

Аддитивная регуляризация тематических моделей в задаче анализа этносоциального дискурса

Мурат Апишев
great-mel@yandex.ru

МГУ им. Ломоносова, Яндекс

April 14, 2016

Тематическое моделирование

Тематическое моделирование — приложение машинного обучения к статистическому анализу текстов.

Тема — терминология предметной области, набор терминов (униграм или n -грам) часто встречающихся вместе в документах.

Тематическая модель исследует скрытую тематическую структуру коллекции текстов:

- *тема* t — это вероятностное распределение $p(w|t)$ над терминами w
- *документ* d — это вероятностное распределение $p(t|d)$ над темами t

Приложения — информационный поиск по длинным текстовым запросам, классификация, категоризация и суммаризация текстов.

Задача тематического моделирования

Дано: W — словарь терминов (униграм или n -биграмм),
 D — коллекция текстовых документов $d \subset W$,
 n_{dw} — счётчик частоты появления слова w в документе d .

Найти: модель $p(w|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$ с параметрами Φ и Θ :
 $\phi_{wt} = p(w|t)$ — вероятности терминов w в каждой теме t ,
 $\theta_{td} = p(t|d)$ — вероятности тем t в каждом документе d .

Критерий максимизация логарифма правдоподобия:

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} \rightarrow \max_{\phi, \theta};$$

$$\phi_{wt} \geq 0; \quad \sum_w \phi_{wt} = 1; \quad \theta_{td} \geq 0; \quad \sum_t \theta_{td} = 1.$$

Проблема: задача стохастического матричного разложения некорректно поставленная: $\Phi\Theta = (\Phi S)(S^{-1}\Theta) = \Phi'\Theta'$.

PLSA и EM-алгоритм

Максимизация логарифма правдоподобия:

$$\sum_{d \in D} \sum_{w \in W} n_{dw} \ln \sum_t \phi_{wt} \theta_{td} \rightarrow \max_{\Phi, \Theta}$$

EM-алгоритм: метод простых итерация для решения системы уравнений

$$\begin{array}{l} \text{E-шаг:} \\ \text{M-шаг:} \end{array} \left\{ \begin{array}{l} p_{tdw} = \operatorname{norm}_{t \in T}(\phi_{wt} \theta_{td}) \\ \phi_{wt} = \operatorname{norm}_{w \in W}(n_{wt}), \quad n_{wt} = \sum_{d \in D} n_{dw} p_{tdw} \\ \theta_{td} = \operatorname{norm}_{t \in T}(n_{td}), \quad n_{td} = \sum_{w \in W} n_{dw} p_{tdw} \end{array} \right.$$

где $\operatorname{norm}_{i \in I} x_i = \frac{\max\{x_i, 0\}}{\sum_{j \in I} \max\{x_j, 0\}}$

ARTM и регуляризованный EM-алгоритм

Максимизация логарифма правдоподобия с **дополнительными аддитивными регуляризаторами R** :

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

EM-алгоритм: метод простых итераций для системы уравнений

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

Примеры регуляризаторов

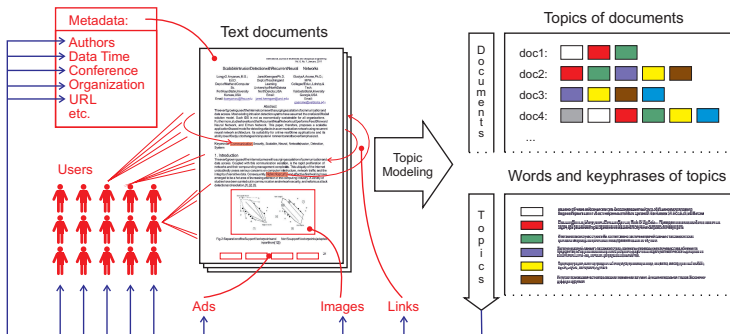
Многие байесовские модели могут быть интерпретированы в терминах ARTM.

Примеры регуляризаторов:

- 1 Сглаживание Φ / Θ (приводит к известной модели LDA)
- 2 Разреживание Φ / Θ
- 3 Декорреляция тем в Φ
- 4 Частичное обучение
- 5 Максимизация когерентности тем
- 6 Отбор тем
- 7 ...

Мультимодальная тематическая модель

Мультимодальная тематическая модель распределения тем на терминах $p(w|t)$, авторах $p(a|t)$, метках времени $p(y|t)$, изображениях $p(o|t)$, связанных документах $p(d'|t)$, рекламных баннерах $p(b|t)$, пользователей $p(u|t)$, и объединяет все эти модальности в одно тематическую модель.



M-ARTM и мультимодальный регуляризованный EM-алгоритм

W^m — словарь терминов m -й модальности, $m \in M$,
 $W = W^1 \sqcup W^m$ как объединение словарей всех модальностей.

