

• Прикладные модели машинного обучения •
**Вероятностное
тематическое моделирование**

Воронцов Константин Вячеславович

k.v.vorontsov@phystech.edu

<http://www.MachineLearning.ru/wiki?title=User:Vokov>

Этот курс доступен на странице вики-ресурса

<http://www.MachineLearning.ru/wiki>

«Машинное обучение (курс лекций, К.В.Воронцов)»

МФТИ • 6 марта 2026

- 1 Вероятностное тематическое моделирование**
 - Цели, приложения, постановка задачи
 - Максимизация на единичных симплексах
 - Аддитивная регуляризация тематических моделей
- 2 Регуляризаторы и модальности**
 - PLSA, LDA, фоновые темы и декоррелирование
 - Мультимодальные тематические модели
 - Комбинирование регуляризаторов
- 3 От «мешка слов» к тематическому вниманию**
 - Нейросетевые тематические модели
 - Тематическая модель локального контекста
 - Аналогия A-ARTM с моделью само-внимания

Эволюция подходов машинного обучения в анализе текстов

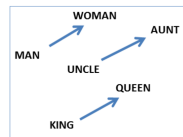
Анализ текстов 15 лет назад: пирамида NLP

- морфологический анализ, лемматизация, опечатки
- синтаксический анализ, выделение терминов, NER
- семантический анализ, выделение фактов, тем



Контекстно независимые эмбединги слов в вероятностных моделях языка на основе матричных разложений

- модели дистрибутивной семантики
word2vec [Mikolov, 2013], FastText [Bojanowski, 2016]
- тематические модели LDA [Blei, 2003], ARTM [2014]



Контекстно зависимые нейросетевые эмбединги

- рекуррентные нейронные сети: LSTM [1997]
- модели внимания и трансформеры: NMT [2015], BERT [2018], GPT-3 [2020], GPT-4 [2023]
- тематические модели внимания?

$$\text{softmax} \left(\frac{\begin{matrix} Q & K^T \\ \begin{matrix} \text{grid} & \times & \text{grid} \end{matrix} \end{matrix}}{\sqrt{d}} \right) \begin{matrix} V \\ \text{grid} \end{matrix}$$

Задача вероятностного тематического моделирования

Дано: коллекция текстовых документов как «мешков-слов»

- n_{dw} — частота слова (терма) $w \in W$ в документе $d \in D$
- $|T|$ — сколько тем хотим определить в коллекции D

Найти: тематическую языковую модель (в.п. $D \times W \times T$)

- $p(w|d) = \sum_{t \in T} p(w|d, t) p(t|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$
- $p(w|t) = \phi_{wt}$ — из каких слов w состоит каждая тема $t \in T$
- $p(t|d) = \theta_{td}$ — из каких тем t состоит каждый документ d

Критерий — log-правдоподобие языковой модели:

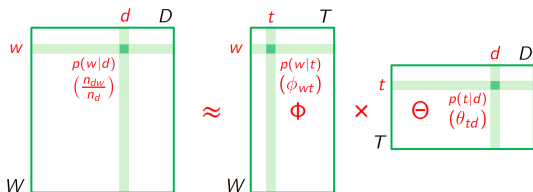
$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} \rightarrow \max_{\Phi, \Theta}$$

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

Hofmann T. Probabilistic Latent Semantic Indexing. ACM SIGIR, 1999.

Три интерпретации задачи тематического моделирования

1. Мягкая би-кластеризация документов и слов по темам
2. Матричное разложение — низкоранговое, стохастическое:



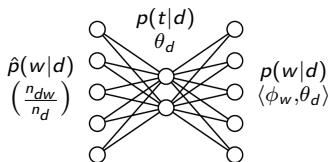
3. Автокодировщик документов в тематические эмбединги:

— кодировщик $f_{\Phi} : \frac{n_{dw}}{n_d} \rightarrow \theta_d$

— декодировщик $g_{\Phi} : \theta_d \rightarrow \Phi \theta_d$

задача реконструкции:

$$\sum_d n_d \sum_w \hat{p}(w|d) \ln p(w|d) \rightarrow \min_{\Phi, \Theta}$$



Примеры приложений тематического моделирования

- разведочный анализ больших текстовых коллекций (сколько в коллекции тем, и о чём они)
- фильтрация и тематизация релевантного контента («поиск и классификация иголок в стоге сена»)
- поиск тематически схожих документов (document-by-document search)
- выявление и отслеживание цепочек событий в новостях (topic detection & tracking — DARPA, 1998)
- поиск тематических сообществ в социальных медиа
- классификация, категоризация, маршрутизация сообщений
- выявление паттернов потребления в банковских данных
- выявление метаболических путей при аннотации генома

J.Boyd-Graber, Yuening Hu, D.Mimno. Applications of Topic Models. 2017.

H.Jelodar et al. Latent Dirichlet allocation (LDA) and topic modeling: models, applications, a survey. 2019.

Цели и не-цели тематического моделирования

Цели:

- выявлять тематическую кластерную структуру коллекции, представляя результат в удобной для человека форме
- получать *интерпретируемые* тематические векторы (эмбединги) слов $p(t|w)$, слов-в-контексте $p(t|d, w)$, документов $p(t|d)$, фрагментов $p(t|s)$, объектов $p(t|x)$
- решать с их помощью задачи поиска, классификации, фильтрации, сегментации, суммаризации текстов

Не-цели:

- угадывать слова по контексту (это слабая модель языка)
- понимать смысл текста (тем не достаточно для этого)
- генерировать осмысленный текст (слабые эмбединги)

Необходимые условия экстремума и метод простых итераций

Операция нормировки вектора: $p_i = \operatorname{norm}_{i \in I}(x_i) = \frac{\max(x_i, 0)}{\sum_k \max(x_k, 0)}$

Лемма. Пусть $f(\Omega)$ непрерывно дифференцируема по Ω .
Если ω_j — вектор локального экстремума нашей задачи
и $\exists i: \omega_{ij} \frac{\partial f}{\partial \omega_{ij}} > 0$, то ω_j удовлетворяет системе уравнений

$$\omega_{ij} = \operatorname{norm}_{i \in I_j} \left(\omega_{ij} \frac{\partial f}{\partial \omega_{ij}} \right).$$

