

Многомодальные тематические модели на гиперграфах

Жариков Илья Николаевич

`zharikov.i.n@yandex.ru`

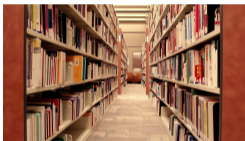
Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем

Научный руководитель: д.ф.-м.н. К. В. Воронцов

Июнь, 2018

Приложения тематического моделирования

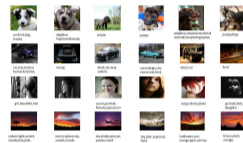
Разведочный поиск в
электронных библиотеках



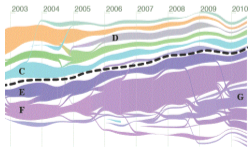
Персонализированный
поиск в соцсетях



Мультимодальный поиск
текстов и изображений



Детектирование и трекинг
новостных сюжетов



Навигация по большим
текстовым коллекциям



Понимание
естественного языка



Цель исследования

Транзакция — взаимодействие некоторого множества объектов.

Проблема

Большинство существующих тематических моделей описывают только попарные взаимодействия между объектами разных типов (модальностей).

Предлагается

Обобщить методы тематического моделирования на случай транзакционных данных, в которых взаимодействуют более двух объектов.

Гиперграф

Примеры
транзакций

- **Данные социальной сети**
 (u, w, d) — пользователь u написал слово w в блоге d
- **Данные сети интернет-рекламы**
 (u, b, d) — 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 b купил товар g у продавца s

Транзакция \leftrightarrow определенное ребро.

Данные \leftrightarrow множество ребер гиперграфа.

Гиперграф — обобщение графа, в котором ребром могут соединяться не только две вершины, но и **любое подмножество вершин**.

Транзакционные данные — наблюдаемые ребра гиперграфа

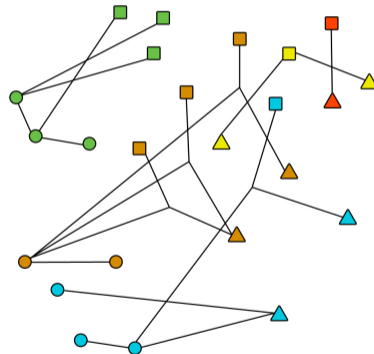
Множество модальностей M :

□ ○ △

Множество типов рёбер K :

d	□	□	○	△	□
x	○	△	○	○	○-△
n_{dx}	3	3	4	2	4

 n_{dx} — число рёбер (d, x) в гиперграфеМножество тем T :



Задача: найти тематические вектора всех вершин

Постановка задачи

Для оптимизации параметров применим принцип максимума правдоподобия:

$$\sum_{k \in K} \tau_k \sum_{(d,x) \in E_k} n_{dx} \ln \sum_{t \in T} \theta_{td} \prod_{v \in X} \varphi_{vtk} + \underbrace{R(\Phi, \Theta)}_{\text{регуляризатор}} \rightarrow \max_{\Phi, \Theta}$$

где τ_k — вес ребёр типа k , $\theta_{td} = p(t | d)$, $\varphi_{vtk} = p_k(v | t)$.

Ограничения:

$$\sum_{v \in V_m} \varphi_{vtk} = 1, \varphi_{vtk} \geq 0$$

$$\sum_{t \in T} \theta_{td} = 1, \theta_{td} \geq 0$$

d — документ

v — вершина

m — модальность

k — тип ребра

t — тема

EM алгоритм

E шаг. Вычисление распределения тем для каждого (d, x) :

$$p_{tdx} = \text{norm}_{t \in T} \left(\theta_{td} \prod_{v \in x} \varphi_{vtk} \right), \text{ где } \text{norm}_{i \in I} a_i = \frac{\max\{a_i, 0\}}{\sum_{j \in I} \max\{a_j, 0\}}.$$

