

# Оптимизация на единичных симплексах для обучения монотонных нейронных сетей, матричных разложений и вероятностных тематических моделей

Воронцов Константин Вячеславович

[k.vorontsov@iai.msu.ru](mailto:k.vorontsov@iai.msu.ru)

д.ф.-м.н., профессор РАН

рук. лаб. машинного обучения и семантического анализа  
Института ИИ МГУ,

зав. каф. математических методов прогнозирования ВМК МГУ,

зав. каф. интеллектуальных систем МФТИ,

г.н.с. ФИЦ «Информатика и управление» РАН

XXIV международная конференция MOTOR  
Новосибирск • 7–11 июля 2025

## 1 Оптимизация на единичных симплексах

- постановка задачи, метод и его сходимость
- монотонные нейронные сети
- неотрицательные матричные разложения

## 2 Вероятностное тематическое моделирование

- постановка задачи
- аддитивная регуляризация тематических моделей
- библиотека BigARTM и прикладные задачи

## 3 На пути к тематической модели внимания

- регуляризация Е-шага
- тематическая модель линейного текста
- тематическая модель локальных контекстов

## Задача максимизации функции на единичных симплексах

Пусть  $\Omega = (\omega_j)_{j \in J}$  — набор нормированных неотрицательных векторов  $\omega_j = (\omega_{ij})_{i \in I_j}$  различных размерностей  $|I_j|$ :

$$\Omega = \left( \begin{array}{c} \text{[yellow]} \\ \text{[blue]} \\ \text{[blue]} \\ \text{[blue]} \\ \text{[purple]} \\ \text{[purple]} \\ \text{[pink]} \\ \text{[pink]} \\ \text{[pink]} \\ \text{[pink]} \\ \text{[green]} \\ \text{[green]} \\ \text{[green]} \\ \text{[green]} \end{array} \right)$$

Задача максимизации функции  $f(\Omega)$  на единичных симплексах:

$$\begin{cases} f(\Omega) \rightarrow \max; \\ \sum_{i \in I_j} \omega_{ij} = 1, \quad \omega_{ij} \geq 0, \quad i \in I_j, \quad j \in J. \end{cases}$$

---

Vorontsov K. V. Rethinking probabilistic topic modeling from the point of view of classical non-Bayesian regularization. 2023.

## Необходимые условия экстремума и метод простых итераций

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

**Лемма.** Пусть  $f(\Omega)$  непрерывно дифференцируема по  $\Omega$ .

Если  $\omega_j$  — вектор локального экстремума нашей задачи и  $\exists i: \omega_{ij} \frac{\partial f}{\partial \omega_{ij}} > 0$ , то  $\omega_j$  удовлетворяет системе уравнений

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

- 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  $\eta$

## 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  $\omega_j$ :

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

Умножим обе части первого равенства на  $\omega_{ij}$ :

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

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

## Теорема о сходимости итерационного процесса

$$\omega_{ij}^{t+1} = \operatorname{norm}_{i \in I_j} \left( \omega_{ij}^t \frac{\partial f(\Omega^t)}{\partial \omega_{ij}^t} \right)$$

**Теорема.** Пусть  $f(\Omega)$  — ограниченная сверху, непрерывно дифференцируемая функция, и все  $\Omega^t$ , начиная с некоторой итерации  $t^0$  обладают свойствами:

- $\forall j \in J \quad \forall i \in I_j \quad \omega_{ij}^t = 0 \rightarrow \omega_{ij}^{t+1} = 0$  (сохранение нулей)
- $\exists \varepsilon > 0 \quad \forall j \in J \quad \forall i \in I_j \quad \omega_{ij}^t \notin (0, \varepsilon)$  (отделимость от нуля)
- $\exists \delta > 0 \quad \forall j \in J \quad \exists i \in I_j \quad \omega_{ij}^t \frac{\partial f(\Omega^t)}{\partial \omega_{ij}} \geq \delta$  (невырожденность)
- $\exists \lambda > 0 \quad f(\Omega^{t+1}) - f(\Omega^t) \geq \lambda H(\Omega^t)$  (монотонный процесс)

Тогда  $|\omega_{ij}^{t+1} - \omega_{ij}^t| \rightarrow 0$  при  $t \rightarrow \infty$ .

## Открытая проблема: неудобное четвёртое условие

**Определение.**  $H(\Omega^t)$  есть линейное приближение приращения функции  $f$  в окрестности точки  $\Omega^t$ :

$$f(\Omega^{t+1}) - f(\Omega^t) = H(\Omega^t) + o(\Delta\Omega^t)$$

**Лемма.** Квадратичное представление функции  $H(\Omega)$ :

$$H(\Omega) = \frac{1}{2} \sum_{j \in J} \sum_{i, k \in I_j} \left( \frac{\partial f(\Omega)}{\partial \omega_{ij}} - \frac{\partial f(\Omega)}{\partial \omega_{kj}} \right)^2 \omega_{ij} \omega_{kj}$$

Следовательно,  $H(\Omega^t) \geq 0$ .

$f(\Omega^{t+1}) - f(\Omega^t) \approx H(\Omega^t)$  — согласно определению;

$f(\Omega^{t+1}) - f(\Omega^t) \geq \lambda H(\Omega^t)$ , начиная с некоторой итерации  $t$  при некотором  $\lambda > 0$  — хотелось бы получить это как результат, а не вводить как предположение. Доказать это пока не удалось.

## Промежуточные итоги и направления исследований

- Метод похож на обычную градиентную оптимизацию, но не требует подбора градиентного шага  $\eta$
- Ограничения неотрицательности и нормировки могут накладываться не на все векторы, а лишь на некоторые
- Операция `norm` может приводить к обнулению части координат, следовательно, к разреживанию векторов  $\omega_j$
- **Приложения:**
  - вероятностное тематическое моделирование
  - неотрицательные матричные разложения
  - монотонные нейронные сети
  - сети для аппроксимации функций распределения
- **Открытая проблема:** упростить четвёртое условие в теореме сходимости (оно представляется избыточным)
- **Открытая проблема:** оценить скорость сходимости

## Монотонная аппроксимация: постановка задачи

**Дано:** выборка  $(x_i, y_i)_{i=1}^m$ ,  $x_i \in \mathbb{R}^n$

**Найти:** предсказательную модель  $a(x, w)$  с параметром  $w$ , обладающую свойством монотонности:

$$x \leqslant x' \Rightarrow a(x, w) \leqslant a(x', w)$$

**Критерий:** минимум эмпирического риска при зашумлённых данных (не обязательно  $x_i \leqslant x_j \Rightarrow y_i \leqslant y_j$ )

$$\sum_{i=1}^m \mathcal{L}(a(x_i, w), y_i) \rightarrow \min_w$$

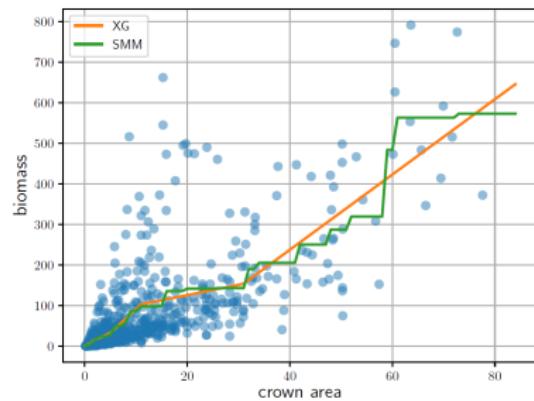
Например, для монотонной регрессии (isotonic regression)

$$\sum_{i=1}^m (a(x_i, w) - y_i)^2 \rightarrow \min_w$$

Линейная модель  $a(x, w) = \langle x, w \rangle$  монотонна  $\Leftrightarrow w \geqslant 0$

## Ограничение монотонности: приложения и интерпретации

- учёт априорных знаний вида «чем больше значение признака, тем выше отклик»
- агрегирование нескольких моделей в ансамбль
- моделирование многомерных функций распределения
- синтез интерпретируемых векторных представлений в глубоких нейронных сетях



### Пример.

Зависимость биомассы  
(и связанного углерода)  
от площади кроны деревьев

*Christian Igel. Smooth monotonic networks. 2023.*

*J.-R. Cano et al. Monotonic classification: An overview on algorithms, performance measures and data sets. Neurocomputing, 2019.*

## Монотонная Min-Max нейронная сеть

Трёхслойная сеть с двумя слоями min и max пулинга:

$$a(x, w) = \min_{k \in K} \max_{h \in H_k} (\langle w_{kh}, x \rangle - b_{kh}), \quad w_{kh} \geq 0$$

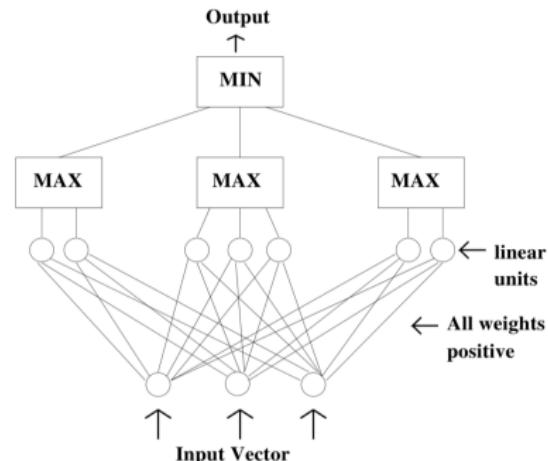
Репараметризация:  $w_{kh} = \exp(z_{kh})$

### Преимущества:

- можно использовать BackProp
- можно встраивать в DeepNN
- доказана аппроксимируемость

### Недостатки:

- кусочно-линейная модель
- затухание градиентов
- недообучение, переусложнение
- чувствительность к инициализации



Joseph Sill. Monotonic networks. NeurIPS, 1997.