- Численное решение системы — методом простых итераций
- Векторы $\omega_j = 0$ отбрасываются как вырожденные решения
- Итерации похожи на градиентную оптимизацию:

$$\omega_{ij} := \omega_{ij} + \eta \frac{\partial f}{\partial \omega_{ij}},$$

но учитывают ограничения и не требуют подбора шага η

Доказательство леммы о максимизации на симплексах

Задача: $f(\Omega) \rightarrow \max_{\Omega}; \quad \sum_{i \in I_j} \omega_{ij} = 1, \quad \omega_{ij} \geq 0, \quad i \in I_j, \quad j \in J.$

Функция Лагранжа:

$$\mathcal{L}(\Omega; \mu, \lambda) = -f(\Omega) + \sum_{j \in J} \lambda_j \left(\sum_{i \in I_j} \omega_{ij} - 1 \right) - \sum_{j \in J} \sum_{i \in I_j} \mu_{ij} \omega_{ij}.$$

Условия Каруша–Куна–Таккера для вектора ω_j :

$$\frac{\partial f(\Omega)}{\partial \omega_{ij}} = \lambda_j - \mu_{ij}, \quad \mu_{ij} \omega_{ij} = 0, \quad \mu_{ij} \geq 0.$$

Умножим обе части первого равенства на ω_{ij} :

$$A_{ij} \equiv \omega_{ij} \frac{\partial f(\Omega)}{\partial \omega_{ij}} = \omega_{ij} \lambda_j.$$

Согласно условию леммы $\exists i: A_{ij} > 0$. Значит, $\lambda_j > 0$.

Если $\frac{\partial f(\Omega)}{\partial \omega_{ij}} < 0$ для некоторого i , то $\mu_{ij} > 0 \Rightarrow \omega_{ij} = 0$.

Тогда $\omega_{ij} \lambda_j = (A_{ij})_+; \quad \lambda_j = \sum_i (A_{ij})_+ \Rightarrow \omega_{ij} = \text{norm}_i(A_{ij}).$

Задачи, некорректно поставленные по Адамару

Задача *корректно поставлена по Адамару*, если её решение

- существует,
- единственно,
- устойчиво.



Жак Саломон Адамар
(1865–1963)

Задача матричного разложения *некорректно поставлена*:
если Φ, Θ — решение, то стохастические Φ', Θ' — тоже решения

- $\Phi'\Theta' = (\Phi S)(S^{-1}\Theta)$, $\text{rank } S = |T|$
- $f(\Phi', \Theta') \approx f(\Phi, \Theta)$

Регуляризация — доопределение решения
путём добавления критерия $+ \tau R(\Phi, \Theta)$

Скаляризация критериев: $+ \sum_i \tau_i R_i(\Phi, \Theta)$



А.Н.Тихонов
(1906–1993)

ARTM: аддитивная регуляризация тематических моделей

Максимизация логарифма правдоподобия с регуляризатором:

$$\sum_{d,w} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}; \quad R(\Phi, \Theta) = \sum_i \tau_i R_i(\Phi, \Theta)$$

Теорема (необходимое условие экстремума). Точка локального экстремума (Φ, Θ) удовлетворяет системе уравнений

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

EM-алгоритм — решение этой системы методом простой итерации

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

Доказательство (по лемме о максимизации на симплексах)

Применим лемму к \log -правдоподобию с регуляризатором:

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

$$\phi_{wt} = \operatorname{norm}_{w \in W} \left(\phi_{wt} \frac{\partial f}{\partial \phi_{wt}} \right) = \operatorname{norm}_{w \in W} \left(\phi_{wt} \sum_{d \in D} n_{dw} \frac{\theta_{td}}{p(w|d)} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right) =$$

$$= \operatorname{norm}_{w \in W} \left(\sum_{d \in D} n_{dw} p_{tdw} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right);$$

$$\theta_{td} = \operatorname{norm}_{t \in T} \left(\theta_{td} \frac{\partial f}{\partial \theta_{td}} \right) = \operatorname{norm}_{t \in T} \left(\theta_{td} \sum_{w \in W} n_{dw} \frac{\phi_{wt}}{p(w|d)} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right) =$$

$$= \operatorname{norm}_{t \in T} \left(\sum_{w \in W} n_{dw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right);$$

где определения вспомогательных переменных $p_{tdw} = \frac{\phi_{wt} \theta_{td}}{p(w|d)}$ выделяются в отдельные уравнения, и в итерационном процессе образуют E-шаг. ■

Свойства алгоритма EM (Expectation–Maximization) для ARTM

E-шаг — это формула Байеса:

$$p(t|d, w) = \frac{p(w, t|d)}{p(w|d)} = \frac{p(w|t)p(t|d)}{p(w|d)} = \frac{\phi_{wt}\theta_{td}}{\sum_s \phi_{ws}\theta_{sd}} = p_{tdw}$$

M-шаг — это оценки частот n_{wt} , n_{td} и условных вероятностей:

$n_{dwt} = n_{dw}p_{dwt}$ — частота тройки (d, w, t) в коллекции

$n_{wt} = \sum_d n_{dwt}$ — частота термина w в теме t

$n_{td} = \sum_w n_{dwt}$ — частота термов темы t в документе d

$n_t = \sum_{d,w} n_{dwt}$ — частота термов темы t в коллекции

$\phi_{wt} = \frac{n_{wt}}{n_t}$ и $\theta_{td} = \frac{n_{td}}{n_d}$ при отсутствии регуляризатора, $R = 0$

Тема t вырождена и исключается из модели (topic selection),

если $n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \leq 0$ для всех $w \in W$

Документ d вырожден и модель не может определить его темы,

если $n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \leq 0$ для всех $t \in T$

PLSA, LDA: первые и самые известные тематические модели

PLSA: probabilistic latent semantic analysis [Hofmann, 1999]
(вероятностный латентный семантический анализ):

$$R(\Phi, \Theta) = 0$$

M-шаг — частотные оценки условных вероятностей:

$$\phi_{wt} = \underset{w}{\text{norm}}(n_{wt}), \quad \theta_{td} = \underset{t}{\text{norm}}(n_{td})$$



Thomas
Hofmann

LDA: latent Dirichlet allocation (латентное размещение Дирихле):

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

M-шаг — частотные оценки со смещением β_w, α_t :

$$\phi_{wt} = \underset{w}{\text{norm}}(n_{wt} + \beta_w), \quad \theta_{td} = \underset{t}{\text{norm}}(n_{td} + \alpha_t)$$



David Blei

Hofmann T. Probabilistic latent semantic indexing. SIGIR 1999.

