Вероятностные тематические модели Лекция 5. Модальности, иерархии и тематический поиск

Kонстантин Вячеславович Воронцов k.v.vorontsov@phystech.edu

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

МФТИ.ФПМИ.ИС.ИАД — ФИЦ ИУ РАН • 2025-10-16

Содержание

- 1 Мультимодальные тематические модели
 - Мультимодальные тематические модели
 - Мультимодальный ЕМ-алгоритм
 - Примеры мультимодальных тематических моделей
- Иерархические тематические модели
 - Регуляризация тематических иерархий
 - Эксперименты с иерархическими моделями
 - Тематические спектры
- 3 Эксперименты с тематическим поиском
 - Методика измерения качества поиска
 - Тематическая модель для документного поиска
 - Оптимизация гиперпараметров

Напоминание. Задача тематического моделирования

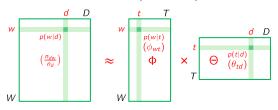
Дано: коллекция текстовых документов, $p(w|d) = rac{n_{dw}}{n_d}$

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:



Hofmann T. Probabilistic Latent Semantic Indexing. ACM SIGIR, 1999. Blei D., Ng A., Jordan M. Latent Dirichlet Allocation. NIPS-2001. JMLR 2003.

Напоминание. ARTM — аддитивная регуляризация

Максимизация \log правдоподобия с регуляризатором R:

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

ЕМ-алгоритм: метод простой итерации для системы уравнений

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

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

Воронцов К. В. Аддитивная регуляризация тематических моделей коллекций текстовых документов. Доклады РАН, 2014.

Hапоминание. Модель локального контекста (Attentive ARTM)

Дано: коллекция текстовых документов, w_1, \ldots, w_n C_i — локальный контекст (окружение) терма w_i α_{ui} — распределение весов термов u в контексте C_i

Найти: параметры ϕ_{wt} , p_t тематической модели

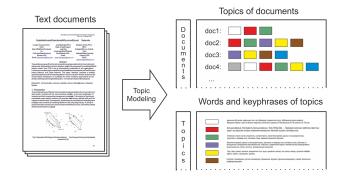
$$p(w|C_i) = \sum_{t \in T} \phi_{wt} \theta_{ti}, \quad \theta_{ti} = \sum_{u \in C_i} \alpha_{ui} \phi'_{tw}, \quad \phi'_{tw} = \underset{t \in T}{\operatorname{norm}} (\phi_{wt} p_t)$$

Критерий: максимум \log правдоподобия с регуляризатором R:

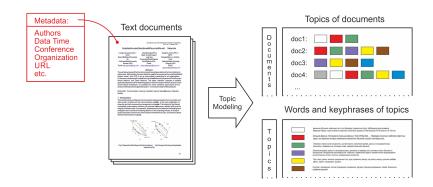
$$\sum_{i=1}^{n} \ln p(w|C_i) + R(\Phi) \to \max_{\Phi}, \quad R(\Phi) = \sum_{i=1}^{k} \tau_i R_i(\Phi)$$

ЕМ-алгоритм:
$$p_{ti} = \underset{t \in T}{\mathsf{norm}} \left(\phi_{w_i t} \theta_{ti}\right), \quad p_t = \frac{1}{n} \sum_{i=1}^n p_{ti}$$
 $\phi_{wt} = \underset{w \in \mathcal{W}}{\mathsf{norm}} \left(\sum_{i=1}^n [w_i = w] p_{ti} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}\right)$

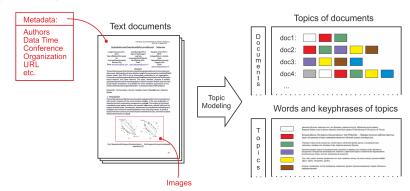
Тема может порождать термы различных модальностей: p(слово|t), p(n-грамма|t),



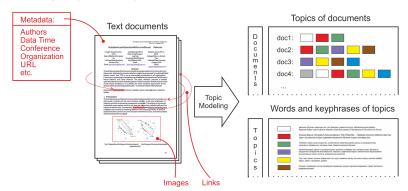
Тема может порождать термы различных модальностей: p(слово|t), p(n-rpamma|t), p(автор|t), p(время|t), p(источник|t),



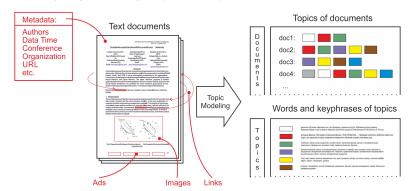
Тема может порождать термы различных модальностей: p(слово|t), p(n-грамма|t), p(автор|t), p(время|t), p(источник|t), p(объект|t),



Тема может порождать термы различных модальностей: $p(\mathsf{словo}|t)$, p(n-грамма|t), $p(\mathsf{автор}|t)$, $p(\mathsf{время}|t)$, $p(\mathsf{источник}|t)$, $p(\mathsf{объект}|t)$, $p(\mathsf{ссылка}|t)$,



