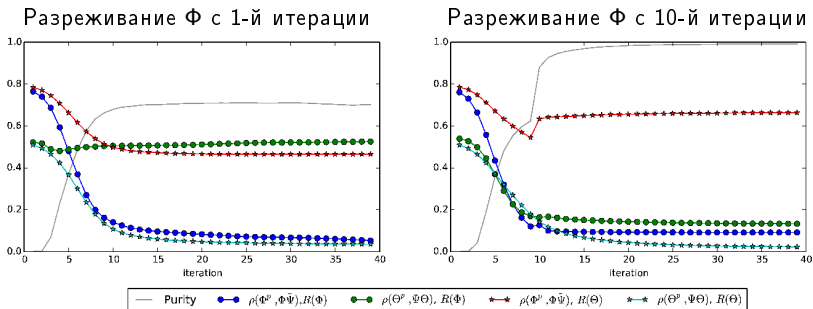
Родительская  $\Theta^p \approx \Psi \Theta$ , отсюда регуляризатор матрицы  $\Theta$ :

$$R(\Theta, \Psi) = \tau \sum_{d \in D} \sum_{t \in T} n_{td} \ln \sum_{s \in S} \psi_{ts} \theta_{sd} \rightarrow \max.$$

Родительские темы  $t$  — модальность с частотами термов  $n_{td}$ .

## Эксперимент на коллекции ММРО-ИОИ

Среднее расстояние Хеллингера  $\rho(\Phi^P, \Phi\tilde{\Psi})$  и  $\rho(\Theta^P, \Psi\Theta)$  для регуляризаторов  $\Phi$  и  $\Theta$  при переходе между уровнями  $1 \rightarrow 2$ :

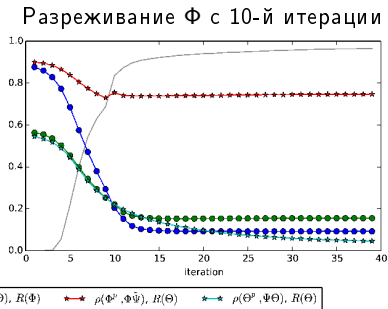
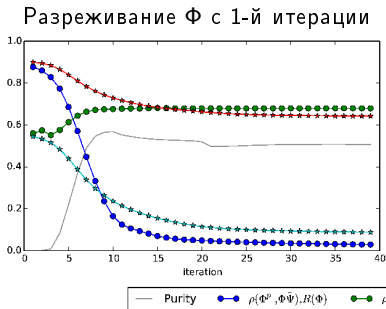


**Вывод.** Регуляризатор  $\Theta$  плохо приближает  $\Phi^P$ .

*Chirkova N. A., Vorontsov K. V. Additive regularization for hierarchical multimodal topic modeling. JMLDA, 2016.*

## Эксперимент на коллекции ММРО-ИОИ

Среднее расстояние Хеллингера  $\rho(\Phi^P, \Phi\tilde{\Psi})$  и  $\rho(\Theta^P, \Psi\Theta)$  для регуляризаторов  $\Phi$  и  $\Theta$  при переходе между уровнями  $2 \rightarrow 3$ :



**Вывод.** Регуляризатор  $\Theta$  плохо приближает  $\Phi^P$ .

Chirkova N. A., Vorontsov K. V. Additive regularization for hierarchical multimodal topic modeling. JMLDA, 2016.



## Выводы

- Регуляризатор  $\Phi$  приближает  $\Phi^P \approx \Phi\tilde{\Psi}$  и  $\Theta^P \approx \Psi\Theta$
- Регуляризатор  $\Theta$  приближает только  $\Theta^P \approx \Psi\Theta$
- Сильное разреживание  $\psi_{ts} \in \{0, 1\}$  даёт иерархию–дерево
- Нельзя допускать вырождения  $\psi_{ts} = p(t|s) \equiv 0$