M шаг. Оценивание параметров модели:

$$\varphi_{kvtk} = \text{norm}_{v \in V_m} \left(n_{vtk} + \varphi_{vtk} \frac{\partial R}{\partial \varphi_{vtk}} \right); \quad n_{vtk} = \sum_{(d,x) \in E_k} [v \in x] \tau_k n_{dx} p_{ktdx};$$

$$\theta_{td} = \text{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad n_{td} = \sum_{k \in K} \sum_{(d,x) \in E_k} \tau_k n_{dx} p_{ktdx};$$

Теорема

Если функция $R(\Phi, \Theta)$ непрерывно дифференцируема и (Φ, Θ) — точка локального максимума рассматриваемой задачи, то выполняется система уравнений, описанная выше, относительно параметров φ_{vtk} , θ_{td} и вспомогательных переменных p_{ktdx} , n_{td} и n_{vtk} .

Цели эксперимента:

- 1 Проверить, способен ли алгоритм восстановить параметры модели, с помощью которой были порождены данные.
- 2 Оценить устойчивость модели относительно инициализации, выбора числа тем, размера данных и разреженности исходной матрицы Θ .

Постановка эксперимента:

- 1 Вычисление матрицы Θ с использованием различных методов:
 - PLSA;
 - MultiARTM;
 - TransARTM.
- 2 Решение задачи классификации документов с использованием распределения тем $p(t | d)$ в качестве признаков.
- 3 Оценка качества восстановления матрицы Θ путем вычисления точности решения задачи классификации документов:

$$\text{Точность} = \frac{1}{N} \sum_{i=1}^N [y_i^{\text{pred}} = y_i^{\text{true}}]$$

Модельные данные:

- 1 Генерация матриц $\Theta = p(t | d)$ и $\Phi_k = p_k(v | t)$ для всех $k \in K$.

Число тем: 50

Число классов: 5

Число документов: 5000

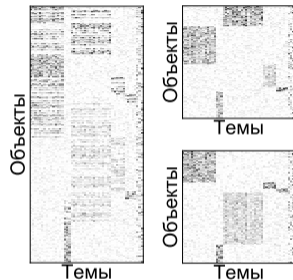
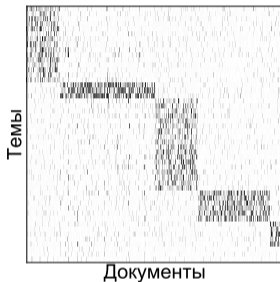
Число объектов: 1000

Число типов ребер: 9

Число модальностей: 3

$$\Theta \in \mathbb{R}^{50 \times 5000}$$

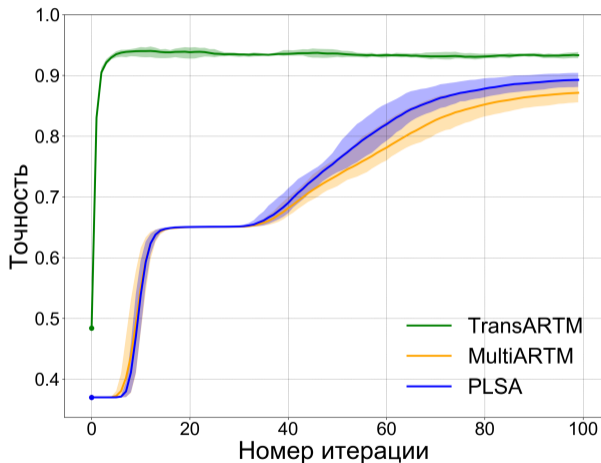
$$\Phi_k \in \mathbb{R}^{1000 \times 50}$$



- 2 Генерация данных (транзакций) на основе матриц Θ , Φ_k .
Общее число транзакций: $\sim 13\,500\,000$.