## Сглаженная монотонная нейронная сеть (Smooth Min-Max)

Введём функцию LogSumExp с параметром  $\beta$ :

$$\text{LSE}_\beta(z_i) = \frac{1}{\beta} \ln \sum_{i \in I} \exp(\beta z_i) \rightarrow \begin{cases} \max_i(z_i), & \beta \rightarrow +\infty \\ \min_i(z_i), & \beta \rightarrow -\infty \end{cases}$$

Min-Max-сеть с заменой  $\min$  и  $\max$  на их сглаженные аналоги:

$$a(x, w) = \underset{k \in K}{\text{LSE}_{-\beta}} \underset{h \in H_k}{\text{LSE}_{+\beta}} (\langle w_{kh}, x \rangle - b_{kh}), \quad w_{kh} \geq 0$$

### Преимущества:

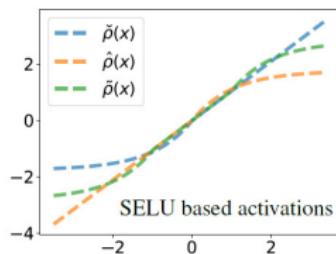
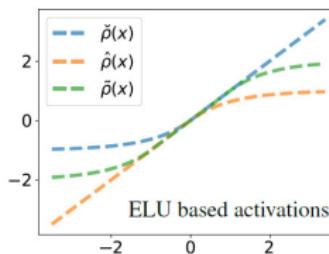
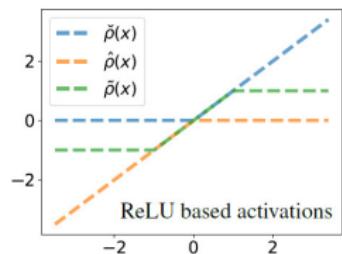
- гладкая аппроксимирующая монотонная функция
- доказаны асимптотические аппроксимационные свойства
- можно использовать BackProp, встраивать в DeepNN
- обобщающая способность существенно выше, чем у Min-Max

## Ограниченнная монотонная нейронная сеть (Constrained MNN)

Двухслойная сеть с монотонными функциями активации  
на скрытом слое: выпуклой, вогнутой и выпукло-вогнутой

$$a(x, w) = \sum_{k \in K} w_k^2 \sigma_k(\langle w_k^1, x \rangle - b_k), \quad w_k^L \geq 0$$

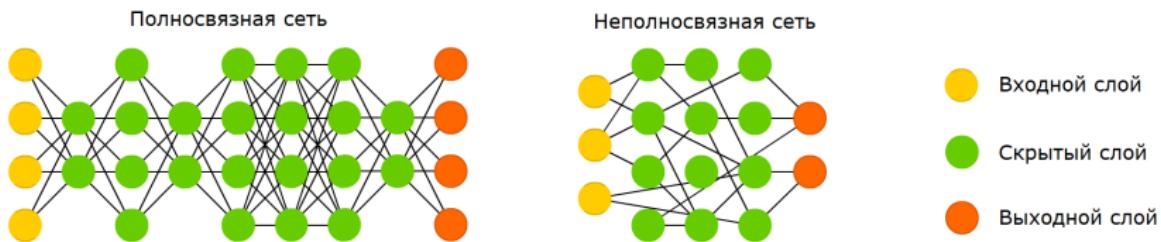
где  $\sigma_k \in \{\check{\rho}, \hat{\rho}, \tilde{\rho}\}$



Преимущества те же, что у нейронной сети Smooth Min-Max,  
структуря решения **разреженная и интерпретируемая**

## Глубокие нейронные сети (Deep Neural Network, DNN)

Вычисление сети:  $x^{\ell+1} = \sigma(W^\ell x^\ell)$ , по слоям  $\ell = 1, 2, \dots, L$



Достаточные условия монотонности DNN:

- неотрицательность весовых коэффициентов  $W^\ell \geq 0$
- монотонность функций активации  $\sigma(z)$

Если предсказывать положительные значения  $y(x)$  при  $W^\ell \geq 0$ ,  
то векторы  $x^\ell$  выучиваются детектировать части объекта  $x$ ,  
веса  $W^\ell$  и векторы  $x^\ell$  — **разреженные и интерпретируемые**

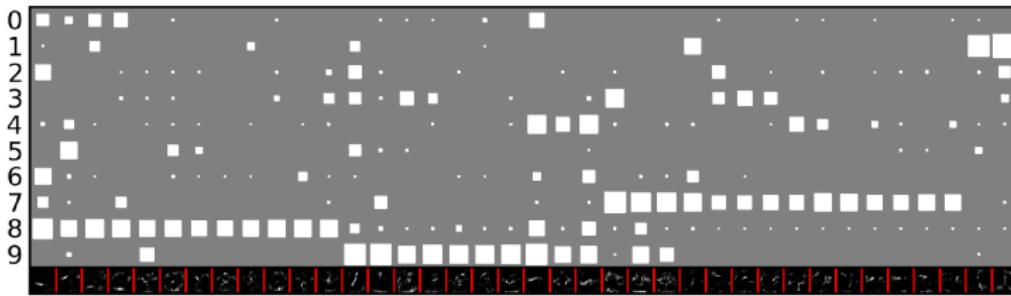
---

J.-R. Cano et al. Monotonic classification: An overview on algorithms, performance measures and data sets. Neurocomputing, 2019.

## Глубокая монотонная нейронная сеть (Constrained MNN)

Двухслойная сеть для распознавания рукописных цифр MNIST

- первый слой выделяет информативные группы пикселей
- второй слой выделяют группы, образующие цифры



Разреженность и интерпретируемость весов  $W^\ell$  и векторов  $x^\ell$  возникает также при обучении автокодировщиков без учителя

---

*J. Chorowski, J.M.Zurada.* Learning understandable neural networks with non-negative weight constraints. 2015

*B.O.Ayinde, J.M.Zurada.* Deep learning of nonnegativity-constrained autoencoders for enhanced understanding of data. 2018

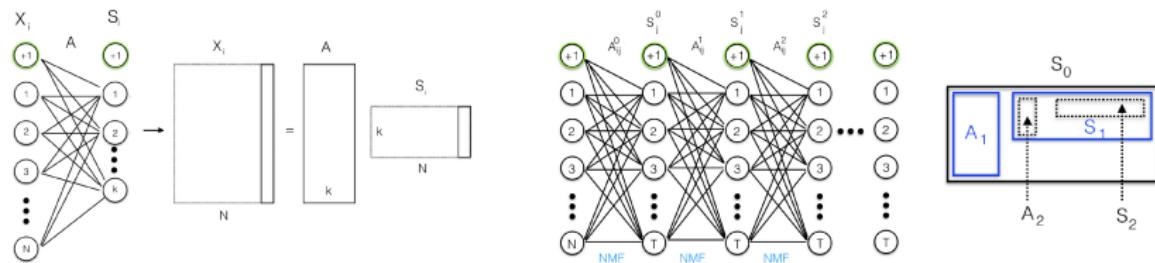
## Неотрицательные матричные разложения и автокодировщики

**Дано:** неотрицательная матрица данных  $X$ ,  $X_{ij} \geq 0$

**Найти:** неотрицательные матрицы  $A, S$ ,  $A_{ik} \geq 0, S_{kj} \geq 0$

**Критерий:**  $\|X - AS\| \rightarrow \min_{A, S}$  — non-negative matrix factorization

Deep NMF — обучаемая иерархия векторных представлений:



J.Flenner, B.Hunter. A deep non-negative matrix factorization neural network. 2017

Zhikui Chen, Shan Jin, Runze Liu, Jianing Zhang. A Deep non-negative matrix factorization model for big data representation learning. 2021

T.Will et al. Neural nonnegative matrix factorization for hierarchical multilayer topic modeling. 2023

## Промежуточные итоги и направления исследований

- Неотрицательность весов в нейросетевых моделях приводит к **разреженности и интерпретируемости**
- Известно, что введение нормировки при оптимизации повышает устойчивость обучения нейронных сетей
- Переход от нормированных векторов (при необходимости) к ненормированным — умножением на коэффициент
- Во всех перечисленных ситуациях может быть использована оптимизация на единичных симплексах
- мультипликативный шаг с нормировкой реализован в pyTorch [Илья Дьяков, курсовая работа, ВМК МГУ, 2024]
- **Открытая проблема:** находить применения этой технике, сравнивать со state-of-the-art на различных задачах