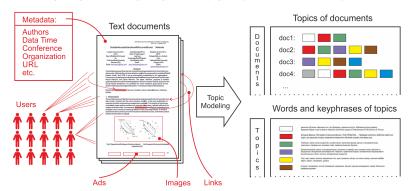
Тема может порождать термы различных модальностей: p(слово|t), p(n-грамма|t), p(автор|t), p(время|t), p(источник|t), p(объект|t), p(ссылка|t), p(баннер|t),



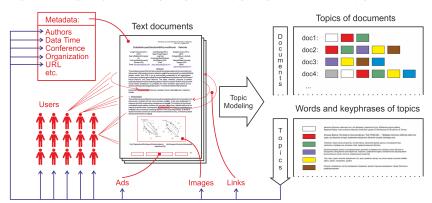
Примеры мультимодальных тематических моделей

Мультимодальная тематическая модель

Тема может порождать термы различных модальностей: p(слово|t), p(n-грамма|t), p(автор|t), p(время|t), p(источник|t), p(объект|t), p(ссылка|t), p(баннер|t), p(пользователь|t)



Тема может порождать термы различных модальностей: p(слово|t), p(n-грамма|t), p(автор|t), p(время|t), p(источник|t), p(объект|t), p(ссылка|t), p(баннер|t), p(пользователь|t)



Мультимодальные тематические модели
Мультимодальный ЕМ-алгоритм
Примеры мультимодальных тематических моделей

EM-алгоритм для мультимодальной ARTM

W_m — словарь термов m-й модальности, $m \in M$

Максимизация суммы log-правдоподобий с регуляризацией:

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

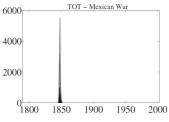
ЕМ-алгоритм: метод простой итерации для системы уравнений

Е-шаг:
$$\begin{cases} p_{tdw} = \underset{t \in T}{\mathsf{norm}} \left(\phi_{wt} \theta_{td} \right) \\ \phi_{wt} = \underset{w \in \mathcal{W}^m}{\mathsf{norm}} \left(\sum_{d \in D} \tau_{m(w)} n_{dw} p_{tdw} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right) \\ \theta_{td} = \underset{t \in T}{\mathsf{norm}} \left(\sum_{w \in d} \tau_{m(w)} n_{dw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right) \end{cases}$$

K. Vorontsov, O. Frei, M. Apishev et al. Non-Bayesian additive regularization for multimodal topic modeling of large collections. CIKM TM workshop, 2015.

Пример. Использование модальностей времени и *п*-грамм

По коллекции выступлений президентов США



800	1850	15	000	1950
1.	mexico		8. to	erritory
2. texas		9. army		
	3. war		10.	peace
4.	mexican		10	l. act
5.	united		12.	policy
6.	country		13.	foreign
7. gc	overnmei	nt	14.	citizens

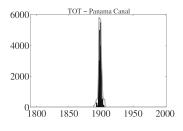
3000	Our Mod	lel – Mexic	an War	
2000				_
1000				
01800	1850	1900	1950	2000

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

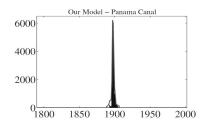
Shoaib Jameel, Wai Lam. An n-gram topic model for time-stamped documents. 2013

Пример. Использование модальностей времени и *п*-грамм

По коллекции выступлений президентов США



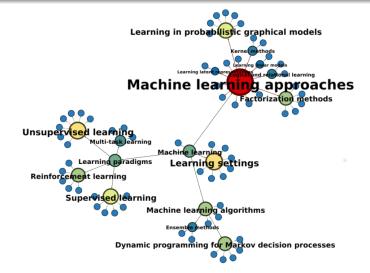
1. government	8. spanish
2. cuba	9. island
3. islands	10. act
4. international	11. commission
5. powers	12. officers
6. gold	13. spain
7. action	14. rico



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	
6. united states	13. american control
7. state panama	14. canal

Shoaib Jameel, Wai Lam. An n-gram topic model for time-stamped documents. 2013

Пример древовидной тематической иерархии



G.Bordea. Domain adaptive extraction of topical hierarchies for expertise mining. 2013.

Стратегии иерархического разделения тем на подтемы

Процесс построения иерархии тем:

- структура: дерево / многодольный граф
- направление: снизу вверх / сверху вниз / одновременно
- наращивание: повершинное / послойное
- обучение: без учителя / по готовым рубрикаторам

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 Ф: 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:

$$\sum_{t \in T} n_t \, \mathsf{KL}_w \Big(p(w|t) \, \Big\| \, \sum_{s \in S} p(w|s) \frac{p(s|t)}{p(s|t)} \Big) \, \to \, \min_{\Phi, \Psi},$$

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

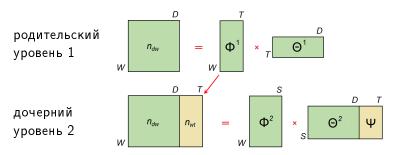
Родительская $\Phi^p \approx \Phi \Psi$, отсюда регуляризатор матрицы Φ :

$$R(\Phi, \Psi) = au \sum_{t \in T} \sum_{w \in W} n_{wt} \ln \sum_{s \in S} \phi_{ws} \psi_{st} \rightarrow \max$$

Родительские темы t- «документы» с частотами термов n_{wt}

Регуляризатор Φ : построение второго уровня с подтемами S

Добавим в коллекцию |T| псевдо-документов родительских тем с частотами термов $n_{wt} = \tau n_t \phi_{wt}, \ t \in T$



Матрица связей тем с подтемами $\Psi = (p(s|t))$ образуется в столбцах матрицы Θ , соответствующих псевдо-документам.