Blei D., Ng A., Jordan M. Latent Dirichlet Allocation. NIPS-2001. JMLR 2003.

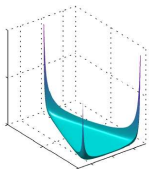
Распределение Дирихле

Гипотеза. Вектор-столбцы $\phi_t = (\phi_{wt})$ и $\theta_d = (\theta_{td})$ порождаются распределениями Дирихле, $\alpha \in \mathbb{R}^{|T|}$, $\beta \in \mathbb{R}^{|W|}$:

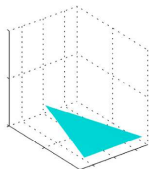
$$\text{Dir}(\phi_t | \beta) = \frac{\Gamma(\beta_0)}{\prod_w \Gamma(\beta_w)} \prod_w \phi_{wt}^{\beta_w - 1}, \quad \phi_{wt} > 0; \quad \beta_0 = \sum_w \beta_w, \quad \beta_w > 0;$$

$$\text{Dir}(\theta_d | \alpha) = \frac{\Gamma(\alpha_0)}{\prod_t \Gamma(\alpha_t)} \prod_t \theta_{td}^{\alpha_t - 1}, \quad \theta_{td} > 0; \quad \alpha_0 = \sum_t \alpha_t, \quad \alpha_t > 0;$$

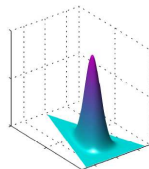
Пример. Распределение $\text{Dir}(\theta | \alpha)$ при $|T| = 3$, $\theta, \alpha \in \mathbb{R}^3$



$$\alpha_1 = \alpha_2 = \alpha_3 = 0.1$$

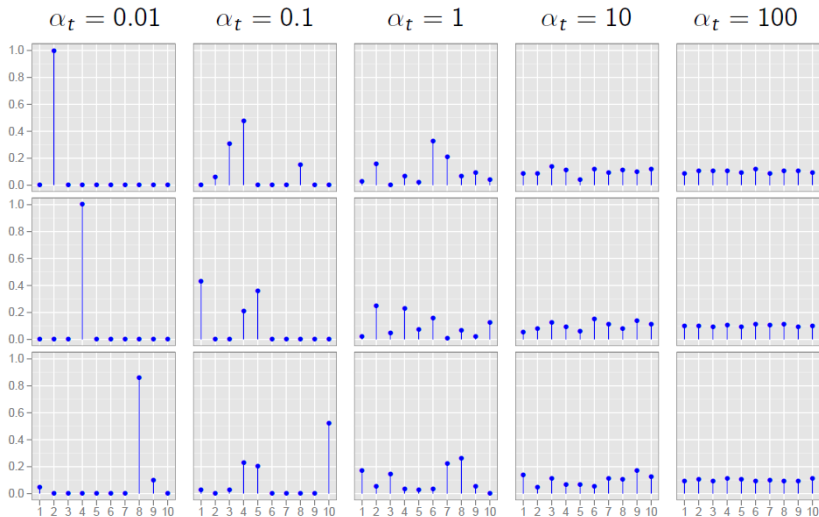


$$\alpha_1 = \alpha_2 = \alpha_3 = 1$$



$$\alpha_1 = \alpha_2 = \alpha_3 = 10$$

Пример. Выборки из трёх 10-мерных векторов $\theta \sim \text{Dir}(\theta|\alpha)$



Максимизация апостериорной вероятности для модели LDA

Совместное правдоподобие данных и модели:

$$\ln \prod_{d \in D} \prod_{w \in d} p(w, d | \Phi, \Theta)^{n_{dw}} \prod_{t \in T} \text{Dir}(\phi_t | \beta) \prod_{d \in D} \text{Dir}(\theta_d | \alpha) \rightarrow \max_{\Phi, \Theta}$$

Регуляризатор — логарифм априорного распределения:

$$R(\Phi, \Theta) = \sum_{t, w} (\beta_w - 1) \ln \phi_{wt} + \sum_{d, t} (\alpha_t - 1) \ln \theta_{td}$$

M-шаг — сглаженные или разреженные частотные оценки:

$$\phi_{wt} = \text{norm}_w(n_{wt} + \beta_w - 1), \quad \theta_{td} = \text{norm}_t(n_{td} + \alpha_t - 1).$$

при $\beta_w > 1$, $\alpha_t > 1$ — сглаживание,

при $0 < \beta_w < 1$, $0 < \alpha_t < 1$ — слабое разреживание,

при $\beta_w = 1$, $\alpha_t = 1$ априорное распределение равномерно, PLSA.