Результаты

Число тем совпадает с заданным при генерации

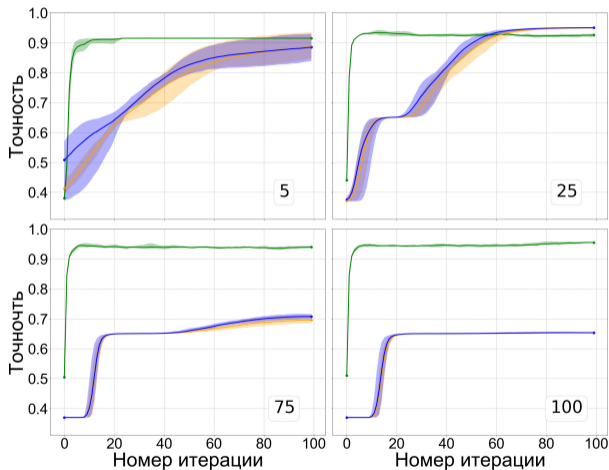


Вывод:
TransARTM достигает высокого качества быстрее других моделей на транзакционных данных.

Результаты

Меняем число тем
от 5 до 100

Число тем в исходной
матрице Θ равно 50.

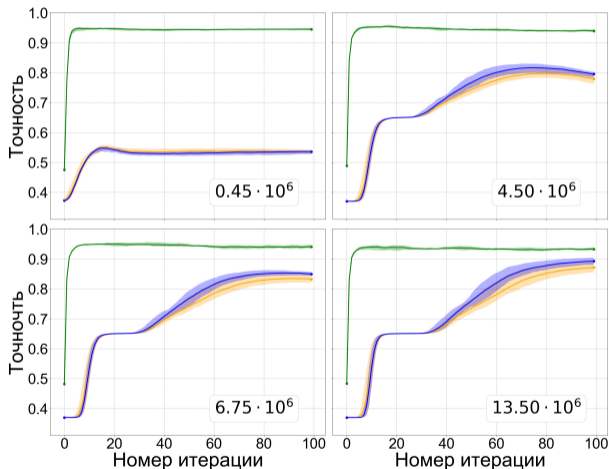


Вывод:

TransARTM наиболее
устойчива относительно
инициализации и выбора
числа тем.

Результаты

Меняем число транзакций от 450 000 до 13 500 000



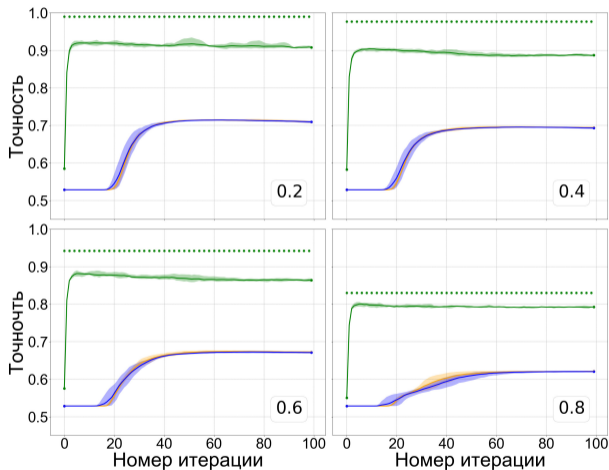
Вывод:

TransARTM

восстанавливает
изначальную структуры
матрицы Θ с высоким
качеством даже при
небольшом числе данных.

Результаты

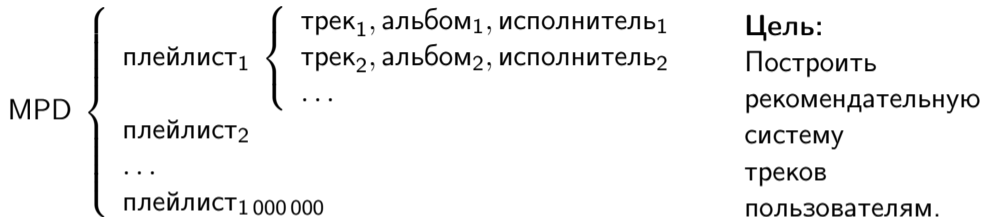
Меняем разреженность
исходной матрицы Θ
от 0.2 до 0.8



Вывод:

TransARTM показывает
близкое к максимальному
качество вне зависимости
от разреженности
матрицы Θ .