## Тематическое моделирование: «о чём все эти тексты?»

**Дано:** коллекция текстовых документов как «мешков-слов»

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

**Найти:** тематическую языковую модель

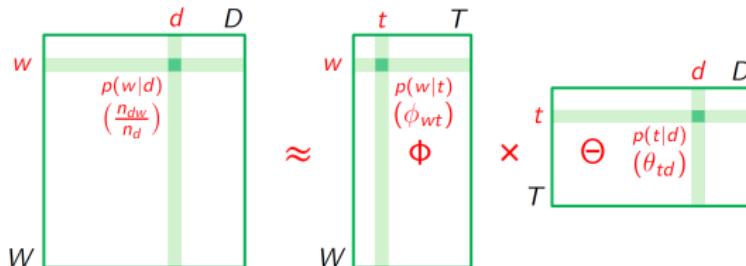
- $p(w|d) = \sum_{t \in T} p(w|\cancel{t}) p(t|d) = \sum_{t \in T} \phi_{wt} \theta_{td}$
- $p(w|t) = \phi_{wt}$  — из каких слов  $w$  состоит каждая тема  $t \in T$
- $p(t|d) = \theta_{td}$  — из каких тем  $t$  состоит каждый документ  $d$

**Критерий:** правдоподобие предсказания слов  $w$  в документах  $d$  с дополнительными критериями-регуляризаторами  $R_i(\Phi, \Theta)$ :

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

## Три интерпретации задачи тематического моделирования

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

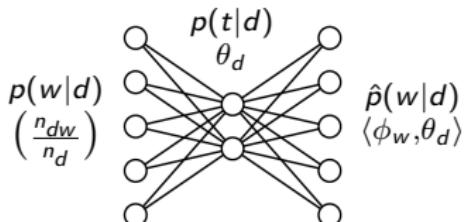


3. Автокодировщик документов в тематические эмбеддинги:

- кодировщик  $f_\Phi: \frac{n_{dw}}{n_d} \rightarrow \theta_d$
- декодировщик  $g_\Phi: \theta_d \rightarrow \Phi \theta_d$

задача реконструкции:

$$\sum_{d,w} n_{dw} \ln \langle \phi_w, \theta_d \rangle \rightarrow \min_{\Phi, \Theta}$$



## Пример 1. Мультиязычная модель Википедии

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

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

Тема №68		Тема №79	
research	4.56	институт	6.03
technology	3.14	университет	3.35
engineering	2.63	программа	3.17
institute	2.37	учебный	2.75
science	1.97	технический	2.70
program	1.60	технология	2.30
education	1.44	научный	1.76
campus	1.43	исследование	1.67
management	1.38	наука	1.64
programs	1.36	образование	1.47

Ассессор оценил 396 тем из 400 как хорошо интерпретируемые.

## Пример 1. Мультиязычная модель Википедии

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

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

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

Ассесор оценил 396 тем из 400 как хорошо интерпретируемые.

## Цели и не-цели тематического моделирования

### Цели:

- выявлять тематическую кластерную структуру текстовой коллекции (сколько в ней тем, и о чём они), 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)
- генерировать связный текст (слабые эмбединги)
- понимать смысл текста (тем не достаточно для этого)

## Некоторые приложения тематического моделирования

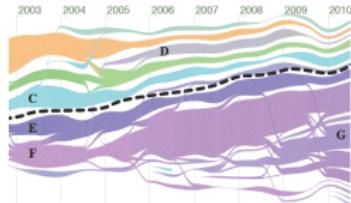
разведочный поиск в электронных библиотеках



поиск тематических сообществ в соцсетях



выявление и отслеживание цепочек новостей



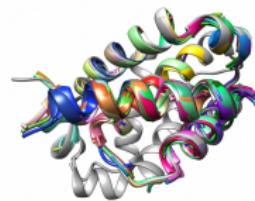
мультимодальный поиск текстов и изображений



анализ банковских транзакционных данных



поиск паттернов в задачах биоинформатики



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.

## Сходства и отличия от LLM

### PTM и LLM — что общего

- языковая модель, которая предсказывает слова в тексте
- автокодировщик, который переводит текст в эмбединги
- обобщаются на мультимодальные, мультиязычные данные
- возможно многозадачное, многокритериальное обучение

### PTM — принципиальные отличия от LLM

- намного более слабая языковая модель
- эмбединги вероятностные, разреженные, интерпретируемые
- простота и скорость матричного разложения

### PTM — дальнейшее развитие навстречу LLM

- отказ от байесовского обучения → алгоритм ближе к SGD
- транзакционные данные → ближе к foundation models
- отказ от мешка слов → ближе к модели внимания

## 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  $\Phi, \Theta$  — решение, то стохастические  $\Phi', \Theta'$  — 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**:

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

EM-алгоритм: метод простой итерации для системы уравнений

E-шаг:  $p_{tdw} \equiv p(t|d, w) = \text{norm}_{t \in T}(\phi_{wt} \theta_{td})$

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

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:

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

$$\begin{aligned} \phi_{wt} &= \underset{w \in W}{\text{norm}} \left( \phi_{wt} \frac{\partial f}{\partial \phi_{wt}} \right) = \underset{w \in W}{\text{norm}} \left( \phi_{wt} \sum_{d \in D} n_{dw} \frac{\theta_{td}}{p(w|d)} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right) = \\ &= \underset{w \in W}{\text{norm}} \left( \sum_{d \in D} n_{dw} p_{tdw} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right); \end{aligned}$$

$$\begin{aligned} \theta_{td} &= \underset{t \in T}{\text{norm}} \left( \theta_{td} \frac{\partial f}{\partial \theta_{td}} \right) = \underset{t \in T}{\text{norm}} \left( \theta_{td} \sum_{w \in W} n_{dw} \frac{\phi_{wt}}{p(w|d)} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right) = \\ &= \underset{t \in T}{\text{norm}} \left( \sum_{w \in d} n_{dw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \end{aligned}$$

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

## PLSA и LDA — две самые известные тематические модели

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

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

М-шаг — частотные оценки условных вероятностей:

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

**LDA**: latent Dirichlet allocation (латентное размещение Дирихле):

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

М-шаг — частотные оценки с поправками  $\beta_w > -1, \alpha_t > -1$ :

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

---

Hofmann T. Probabilistic latent semantic indexing. SIGIR 1999.

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

## От байесовского вывода к аддитивной регуляризации

$X$  — исходные данные,  $\Omega$  — параметры порождающей модели

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

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

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

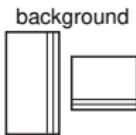
**Максимизация апостериорной вероятности** (MAP)  
даёт точечную оценку  $\Omega$  напрямую, без вывода Posterior:

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

**Многокритериальная аддитивная регуляризация** (ARTM)  
обобщает MAP на любые регуляризаторы и их комбинации:

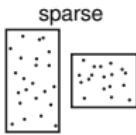
$$\Omega := \arg \max_{\Omega} (\ln p(X|\Omega) + \sum_{i=1} \tau_i R_i(\Omega))$$

## Регуляризаторы для улучшения интерпретируемости тем



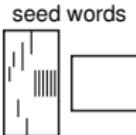
Сглаживание фоновых тем  $B \subset T$ :

$$R(\Phi, \Theta) = \beta_0 \sum_{t \in B} \sum_w \beta_w \ln \phi_{wt} + \alpha_0 \sum_d \sum_{t \in B} \alpha_t \ln \theta_{td}$$

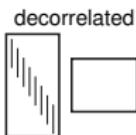


Разреживание предметных тем  $S = T \setminus B$ :

$$R(\Phi, \Theta) = -\beta_0 \sum_{t \in S} \sum_w \beta_w \ln \phi_{wt} - \alpha_0 \sum_d \sum_{t \in S} \alpha_t \ln \theta_{td}$$

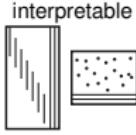


Сглаживание для выделения релевантных тем  
с помощью словаря «затравочных» ключевых слов



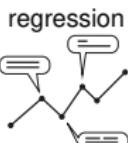
Декоррелирование для повышения различности тем:

$$R(\Phi) = -\frac{\tau}{2} \sum_{t,s} \sum_w \phi_{wt} \phi_{ws}$$



Сглаживание + разреживание + декоррелирование  
для улучшения интерпретируемости тем

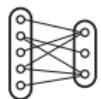
## Регуляризаторы для учёта дополнительной информации



biterm



relational



hierarchy



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

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

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

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

Связи или ссылки между документами:

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

Связи родительских тем  $t$  с дочерними подтемами  $s$ :

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

## Мультимодальная тематическая модель ARTM

$W_m$  — словарь термов  $m$ -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

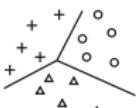
E-шаг:  $p_{tdw} = \text{norm}_{t \in T}(\phi_{wt} \theta_{td})$

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

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

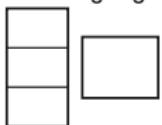
## Регуляризаторы для мультимодальных тематических моделей

supervised



Модальности меток классов или категорий для задач классификации и категоризации текстов.