От байесовского обучения к аддитивной регуляризации

X — исходные данные, $\Omega = (\Phi, \Theta)$ — параметры модели

Байесовский вывод апостериорного распределения $p(\Omega|X)$
(громоздкий, приближённый) **только** ради точечной оценки Ω :

$$\text{Posterior}(\Omega|X, \gamma) = \frac{p(X|\Omega) \text{Prior}(\Omega|\gamma)}{\int p(X|\Omega) \text{Prior}(\Omega|\gamma) d\Omega}$$

$$\Omega = \arg \max_{\Omega} \text{Posterior}(\Omega|X, \gamma)$$

Максимизация апостериорной вероятности (MAP)

даёт точечную оценку Ω напрямую, без вывода Posterior:

$$\Omega = \arg \max_{\Omega} (\ln p(X|\Omega) + \ln \text{Prior}(\Omega|\gamma))$$

Многокритериальная аддитивная регуляризация (ARTM)

обобщает MAP на любые регуляризаторы и их комбинации:

$$\Omega = \arg \max_{\Omega} (\ln p(X|\Omega) + \sum_{i=1} \tau_i R_i(\Omega))$$

Обобщённая модель LDA (без ограничений на параметры)

Сглаживание ($\beta_{wt} > 0, \alpha_{td} > 0$) и разреживание ($\beta_{wt} < 0, \alpha_{td} < 0$):

$$R(\Phi, \Theta) = \sum_{t \in T} \sum_{w \in W} \beta_{wt} \ln \phi_{wt} + \sum_{d \in D} \sum_{t \in T} \alpha_{td} \ln \theta_{td}$$

Сглаживание фоновой темы t_ϕ с общей лексикой языка:

- $\beta_{wt_\phi} = \beta_0 p_\phi(w)$ — тема t_ϕ похожа на заданное $p_\phi(w)$
- $\alpha_{t_\phi d} = \alpha_0$ — общая лексика есть в каждом документе d

Сглаживание по «белым спискам» (seed words, seed topics):

- $\beta_{wt} = \beta_0 [w \in W_t]$ — термы из W_t должны быть в t
- $\alpha_{td} = \alpha_0 [t \in T_d]$ — темы из T_d должны быть в d

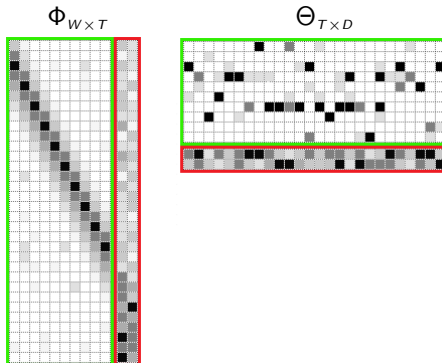
Разреживание по «чёрным спискам»:

- $\beta_{wt} = -\beta_0 [w \in W_t]$ — термов из W_t не должно быть в t
- $\alpha_{td} = -\alpha_0 [t \in T_d]$ — тем из T_d не должно быть в d

Разделение тем на предметные и фоновые

Предметные темы S содержат термины предметной области, $p(w|t)$, $p(t|d)$, $t \in S$ — разреженные, существенно различные

Фоновые темы B содержат слова общей лексики, $p(w|t)$, $p(t|d)$, $t \in B$ — существенно отличные от нуля



Регуляризатор декоррелирования тем

Цель: усилить различность тем; выделить в каждой теме лексическое ядро, отличающее её от других тем; способствовать переходу общей лексики в фоновые темы.

Минимизируем ковариации между вектор-столбцами ϕ_t :

$$R(\Phi) = -\frac{\tau}{2} \sum_{t \in T} \sum_{s \in T \setminus t} \sum_{w \in W} \phi_{wt} \phi_{ws} \rightarrow \max.$$

Подставляем в формулы M-шага, получаем ещё один вариант разреживания — контрастирование строк матрицы Φ (малые вероятности ϕ_{wt} в строке становятся ещё меньше):

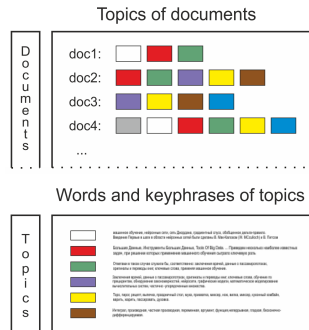
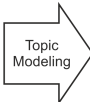
$$\phi_{wt} = \operatorname{norm}_w \left(n_{wt} - \tau \phi_{wt} \sum_{s \in T \setminus t} \phi_{ws} \right).$$

Мультимодальная тематическая модель

Тема может порождать термины различных модальностей:

$$p(\text{слово} | t), p(n\text{-грамма} | t), p(\text{автор} | t), p(\text{время} | t), p(\text{источник} | t),$$

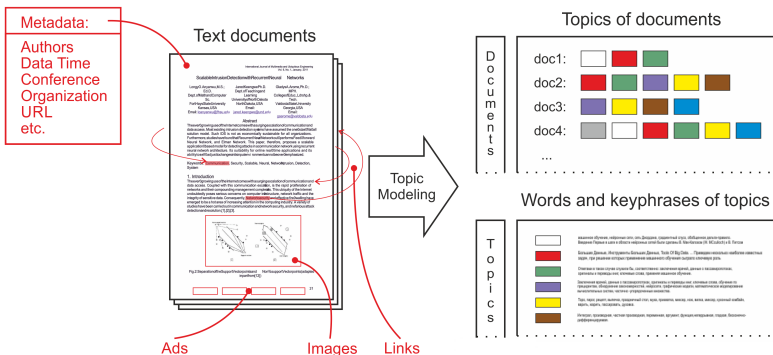
Metadata:
 Authors
 Data Time
 Conference
 Organization
 URL
 etc.



Мультимодальная тематическая модель

Тема может порождать термины различных модальностей:

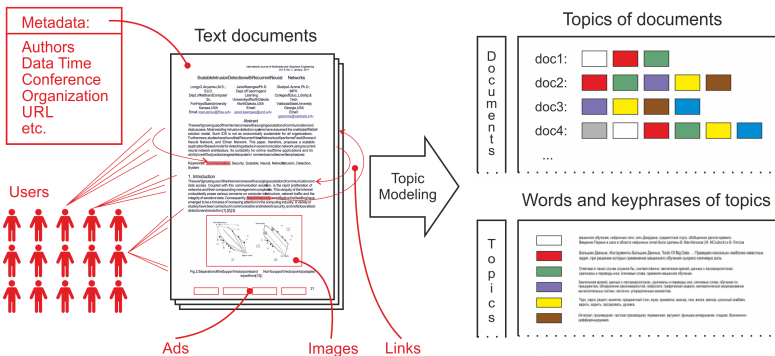
$p(\text{слово} | t)$, $p(n\text{-грамма} | t)$, $p(\text{автор} | t)$, $p(\text{время} | t)$, $p(\text{источник} | t)$,
 $p(\text{объект} | t)$, $p(\text{ссылка} | t)$, $p(\text{баннер} | t)$,