The Million Playlist Dataset (MPD)



	MPD	Train	Test	Valid
Число плейлистов:	1 000 000	100 000	1 000	1 000
Число треков:	66 346 428	9 875 306	232 613	232 808
Число уникальных треков:	2 262 292	296 882	39 368	38 641
Число уникальных альбомов:	734 684	140 983	20 690	20 483
Число уникальных артистов:	295 860	69 280	10 081	10 008

Цели эксперимента:

- 1 Применить TransARTM к задаче построения рекомендаций.
- 2 Проанализировать различные гипотезы порождения плейлиста.
- 3 Сравнить результаты с другими моделями.

Постановка эксперимента:

- 1 Для нахождения параметров моделей используем обучающую выборку, состоящую из 100 000 плейлистов.
- 2 Настраиваем коэффициенты регуляризации по сетке на валидационной выборке.
- 3 Предсказываем ранжированный список треков для каждого плейлиста тестовой выборки (последние 70 каждого плейлиста используются для оценки качества).
- 4 Используем следующие метрики для оценки качества:
 - **precision**;
 - **recall**;
 - **fscore**;
 - **ndcg** — normalized discounted cumulative gain.

Модель	Рассматриваемые взаимодействия	Метрики, @500			
		precision	recall	fscore	ndcg
TopTracks	-	0.0230	0.1646	0.0404	0.1152
PLSA	(Pl, Tr)	0.0592	0.4228	0.1038	0.3025
LDA	(Pl, Tr)	0.0583	0.4162	0.1022	0.2988
MultiARTM	(Pl, Al), (Pl, Tr)	0.0594	0.4245	0.1043	0.3029
	(Pl, Ar), (Pl, Tr)	0.0608	0.4343	0.1067	0.3110
	(Pl, Ar), (Pl, Al), (Pl, Tr)	0.0605	0.4321	0.1061	0.3098
TransARTM	(Pl, Al, Tr)	0.0490	0.3497	0.0859	0.2484
	(Pl, Ar, Tr)	0.0504	0.3603	0.0885	0.2555
	(Pl, Al, Tr), (Pl, Ar, Tr)	0.0502	0.3587	0.0879	0.2548
	(Pl, Ar, Al, Tr)	0.0476	0.3398	0.0835	0.2374

Pl aylist
Ar tist
Al bum
Tr ack

Вывод: TransARTM не превосходит рассматриваемые тематические модели на данных, которые могут быть адекватно описаны парными взаимодействиями.

Интерпретируемость

Представление трех различных тем в виде 10 наиболее вероятных артистов в каждой теме (в убывающем порядке):

1	Nicki Minaj	Lil Jon	The Beatles
2	Beyonce	50 Cent	John Lennon
3	Rihanna	Snoop Dogg	George Harrison
4	Tinashe	J-Kwon	The Beach Boys
5	Omarion	Nelly	Elvis Presley
6	Jeremih	Usher	Paul McCartney
7	Trey Songz	Kanye West	David Bowie
8	Chris Brown	R. Kelly	Jim Sturgess
9	Big Sean	Youngbloodz	The Mamas & The Papas
10	Sage The Gemini	Bubba Sparxxx	The Turtles

Результаты

- 1 Предложено обобщение методов тематического моделирования на случай, когда исходные данные представимы в виде гиперграфа.
- 2 Проведены эксперименты на модельных данных, демонстрирующие корректность предложенного метода и преимущество его использования для сложноструктурированных данных.
- 3 Продемонстрировано применение гиперграфовой многомодальной тематической модели для построения рекомендательной системы.

Доклады:

- 1 Международная научная конференция «Ломоносов-2018», «Многомодальные тематические модели на гиперграфах».
- 2 Data Fest⁵, «Гиперграфовые тематические модели для анализа транзакционных данных».