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

Регуляризатор Θ : 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:

$$\sum_{d \in D} n_d \operatorname{\mathsf{KL}}_t \Big(p(t|d) \, \Big\| \, \sum_{s \in S} \underbrace{p(t|s)} p(s|d) \Big) \, \, \to \, \, \min_{\Theta, \Psi},$$

где $\Psi = (\psi_{ts})_{T imes S} - ($ другая!) матрица связей, $\psi_{ts} = p(t|s)$

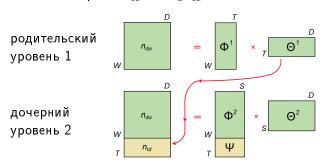
Родительская $\Theta^p pprox \Psi\Theta$, отсюда регуляризатор матрицы Θ :

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

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

Регуляризатор Θ : построение второго уровня с подтемами S

Добавим в каждый документ модальность родительских тем с частотами термов $n_{td}= au n_d heta_{td},\ t\in T$

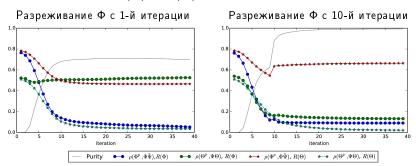


Матрица связей тем с подтемами $\Psi = (p(t|s))$ образуется в строках матрицы Φ , соответствующих родительским темам.

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 $R(\Phi)$ и $R(\Theta)$ при переходе с уровня 1 на 2:

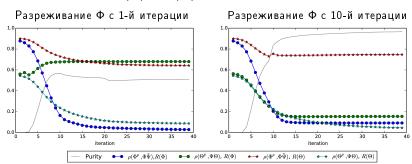


Выводы. $R(\Theta)$ плохо приближает Φ^p . При разреживании Φ с 10-й итерации $R(\Phi)$ хорошо приближает Φ^p и Θ^p

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 $R(\Phi)$ и $R(\Theta)$ при переходе с уровня 2 на 3:



Выводы. $R(\Theta)$ плохо приближает Φ^p . При разреживании Φ с 10-й итерации $R(\Phi)$ хорошо приближает Φ^p и Θ^p

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

Выводы

- $R(\Phi)$ лучше $R(\Theta)$, т.к. добавлять псевдо-документы удобнее, чем вставлять модальности в каждый документ
- ullet $R(\Phi)$ хорошо приближает $\Phi^p pprox \Phi ilde{\Psi}$ и $\Theta^p pprox \Psi \Theta$ при осторожном (с 10-й итерации) разреживании Φ
- ullet $R(\Theta)$ приближает только $\Theta^p pprox \Psi\Theta$
- ullet сильное разреживание $\psi_{ts} \in \{0,1\}$ даёт иерархию-дерево
- ullet нельзя допускать вырождения $\psi_{ts}= {\it p}(t|s)\equiv 0$

Трудные и/или открытые проблемы:

- тематические иерархии с ветвлением различной глубины
- автоматическое оценивание качества иерархии
- автоматическое именование подтем с учётом родительской
- определение типа документа по его следу в иерархии

Визуализация тематической иерархии

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



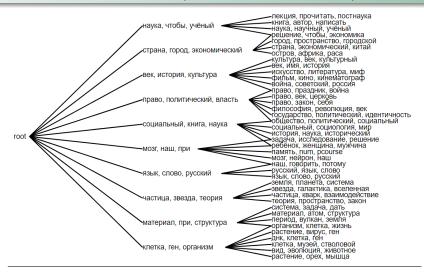


Для именования темы используются три топовых слова темы

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.

Иерархический спектр тем (коллекция postnauka.ru)



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

Построение спектра тем. Постановка задачи

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

$$\sum_{i=2}^{|T|}
ho(t_i,t_{i-1}) o \mathsf{min}$$

 Φ ункция расстояния ho(t,t') между темами, примеры:

- ullet Манхэттенское: $ho(t,t') = \sum\limits_{w \in W} \left| \phi_{wt} \phi_{wt'} \right|$
- ullet Хеллингера: $ho^2(t,t')=rac{1}{2}\sum_{w\in W}\left(\sqrt{\phi_{wt}}-\sqrt{\phi_{wt'}}
 ight)^2$
- ullet Жаккара: $ho(t,t')=1-rac{|W_t\cap W_{t'}|}{|W_t\cup W_{t'}|},\;\;W_t=ig\{w\colon \phi_{wt}>rac{1}{|W|}ig\}$

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

Построение спектра тем — это задача коммивояжёра

Задача TSP (traveling salesman problem)

Найти путь минимальной суммарной стоимости, соединяющий $\mathcal T$ 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 алгоритм — $O(T^{2.2})$.