Максимизация логарифма мультимодального правдоподобия с аддитивными регуляризаторами R :

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

EM-алгоритм: метод простых итерация для системы уравнений

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

Проект BigARTM

Особенности BigARTM:

- Быстрая¹ параллельная и онлайн-обработка данных;
- Поддержка мультимодальных регуляризованных тематических моделей;
- Встроенная расширяемая библиотека регуляризаторов и метрик качества;

Сообщество BigARTM:

- Открытый репозиторий <https://github.com/bigartm>
- Описание и документация <http://bigartm.org>

Лицензия BigARTM и программные особенности:

- Бесплатное коммерческое использование (BSD 3-Clause license)
- Кроссплатформенная — Windows, Linux, Mac OS X (32 bit, 64 bit)
- Программные API: command line, C++, Python

¹Vorontsov K., Frei O., Apishev M., Romov P., Dudarenko M. BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections Analysis of Images, Social Networks and Texts. 2015

Поиск этнорелевантного контента в блогосфере

Создание концепции и методологии для мониторинга остояния межэтнических отношений по данным социальных медиа.

Задачи тематического моделирования в проекте:

- 1 Извлечение этно-релевантных тем из данных социальных медиа
- 2 Распознавание событийных и непрерывных во времени тем
- 3 Сантимент-анализ этнического дискурса

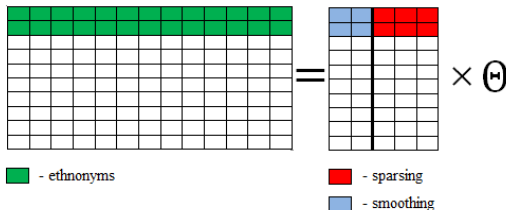
Грант Российского Научного Фонда 15-18-00091 (2015–2017)
(Высшая Школа Экономики, С.-Петербургская Школа социальных и общественных наук, Лаборатория интрнет-исследований ЛИНИС)

Примеры этнонимов для частичного обучения

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

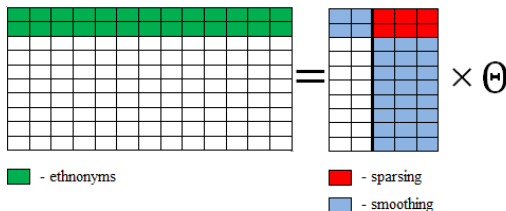
Регуляризация для поиска этнических тем

- сглаживание этнонимов в этнических темах
- разреживание этнонимов в общих темах
-
-
-



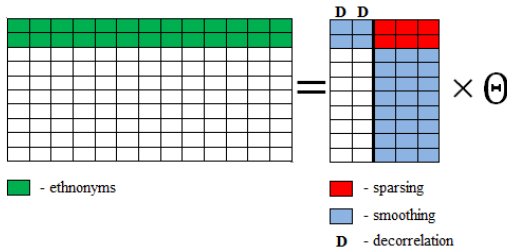
Регуляризация для поиска этнических тем

- сглаживание этнонимов в этнических темах
- разреживание этнонимов в общих темах
- сглаживание не-этнонимов в общих темах
-
-



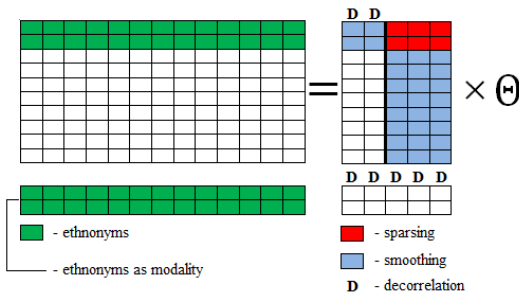
Регуляризация для поиска этнических тем

- сглаживание этнонимов в этнических темах
- разреживание этнонимов в общих темах
- сглаживание не-этнонимов в общих темах
- декоррелирование этнических тем
-



Регуляризация для поиска этнических тем

- сглаживание этнонимов в этнических темах
- разреживание этнонимов в общих темах
- сглаживание не-этнонимов в общих темах
- декоррелирование этнических тем
- дублирование этнонимов в качестве новой модальности и декоррелирование тем в ней



Experiment

- Коллекция Живого Журнала: 1.58М документов
- 860К слов в словаре после лемматизации
- 90К слов после фильтрации
 - коротких слов с длиной ≤ 2 ,
 - редких слов со встречаемости в коллекции $n_w < 20$,
 - нерусских слов
- 250 этнонимов в словаре

ARTM с частичным обучением для поиска этнорелевантных тем

Число и качество этнических тем, найденных моделью:

модель	ethnic $ S $	background $ B $	++	+-	-+	coh ₂₀ ²	tfidf ₂₀
PLSA		400	12	15	17	-1447	-1012
LDA		400	12	15	17	-1540	-1121
ARTM-4	250	150	21	27	20	-1651	-1296
ARTM-5	250	150	38	42	30	-1342	-908

- ARTM-4:

- этнические темы: декорреляция, сглаживание этнонимов
- фоновые темы: сглаживание, разреживание этнонимов

- ARTM-5:

- ARTM-4 + декоррелируемая модальность этнонимов

²Когерентность и TF-IDF когерентность — метрики, коррелирующие с человеческими оценками интерпретируемости. Тема тем лучше, чем выше её когерентность.

Ethnic topics examples

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

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

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

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

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

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

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

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

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

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

Заключение

- BigARTM — открытая библиотека с поддержкой ARTM и мультимодальных моделей.
- Комбинация восьми регуляризаторов в задаче поиска этнорелевантных тем показала свое превосходство над LDA.
- Дальнейшие исследования предполагают работу с короткими текстами коллекции ВКонтакте, использование биграмм этнонимов и модальностей меток времени и геотегов для мониторинга изменения тем во времени и по регионам.

Contacts: bigartm.org, great-mel@yandex.ru