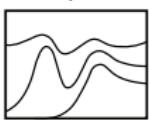
multilanguage



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

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

temporal



Темпоральные модели с модальностью времени  $i$ :

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

geospatial



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

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

## Трёхматричная тематическая модель ARTM

Темы порождаются модальностью  $C$  (категории, авторы):

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \sum_{c \in C} \phi_{wt} \psi_{tc} \theta_{cd} + R(\Phi, \Psi, \Theta) \rightarrow \max_{\Phi, \Psi, \Theta};$$

EM-алгоритм: метод простой итерации для системы уравнений

Е-шаг:  $p_{tcdw} \equiv p(t, c | d, w) = \text{norm}_{(t,c) \in T \times C} (\phi_{wt} \psi_{tc} \theta_{cd});$

М-шаг:  $\left\{ \begin{array}{l} \phi_{wt} = \text{norm}_{w \in W_m} \left( n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right); \quad n_{wt} = \sum_{d,c} n_{dw} p_{tcdw} \\ \psi_{tc} = \text{norm}_{t \in T} \left( n_{tc} + \psi_{tc} \frac{\partial R}{\partial \psi_{tc}} \right); \quad n_{tc} = \sum_{d,w} n_{dw} p_{tcdw} \\ \theta_{cd} = \text{norm}_{c \in C} \left( n_{cd} + \theta_{cd} \frac{\partial R}{\partial \theta_{cd}} \right); \quad n_{cd} = \sum_{w,t} n_{dw} p_{tcdw} \end{array} \right.$

*M. Rosen-Zvi et al. The author-topic model for authors and documents. 2004.*

## Транзакционные данные

Текст — это двудольный граф с рёбрами вида  $(d, w)$ ,  $w \in W_m$ .

Когда данные содержат  $n$ -ки термов разных модальностей:

- **данные социальной сети:**

$(d, u, w)$  — пользователь  $u$  записал слово  $w$  в блоге  $d$

- **данные сети интернет-рекламы:**

$(u, d, b)$  — пользователь  $u$  кликнул баннер  $b$  на странице  $d$

- **данные рекомендательной системы:**

$(u, f, s)$  — пользователь  $u$  оценил фильм  $f$  в ситуации  $s$

- **данные финансовых организаций:**

$(b, s, g)$  — покупатель  $u$  купил у продавца  $s$  товар  $g$

- **данные о пассажирских авиаперелётах:**

$(u, a, b, c)$  — перелёт клиента  $u$  из  $a$  в  $b$  авиакомпанией  $c$

**Задача:** по наблюдаемой выборке рёбер гиперграфа найти латентные тематические векторные представления его вершин.

## Гиперграфовая тематическая модель транзакционных данных

**Дано:**  $E_k$  — выборка транзакций (ребер гиперграфа) типа  $k$   
 $n_{kdx}$  — число вхождений ребра  $(d, x)$  в выборку  $E_k$

$$\sum_{k \in K} \tau_k \sum_{(d,x) \in E_k} n_{kdx} \ln \sum_{t \in T} \theta_{td} \prod_{v \in x} \phi_{vt} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

EM-алгоритм: метод простой итерации для системы уравнений со вспомогательными переменными  $p_{tdx} = p(t|d, x)$ :

$$\left. \begin{array}{l} \text{E-шаг: } p_{tdx} = \text{norm} \left( \theta_{td} \prod_{v \in x} \phi_{vt} \right) \\ \text{M-шаг: } \begin{cases} \phi_{vt} = \text{norm} \left( \sum_{k \in K} \tau_k \sum_{(d,x) \in E_k} [v \in x] n_{kdx} p_{tdx} + \phi_{vt} \frac{\partial R}{\partial \phi_{vt}} \right) \\ \theta_{td} = \text{norm} \left( \sum_{k \in K} \tau_k \sum_{(d,x) \in E_k} n_{kdx} p_{tdx} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right) \end{cases} \end{array} \right.$$

## Частный случай: тематическая модель предложений

$S_d$  — множество предложений документа  $d$

$n_{sw}$  — сколько раз терм  $w$  встречается в предложении  $s$

**Критерий:** максимум регуляризованного log-правдоподобия

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

**Свобода выбора** гипер-ребер сегментоидов — подмножеств термов, связанных по смыслу и порождаемых общей темой:

- предложение / фраза / синтагма / именная группа
- факт «объект, субъект, действие»
- лексическая цепочка: синонимы, гипонимы, меронимы
- текст комментария, дата–время, автор

---

Wayne Xin Zhao et al. Comparing Twitter and traditional media using topic models. 2011.  
G.Balikas, M.-R.Amini, M.Clausel. On a topic model for sentences. SIGIR 2016.

## Почему отказ от байесовского обучения даёт больше свободы

В байесовском подходе регуляризаторы являются априорным распределением — частью генеративной модели.

В подходе ARTM это инструмент управления сходимостью итерационного процесса в некорректно поставленной задаче.

**ARTM позволяет свободно:**

- подбирать комбинацию регуляризаторов под задачу
- менять коэффициенты регуляризации в ходе итераций
- включать и отключать регуляризаторы в ходе итераций
- ставить задачу управления траекторией регуляризации
- отключать все регуляризаторы в конце итераций, чтобы получать несмешённые оценки параметров

---

Vorontsov K.V. Rethinking probabilistic topic modeling from the point of view of classical non-Bayesian regularization. 2023.

## Модульный подход к синтезу моделей с заданными свойствами

Для построения композитных моделей в BigARTM не нужны ни математические выкладки, ни программирование «с нуля»

Этапы моделирования	Bayesian TM	ARTM
Формализация:	Анализ требований	Анализ требований
Алгоритмизация:	Вероятностная модель порождения данных	Стандартные критерии Свои критерии
Реализация:	Байесовский вывод для данной порождающей модели (VI, GS, EP)	Единый регуляризованный ЕМ-алгоритм для любых моделей и их композиций
Оценивание:	Исследовательский код (Matlab, Python, R)	Промышленный код BigARTM (C++, Python API)
	Исследовательские метрики, исследовательский код	Стандартные метрики Свои метрики
	Внедрение	Внедрение

-- нестандартизуемые этапы, уникальная разработка для каждой задачи

-- стандартизуемые этапы

# BigARTM: библиотека тематического моделирования

## Ключевые возможности:

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

## 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.

## Качество и скорость: BigARTM vs Gensim и Vowpal Wabbit

3.7М статей Википедии, 100К слов: время min (перплексия)

проц.	$ T $	Gensim	Vowpal Wabbit	BigARTM	BigARTM асинхрон
1	50	142m (4945)	50m (5413)	42m (5117)	25m (5131)
1	100	287m (3969)	91m (4592)	52m (4093)	32m (4133)
1	200	637m (3241)	154m (3960)	83m (3347)	53m (3362)
2	50	89m (5056)		22m (5092)	13m (5160)
2	100	143m (4012)		29m (4107)	19m (4144)
2	200	325m (3297)		47m (3347)	28m (3380)
4	50	88m (5311)		12m (5216)	7m (5353)
4	100	104m (4338)		16m (4233)	10m (4357)
4	200	315m (3583)		26m (3520)	16m (3634)
8	50	88m (6344)		8m (5648)	5m (6220)
8	100	107m (5380)		10m (4660)	6m (5119)
8	200	288m (4263)		15m (3929)	10m (4309)

D.Kochedykov, M.Apishev, L.Golitsyn, K.Vorontsov. Fast and modular regularized topic modelling. FRUCT ISMW, 2017.

## Разведочный поиск в технологических блогах

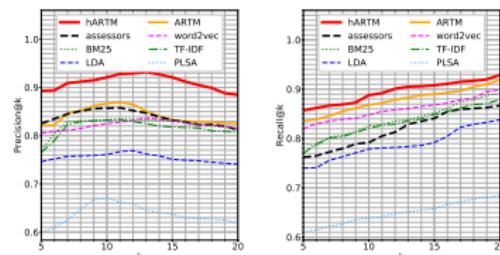
**Цель:** поиск документов  
по длинным текстовым запросам  
— Habr.ru (175K документов),  
— TechCrunch.com (760K док.).

**Регуляризаторы:**

$$\mathcal{L}\left(\begin{array}{|c|c|}\hline \Phi & \Theta \\ \hline\end{array}\right) + R\left(\text{hierarchy}\right) + R\left(\text{interpretable}\right) + R\left(\text{multimodal}\right) + R\left(\text{n-gram}\right) \rightarrow \max$$

**Результаты:**

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



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

# 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:**

$$\mathcal{L} \left( \begin{array}{c} \text{PLSA} \\ \boxed{\Phi} \quad \boxed{\Theta} \end{array} \right) + R \left( \begin{array}{c} \text{seed words} \\ \boxed{\text{---}} \quad \boxed{\square} \end{array} \right) + R \left( \begin{array}{c} \text{interpretable} \\ \boxed{\text{---}} \quad \boxed{\text{---}} \end{array} \right) + R \left( \begin{array}{c} \text{multimodal} \\ \boxed{\text{---}} \quad \boxed{\square} \end{array} \right) \\ + R \left( \begin{array}{c} \text{temporal} \\ \boxed{\text{---}} \quad \boxed{\text{---}} \end{array} \right) + R \left( \begin{array}{c} \text{geospatial} \\ \boxed{\text{---}} \quad \boxed{\text{---}} \end{array} \right) + R \left( \begin{array}{c} \text{sentiment} \\ \boxed{\text{---}} \quad \boxed{\text{---}} \end{array} \right) \rightarrow \max$$

**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.

**(японцы):** японский, япония, корея, китайский, жилища, авария, фукусиму, цунами, сообщество, океан, станция, катако, район, правительство, атомный.