### Дальнейшие задачи:

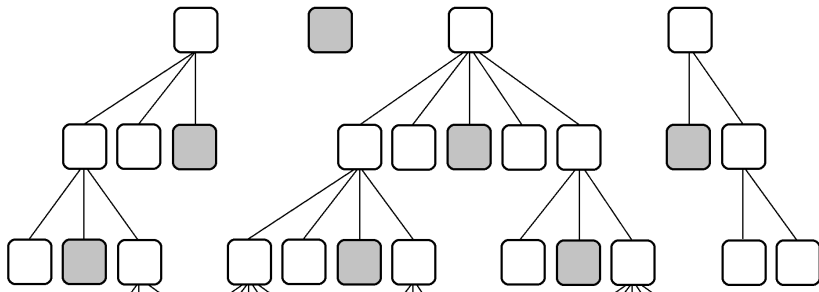
- Согласованная регуляризация:  $\tilde{\psi}_{st}p(t) = \psi_{ts}p(s)$

$$\tau_1 \sum_{t,w} n_{wt} \ln \sum_s \phi_{ws} \psi_{ts} \frac{n_s}{n_t} + \tau_2 \sum_{d,t} n_{td} \ln \sum_s \psi_{ts} \theta_{sd} \rightarrow \max_{\Phi, \Psi, \Theta}$$

- Иерархии с темами различной глубины:
  - наращивание уровня для подмножества  $T' \subseteq T$
  - критерий неоднородности темы для включения её в  $T'$

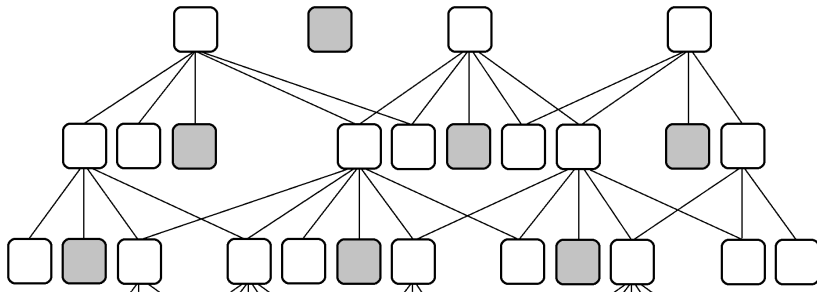
## Иерархии с темами различной глубины

- На каждом уровне расщепляются не все темы (допускается вырожденность:  $p(s|t) \equiv 0$  для некоторых  $t$ )
- Расщепляемая тема может иметь дочернюю фоновую, в которой собирается общая лексика родительской темы
- При максимальном разреживании  $p(t|s) \in \{0, 1\}$  иерархия является деревом (корень не показан)



## Иерархии с темами различной глубины

- На каждом уровне расщепляются не все темы (допускается вырожденность:  $p(s|t) \equiv 0$  для некоторых  $t$ )
- Расщепляемая тема может иметь дочернюю фоновую, в которой собирается общая лексика родительской темы
- При умеренном разреживании  $p(t|s)$  у вершины может быть несколько родителей (корень не показан)



## Иерархии с темами различной глубины

След документа в тематической иерархии определяет степень его специализации, назначение, аудиторию



узко специализированный,  
для профессионалов



междисциплинарное исследование,  
для профессионалов



обзорный,  
для ознакомления с предметной областью



популярный или энциклопедический,  
для расширения кругозора

## Способы оценивания качества тематических иерархий

- *Перплексия* или правдоподобие: приводит ли постепенное дробление тем к более точному разложению
- *Полезность*: сколько шагов делает пользователь, чтобы найти документ по иерархии
- *Когерентность*: как часто слова темы и её подтемы совместно встречаются рядом в тексте
- *Метод интрузий*: правильно ли ассессоры определяют чужую тему, внедрённую в список дочерних тем
- *Сравнение с «золотым стандартом»*: насколько иерархия похожа на имеющуюся категоризацию документов

## Что такое «спектр тем» и зачем он нужен

Визуализация иерархии тем во времени (концепт):



- Интерпретируемые оси «время–темы»
- Близкие темы должны находиться рядом
- *Тематический спектр* — одномерная линейная проекция (например, науки: гуманитарные → естественные → точные)