Мультимодальная тематическая модель

Тема может порождать термины различных модальностей:

$p(\text{слово} | t)$, $p(n\text{-грамма} | t)$, $p(\text{автор} | t)$, $p(\text{время} | t)$, $p(\text{источник} | t)$,
 $p(\text{объект} | t)$, $p(\text{ссылка} | t)$, $p(\text{баннер} | t)$, $p(\text{пользователь} | t)$

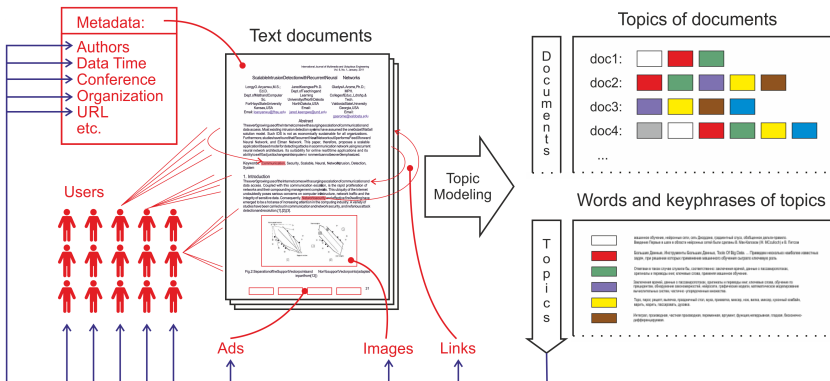


Мультимодальная тематическая модель

Тема может порождать термины различных модальностей:

$$p(\text{слово} | t), p(n\text{-грамма} | t), p(\text{автор} | t), p(\text{время} | t), p(\text{источник} | t),$$

$$p(\text{объект} | t), p(\text{ссылка} | t), p(\text{баннер} | t), p(\text{пользователь} | t)$$



Мультимодальная ARTM

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

Максимизация суммы \log правдоподобий с регуляризацией:

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

EM-алгоритм: метод простой итерации для системы уравнений

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

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

Пример 1. Модальности языков в мультязычных моделях

216 175 русско-английских пар статей Википедии, $|T| = 400$

Первые 10 слов и их частоты $p(w|t)$ в %:

Тема №68				Тема №79			
research	4.56	институт	6.03	goals	4.48	матч	6.02
technology	3.14	университет	3.35	league	3.99	игрок	5.56
engineering	2.63	программа	3.17	club	3.76	сборная	4.51
institute	2.37	учебный	2.75	season	3.49	фк	3.25
science	1.97	технический	2.70	scored	2.72	против	3.20
program	1.60	технология	2.30	cup	2.57	клуб	3.14
education	1.44	научный	1.76	goal	2.48	футболист	2.67
campus	1.43	исследование	1.67	apps	1.74	гол	2.65
management	1.38	наука	1.64	debut	1.69	забивать	2.53
programs	1.36	образование	1.47	match	1.67	команда	2.14

Ассессор оценил 396 тем из 400 как хорошо интерпретируемые.

Vorontsov, Frei, Apishev, Romov, Suvorova. BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections. AIST-2015.

Пример 1. Модальности языков в мультязычных моделях

216 175 русско-английских пар статей Википедии, $|T| = 400$

Первые 10 слов и их частоты $p(w|t)$ в %:

Тема №88				Тема №251			
opera	7.36	опера	7.82	windows	8.00	windows	6.05
conductor	1.69	оперный	3.13	microsoft	4.03	microsoft	3.76
orchestra	1.14	дирижер	2.82	server	2.93	версия	1.86
wagner	0.97	певец	1.65	software	1.38	приложение	1.86
soprano	0.78	певица	1.51	user	1.03	сервер	1.63
performance	0.78	театр	1.14	security	0.92	server	1.54
mozart	0.74	партия	1.05	mitchell	0.82	программный	1.08
sang	0.70	сопрано	0.97	oracle	0.82	пользователь	1.04
singing	0.69	вагнер	0.90	enterprise	0.78	обеспечение	1.02
operas	0.68	оркестр	0.82	users	0.78	система	0.96

Ассессор оценил 396 тем из 400 как хорошо интерпретируемые.

Vorontsov, Frei, Apishev, Romov, Suvorova. BigARTM: Open Source Library for Regularized Multimodal Topic Modeling of Large Collections. AIST-2015.

Пример 2. Модальности униграмм и биграмм

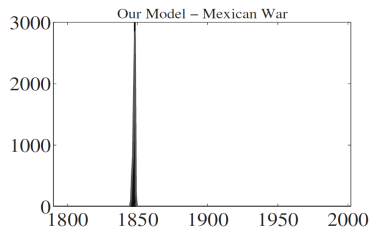
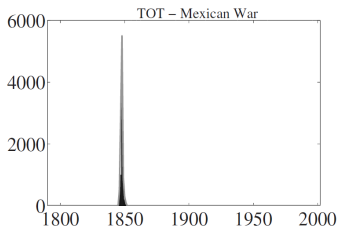
Коллекция 1000 статей конференций ММРО, ИОИ на русском

распознавание образов в биоинформатике		теория вычислительной сложности	
униграммы	биграмммы	униграммы	биграмммы
объект	задача распознавания	задача	разделять множества
задача	множество мотивов	множество	конечное множество
множество	система масок	подмножество	условие задачи
мотив	вторичная структура	условие	задача о покрытии
разрешимость	структура белка	класс	покрытие множества
выборка	распознавание вторичной	решение	сильный смысл
маска	состояние объекта	конечный	разделяющий комитет
распознавание	обучающая выборка	число	минимальный аффинный
информативность	оценка информативности	аффинный	аффинный комитет
состояние	множество объектов	случай	аффинный разделяющий
закономерность	разрешимость задачи	покрытие	общее положение
система	критерий разрешимости	общий	множество точек
структура	информативность мотива	пространство	случай задачи
значение	первичная структура	схема	общий случай
регулярность	тупиковое множество	комитет	задача MASC