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.

Иерархическая тематизация коллекции научных публикаций

Гипотеза. След документа в глубокой тематической иерархии определяет его научный жанр (специализацию, назначение):



узко специализированный, для профессионалов



междисциплинарное исследование, для профессионалов



обзорный, для ознакомления с предметной областью

популярный или энциклопедический, для самообразования, расширения кругозора

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)
- ullet удалена пунктуация, $\ddot{e} \rightarrow e$, лемматизация pymorphy2

Анастасия Янина. Тематические и нейросетевые модели языка для разведочного информационного поиска. Диссертация к.ф.-м.н., МФТИ. 2022.

Методика измерения качества поиска

Тематическая модель для документного поиска Оптимизация гиперпараметров

Методика оценивания качества разведочного поиска

Поисковый запрос

набор ключевых слов или фрагментов текста, около одной страницы А4

Поисковая выдача

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

Два задания асессорам

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

op MapReduce

Набор МарКофис — программия модели (фициароді) заполитивих распределенням вагоплений для большо офедура, даннах в разоках порадилная предстабор, представляющих сооби выбор Зати-милосов и исполительно, учинит для соц

- Основные концепции <u>Набоор Марбафор</u>д можно сформулировата нак обработна вычисление больших объемов даниых,
- SECONDARY OF SECONDARY
- работа на ненадежном оборудования;
- ватоматическая обработна отначие выполнения надаемей.
 Цебовр популнувная программена платформа (собщае бългетоск) построения распрадлениям приложений для высосно-прадлениям обработиям (макадет разлібе) пределениям МТРУ деятеми.
- Наборр включает в себе спедующие компоненты: 1. HDFS – распределенных файтовах система;

 Наборо МарЯефоре — программия моделя (брацепрорк) выполнения распределенных выполнений для больших обредор, данных в рамках перадилых цар induce.

учесь. «ИК пераменями налини», сольная саказность форфирора расправлениях качислений и коментские бейшения, реализующие расправлениям атмурить. Хах соверствая: Откустения подержих автерениямий регурациямий можем выполнения расправлениями канислений в будбор v1.0 подверживается только можем важисления другофобо.

Маличи «присован гочк отказа и дах спедства», вмониковность испланования о орган с выголания гребованнями вырачности. Проблеми версиосной совмествности: гребования по единовременнями общениемы всех вычастительных учло клистера при общениемы или бромы Ценеро установа и посей версии или пакта обмествений;

Пример запроса для разведочного поиска

Пример: фрагмент запроса «Система IBM Watson»

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

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

...

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

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

Тематика запросов разведочного поиска

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

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

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

Векторный поиск тематически близких документов

$$heta_{tq} = p(t|q)$$
 — тематический вектор запроса q $heta_{td} = p(t|d)$ — тематические векторы документов $d \in D$

Косинусная мера близости документа d и запроса q:

$$sim(q,d) = \frac{\sum_{t} \theta_{tq} \theta_{td}}{\left(\sum_{t} \theta_{tq}^{2}\right)^{1/2} \left(\sum_{t} \theta_{td}^{2}\right)^{1/2}}.$$

Ранжируем документы коллекции $d \in D$ по убыванию $\mathrm{sim}(q,d)$ Выдача тематического поиска — k первых документов.

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

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

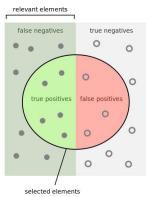
A.lanina, K. Vorontsov. Regularized multimodal hierarchical topic model for document-by-document exploratory search. FRUCT-ISMW, 2019.

Оценивание качества поиска

Precision — доля релевантных среди найденных Recall — доля найденных среди релевантных

$$P=rac{ ext{TP}}{ ext{TP}+ ext{FP}}$$
 — точность (precision) $R=rac{ ext{TP}}{ ext{TP}+ ext{FN}}$ — полнота, (recall) $F_1=rac{2PR}{P+R}$ — F1-мера

TP (true positive) — найденные релевантные FP (false positive) — найденные нерелевантные FN (false negative) — ненайденные релевантные





Какие модели поиска сравнивались

- assessors: результаты поиска, выполненного асессорами
- TF-IDF, BM25: сравнение документов по частотам слов
- word2vec: нетематические векторные представления слов
- PLSA: Probabilistic Latent Semantic Analysis (1999)
- LDA: Latent Dirichlet Allocation (2001)
- ARTM: тематическая модель с тремя регуляризаторами
- hARTM: иерархические модели ARTM 2х и 3х уровней

Задачи регуляризаторов в ARTM и hARTM:

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

Стратегия регуляризации

Последовательное применение трёх регуляризаторов

• декоррелирование тем:

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

 $oldsymbol{2}$ разреживание распределений p(t|d):