## Построение спектра тем. Постановка задачи

*Тематический спектр* — такая перестановка тем  $t_1, \dots, t_{|T|}$ , что сумма расстояний между соседними темами минимальна:

$$\sum_{i=2}^{|T|} \rho(t_i, t_{i-1}) \rightarrow \min$$

*Функция расстояния*  $\rho(t, t')$  между темами, примеры:

- Манхэттенское:  $\rho(t, t') = \sum_{w \in W} |\phi_{wt} - \phi_{wt'}|$
- Хеллингера:  $\rho^2(t, t') = \frac{1}{2} \sum_{w \in W} (\sqrt{\phi_{wt}} - \sqrt{\phi_{wt'}})^2$
- Жаккара:  $\rho(t, t') = 1 - \frac{|W_t \cap W_{t'}|}{|W_t \cup W_{t'}|}$ ,  $W_t = \{w : \phi_{wt} > \frac{1}{|W|}\}$

## Построение спектра тем — это задача коммивояжёра

### Задача TSP (traveling salesman problem)

Найти путь минимальной суммарной стоимости, соединяющий  $T$  городов так, чтобы в каждом городе побывать один раз.

Алгоритм Лина–Кернигана в реализации Хельсгауна — лучший для решения задачи TSP, по данным *Encyclopedia of operations research* на 2013 год.

Вычислительная сложность  $T^{2.2}$ .

Другие алгоритмы оказались не только медленнее, но и хуже по качеству тематических спектров.

---

*Keld Helsgaun*. An effective implementation of the Lin–Kernighan traveling salesman heuristic. EJOR, 2000.

*Дмитрий Федоряка*. Технология интерактивной визуализации тематических моделей. Бакалаврская диссертация. МФТИ, 2017.



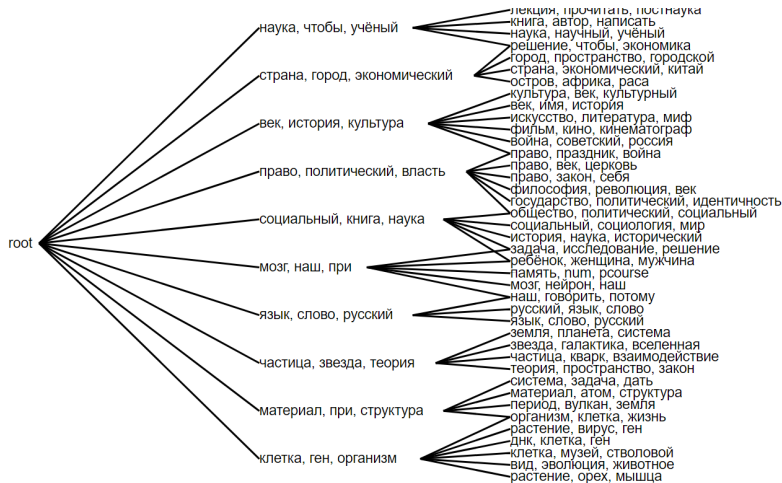
## Пример спектра (коллекция postnauka.ru)