**(корейцы):** дети, ребенок, родиться, детский, семья, воспитанный, право, мораль, фольклор, языковой, нормативный, родительский, родить, малыши, взрослый, опека, сын.

**(венесуэльцы):** куба, карто, венесуэла, чавес, президент, уго, мадуро, боливия, фидель, глава, латинский, венесуэльский, лидер, боливарианская,

президентский, элиэнде, гевара.

**(китайцы):** китайский, россия, производство, китай, продукция, страна, предприятия, компании, технология, военный, регион, производить,

промышленный, промышленность, российский, экономический, кир,

**(азербайджанцы):** русский, азербайджан, азербайджанцы, россия, азербайджан, ташкент, дистректор, анала, народ, москва, страна, армянин,

слово, рынок.

**(грузины):** грузинский, спасиас, военный, август, батиашвила, российский,

специаловец, мигранты, операция, румыны, бригада, миротворческий, абхазия,

группа, войска, русский, цинхана.

**(осетины):** конституция, осетия, аминат, русский, осетинский, южный, северный,

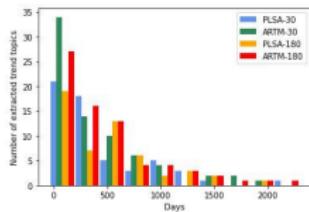
россия, война, республика, вопрос, аланай, российский, население, конфликт,

**(цыгане):** наркотик, цыган, цыганка, хороший, место, страна, деньги, время,

работать, жизнь, жить, рука, дом, цыганский, наркоманка.

## Выявление трендов в коллекции научных публикаций

**Цель:** раннее обнаружение трендовых тем с начальным экспоненциальным ростом в области AI/ML 2009–2021 гг.



**Регуляризаторы:**

$$\mathcal{L}\left(\begin{array}{c} \text{PLSA} \\ \Phi \quad \Theta \end{array}\right) + R\left(\begin{array}{c} \text{interpretable} \\ \text{grid} \end{array}\right) + R\left(\begin{array}{c} \text{dynamic} \\ \text{wavy line} \end{array}\right) + R\left(\begin{array}{c} \text{multimodal} \\ \text{bar chart} \end{array}\right) + R\left(\begin{array}{c} \text{n-gram} \\ \text{matrix} \end{array}\right) \rightarrow \max$$

**Результаты:**

- выделение 90 из 91 тренда в области машинного обучения
- 63% тем выделяется за год, 79% за два года

---

N.Gerasimenko, A.Chernyavskiy, M.Nikiforova, M.Nikitin, K.Vorontsov. Incremental topic modeling for scientific trend detection Doklady RAS, 2022.

## Выделение поляризованных мнений в политических новостях

**Цель:** найти признаки, по которым событийная тема разделяется на кластеры-мнения

**Регуляризаторы:**

Modalities	Pr	Rec	F1
TF-IDF	0.51	0.95	0.67
SPO	0.59	0.7	0.64
FR	0.86	0.49	0.65
Sent	0.69	0.57	0.66
SPO+FR	0.86	0.68	0.76
SPO+Sent	0.83	0.78	0.81
FR+Sent	0.9	0.52	0.67
All	<b>0.77</b>	<b>0.97</b>	<b>0.86</b>

$$\mathcal{L} \left( \begin{array}{|c|c|} \hline \text{PLSA} & \\ \hline \Phi & \Theta \\ \hline \end{array} \right) + R \left( \begin{array}{|c|c|} \hline \text{interpretable} & \\ \hline \text{grid} & \text{matrix} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|c|} \hline \text{multimodal} & \\ \hline \text{grid} & \square \\ \hline \end{array} \right) + R \left( \begin{array}{|c|c|} \hline \text{n-gram} & \\ \hline \text{grid} & \text{matrix} \\ \hline \end{array} \right) + R \left( \begin{array}{|c|c|} \hline \text{syntax} & \\ \hline \text{tree} & \text{grid} \\ \hline \end{array} \right) \rightarrow \max$$

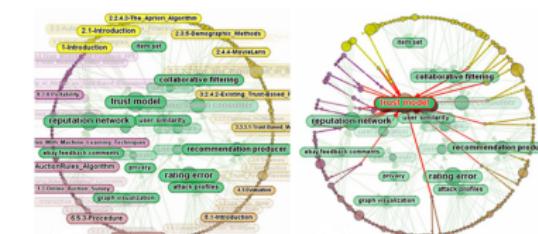
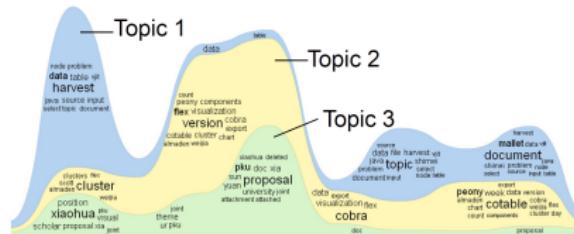
**Результаты:**

- выделение мнений внутри тем: F1-мера = 0.86%
- совместное использование трёх модальностей:
  - SPO — факты как триплеты «субъект–предикат–объект»
  - FR — семантические роли слов по Филлмору
  - Sent — тональности именованных сущностей

D.Feldman, T.Sadekova, K.Vorontsov. Combining facts, semantic roles and sentiment lexicon in a generative model for opinion mining. Dialogue 2020.

## Мотивации. Что хотим от PTMs глядя на LLM (BERT, GPT)

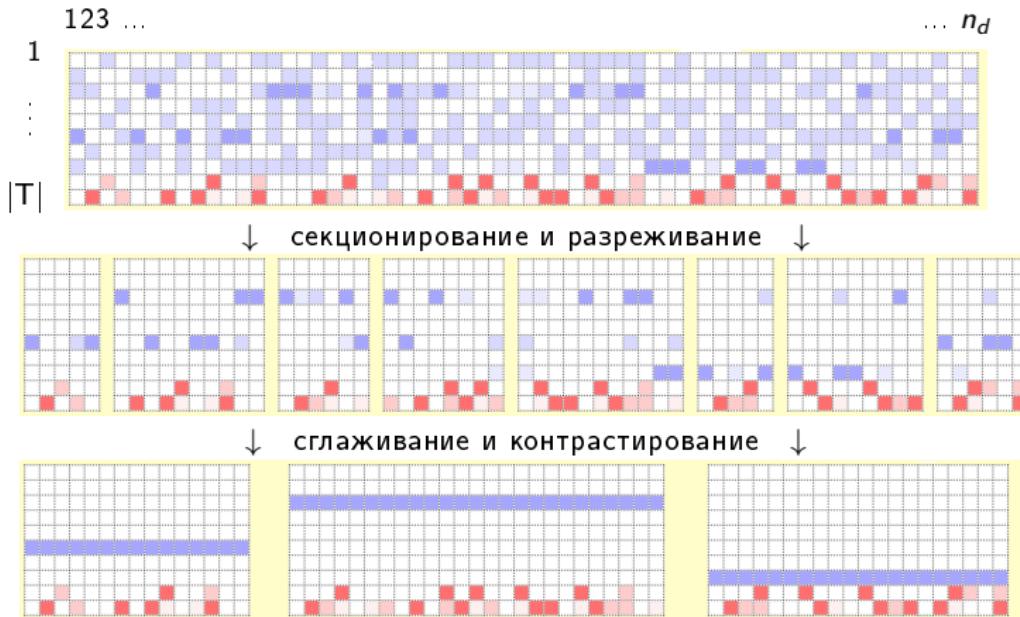
- вместо «мешка слов» — последовательность  $w_1, \dots, w_n$
- вместо документов — локальные контексты слов
- определять тематику любого фрагмента текста
- быстро находить фрагменты, относящиеся к данной теме
- в том числе фразы для суммаризации документа или темы
- разделять документ на тематически однородные сегменты
- визуализировать тематическую структуру документа



## Сегментная структура текста и пост-обработка Е-шага

Документ  $d = \{w_1, \dots, w_{n_d}\}$ ,  $n_d$  — длина документа  $d$

Тематика термов в документе  $p(t|d, w_i)$  — матрица  $T \times n_d$ :



## Регуляризация Е-шага

Трёхмерная матрица  $\Pi = (p_{tdw} = p(t|d, w))_{T \times D \times W}$

Регуляризатор Е-шага:  $\tilde{R}(\Phi, \Theta) = R(\Pi(\Phi, \Theta), \Phi, \Theta)$ :

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

EM-алгоритм: метод простой итерации для системы уравнений

Е-шаг: 
$$\begin{cases} p_{tdw} = \text{norm}_{t \in T}(\phi_{wt} \theta_{td}) \\ \tilde{p}_{tdw} = p_{tdw} \left( 1 + \frac{1}{n_{dw}} \left( \frac{\partial R}{\partial p_{tdw}} - \sum_{z \in T} p_{zdw} \frac{\partial R}{\partial p_{zdw}} \right) \right) \end{cases} \quad (*)$$

М-шаг: 
$$\begin{cases} \phi_{wt} = \text{norm} \left( \sum_{d \in D} n_{dw} \tilde{p}_{tdw} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right) \\ \theta_{td} = \text{norm} \left( \sum_{w \in d} n_{dw} \tilde{p}_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right) \end{cases}$$