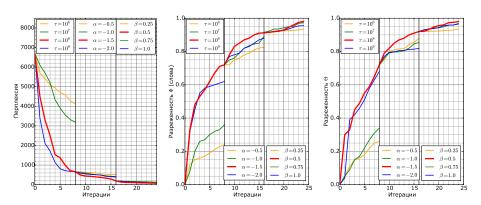
$$R(\Theta) = -\alpha \sum_{d \ t} \ln \theta_{td}$$

 \odot сглаживание распределений p(w|t):

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

Последовательный подбор коэффициентов регуляризации

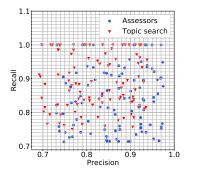
- ullet декоррелирование распределений термов в темах (au),
- ullet разреживание распределений тем в документах (lpha),
- ullet сглаживание распределений термов в темах (eta).



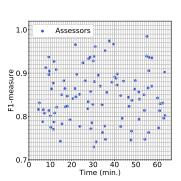
Результаты измерения точности и полноты по запросам

100 запросов, 3 асессора на запрос

точность и полнота поиска



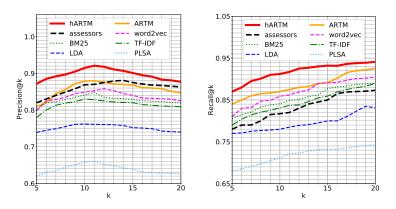
время и F_1 -мера (асессоры)



- среднее время обработки запроса асессором 30 минут
- точность выше у асессоров, полнота у поисковика

Сравнение с асессорами по качеству поиска

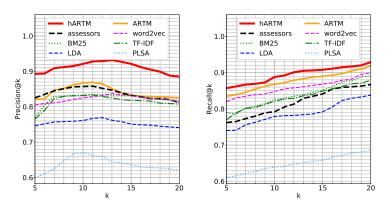
Точность и полнота по первым k позициям поисковой выдачи (коллекция Habrahabr.ru)



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

Сравнение с асессорами по качеству поиска

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



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

Влияние числа тем на качество поиска

Все регуляризаторы и модальности, плоская модель

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

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

Методика измерения качества поиска Тематическая модель для документного поиска Оптимизация гиперпараметров

Влияние числа тем на качество поиска

Habrahabr. Все регуляризаторы и модальности, два уровня

$ T_1 $	2	0		25							30
$ T_2 $	150	200	2	50		275		30	00	400	450
Pr@5	0.621	0.742	0.839	0.850	0.865	0.869	0.869	0.803	0.769	0.701	0.670
Pr@10	0.645	0.749	0.850	0.861	0.879	0.911	0.895	0.809	0.796	0.719	0.689
Pr@15	0.635	0.751	0.848	0.869	0.873	0.893	0.887	0.807	0.781	0.721	0.701
Pr@20	0.630	0.745	0.841	0.855	0.864	0.874	0.875	0.800	0.775	0.709	0.675
R@5	0.628	0.773	0.843	0.865	0.881	0.881	0.868	0.849	0.839	0.715	0.691
R@10	0.652	0.782	0.855	0.871	0.902	0.918	0.877	0.871	0.845	0.745	0.699
R@15	0.671	0.801	0.870	0.889	0.929	0.939	0.901	0.883	0.861	0.781	0.722
R@20	0.680	0.819	0.886	0.892	0.955	0.955	0.907	0.901	0.872	0.801	0.729

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

Влияние числа тем на качество поиска

Habrahabr. Все регуляризаторы и модальности, три уровня

$ T_1 $	2	0		25							30
$ T_2 $	150	200	2!	50		275		30	00	400	450
$ T_3 $	750	800	1200	1300	1300	1400	1500	1500	1600	3000	3500
Pr@5	0.625	0.743	0.840	0.852	0.869	0.872	0.870	0.805	0.771	0.705	0.672
Pr@10	0.648	0.754	0.851	0.867	0.882	0.915	0.901	0.811	0.799	0.722	0.694
Pr@15	0.632	0.752	0.850	0.872	0.878	0.895	0.889	0.809	0.785	0.729	0.703
Pr@20	0.629	0.745	0.845	0.861	0.871	0.877	0.882	0.803	0.778	0.710	0.681
R@5	0.632	0.780	0.845	0.869	0.883	0.889	0.872	0.851	0.841	0.721	0.695
R@10	0.654	0.792	0.859	0.873	0.905	0.922	0.881	0.873	0.850	0.749	0.703
R@15	0.675	0.805	0.874	0.892	0.932	0.942	0.905	0.889	0.863	0.787	0.725
R@20	0.684	0.824	0.889	0.901	0.958	0.961	0.912	0.904	0.878	0.805	0.734

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

Методика измерения качества поиска Тематическая модель для документного поиска Оптимизация гиперпараметров