1. остров, земля, период, там, территория, океан, где, более, вид, найти, вулкан, находится, южный
2. растение, япония, раса, при, более, чем, например, исследование, вид, страна, население
3. вид, эволюция, самец, мозг, самка, животное, отбор, ген, более, птица, наш, между, чтобы, чем, друг
4. мозг, нейрон, при, заболевание, наш, пациент, состояние, система, болезнь, сон, исследование
5. клетка, музей, стволовой, ткань, организм, чтобы, опухоль, система, использовать, технология
6. клетка, ген, днк, организм, молекула, геном, белок, белка, бактерия, система, процесс, жизнь
7. система, материал, задача, структура, метод, компьютер, дать, при, химический, область, химия
8. квантовый, свет, волна, атом, информация, фотон, сигнал, использовать, два, при, частота, состояние
9. частица, энергия, кварк, взаимодействие, магнитный, электрон, масса, физика, бозон, протон, модель
10. звезда, галактика, земля, планета, вселенная, дыра, чёрный, объект, солнце, масса, наш, система
11. теория, пространство, вселенная, закон, физика, математический, уравнение, число, два, мир, система
12. наш, сеть, информация, дать, объект, культура, задача, например, образ, память, слово, разный
13. язык, слово, русский, например, говорить, словарь, речь, разный, языковой, текст, два, лингвист
14. наука, учёный, научный, потому, чтобы, лекция, хороший, университет, сейчас, наш, заниматься
15. экономический, экономика, страна, чтобы, более, рынок, компания, цена, решение, деньга, работа, чем
16. страна, война, государство, политический, россия, советский, власть, политика, германия, статья
17. ребёнок, женщина, мужчина, жизнь, культура, общество, себя, семья, социальный, советский, женский
18. город, пространство, социальный, городской, общество, место, культурный, жизнь, более, современный
19. исследование, социальный, поведение, группа, решение, and, the, теория, проблема, наука
20. социальный, социология, мир, теория, объект, социологический, действие, событие, социолог, наука
21. политический, философия, идея, наука, свобода, понятие, революция, история, философ, век, себя
22. право, власть, закон, король, век, римский, бог, себя, церковь, правовой, политический, суд, два
23. век, история, русский, исторический, имя, традиция, христианский, культура, историк, текст, уже
24. себя, искусство, литература, говорить, потому, мир, сам, миф, жизнь, слово, текст, роман, век
25. книга, фильм, автор, кино, rcourse, num, читатель, посвятить, тема, история, исследование, работа

## Пример спектра (коллекция lenta.ru)

1. спортсмен, допинг, олимпиада, рию, де, россия, проба, жанейро, wada, олимпийский\_игра, соревнование
2. команда, матч, счёт, клуб, победа, чемпионат, турнир, минута, футболист, встреча, летний, футбол
3. евро, евровидение, страна, россия, конкурс, франция, болельщик, англяя, украина, футбол, певец
4. пройти, мероприятие, россия, акция, фестиваль, москва, фильм, участник, картина, театр, музей
5. фильм, сериал, продукт, актёр, компания, продукция, процент, россия, книга, товар, картина, сезон
6. россия, москва, турист, процент, россиянин, страна, отель, рейс, путешественник, город, тысяча
7. процент, доллар, рубль, нефть, цена, россия, баррель, страна, уровень, вырасти, рынок, рост
8. компания, миллиард\_рубль, процент, миллиард\_доллар, россия, сумма, миллион\_доллар, банк, банка
9. закон, законопроект, документ, реклама, использование, деятельность, поправка, внести, организация
10. россия, страна, керченский\_пролив, российский, боинг, работа, чайка, ряд, гражданин, аэропорт
11. партия, кандидат, журналист, праймериза, выбор, единый\_россия, госдума, выборы
12. россия, украина, крым, решение, киев, депутат, вопрос, отношение, страна, мнение, право, москва
13. россия, страна, турция, сша, ес, евросоюз, москва, санкция, отношение, украина, вопрос, государство
14. россия, сирия, исламский\_государство, сша, нато, иго, запретить, террорист, страна, боевик
15. ракета, путин, россия, запуск, глава\_государство, союз, спутник, президент
16. учёный, клетка, исследование, исследователь, ген, университет, оказать, процент, помощь, организм
17. земля, животное, учёный, животный, тысяча, звезда, планета, обнаружить, кошка, территория, жизнь
18. самолёт, километр, машина, борт, пассажир, вертолёт, погибнуть, лайнер, пилот, час, район, яхта
19. полицейский, полиция, мужчина, задержать, автомобиль, улица, москва, пострадать, life
20. статья, убийство, задержать, суд, отношение, ук\_рф, подозревать, следствие, обвинять, трамп, часть
21. ребёнок, женщина, мужчина, летний, дом, сын, семья, мальчик, жена, полиция, дочь, школа, врач
22. видео, youtube, ролик, фото, фотография, канал, снимка, auto, instagram, девушка, страница, группа
23. facebook, пользователь, интернет, страница, twitter, пост, написать, соцсеть, вконтакте, аккаунт
24. устройство, смартфон, компания, мотоциклист, игра, байкер, видео, миллион\_доллар, робот, молодая
25. бренд, модель, компания, обувь, основать, одежда, релиз, коллекция, редакция, часы, поступить