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

Пример 3. Модальности времени и n -грамм

Коллекция еженедельных выступлений президентов США



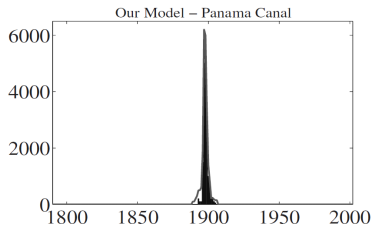
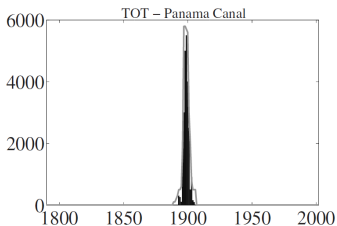
1. mexico	8. territory
2. texas	9. army
3. war	10. peace
4. mexican	11. act
5. united	12. policy
6. country	13. foreign
7. government	14. citizens

1. east bank	8. military
2. american coins	9. general herrera
3. mexican flag	10. foreign coin
4. separate independent	11. military usurper
5. american commonwealth	12. mexican treasury
6. mexican population	13. invaded texas
7. texan troops	14. veteran troops

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

Пример 3. Модальности времени и n -грамм

Коллекция еженедельных выступлений президентов США



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

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

Shoib Jameel, Wai Lam. An N -gram topic model for time-stamped documents. 2013.

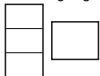
Регуляризаторы для мультимодальных тематических моделей

supervised



Модальности меток классов или категорий для задач классификации и категоризации текстов

multilanguage



Модальность языков и регуляризация со словарём $\pi_{uwt} = p(u|w, t)$ переводов с языка k на ℓ :

$$R(\Phi, \Pi) = \tau \sum_{u \in W^k} \sum_{t \in T} n_{ut} \ln \sum_{w \in W^\ell} \pi_{uwt} \phi_{wt}$$

temporal



Модальность интервалов времени i , сглаживание тем:

$$R(\Phi) = -\tau \sum_{i \in I} \sum_{t \in T} |\phi_{it} - \phi_{i-1,t}|$$

geospatial



Модальность геолокаций g с близостью $S_{gg'}$:

$$R(\Phi) = -\frac{\tau}{2} \sum_{g, g' \in G} S_{gg'} \sum_{t \in T} n_t^2 \left(\frac{\phi_{gt}}{n_g} - \frac{\phi_{g't}}{n_{g'}} \right)^2$$

Регуляризаторы для учёта взаимосвязей и зависимостей

regression



Линейная модель регрессии $\hat{y}_d = \langle v, \theta_d \rangle$ документов:

$$R(\Theta, v) = -\tau \sum_{d \in D} \left(y_d - \sum_{t \in T} v_t \theta_{td} \right)^2$$

biterm



Связи сочетаемости слов (n_{uv} — частота битерма):

$$R(\Phi) = \tau \sum_{u \in W} \sum_{v \in W} n_{uv} \ln \sum_{t \in T} n_t \phi_{ut} \phi_{vt}$$

relational



Связи между документами (n_{dc} — число ссылок):

$$R(\Theta) = \tau \sum_{d, c \in D} n_{dc} \sum_{t \in T} \theta_{td} \theta_{tc}$$

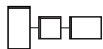
hierarchy



Связи родительских тем t с дочерними подтемами s , $\psi_{st} = p(s|t)$ — псевдодокумент родительской темы:

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

Тематические модели сложно структурированных данных



Модель трёхматричного разложения «авторы-темы»:

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \sum_{a \in A} \phi_{wt} \theta_{ta} \psi_{ad} + R(\Phi, \Theta, \Psi) \rightarrow \max_{\Phi, \Theta, \Psi}$$

hypergraph



Гиперграфовая модель транзакционных данных:

$$\sum_{k \in K} \tau_k \sum_{x \in E_k} n_{kx} \ln \left(\sum_{t \in T} \pi_t \prod_{v \in X} \phi_{vt} \right) + R(\Phi, \pi) \rightarrow \max_{\Phi, \pi}$$

sentence



Тематическая модель предложений:

$$\sum_{d \in D} \sum_{s \in S_d} \ln \left(\sum_{t \in T} \theta_{td} \prod_{w \in s} \phi_{wt}^{n_{sw}} \right) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

segmentation

Регуляризация E-шага, $\Pi = (p_{tdw} = p(t|d, w))_{T \times D \times W}$:

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

Воронцов К. В. Вероятностное тематическое моделирование: теория регуляризации ARTM и библиотека с открытым кодом BigARTM. 2025. 224 с.

Библиотека BigARTM

Ключевые возможности:

- Встроенная библиотека регуляризаторов и метрик качества
- Большие данные: коллекция не хранится в памяти
- Онлайн-параллельный мультимодальный ARTM

Сообщество:

- Открытый код <https://github.com/bigartm>
(discussion group, issue tracker, pull requests)
- Документация <http://bigartm.org>



Лицензия и среда разработки:

- Свободная коммерческая лицензия (BSD 3-Clause)
- Кросс-платформенность: Windows, Linux, MacOS (32/64 bit)
- Интерфейсы API: command-line, C++, and Python

K. Vorontsov, O. Frei, M. Apishev, P. Romov, M. Suvorova. BigARTM: open source library for regularized multimodal topic modeling of large collections. 2015.

Разведочный поиск в технологических блогах

Цель: поиск документов

по длинным текстовым запросам

— Habr.ru (175К документов),

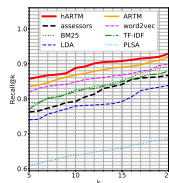
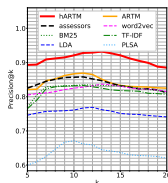
— TechCrunch.com (760К док.).