Литература

- ① Daud A. et al. Knowledge discovery through directed probabilistic topic models: a survey // *Frontiers of computer science in China*. – 2010. – Т. 4. – №. 2. – С. 280-301.
- ② Vorontsov K. et al. Non-Bayesian additive regularization for multimodal topic modeling of large collections // *Proceedings of the 2015 Workshop on Topic Models: Post-Processing and Applications*. – ACM, 2015. – С. 29-37.
- ③ Li L., Li T. News recommendation via hypergraph learning: encapsulation of user behavior and news content // *Proceedings of the sixth ACM international conference on Web search and data mining*. – ACM, 2013. – С. 305-314.
- ④ Zhou D., Huang J., Schölkopf B. Learning with hypergraphs: Clustering, classification, and embedding // *Advances in neural information processing systems*. – 2007. – С. 1601-1608.

Каждая вершина $v \in V$ имеет **модальность** $m = \mu(v) \in M$:

$$V = \bigsqcup_{m \in M} V_m, \text{ where } V_m = \{v \in V : \mu(v) = m\}.$$

Каждое транзакция $e \in E$ имеет **тип** $k = \varkappa(e) \in K$:

$$E = \bigsqcup_{k \in K} E_k, \text{ where } E_k = \{e \in E : \varkappa(e) = k\}.$$

Каждому типу рёбер k соответствует дискретное вероятностное пространство:

$$\Omega_k \subseteq 2^V \times T$$

с функцией вероятности $p_k : \Omega_k \rightarrow [0, 1]$.

Обозначения

Будем полагать, что ребра гиперграфа $e \in E_k$ это выборка независимых наблюдений $(e, t) \in \Omega_k$, и каждое ребро входит в выборку n_e раз, и с каждым вхождением ребра связана своя латентная тема $t \in T$.

Каждая вершина v связана с латентными темами:

$$p_k(v, t) = p_k(v | t)p_k(t) = p_k(t | v)p_k(v).$$

Вероятностные распределения нормируются внутри каждой модальности:

$$\sum_{v \in V_m} p_k(v) = 1; \quad \sum_{v \in V_m} p_k(v | t) = 1.$$

Предположения

- 1 Для каждого типа рёбер k первая модальность является *контейнером*.

D — множество всех вершин-контейнеров в гиперграфе.

$(d, x) \in E_k$ произвольное ребро типа k , где x — множество всех остальных вершин ребра.

- 2 Распределение тем в вершине-контейнере d не зависит от типа ребра:

$$p_k(t | d) = p(t | d) \text{ for all } k \in K.$$

- 3 Гипотеза условной независимости вершин в ребрах:

$$p_k(x | t, d) = p_k(x | t) = \prod_{v \in x} p_k(v | t).$$

Гиперграфовая тематическая модель

Вероятности появления рёбер гиперграфа:

$$\begin{aligned}
 p_k(d, x) &= p_k(d) p_k(x | d) = p_k(d) \sum_{t \in T} p_k(x | d, t) p_k(t | d) = \\
 &= p_k(d) \sum_{t \in T} p(t | d) \prod_{v \in X} p_k(v | t) = p_k(d) \sum_{t \in T} \theta_{td} \prod_{v \in X} \varphi_{vtk}.
 \end{aligned}$$

Параметры модели:

- 1 Распределения вершин в темах: $\varphi_{vtk} = p_k(v | t)$;
- 2 Распределение тем в вершинах-контейнерах $\theta_{td} = p(t | d)$;
- 3 Вероятности $p_k(d)$ оцениваются по наблюдаемым данным:

$$p_k(d) = \frac{\sum_{(d,x) \in E_k} n_{dx}}{\sum_{e \in E_k} n_e}.$$

Гиперграфовая тематическая модель (TransARTM)

Гиперграфовая тематическая модель определяется:

- 1 ориентированным гиперграфом $\Gamma = \langle V, E \rangle$,
- 2 множеством модальностей M ,
- 3 разбиением множества вершин по модальностям $\mu: V \rightarrow M$,
- 4 множеством типов рёбер K ,
- 5 разбиением множества рёбер по типам $\varkappa: E \rightarrow K$,
- 6 множеством тем T ,
- 7 вероятностным пространством Ω_k с распределением p_k для каждого $k \in K$,
- 8 параметрами модели $\varphi_{vtk} = p_k(v | t)$ и $\theta_{td} = p(t | d)$.