## Иерархический спектр (коллекция postnauka.ru)



## Иерархический спектр (коллекция lenta.ru)



## Две коллекции новостей про технологии

### Habrahabr.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

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

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

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

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

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

Выдача тематического поиска —  $k$  первых документов.

Реализация: *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

### Поиск MapReduce

**Поиск MapReduce** – программа поиска (библиотека) написанная распределенно: выделены для больших объемов данных и разная параллельная обработка, представляющая собой набор Java-классов и исполняемых утилит для создания и обработки данных на распределенной обработке.

**Основные компоненты Поиска MapReduce** можно сформулировать как:

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

**Поиск** – популярная программная платформа (язык Java, библиотека) построена распределенными приложениями для массово-параллельной обработки (разные работы, ресурсы, CPU) данных.

**Поиск** включает в себе следующие компоненты:

1. HDFS – распределенная файловая система;

2. **Поиск MapReduce** – программная модель (библиотека) написанная распределенно: выделены для больших объемов данных и разная параллельная обработка;

Ключевые, объектные и архитектура **Поиска MapReduce** и структура HDFS, стали примером того, как можно сделать в своем компоненте, в том числе и единичные точки отказа. Это, в конечном итоге, определило ограничение платформ **Поиск** и целью. К последним можно отнести:

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

Сильная связность **Поиска** распределенно вычислений и элементов выделены распределенно распределенной алгоритмы. Как следствие:

Отсутствие поддержки альтернативной программы выделены распределенно вычислений: в **Поиск** v1.0 поддерживается только модель выделены распределенно;

Многие выделены точки отказа и как следствие, необходимость выделены в средстве с выделены требования к надежности;

Проблема совместности требований по единичному выделены выделены всех выделены узлов кластера при обновлении платформ **Поиск** (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

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

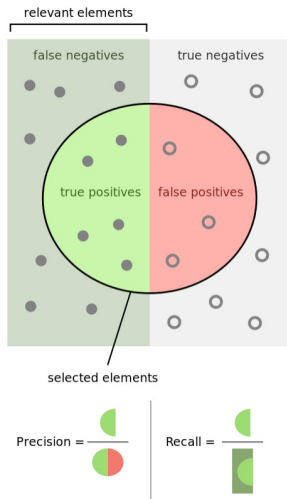
$$R = \frac{TP}{TP + FN} \text{ — полнота, (recall)}$$

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

TP (true positive) — найденные релевантные

FP (false positive) — найденные нерелевантные

FN (false negative) — не найденные релевантные



## Какие модели поиска сравнивались

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

Задачи регуляризаторов в ARTM и hARTM:

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

## Стратегия регуляризации

Последовательное применение трёх регуляризаторов

- 1 декоррелирование тем:

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

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

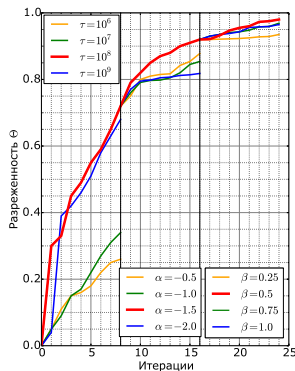
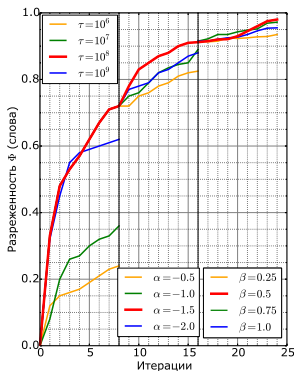
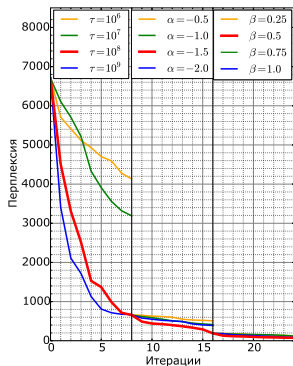
$$R(\Theta) = -\alpha \sum_{d,t} \ln \theta_{td}$$