Регуляризаторы:

$$\mathcal{L} \left(\begin{array}{c} \text{PLSA} \\ \Phi \quad \Theta \end{array} \right) + R \left(\begin{array}{c} \text{hierarchy} \\ \text{graph} \end{array} \right) + R \left(\begin{array}{c} \text{interpretable} \\ \text{matrix} \end{array} \right) + R \left(\begin{array}{c} \text{multimodal} \\ \text{stack} \end{array} \right) + R \left(\begin{array}{c} \text{n-gram} \\ \text{grid} \end{array} \right) \rightarrow \max$$

Результаты:

- Точность и полнота **93%**, превосходит ассессоров и другие методы (tf-idf, BM25, word2vec, PLSA, LDA, ARTM).
- Увеличилась оптимальная размерность векторов:
200 → 1400 (Habr.ru), 475 → 2800 (TechCrunch.com).



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

Поиск и классификация этно-релевантных тем в соцсетях

Цель: выявление как можно большего числа тем о национальностях и межнациональных отношениях (затравка — словарь 300 этнонимов).

Регуляризаторы:

$$\mathcal{L} \left(\begin{array}{|c|} \hline \text{PLSA} \\ \hline \Phi \quad \Theta \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{seed words} \\ \hline \text{[bar chart]} \quad \square \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{interpretable} \\ \hline \text{[bar chart]} \quad \text{[scatter plot]} \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{multimodal} \\ \hline \text{[stacked bars]} \quad \square \\ \hline \end{array} \right) \\ + R \left(\begin{array}{|c|} \hline \text{temporal} \\ \hline \text{[line graph]} \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{geospatial} \\ \hline \text{[map]} \\ \hline \end{array} \right) + R \left(\begin{array}{|c|} \hline \text{sentiment} \\ \hline \text{[sentiment icons]} \\ \hline \end{array} \right) \rightarrow \max$$

(японцы): японский, япония, корей, китайский, жилища, авария, фукусима, цунами, сообщать, омега, станция, хатико, район, правительство, атомный,
(норвежцы): дитя, ребенок, родился, детский, семья, воспитанный, право, возраст, отец, воспитание, норвежский, родительский, родить, мальчик, взрослый, омега, сын,
(венесуэльцы): куба, кастро, венесуэла, часеб, президент, уго, мадура, боливия, фидель, глава, латиноский, венесуэльский, лидер, боливарианский, президентский, альенде, гевару,
(китайцы): китайский, россия, производство, китай, продукция, страна, предприятие, компания, технология, военный, регион, производить, производственный, промышленность, российский, экономический, кыр,
(азербайджанцы): русский, азербайжан, азербайджанец, россия, азербайджанский, таксист, диаспора, анапа, народ, москва, страна, армянин, слово, рынок,
(грузины): грузинский, спецназ, военной, август, баташева, российский, спецназовца, миротворец, операция, дурын, бригада, миротворческой, абхазия, группа, войска, русский, цхинвалс,
(осетины): конституция, осетия, аминат, русский, осетинский, южный, северный, россия, война, республика, вопрос, алакай, российский, население, конфликт,
(цыгане): народный, цыган, цыганка, хордовый, место, страна, денга, время, работать, жилье, жить, рука, дом, цыганский, наркоманка.

Результаты: число релевантных тем: 45 (LDA) \rightarrow 83 (ARTM).

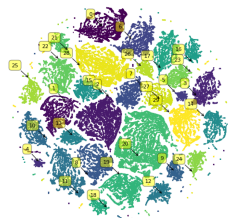
M. Apishev, S. Koltcov, O. Koltsova, S. Nikolenko, K. Vorontsov. Additive regularization for topic modeling in sociological studies of user-generated text content. MICAI, 2016.

–, –, –, –, –. Mining ethnic content online with additively regularized topic models. 2016.

Тематическая модель банковских транзакционных данных

Цель: Выявление паттернов потребительского поведения клиентов банка, причём

- документы \rightarrow клиенты,
- слова \rightarrow MCC-коды продавцов.



Регуляризаторы:

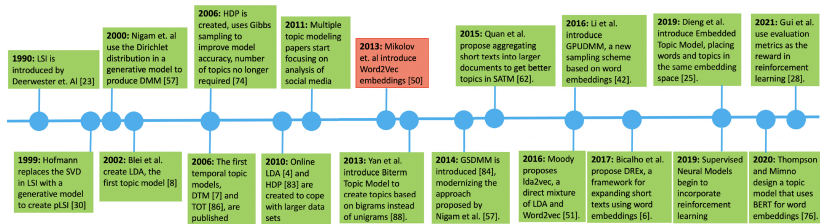
$$\mathcal{L}\left(\begin{array}{c} \text{PLSA} \\ \Phi \quad \Theta \end{array}\right) + R\left(\begin{array}{c} \text{interpretable} \\ \begin{array}{|c|} \hline \text{[Bar Chart Icon]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Scatter Plot Icon]} \\ \hline \end{array} \end{array}\right) + R\left(\begin{array}{c} \text{multimodal} \\ \begin{array}{|c|} \hline \text{[Stacked Bars Icon]} \\ \hline \end{array} \quad \begin{array}{|c|} \hline \text{[Box Icon]} \\ \hline \end{array} \end{array}\right) + R\left(\begin{array}{c} \text{supervised} \\ \begin{array}{|c|} \hline \text{[Decision Tree Icon]} \\ \hline \end{array} \end{array}\right) \rightarrow \max$$

Результаты:

- темы — паттерны потребительского поведения
- предсказание пола, возраста, достатка клиентов

E.Egorov, F.Nikitin, A.Goncharov, V.Alekseev, K.Vorontsov. Topic modelling for extracting behavioral patterns from transactions data. 2019.

Neural Topic Models — эволюция PTM в сторону LLM



Как «объединить лучшее от двух миров»?

- **Neural:** общность, качество, предобучение, генерация
- **Topics:** интерпретируемость, полнота, простота, скорость

Что объединяет PTM и LLM, и что их разобщает:

- ⊕ обе — вероятностные языковые модели,
- ⊕ обе — автокодировщики, векторные представления текста
- ⊖ **PTM:** байесовское обучение, архитектура MF, мешок слов

Rob Churchill, Lisa Singh. The Evolution of Topic Modeling. 2022.