Постановка задачи

Для оптимизации параметров модели применим *принцип максимума правдоподобия* для каждого типа рёбер k :

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

где τ_k — веса ребер типа k , $\theta_{td} = p(t | d)$, $\varphi_{vtk} = p_k(v | t)$.

Ограничения:

$$\sum_{v \in V_m} \varphi_{vtk} \in \{0, 1\}, \quad \varphi_{vtk} \geq 0, \quad k \in K, m \in M, t \in T.$$

$$\sum_{t \in T} \theta_{td} \in \{0, 1\}, \quad \theta_{td} \geq 0, \quad t \in T, d \in D;$$

Теоретическое обоснование

Теорема 1. Если функция $R(\Phi, \Theta)$ непрерывно дифференцируема и (Φ, Θ) — точка локального максимума рассматриваемой задачи, то выполняется система уравнений, приведенная ниже, относительно параметров φ_{vtk} , θ_{td} и вспомогательных переменных p_{ktdx} , n_{td} и n_{vtk} :

$$p_{ktdx} = \operatorname{norm}_{t \in T} \left(\theta_{td} \prod_{v \in X} \varphi_{vtk} \right),$$

$$\varphi_{vtk} = \operatorname{norm}_{v \in V_m} \left(n_{vtk} + \varphi_{vtk} \frac{\partial R}{\partial \varphi_{vtk}} \right), \quad n_{vtk} = \sum_{(d,x) \in E_k} [v \in X] \tau_k n_{dx} p_{ktdx},$$

$$\theta_{td} = \operatorname{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right), \quad n_{td} = \sum_{k \in K} \sum_{(d,x) \in E_k} \tau_k n_{dx} p_{ktdx}.$$

Доказательство

Сначала выведем уравнение для p_{ktdx} используя формулу Байеса:

$$\begin{aligned}
 p_{ktdx} = p_k(t | d, x) &= \frac{p_k(t, d, x)}{p_k(d, x)} = \frac{p_k(x | d, t) p_k(t | d)}{p_k(x | d)} = \frac{p_k(x | d, t) p_k(t | d)}{\sum_{t \in T} p_k(x | d, t) p_k(t | d)} = \\
 &= \operatorname{norm}_{t \in T} (p_k(x | d, t) p_k(t | d)) = \operatorname{norm}_{t \in T} (p_k(x | t) p(t | d)) = \operatorname{norm}_{t \in T} \left(\theta_{td} \prod_{v \in X} \varphi_{vtk} \right).
 \end{aligned}$$

Воспользуемся условиями Каруша–Куна–Таккера и запишем лагранжиан оптимизационной задачи:

$$\begin{aligned}
 \mathcal{L}(\Phi, \Theta) &= \sum_{k \in K} \tau_k \sum_{(d, x) \in E_k} n_{dx} \ln \sum_{t \in T} \theta_{td} \prod_{v \in X} \varphi_{vtk} + R(\Phi, \Theta) - \sum_{k \in K} \sum_{m \in M} \sum_{t \in T} \lambda_{kmt} \left(\sum_{v \in V_m} \varphi_{vtk} - 1 \right) - \\
 &- \sum_{k \in K} \sum_{m \in M} \sum_{v \in V_m} \sum_{t \in T} \lambda_{kmvt} \varphi_{vtk} - \sum_{d \in D} \mu_d \left(\sum_{t \in T} \theta_{td} - 1 \right) - \sum_{d \in D} \sum_{t \in T} \mu_{td} \theta_{td}.
 \end{aligned}$$