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

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

## Последовательный подбор коэффициентов регуляризации

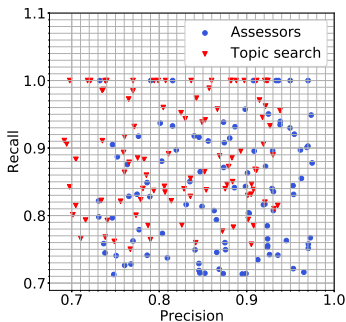
- декоррелирование распределений терминов в темах ( $\tau$ ),
- разреживание распределений тем в документах ( $\alpha$ ),
- сглаживание распределений терминов в темах ( $\beta$ ).



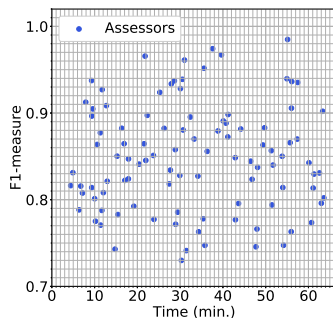
## Результаты измерения точности и полноты по запросам

100 запросов, 3 ассессора на запрос

точность и полнота поиска



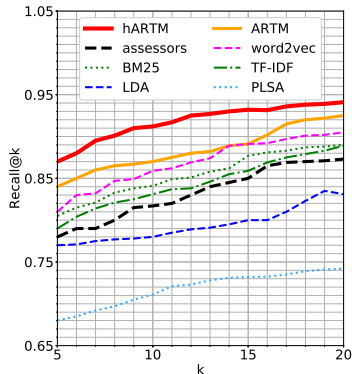
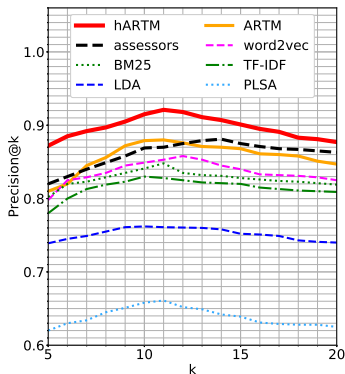
время и  $F_1$ -мера (ассессоры)



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

## Сравнение с ассессорами по качеству поиска

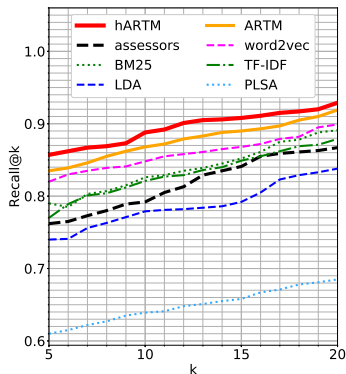
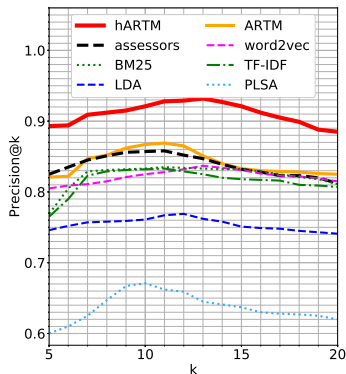
Точность и полнота по первым  $k$  позициям поисковой выдачи (коллекция Habrahabr.ru)