## Набросок доказательства: три шага

1. Для функции  $p_{tdw}(\Phi, \Theta) = \frac{\phi_{wt}\theta_{td}}{\sum_z \phi_{wz}\theta_{zd}}$  и любого  $z \in T$

$$\phi_{wt} \frac{\partial p_{zdw}}{\partial \phi_{wt}} = \theta_{td} \frac{\partial p_{zdw}}{\partial \theta_{td}} = p_{tdw}([z=t] - p_{zdw}).$$

2. Введём вспомогательную функцию от переменных  $\Pi, \Phi, \Theta$ :

$$Q_{tdw}(\Pi, \Phi, \Theta) = \frac{\partial R(\Pi, \Phi, \Theta)}{\partial p_{tdw}} - \sum_{z \in T} p_{zdw} \frac{\partial R(\Pi, \Phi, \Theta)}{\partial p_{zdw}}.$$

Если  $R(\Pi, \Phi, \Theta)$  не зависит от  $p_{tdw}$  при  $w \notin d$ , то

$$\phi_{wt} \frac{\partial \tilde{R}}{\partial \phi_{wt}} = \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} + \sum_{d \in D} p_{tdw} Q_{tdw}; \quad \theta_{td} \frac{\partial \tilde{R}}{\partial \theta_{td}} = \theta_{td} \frac{\partial R}{\partial \theta_{td}} + \sum_{w \in d} p_{tdw} Q_{tdw}.$$

3. Подставляем это в формулы М-шага:

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

$$\theta_{td} = \text{norm} \left( \sum_{w \in d} n_{dw} p_{tdw} + \sum_{w \in d} Q_{tdw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right).$$



## Любая пост-обработка Е-шага — это регуляризатор $R(\Pi)$

Итак, произвольному гладкому регуляризатору  $R(\Pi, \Phi, \Theta)$  однозначно соответствует пост-обработка  $p_{tdw} \rightarrow \tilde{p}_{tdw}$ .

Верно и обратное:

**Теорема.** Если на  $k$ -й итерации ЕМ-алгоритма для каждого  $(d, w)$ :  $n_{dw} > 0$  в формулах М-шага вместо вектора  $(p_{tdw}^k)_{t \in T}$  подставить вектор  $(\tilde{p}_{tdw}^k)_{t \in T}$ , удовлетворяющий условию нормировки  $\sum_t \tilde{p}_{tdw}^k = 1$ , то это эквивалентно добавлению регуляризатора сглаживания–разреживания

$$R(\Pi) = \sum_{d \in D} \sum_{w \in d} n_{dw} \sum_{t \in T} (\tilde{p}_{tdw}^k - p_{tdw}^k) \ln p_{tdw}.$$

$p(t|d, w)$  можно подвергать любой разумной пост-обработке!

---

Vorontsov K.V. Rethinking probabilistic topic modeling from the point of view of classical non-Bayesian regularization. 2023.

## Доказательство

В системе (\*) дифф. уравнений относительно  $R$  введём переменные  $x_{tdw}$ :

$$\underbrace{p_{tdw}^k \frac{\partial R}{\partial p_{tdw}}}_{x_{tdw}} = n_{dw} (\tilde{p}_{tdw}^k - p_{tdw}^k) + p_{tdw}^k \sum_{z \in T} \underbrace{p_{zdw}^k \frac{\partial R}{\partial p_{zdw}}}_{x_{zdw}}, \quad t \in T.$$

Для любой пары  $(d, w)$  такой, что  $n_{dw} > 0$ , это система  $|T|$  линейных уравнений относительно  $|T|$  переменных  $x_{tdw}$ ,  $t \in T$ .

Подстановкой убеждаемся, что  $x_{tdw} = n_{dw} (\tilde{p}_{tdw}^k - p_{tdw}^k)$  — решение системы. Взяв это решение, получим систему дифф. уравнений относительно  $R$ :

$$\frac{\partial R}{\partial p_{tdw}} = \frac{x_{tdw}}{p_{tdw}}, \quad d \in D, w \in d, t \in T.$$

Система декомпозируется по переменным  $p_{tdw}$ : каждой тройке  $(d, w, t)$  соответствует частное решение  $R(\Pi) = x_{tdw} \ln p_{tdw} + C$ . Общее решение:

$$R(\Pi) = \sum_{d \in D} \sum_{w \in d} \sum_{t \in T} x_{tdw} \ln p_{tdw} + C.$$

Подставляя сюда найденное решение  $x_{tdw}$ , получаем требуемое. ■

## Идея тематизации текста за один проход

**Дано:**  $s$  — фрагмент текста  $d$ ,  $\Phi$  — тематическая модель

**Найти:**  $p(t|s)$  — тематический вектор фрагмента текста

### Проблемы:

- как не переобучить вектор  $p(t|s)$ , если текст короткий?
- как согласовать  $p(t|s)$  с объемлющим контекстом  $p(t|d)$ ?
- как согласовать  $p(t|s)$  с  $p(t|w) = \phi_{wt} \frac{p(t)}{p(w)}$  термов  $w \in s$ ?

### Наводящие соображения:

- первая итерация EM-алгоритма с инициализацией  $\theta_{td}^0 = \frac{1}{|T|}$ :

$$\theta_{td}(\Phi) = \text{norm} \left( \sum_{w \in W} n_{dw} p_{tdw} \right) = \sum_{w \in d} \frac{n_{dw}}{n_d} \text{norm} \left( \phi_{wt} \theta_{td}^0 \right)$$

- формула полной вероятности + гипотеза усл. независ.:

$$\theta_{td}(\Phi) = \sum_{w \in d} p(w|d) p(t|w) \cancel{\theta_{td}} = \sum_{w \in d} \frac{n_{dw}}{n_d} \text{norm} \left( \phi_{wt} \cancel{\theta_{td}} \right)$$

## EM-алгоритм для ARTM с явным выражением $\Theta$ через $\Phi$

Максимизация логарифма правдоподобия:

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

EM-алгоритм: метод простой итерации для системы уравнений

$$p_{tdw} \equiv p(t|d, w) = \text{norm}_{t \in T}(\phi_{wt} \theta_{td}); \quad n_{td} = \sum_{w \in d} n_{dw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}$$

$$p'_{tdw} = p_{tdw} + \frac{1}{n_{dw}} \sum_{s \in T} \frac{n_{sd}}{\theta_{sd}} \phi_{wt} \frac{\partial \theta_{sd}}{\partial \phi_{wt}}$$

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

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

## Доказательство (по Лемме о максимизации на симплексах)

Оптимационная задача М-шага относительно  $\Phi$  и  $\Theta(\Phi)$ :

$$Q(\Phi) = \sum_{d \in D} \sum_{u \in W} \sum_{s \in T} n_{du} p_{sdu} \ln(\phi_{us} \theta_{sd}(\Phi)) + R(\Phi, \Theta(\Phi)) \rightarrow \max_{\Phi}$$

Применим Лемму к регуляризованному log-правдоподобию  $Q$ :

$$\begin{aligned} \phi_{wt} \frac{\partial Q}{\partial \phi_{wt}} &= \sum_{d \in D} n_{dw} p_{tdw} + \sum_{d, s, u} n_{du} p_{sdu} \frac{\phi_{wt}}{\theta_{sd}} \frac{\partial \theta_{sd}}{\partial \phi_{wt}} + \phi_{wt} \sum_{d, s} \frac{\partial R}{\partial \theta_{sd}} \frac{\partial \theta_{sd}}{\partial \phi_{wt}} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} = \\ &= \sum_{d \in D} n_{dw} \left( p_{tdw} + \frac{1}{n_{dw}} \sum_{s \in T} \frac{\phi_{wt}}{\theta_{sd}} \underbrace{\left( \sum_{u \in d} n_{du} p_{sdu} + \theta_{sd} \frac{\partial R}{\partial \theta_{sd}} \right)}_{n_{sd}} \frac{\partial \theta_{sd}}{\partial \phi_{wt}} \right) + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} = \\ &= \sum_{d \in D} n_{dw} \underbrace{\left( p_{tdw} + \frac{1}{n_{dw}} \sum_{s \in T} \frac{n_{sd}}{\theta_{sd}} \phi_{wt} \frac{\partial \theta_{sd}}{\partial \phi_{wt}} \right)}_{p'_{tdw}} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}. \end{aligned}$$

■

## EM-алгоритм для ARTM с линейной тематизацией документов

$$\theta_{td}(\phi) = \sum_{w \in d} \frac{n_{dw}}{n_d} \underset{t \in T}{\text{norm}}(\phi_{wt} p_t) \Rightarrow \phi_{wt} \frac{\partial \theta_{sd}}{\partial \phi_{wt}} = \frac{n_{dw}}{n_d} \phi'_{tw} (\delta_{st} - \phi'_{sw})$$

EM-алгоритм: метод простой итерации для системы уравнений

$$\phi'_{tw} \equiv p(t|w) = \underset{t \in T}{\text{norm}}(\phi_{wt} n_t); \quad \theta_{td} = \sum_{w \in d} \frac{n_{dw}}{n_d} \phi'_{tw}$$

$$p_{tdw} \equiv p(t|d, w) = \underset{t \in T}{\text{norm}}(\phi_{wt} \theta_{td}); \quad n_t = \sum_{d \in D} \sum_{w \in d} n_{dw} p_{tdw}$$

$$n_{td} = \sum_{w \in d} n_{dw} p_{tdw} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}$$