Доказательство

Приравняем нулю производные лагранжиана по параметрам модели:

$$\frac{\partial \mathcal{L}}{\partial \varphi_{vtk}} = \sum_{(d,x) \in E_k} [v \in x] \tau_k n_{dx} \frac{\theta_{td} \prod_{u \in x \setminus v} \varphi_{utk}}{p_k(x|d)} + \frac{\partial R}{\partial \varphi_{vtk}} - \lambda_{k\mu(v)t} - \lambda_{k\mu(v)vt} = 0;$$

$$\frac{\partial \mathcal{L}}{\partial \theta_{td}} = \sum_{k \in K} \sum_{(d,x) \in E_k} \tau_k n_{dx} \frac{\prod_{v \in x} \varphi_{vtk}}{p_k(x|d)} + \frac{\partial R}{\partial \theta_{td}} - \mu_d - \mu_{td} = 0.$$

Умножим левую и правую части первого равенства на φ_{vtk} , второго — на θ_{td} :

$$\sum_{(d,x) \in E_k} [v \in x] \tau_k n_{dx} \underbrace{\frac{\theta_{td} \prod_{u \in x} \varphi_{utk}}{p_k(x|d)}}_{p_{ktdx} = p_k(t|d,x)} + \varphi_{vtk} \frac{\partial R}{\partial \varphi_{vtk}} = \lambda_{k\mu(v)t} \varphi_{vtk};$$

$$\sum_{k \in K} \sum_{(d,x) \in E_k} \tau_k n_{dx} \underbrace{\frac{\theta_{td} \prod_{v \in x} \varphi_{vtk}}{p_k(x|d)}}_{p_{ktdx} = p_k(t|d,x)} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} = \mu_d \theta_{td}.$$

Доказательство

Запишем эти уравнения короче через переменные n_{vtk} и n_{td} :

$$\begin{aligned}\varphi_{vtk} \lambda_{k\mu(v)t} &= n_{vtk} + \varphi_{vtk} \frac{\partial R}{\partial \varphi_{vtk}}. \\ \theta_{td} \mu_d &= n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}};\end{aligned}$$

Если $\lambda_{kmt} \leq 0$, то тема t вырождена: $\varphi_{vtk} = 0$ для всех $v \in V_m$. Если $\lambda_{kmt} > 0$, то либо обе части равенства выше положительны, либо $\varphi_{vtk} = 0$. Объединяя оба случая в одну формулу, получим:

$$\varphi_{vtk} \lambda_{k\mu(v)t} = \left(n_{vtk} + \varphi_{vtk} \frac{\partial R}{\partial \varphi_{vtk}} \right)_+.$$

Если $\mu_d \leq 0$, то вершина-контейнер d вырождена: $\theta_{td} = 0$ для всех $t \in T$. Если $\mu_d > 0$, то либо обе части равенства выше положительны, либо $\theta_{td} = 0$. Объединяя оба случая в одну формулу, получим:

$$\theta_{td} \mu_d = \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right)_+.$$

Просуммируем равенства, полученные на предыдущем слайде, по $v \in V_m$ и по $t \in T$, применим условия нормировки и выразим двойственные переменные:

$$\lambda_{kmt} = \sum_{v \in V_m} \left(n_{vtk} + \varphi_{vtk} \frac{\partial R}{\partial \varphi_{vtk}} \right)_+ ;$$

$$\mu_d = \sum_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right)_+ .$$

Подставив λ_{kmt} и μ_d , получим искомые выражения. Теорема доказана. ■

Онлайновый EM-алгоритм для TransARTM

Вход: коллекция $\bigcup_{k \in K} D_k$, разделенная на части $D_b, b = 1, \dots, B$;

Выход: φ_{vtk} для всех $v \in V, t \in T, k \in K$;