Влияние числа тем на качество поиска

TechCrunch. Все регуляризаторы и модальности, два уровня

$ T_1 $	8	0		100						1	20
$ T_2 $	300	350	50	00	550		600		700	750	
Pr@5	0.651	0.701	0.749	0.789	0.883	0.889	0.889	0.785	0.721	0.701	0.675
Pr@10	0.675	0.709	0.771	0.821	0.891	0.918	0.902	0.803	0.738	0.718	0.691
Pr@15	0.687	0.712	0.773	0.827	0.899	0.919	0.905	0.817	0.741	0.721	0.701
Pr@20	0.683	0.707	0.759	0.815	0.885	0.888	0.895	0.805	0.732	0.716	0.679
R@5	0.749	0.791	0.801	0.854	0.868	0.875	0.861	0.849	0.829	0.731	0.701
R@10	0.765	0.809	0.823	0.873	0.890	0.904	0.875	0.867	0.835	0.745	0.708
R@15	0.771	0.820	0.841	0.882	0.909	0.921	0.895	0.890	0.848	0.769	0.717
R@20	0.778	0.825	0.851	0.887	0.928	0.942	0.929	0.901	0.869	0.785	0.728

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

Методика измерения качества поиска
Тематическая модель для документного поиска
Оптимизация гиперпараметров

Влияние числа тем на качество поиска

TechCrunch. Все регуляризаторы и модальности, три уровня

$\overline{ T_1 }$	8	0		100						1	20
$ T_2 $	300	350	50	00		550		60	00	700	750
$ T_3 $	1500	1700	2500	2600	2600	2800	3000	3000	3200	4500	4700
Pr@5	0.655	0.707	0.751	0.792	0.887	0.893	0.890	0.789	0.722	0.703	0.678
Pr@10	0.678	0.712	0.773	0.823	0.895	0.922	0.905	0.805	0.741	0.722	0.692
Pr@15	0.692	0.715	0.775	0.831	0.902	0.921	0.907	0.821	0.743	0.725	0.703
Pr@20	0.687	0.709	0.761	0.819	0.889	0.885	0.898	0.809	0.736	0.719	0.683
R@5	0.751	0.795	0.802	0.856	0.871	0.877	0.863	0.852	0.831	0.738	0.705
R@10	0.767	0.812	0.825	0.875	0.892	0.908	0.879	0.871	0.842	0.751	0.711
R@15	0.772	0.824	0.841	0.887	0.912	0.927	0.901	0.893	0.854	0.772	0.721
R@20	0.783	0.830	0.854	0.892	0.931	0.949	0.935	0.905	0.871	0.790	0.732

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

Влияние модальностей на качество поиска

Все регуляризаторы и модальности, 3 уровня, оптимальное |T| Модальности: Words, Bigrams, Authors, Comments, Tags, Hubs, Categories

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

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

Влияние регуляризаторов на качество поиска

Все регуляризаторы и модальности, 3 уровня, оптимальное |T| Регуляризаторы: Decorrelation, $\underline{\Theta}$ -sparsing, $\underline{\Phi}$ -smoothing, \underline{H} ierarchy

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

- лучше использовать все регуляризаторы
- модели со слабой регуляризацией (PLSA, LDA) слабы

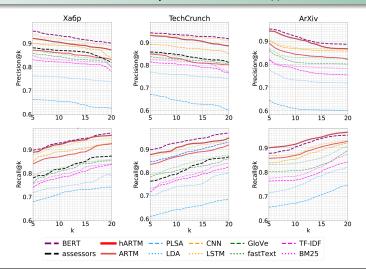
Влияние функции близости на качество поиска

Все регуляризаторы и модальности, 3 уровня, оптимальное |T| Функции близости: Euclidean, Cosine, Manhattan, Hellinger, KL-div

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

• косинусная функция близости уверенно лидирует

Сравнение с поиском по нейросетевым эмбедингам



A.lanina, K. Vorontsov. Hierarchical interpretable topical embeddings for exploratory search and real-time document tracking. 2020.

Выводы по результатам экспериментов

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

А. О. Янина. Тематические и нейросетевые модели языка для разведочного информационного поиска. Диссертация к.ф.-м.н., МФТИ. 2022.

Тематическая модель для научного поиска должна быть...

- Интерпретируемая: объяснять смысл каждой темы
- Иерархическая: разделять тем на подтемы
- Хронологическая: прослеживать темы во времени
- 🚇 Мультимодальная: слова, авторы, категории, связи, теги,...
- Мультиграммная: слова, термины-словосочетания
- Мультиязычная для кросс- и много-языкового поиска
- Сегментирующая документ на тематические блоки
- Обучаемая по обратной связи с пользователями
- Определяющая число тем автоматически
- Создающая и именующая новые темы автоматически
- Онлайновая: обрабатывать новые документы в потоке
- 🚇 Параллельная, распределённая для больших данных

Резюме

Разведочный информационный поиск (exploratory search):

- это поиск по смыслу, а не по ключевым словам
- строится на векторных представлениях текста (тематических или нейросетевых эмбедингах текста)
- требует от тематических моделей многофункциональности
- является одной из главных мотиваций для ARTM,
- в том числе для мультимодальных и иерархических ARTM

Открытые проблемы:

- тематизация подборок с дисбалансом тем
- автоматическое именование и суммаризация тем
- эффективные методы визуализации (картирования)

Задания по курсу

Задача-минимум: научиться решать задачи NLP с использованием тематического моделирования в BigARTM

Задача-максимум: сделать полезное мини-исследование

виды деятельности	оценка
теоретические задания	$\sum_{i} X_{i}$
решение прикладной задачи	5 <i>X</i>
обзор по последним NeuralTM	5 <i>X</i>
интеграция ARTM в pyTorch	5 <i>X</i>
участие в одном из проектов	10 <i>X</i>
работа над открытой проблемой	10 <i>X</i>

где X — оценка за вид деятельности по 5-балльной шкале. score — суммарная оценка по всем видам деятельности.

Итоговая оценка: $\min(10, \lfloor score/5 \rfloor)$ по 10-балльной шкале.

Теоретическое задание к лекции 1

Упражнения на принцип максимума правдоподобия:

- 1. Униграммная модель документов: $p(w|d)=\xi_{dw}$ Найти параметры модели ξ_{dw} .
- 2. Униграммная модель коллекции: $p(w|d) = \xi_w$ для всех d Найти параметры модели ξ_w .

Подсказка: применить условия ККТ или основную лемму.

- 3. Творческое задание (возможны разные решения) Предложите модель, определяющую роли слов в текстах:
- тематические слова
- специфичные слова документа (шум)
- слова общей лексики (фон)
- Подсказка 1: искать распределение ролей слов p(r|w), $r \in \{\tau, \mu, \phi\}$.
- Подсказка 2: можно разреживать p(r|w) для жёсткого определения ролей. Подсказка 3: можно использовать документную частоту слов.
- подсказка э: можно использовать документную частоту слов.

Теоретическое задание к лекции 2

- 4. Запишите критерий логарифма правдоподобия с регуляризацией для тематической модели $p(w|d) = \sum_t \phi_{wt} \theta_{td}$, используя исходные данные $(d_i, w_i)_{i=1}^n$ вместо счётчиков n_{dw} . Выведете из него EM-алгоритм, докажите его эквивалентность обычному EM-алгоритму для ARTM.
- **5.** Запишите критерий логарифма правдоподобия для локализованной тематической модели $p(w|C_i) = \sum_t \phi_{wt} p(t|C_i)$. Выведете из него EM-алгоритм с локализованным E-шагом.

Какие приближения пришлось сделать в процессе вывода? Какие переменные удобнее оставить в модели, ϕ_{wt} или ϕ_{tw}' ?

6. Творческое задание (возможны разные решения) Предложите «какую-нибудь разумную» параметризацию для тематической модели внимания. Используя «основную лемму», получите уравнения для новых параметров модели.

Исследовательское задание к лекции 2

Открытая проблема. Продолжить исследование Ильи Ирхина:

- Освоить код: https://github.com/ilirhin/python_artm
- Реализовать локализованный Е-шаг

Исследовать зависимость метрик качества от параметров (перплексия, разреженность, различность, когерентность):

- L число проходов
- ullet $\vec{\gamma}_i,\; \dot{\overline{\gamma}}_i$ длина скользящего среднего
- ullet $ec{\gamma}_i, \ \dot{\overline{\gamma}}_i$ асимметричность левого и правого контекста
- ullet $\vec{\gamma}_i$, $\dot{\gamma}_i$ учёт границ предложений, абзацев, глав
- ullet β баланса левого и правого контекста
- ullet α , δ параметры онлайнового EM-алгоритма
- ullet опция «подставлять p_{ti}/n_t вместо ϕ_{w_it} на $\hbox{E-шаге}$ »
- ullet опция «исключать p_{ti} позиции i из контекстов $\stackrel{
 ightarrow}{ heta}_{ti}$ $\stackrel{
 ightarrow}{ heta}_{ti}$ »

Теоретическое задание к лекции 3

7. Выведете формулы ЕМ-алгоритма в случае, когда логарифм в функции потерь заменяется гладкой монотонно возрастающей функцией ℓ :

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ell \left(\sum_{t \in T} \phi_{wt} \theta_{td} \right) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

Подумайте, какие замены логарифма полезны, и почему.

8. Замените In гладкой монотонно возрастающей функцией μ в регуляризаторе сглаживания—разреживания (модель LDA):

$$R(\Phi, \Theta) = \sum_{t \in T} \sum_{w \in W} \beta_w \mu(\phi_{wt}) + \sum_{d \in D} \sum_{t \in T} \alpha_t \mu(\theta_{td}).$$

Как изменится М-шаг и воздействие регуляризатора на модель?

9. Какому регуляризатору соответствует формула М-шага

$$\phi_{wt} = \operatorname{norm}_{w} \left(n_{wt} [n_{wt} > \gamma n_t] \right)$$

Теоретическое задание к лекции 3

Аналитик построил тематическую модель Φ^0 , Θ^0 и отметил среди столбцов матрицы Φ^0 темы двух типов: удачные $\mathcal{T}_+ \subset \mathcal{T}$ и неудачные $\mathcal{T}_- \subset \mathcal{T}$.