Нейросетевая тематическая модель Contextual-Top2Vec

Вместо РТМ — сумма технологий:

- 1 векторизация токенов (Sentence-BERT)
- 2 векторизация предложений скользящим окном в 50 токенов (mean pooling)
- 3 понижение размерности векторов (UMAP)
- 4 иерархическая кластеризация (hDbSCAN) с автоматическим определением числа тем
- 5 иерархическое укрупнение тем слиянием мелких кластеров с ближайшими соседями (Top2Vec)
- 6 разбиение документа на монотематические сегменты
- 7 $p(t|d)$ = доля векторов данной темы в документе
- 8 именованые тем: поиск фраз, ближайших к центроиду темы



Dimo Angelov. Top2vec: Distributed representations of topics. 2020.

D. Angelov, D. Inkpen. Topic modeling: contextual token embeddings are all you need. 2024.

Нейросетевая тематическая модель Contextual-Top2Vec

Достоинства:

- модель BERT предобучена по большим внешним данным, поэтому качество тем не зависит от размера коллекции
- документ разбивается на монотематические сегменты
- тема описывается фразами, а не отдельными словами

Недостатки:

- долго-дорого, особенно на больших коллекциях
- инкрементное добавление документов не предполагается

Идеи новой модели тематического внимания в ARTM:

- вместо «мешка слов» — последовательность w_1, \dots, w_n
- вместо документов — локальные контексты как в BERT
- вместо $p(w|t)$ — тематические эмбединги $\phi_{tw} = p(t|w)$

Воронцов К. В. Вероятностное тематическое моделирование: теория регуляризации ARTM и библиотека с открытым кодом BigARTM. 2025. 224 с.

Контекстная тематическая модель Attentive ARTM

Дано: коллекция текстовых документов, w_1, \dots, w_n
 $C_i \subset \{1, \dots, n\}$ — локальный контекст (окружение) термина w_i
 α_{ci} — коэффициент внимания, вес термина w_c из C_i для w_i

Найти: $\phi_{tw} = p(t|w)$ — параметры тематической модели

$$p(w|C_i) = \sum_{t \in T} p(w|t)p(t|C_i) = \sum_{t \in T} p(t|w) \frac{p(w)}{p(t)} p(t|C_i)$$

$$p(t|C_i) \equiv \theta_{ti} = \sum_{c \in C_i} \alpha_{ci} p(t|w_c), \quad \sum_{c \in C_i} \alpha_{ci} = 1, \quad \alpha_{ci} \geq 0$$

Критерий: максимум \log правдоподобия с регуляризатором R :

$$\sum_{i=1}^n \ln \left(\sum_{t \in T} \phi_{tw_i} \frac{p(w_i)}{p(t)} \sum_{c \in C_i} \alpha_{ci} \phi_{tw_c} \right) + R(\Phi) \rightarrow \max_{\Phi}$$

EM-алгоритм для модели Attentive ARTM

EM-алгоритм: метод простой итерации для системы уравнений

$$p_{ti} \equiv p(t|C_i, w_i) = \operatorname{norm}_{t \in T}(\phi_{tw_i} \theta_{ti} / n_t), \quad \theta_{ti} = \sum_{c \in C_i} \alpha_{ci} \phi_{tw_c}$$

$$\phi_{tw} = \operatorname{norm}_{t \in T} \left(n_{tw} + \phi_{tw} \frac{\partial R}{\partial \phi_{tw}} \right), \quad n_{tw} = \sum_{i=1}^n p_{ti}[w_i = w], \quad n_t = \sum_{i=1}^n p_{ti}$$

Отличия от обычных $\Phi\Theta$ -моделей:

- это не матричное разложение, тут нет матрицы Θ
- похоже на модель внимания без обучаемых параметров
- θ_{ti} можно вычислять скользящим средним за $O(n)$
- при $C_i = \text{документ}$, $\alpha_{ci} = \frac{1}{n_d}$, это модель из [Ирхин, 2020]

И.А.Ирхин, В.Г.Булатов, К.В.Воронцов. Аддитивная регуляризация тематических моделей с быстрой векторизацией текста, 2020.

Аналогия Attentive ARTM с моделью само-внимания

Контекстный тематический вектор на выходе E-шага:

$$p(t|C_i, w_i) \equiv p_{ti} = \operatorname{norm}_{t \in T} \left(\sum_{c \in C_i} \phi_{twc} \alpha_{ci} \frac{1}{p(t)} \phi_{tw_i} \right)$$

Контекстный вектор на выходе модели само-внимания:

$$h_i = \sum_{c \in C_i} W_v x_c \alpha_{ci} = \sum_{c \in C_i} W_v x_c \operatorname{SoftMax}_{c \in C_i} \langle W_k x_c, W_q x_i \rangle$$

Сходство:

- вектор термина w_i трансформируется в контекстный вектор
- путём усреднения векторов термов w_c из его контекста,
- наиболее схожих с ним по тематике

Отличия локализованного E-шага:

- адамарово умножение вектора ϕ_{w_c} на вектор-фильтр ϕ_{w_i}
- нет обучаемых матриц W_q, W_k, W_v как у модели внимания
- проецирование итогового вектора на единичный симплекс

- Тематическое моделирование — мягкая кластеризация, автокодировщик, стохастическое матричное разложение, но при этом весьма посредственная языковая модель.
- Самые известные модели — PLSA [1999] и LDA [2001]. Проблема — неудобства байесовского обучения.
- *Аддитивная регуляризация ARTM* — многокритериальная оптимизация с возможностью комбинирования моделей и построения моделей с заданными свойствами. Проблема — подбор коэффициентов регуляризации.
- *Лемма о максимизации на единичных симплексах* радикально упрощает как теорию, так и практику РТМ. Эта лемма применима далеко за пределами РТМ, например, для релаксации задач дискретной оптимизации.
- Развитие ARTM — тематические модели внимания.