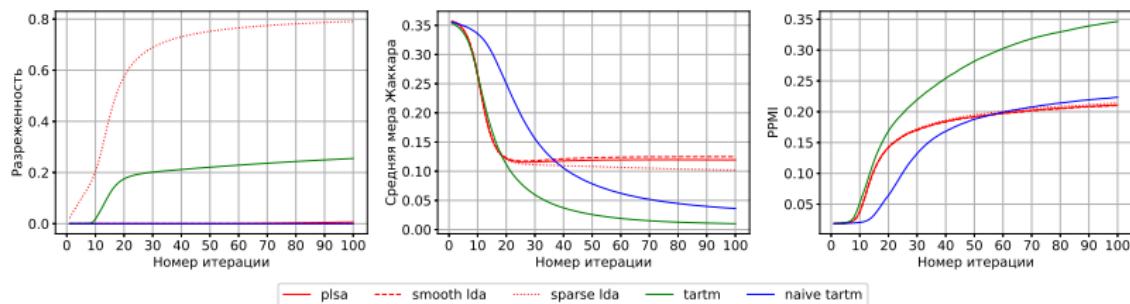
$$p'_{tdw} = p_{tdw} + \frac{\phi'_{tw}}{n_d} \left( \frac{n_{td}}{\theta_{td}} - \sum_{s \in T} \phi'_{sw} \frac{n_{sd}}{\theta_{sd}} \right)$$

$$\phi_{wt} = \underset{w \in W}{\text{norm}} \left( \sum_{d \in D} n_{dw} p'_{tdw} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right)$$

## Эксперимент. Проверка модифицированного EM-алгоритма

Коллекция NIPS,  $|T| = 50$ , модели:

- TARTM ( $\Theta$ less ARTM) — модифицированный EM-алгоритм
- naive TARTM — одна итерация обычного EM-алгоритма



- TARTM очищает темы от общеупотребительных слов,
- улучшает разреженность, различность и когерентность тем

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

[https://github.com/ilirhin/python\\_artm](https://github.com/ilirhin/python_artm)

## Упрощение EM-алгоритма для линейной тематизации

- Нет регуляризации по  $\Theta$ , следовательно,  $\frac{\partial R}{\partial \theta_{td}} = 0$
- Значение отношения  $\frac{n_{td}}{\theta_{td}} \approx n_d$  не зависит от  $t$ , подстановка в формулу М-шага приводит к упрощению:  $p'_{tdw} = p_{tdw}$

EM-алгоритм: метод простой итерации для системы уравнений

$$\phi'_{tw} = \underset{t \in T}{\text{norm}}(\phi_{wt} n_t); \quad \theta_{td} = \sum_{w \in d} \frac{n_{dw}}{n_d} \phi'_{tw};$$

$$p_{tdw} = \underset{t \in T}{\text{norm}}(\phi_{wt} \theta_{td}); \quad n_t = \sum_{d \in D} \sum_{w \in d} n_{dw} p_{tdw};$$

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

Это обычный EM-алгоритм, только с однопроходным Е-шагом!

---

Vorontsov K. V. Rethinking probabilistic topic modeling from the point of view of classical non-Bayesian regularization. 2023.

## Линейная тематизация: от документа к локальным контекстам

Тематизация документа  $d = (w_1, \dots, w_{n_d})$  за один проход:

$$\theta_{td}(\Phi) \equiv p(t|d) = \frac{1}{n_d} \sum_{i=1}^{n_d} p(t|w_i) = \frac{1}{n_d} \sum_{i=1}^{n_d} \phi'_{tw_i}$$

Тематизация локального контекста  $C_i = (\dots, w_i, \dots)$  терма  $w_i$ :

$$\theta_{ti}(\Phi) \equiv p(t|C_i) = \frac{1}{|C_i|} \sum_{u \in C_i} p(t|u) = \frac{1}{|C_i|} \sum_{u \in C_i} \phi'_{tu}$$

Тематизация локального контекста с распределением весов  $\alpha_{ui}$ :

$$\theta_{ti}(\Phi) \equiv p(t|C_i) = \sum_{u \in C_i} \alpha_{ui} \phi'_{tu}, \quad \sum_{u \in C_i} \alpha_{ui} = 1, \quad \alpha_{ui} \geq 0$$

Локализованная тематическая модель:

$$p(w|C_i) = \sum_{t \in T} p(w|t) p(t|C_i) = \sum_{t \in T} \phi_{wt} \sum_{u \in C_i} \phi'_{tu} \alpha_{ui}$$

## ЕМ-алгоритм с локализованным Е-шагом

$w_1, \dots, w_n$  — сквозная нумерация термов во всей коллекции

$C_i$  — локальный контекст (окружение) терма  $w_i$

$\alpha_{ui}$  — распределение весов термов  $u \in C_i$  около терма  $w_i$

- отказались от гипотезы «мешка слов»
- отказались от разбиения коллекции на документы

ЕМ-алгоритм: метод простой итерации для системы уравнений

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

$$p_{ti} \equiv p(t|C_i, w_i) = \operatorname{norm}_{t \in T}(\phi_{w_i t} \theta_{ti}); \quad p_t \equiv p(t) = \frac{1}{n} \sum_{i=1}^n p_{ti};$$

$$\phi_{wt} = \operatorname{norm}_{w \in W} \left( \sum_{i=1}^n [w_i = w] p_{ti} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right).$$

## Быстрое вычисление двунаправленных векторов контекста

Два прохода по тексту — «слева направо» и «справа налево» для вычисления экспоненциальных скользящих средних (ЭСС):

$$\vec{p}(t|i) = \vec{\gamma}_i p(t|w_i) + (1 - \vec{\gamma}_i) \vec{p}(t|i-1), \quad i = 1, \dots, n, \quad \vec{\gamma}_1 = 1$$

$$\bar{p}(t|i) = \bar{\gamma}_i p(t|w_i) + (1 - \bar{\gamma}_i) \bar{p}(t|i+1), \quad i = n, \dots, 1, \quad \bar{\gamma}_n = 1$$

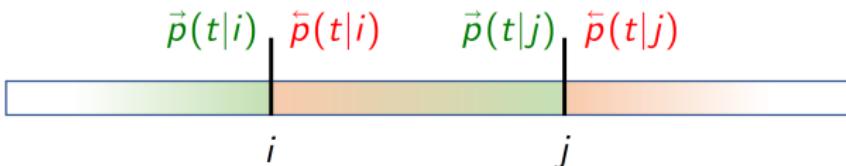
где  $\vec{\gamma}_i$ ,  $\bar{\gamma}_i$  — коэффициенты сглаживания в позиции  $i$

**Основное свойство:** если  $\gamma_i = \gamma$ , то  $\alpha_{w_k} i = \gamma(1 - \gamma)^{|i-k|}$

**Несколько соображений**, как распоряжаться выбором  $\vec{\gamma}_i$ ,  $\bar{\gamma}_i$ :

- $\gamma_i \approx \frac{1}{h}$ , где  $h$  — ширина окна, размер контекста
- $\gamma_i = 1$ , если надо забыть контекст, сменить документ
- $\gamma_i = 0$ , если надо проигнорировать терм
- $\gamma_i$  можно умножать на оценку важности терма

## Использование двунаправленных векторов контекста



Через *дву направленные тематические векторы* определяется:

- $\vec{p}(t|i)$  — тематика левого контекста терма  $w_i$ ;
- $\bar{p}(t|i)$  — тематика правого контекста терма  $w_i$ ;
- $\frac{1}{2}(\vec{p}(t|i) + \bar{p}(t|i))$  — тематика двустороннего контекста  $w_i$ ;
- $p(t|i \dots j) = \frac{1}{2}(\bar{p}(t|i) + \vec{p}(t|j))$  — тематика сегмента  $[i \dots j]$
- $\bar{p}(t|i) \approx \vec{p}(t|j)$  — однородность тематики сегмента  $[i \dots j]$
- $\max_i \|\vec{p}(t|i) - \bar{p}(t|i)\|$  — граница  $i$  между сегментами
- при различных  $\gamma_i$  — короткие и длинные контексты

**Аналогия** с моделями языка GCNN, Attention, Transformer

## Свёрточная нейросеть GCNN (Gated Convolutional Network)

Входные векторы слов (эмбеддинги)

$$X = (x_1, \dots, x_n) \in \mathbb{R}^T$$

трансформируются в векторы слов, зависящие от контекстов  $C_i$ :

$$H = (h_1, \dots, h_n) \in \mathbb{R}^d$$

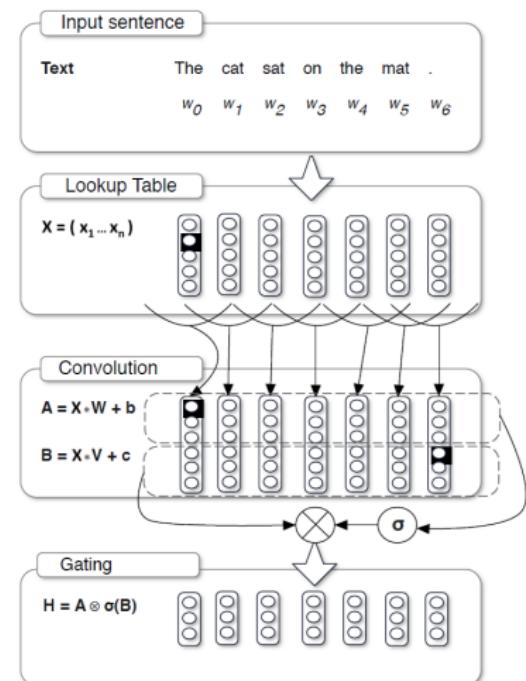
через адамарово произведение:

$$h_i = a_i \otimes \sigma(b_i), \text{ где}$$

$$a_i = \sum_{u \in C_i} W_u x_u \text{ — свёртка-контекст,}$$

$$b_i = \sum_{u \in C_i} V_u x_u \text{ — свёртка-фильтр,}$$

$W_u, V_u$  — матрицы размера  $d \times T$ , параметры модели GCNN.