- 1: инициализировать φ_{vtk} для всех $v \in V, t \in T, k \in K$;
- 2: $n_{vtk} := 0, \tilde{n}_{vtk} := 0$ для всех $v \in V, t \in T, k \in K$;
- 3: **for all** $D_b, b = 1, \dots, B$ **do**
- 4: обработать каждый $d \in D_b$ при неизменной матрице Φ :
 $(\tilde{n}_{vtk}) := (\tilde{n}_{vtk}) + \text{ProcessBatch}(D_b, \Phi)$;
- 5: **if** synchronize **then**
- 6: $n_{vtk} := n_{vtk} + \tilde{n}_{vtk}$ для всех $v \in V, t \in T, k \in K$;
- 7: $\varphi_{vtk} = \text{norm}_{v \in V_m} \left(n_{vtk} + \varphi_{vtk} \frac{\partial R}{\partial \varphi_{vtk}} \right)$ для всех $v \in V_m, m \in M, t \in T, k \in K$;
- 8: $\tilde{n}_{vtk} := 0$ для всех $v \in V, t \in T, k \in K$;

Онлайновый EM-алгоритм для TransARTM

ProcessBatch обрабатывает $d \in D_b$ не изменяя матрицу Φ .

Вход: множество вершин-контейнеров D_b , матрица Φ ;

Выход: матрица (\tilde{n}_{vtk}) ;

- 1: $\tilde{n}_{vtk} := 0$ для всех $v \in V, t \in T, k \in K$;
- 2: **for all** $d \in D_b$ **do**
- 3: инициализировать $\theta_{td} := \frac{1}{|T|}$ для всех $t \in T$;
- 4: **repeat**
- 5: $p_{tdx} = \mathop{\text{norm}}_{t \in T} \left(\theta_{td} \prod_{v \in x} \varphi_{vtk} \right)$ для всех $t \in T, k \in K, (d, x) \in E_k$;
- 6: $n_{td} = \sum_{k \in K} \sum_{(d, x) \in E_k} \tau_k n_{dx} p_{tdx}$ для всех $t \in T$;
- 7: $\theta_{td} = \mathop{\text{norm}}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right)$ для всех $t \in T$;
- 8: **until** θ_{td} не сойдутся;
- 9: $\tilde{n}_{vtk} := \tilde{n}_{vtk} + \sum_{(d, x) \in E_k} [v \in x] \tau_k n_{dx} p_{tdx}$ для всех $v \in V, t \in T, k \in K$;

Вероятностный процесс порождения данных:

Вход: K , распределения $p(t | d)$, $p_k(v | t)$, для всех $k \in K$;

Выход: ребра гиперграфа (транзакции);

1: **for all** $d \in D$ **do**

▷ D — множество вершин-контейнеров

2: определить $K' \subset K$;

▷ K — множество типов ребер

3: **for all** $k \in K'$ **do**

4: определить число транзакций — n_{dk} ;

5: **for all** $i = 1, \dots, n_{dk}$ **do**

6: $d_i := d$;

7: выбрать случайную тему t_i из $p(t | d_i)$;

8: **for all** $j = 2, \dots, h(k)$ **do**

▷ $h(k) = |e|$, $e \in E_k$

9: выбрать случайный объект v_i из $p_k(v | t_i)$;

The Million Playlist Dataset (MPD)

Датасет содержит 1 000 000 плейлистов:

- pid** – playlist id;
- name** – the name of the playlist;
- description** – the description given to the playlist;
- modified_at** – timestamp when this playlist was last updated;
- num_artists** – the total number of unique artists;
- num_albums** – the number of unique albums;
- num_tracks** – the number of tracks in the playlist;
- num_followers** – the number of followers this playlist;
- num_edits** – the number of separate editing sessions;
- duration_ms** – the total duration of all the tracks in the playlist;
- collaborative** – if true, the playlist is a collaborative playlist;
- tracks** – an array of information about each track in the playlist:
 - track_name** – the name of the track;
 - track_uri** – the Spotify URI of the track;
 - album_name** – the name of the track's album;
 - album_uri** – the Spotify URI of the album;
 - artist_name** – the name of the track's primary artist;
 - artist_uri** – the Spotify URI of track's primary artist;
 - duration_ms** – the duration of the track in milliseconds;
 - pos** – the position of the track in the playlist (zero-based).