A. Ianina, K. Vorontsov. 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	<b>0.810</b>	0.682	0.709	0.721	0.635	<b>0.819</b>	0.673	0.732	0.715
Prec@10	0.657	<b>0.879</b>	0.697	0.735	0.749	0.665	<b>0.867</b>	0.683	0.752	0.732
Prec@15	0.627	<b>0.868</b>	0.635	0.727	0.711	0.643	<b>0.833</b>	0.642	0.742	0.724
Prec@20	0.619	<b>0.847</b>	0.627	0.728	0.707	0.638	<b>0.825</b>	0.638	0.729	0.708
Recall@5	0.672	<b>0.840</b>	0.692	0.721	0.803	0.658	<b>0.835</b>	0.669	0.733	0.775
Recall@10	0.682	<b>0.870</b>	0.707	0.775	0.856	0.671	<b>0.868</b>	0.682	0.753	0.787
Recall@15	0.705	<b>0.891</b>	0.725	0.791	0.878	0.715	<b>0.890</b>	0.708	0.785	0.809
Recall@20	0.703	<b>0.925</b>	0.732	0.812	0.888	0.712	<b>0.919</b>	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	<b>0.810</b>	0.652	0.775	0.779	<b>0.819</b>
Prec@10	0.653	0.776	0.812	<b>0.879</b>	0.679	0.787	0.819	<b>0.867</b>
Prec@15	0.642	0.765	0.792	<b>0.868</b>	0.669	0.773	0.798	<b>0.833</b>
Prec@20	0.643	0.759	0.783	<b>0.847</b>	0.673	0.777	0.792	<b>0.825</b>
Recall@5	0.692	0.784	0.805	<b>0.840</b>	0.673	0.812	0.812	<b>0.835</b>
Recall@10	0.714	0.814	0.834	<b>0.870</b>	0.685	0.821	0.845	<b>0.868</b>
Recall@15	0.725	0.835	0.867	<b>0.891</b>	0.712	0.859	0.869	<b>0.890</b>
Recall@20	0.735	0.862	0.891	<b>0.925</b>	0.723	0.882	0.895	<b>0.919</b>

- 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	<b>0.810</b>
Prec@10	0.869	0.635	0.568	0.701	0.752	<b>0.879</b>
Prec@15	0.875	0.625	0.532	0.685	0.682	<b>0.868</b>
Prec@20	0.863	0.616	0.533	0.682	0.687	<b>0.847</b>
Recall@5	0.780	0.722	0.636	0.797	0.827	<b>0.840</b>
Recall@10	0.817	0.744	0.648	0.812	0.875	<b>0.870</b>
Recall@15	0.850	0.778	0.677	0.842	0.893	<b>0.891</b>
Recall@20	0.873	0.803	0.685	0.852	0.898	<b>0.925</b>

- 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.

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

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

## Влияние числа тем на качество поиска

### Коллекция Nabrhabr.ru

Используем все 5 модальностей, меняем  $|T|$

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

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

## Влияние числа тем на качество поиска

### Коллекция TechCrunch.com

Используем все 4 модальности, меняем  $|T|$

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

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

## Выводы по результатам экспериментов

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

## Резюме

Разведочный информационный поиск (exploratory search):

- это поиск по смыслу, а не по ключевым словам
- может быть построен на тематическом моделировании
- требует многофункциональности от тематических моделей
- является одной из главных мотиваций для ARTM
- и, в частности, для иерархических моделей

### Открытые проблемы

- построение разноуровневых иерархий в ARTM
- оценивание качества тематических иерархий
- оптимизация числа тем на каждом уровне иерархии



- Научно-популярные статьи: ПостНаука, Элементы, Хабр
- Википедия
- Вики-227
- Новостной поток (RSS lenta.ru / нефильТРованный поток)
- Акты арбитражных судов РФ
- TechCrunch (английский)
- Открытые датасеты (английский): 20 newsgroups, NIPS, KOS
- Транзакции клиентов Sberbank DSD 2016

**Задача-минимум:** научиться решать задачи NLP и строить тематические модели в BigARTM

**Задача-максимум:** решить открытую проблему

- Несбалансированность и семантическая однородность тем
- Агрегирование гетерогенных коллекций
- Создание новых тем при расширении коллекции
- В том числе, создание новых тем в иерархиях
- Прослеживание тем в новостных потоках
- Тематическая сегментация и посегментный поиск
- Предобученные тематические векторные представления слов
- Визуализация «карт знаний» (D3.js)