## Аналогия локализованного Е-шага с моделью GCNN

Контекстный тематический вектор на выходе Е-шага:

$$p(t|C_i, w_i) \equiv p_{ti} = \operatorname{norm}_{t \in T}(\phi_{w_i t} \theta_{ti}) = \operatorname{norm}_{t \in T}\left(\sum_{u \in C_i} \alpha_{ui} \phi'_{tu} \phi_{w_i t}\right)$$

Контекстный вектор на выходе модели GCNN:

$$h_i = \left( \sum_{u \in C_i} W_u x_u \right) \otimes \sigma \left( \sum_{u \in C_i} V_u x_u \right)$$

### Сходство:

- вектор терма  $w_i$  трансформируется в контекстный вектор
- путём усреднения векторов  $\phi'_u$  его контекста,
- семантически схожих с вектором терма  $w_i$ , фильтруемых адамаровым умножением на неотрицательный вектор

### Отличия:

- нет обучаемых матриц  $W_u, V_u$  как у модели GCNN
- вектор-фильтр  $\phi_{w_i}$  без усреднения по контексту  $C_i$
- проецирование итогового вектора на единичный симплекс

## Модель внимания (self-attention) Query–Key–Value

Входные векторы слов (эмбединги)

$$X = (x_1, \dots, x_n) \in \mathbb{R}^T$$

трансформируются в векторы слов,  
зависящие от контекстов  $C_i$ :

$$H = (h_1, \dots, h_n) \in \mathbb{R}^d$$

Модель внимания (self-attention):

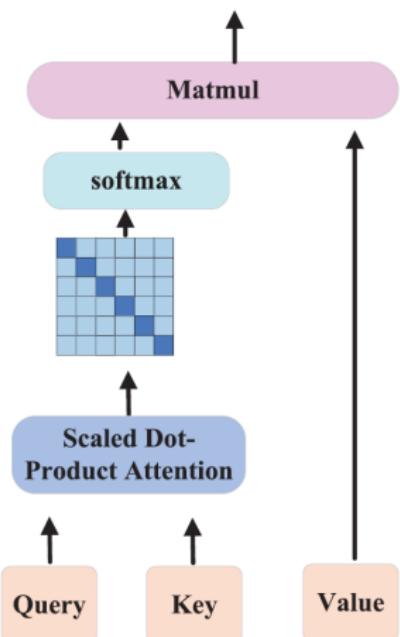
$$h_i = \sum_{u \in C_i} W_v x_u \operatorname{SoftMax}_{u \in C_i} \langle W_k x_u, W_q x_i \rangle$$

$W_v x_u$  — вектор-значение (value)

$W_k x_u$  — вектор-ключ (key)

$W_q x_i$  — вектор-запрос (query)

$W_q, W_k, W_v$  — матрицы параметров



## Аналогия локализованного Е-шага с моделью само-внимания

Контекстный тематический вектор на выходе Е-шага:

$$p(t|C_i, w_i) \equiv p_{ti} = \operatorname{norm}_{t \in T}(\phi_{w_i t} \theta_{ti}) = \operatorname{norm}_{t \in T}\left(\sum_{u \in C_i} \phi'_{tu} \alpha_{ui} \phi_{w_i t}\right)$$

Контекстный вектор на выходе модели само-внимания:

$$h_i = \sum_{u \in C_i} W_v x_u \alpha_{ui} = \sum_{u \in C_i} W_v x_u \operatorname{SoftMax}_{u \in C_i} \langle W_k x_u, W_q x_i \rangle$$

### Сходство:

- вектор терма  $w_i$  трансформируется в контекстный вектор
- путём усреднения векторов  $\phi'_u$  из контекста терма  $w_i$ ,
- наиболее (семантически) схожих с вектором терма  $w_i$

### Отличия:

- адамарово умножение вектора  $\phi'_u$  на вектор-фильтр  $\phi_{w_i}$
- нет обучаемых матриц  $W_q, W_k, W_v$  как у модели внимания
- проецирование итогового вектора на единичный симплекс

## Аналогия локализованного Е-шага с моделью трансформера

Один проход документа аналогичен модели внимания:

- для каждого  $d \in D$ , для каждой позиции  $i = 1, \dots, n_d$  вычисляются 5 тематических векторов, связанных с термом  $w_i$ :

$\phi'_{tw_i} = \text{norm}_t(\phi_{w_i t} p_t)$  — бесконтекстный вектор терма  $p(t|w_i)$

$\vec{p}(t|i) = \vec{\theta}_{ti}$ ,  $\hat{p}(t|i) = \hat{\theta}_{ti}$  — векторы левого и правого контекста

$\theta_{ti} = \beta \vec{\theta}_{ti} + (1 - \beta) \hat{\theta}_{ti}$  — вектор двустороннего контекста

$p_{ti} = \text{norm}_t(\phi_{w_i t} \theta_{ti})$  — контекстный вектор терма  $p(t|C_i, w_i)$

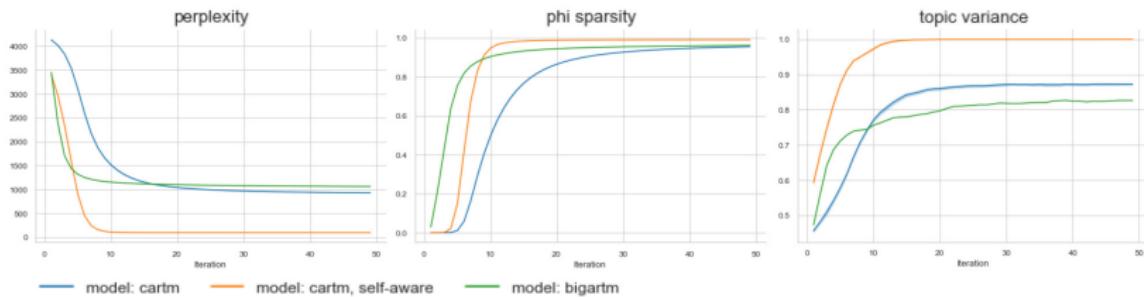
Несколько таких проходов аналогичны трансформеру:

контекстный вектор терма  $p_{ti} = p(t|C_i, w_i)$  с предыдущего прохода используется вместо его бесконтекстного вектора  $\phi'_{tw_i} = p(t|w_i)$

$L$  итераций аналогичны  $L$  необучаемым блокам внимания

## Первые эксперименты с реализацией Context-ARTM

Коллекция «20 Newsgroups»:  $|D| = 18846$ ,  $|W| = 107672$ .



— улучшилась перплексия, разреженность  $\Phi$ , различность тем

model	time@10 topics	time@30 topics	time@70 topics	time@100 topics
CARTM@CPU	4min 55s $\pm$ 1.7s	9min 20s $\pm$ 9.52s	20min 37s $\pm$ 2.36s	25min 52s $\pm$ 4.63s
BigARTM	1min 30s $\pm$ 1.6s	3min 5s $\pm$ 1.98s	4min 55s $\pm$ 653ms	6min 22s $\pm$ 5.81s

— время хуже в несколько раз, при этом реализация CARTM на Python/JAX, тогда как ядро BigARTM на C++

---

Дьяков И.А. Тематические модели внимания для анализа связного текста.  
ВКР бакалавра, ВМК МГУ, 2025.

## Выводы

- Лемма о максимизации на единичных симплексах:
  - позволила переосмыслить и упростить теорию PTMs
  - приводит к разреженным интерпретируемым моделям
  - имеет много потенциальных применений, но...
  - видимо, не известна в сообществах DS/AI/ML
  - могла бы послужить основой для Neural Topic Models
- Больше свободы выбора для моделирования:
  - аддитивная регуляризация, траектории регуляризации
  - модульная технология, реализованная в BigARTM
  - модальности, транзакционные данные, иерархии тем
  - эвристики пост-обработки Е-шага — тоже регуляризация
  - большая вариативность тематических моделей внимания

---

Potapenko A.A., Vorontsov K.V. Robust PLSA performs better than LDA. ECIR-2013.  
Vorontsov K.V. Additive regularization for topic models of text collections. 2014.  
Vorontsov K.V. Rethinking probabilistic topic modeling from the point of view of classical non-Bayesian regularization. 2023.