Теперь он хочет построить модель ещё раз так, чтобы

- удачные темы остались в матрице Ф;
- ullet остальные темы построились по-другому и были не похожи на каждую из неудачных тем $t\in \mathcal{T}_-$.
- 10. Предложите регуляризаторы для этого.
- 11. Не получится ли так, что новые темы будут отдаляться от суммы неудачных тем $\sum_{t\in T_-}\phi_{wt}^0$ вместо того, чтобы отдаляться от каждой из неудачных тем по отдельности? Почему это плохо и как этого избежать?
- 12. Предложите способ инициализации Ф для новой модели.

Исследовательские задания к лекции 4

- Проблема несбалансированности тем
 - генераторы синтетических несбалансированных коллекций
 - модели локального контекста лишены этой проблемы?
 - регуляризаторы декоррелирования + семантической однородности
- Семейство средневзвешенных статистик
 - генераторы синтетических коллекций, удовлетворяющих гипотезе условной независимости
 - как (и нужно ли) определять пороги для построения статистических тестов условной независимости?
 - как ослабить проверку гипотезы условной независимости в модели локального контекста?
 - как перестраивать несогласованные темы?
- Критерий внутритекстовой когерентности
 - найти лучший вариант критерия с помощью калибровки по размеченным тематическим цепочкам
 - вычисление критерия должно естественным образом встраиваться в модель локального контекста

Теоретическое задание к лекции 5

- 13. Для иерархической тематической модели с рег. $R(\Phi, \Psi)$ предложите способ разреживания матрицы связей $\Psi = (p(s|t))$, гарантирующий, что
- 1) у каждой родительской темы будет хотя бы одна дочерняя; 2) у каждой дочерней темы будет хотя бы одна родительская.

Подсказка: можно придумывать критерий регуляризации, а можно — формулу М-шага для матрицы Ψ .

- 14. Предложите способ гарантировать, что если родительская тема t получает только одну дочернюю s, то она переходит в неё целиком и как распределение: p(w|s) = p(w|t), то есть тема t на данном уровне не расщепляется на подтемы.
- 15. Предложите способ согласования вероятностных смесей $p(w|t) \approx \sum_{s \in S} p(w|s)p(s|t)$ и $p(t|d) \approx \sum_{s \in S} p(t|s)p(s|d)$ с учётом тождества p(s|t)p(t) = p(t|s)p(s).

Исследовательское задание к лекции 5

Участие в проекте «Мастерская знаний»

Дано:

- подборки, сгенерированные SciRus по одной статье
- асессорская разметка статей подборки по релевантности
- несколько вариантов токенизации
 - в том числе с автоматическим выделением терминов

Найти:

- тематическую модель
- модель ранжирования подборки по релевантности
- оптимальные: токенизацию, число тем, регуляризаторы
- распределение терминов по тематичности

Критерий:

- качество ранжирования
- (визуально) интерпретируемость тем
 - в том числе автоматического именования тем

Примеры датасетов для практических заданий по курсу

- Открытые датасеты (английский): 20 newsgroups, NIPS, KOS
- Научные статьи: eLibrary, Semantic Scholar, arXiv, PubMed
- Научно-популярные статьи: ПостНаука, Элементы, Хабр,...
- TechCrunch (английский)
- Данные социальных сетей: VK, Twitter, Telegram,...
- Википедия
- Новостной поток (20 источников на русском языке)
- Данные кадровых агентств: резюме + вакансии
- Транзакции клиентов Sberbank DSD 2016
- Акты арбитражных судов РФ

Проекты

- «Мастерская знаний» для научного поиска
 - пользователь строит тематические подборки статей,
 - поисковая выдача формируется моделью SciRus.
 - задача: показать пользователю тематику подборки
 - понадобится автоматическое выделение терминов,
 - выделение тематических фраз из документов,
 - автоматическое именование и суммаризация тем
 - конечная цель: ускорить понимание предметной области
- «Тематизатор» для социо-гуманитарных исследований
 - пользователь задаёт грубый фильтр текстового потока
 - задача: «классифицировать иголки в стоге сена»,
 - разделив темы на информативные и мусорные,
 - выделив аспекты и тональности в каждой теме
 - конечная цель: q&q аналитика проблемной среды

Открытые проблемы тематического моделирования

- 💶 Проблема несбалансированности тем в коллекции
- Обеспечение 100%-й интерпретируемости тем
- Тематические модели внимания последовательного текста
- Обнаружение новых тем или трендов в потоке текстов
- Автоматическое именование и аннотирование тем
- Обзор подходов в нейросетевых тематических моделях
- Обеспечение полноты и устойчивости множества тем
- Автоматический подбор гиперпараметров, AutoML
- Оптимизация гиперпараметров в потоковом режиме
- 💿 Проблема несбалансированности текстов по длине
- 🚇 Бережное слияние моделей нескольких коллекций
- Гиперграфовые тематические модели в RecSys