

МИНОБРНАУКИ РОССИИ  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования

**«САРАТОВСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
ГОСУДАРСТВЕННЫЙ УНИВЕРСИТЕТ  
ИМЕНИ Н. Г. ЧЕРНЫШЕВСКОГО»**

Кафедра математической кибернетики и компьютерных наук

**СОЗДАНИЕ ФРЕЙМВОРКА ДЛЯ РАБОТЫ С  
АВТОМАТИЗИРОВАННЫМ ТЕСТИРОВАНИЕМ**

КУРСОВАЯ РАБОТА

Студента 3 курса 351 группы  
направления 09.03.04 — Программная инженерия  
факультета КНиИТ  
Кондрашова Даниила Владиславовича

Научный руководитель  
зав.каф.техн.пр.,  
доцент, к. ф.-м. н.

\_\_\_\_\_

С. В. Папшев

Заведующий кафедрой  
к. ф.-м. н.

\_\_\_\_\_

С. В. Миронов

Саратов 2024

# **СОДЕРЖАНИЕ**

## **ВВЕДЕНИЕ**

С ростом объёмов информации в наше время для её эффективного поиска и изучения стало необходимо уметь её классифицировать и структурировать. Сегодня физически невозможно найти нужные данные просто перебирая все ресурсы подряд, появилась острая потребность в поиске по темам, в классификации данных.

Данную проблему призвано решить тематическое моделирование. Оно способно быстро и эффективно автоматически разбить по темам огромные объёмы информации.

# 1 Математические основы тематического моделирования

## 1.1 Основная гипотеза тематического моделирования

Тематическое моделирование — это подход анализа текстовых данных, направленный на выявление семантических структур в коллекции документов.

Само тематическое моделирование основывается на предположении, что слова в тексте связаны не с документом, а с темой. Кроме того первично текст разбивается на темы, затем каждая из них порождает слово для соответствующих позиций в документе. Таким образом, сначала порождается тема, а потом термины.

Благодаря этой гипотезе можно по частоте и встречаемости слов производить тематическую классификацию текстов.

## 1.2 Аксиоматика тематического моделирования

Каждый текст можно количественно охарактеризовать. Вот основные количественные характеристики, используемые при тематическом моделировании:

- $W$  — конечное множество термов;
- $D$  — конечное множество текстовых документов;
- $T$  — конечное множество тем;
- $D \times W \times T$  — дискретное вероятностное пространство;
- коллекция — i.i.d выборка  $(d_i, w_i, t_i)_{i=1}^n$ ;
- $n_{dwt} = \sum_{i=1}^n [d_i = d][w_i = w][t_i = t]$  — частота  $(d, w, t)$  в коллекции;
- $n_{wt} = \sum_d n_{dwt}$  — частота термина  $w$  в документе  $d$ ;
- $n_{td} = \sum_w n_{dwt}$  — частота термов темы  $t$  в документе  $d$ ;
- $n_t = \sum_{d,w} n_{dwt}$  — частота термов темы  $t$  в коллекции;
- $n_{dw} = \sum_t n_{dwt}$  — частота термина  $w$  в документе  $d$ ;
- $n_W = \sum_d n_{dw}$  — частота термина  $w$  в коллекции;
- $n_d = \sum_w n_{dw}$  — длина документа  $d$ ;
- $n = \sum_{d,w} n_{dw}$  — длина коллекции.

Также в тематическом моделировании используются следующие гипотезы и аксиомы:

- Независимость слов от порядка в документе: порядок слов в документе не важен;
- Независимость от порядка документов в коллекции: порядок документов

в коллекции не важен;

- Зависимость терма от темы: каждый терм связан с соответствующей темой и порождается ей;
- Гипотеза условной независимости:  $p(w|d, t) = p(w|t)$ .

Вышеперечисленные характеристики, гипотезы и аксиомы являются основой тематического моделирования, являющейся достаточной для построения тематической модели.

### 1.3 Задача тематического моделирования

Как уже говорилось ранее, документ порождается следующим образом:

1. для каждой позиции в документе генерируется тема  $p(t|d)$ ;
2. для каждой сгенерированной темы в соответствующей позиции генерируем терм  $p(w|d, t)$ .

Тогда вероятность появления слова в документе можно описать по формуле полной вероятности:

$$p(w|d) = \sum_{t \in T} p(w|d, t)p(t|d) = \sum_{t \in T} p(w|t)p(t|d) \quad (1)$$

Такой алгоритм является прямой задачей порождения текста. Тематическое моделирование призвано решить обратную задачу:

1. для каждого терма  $w$  в тексте найти вероятность появления в теме  $t$  (найти  $p(w|t) = \phi_{wt}$ );
2. для каждой темы  $t$  найти вероятность появления в документе  $d$  (найти  $p(t|d) = \theta_{td}$ ).

Обратную задачу можно представить в виде стохастического матричного разложения.

Таким образом, тематическое моделирование ищет величину  $p(w|d)$ .

### 1.4 Решение обратной задачи

Для решения задачи тематического моделирования необходимо найти величину  $p(w|d)$ , сделать это можно с помощью метода максимального правдоподобия.

#### 1.4.1 Лемма о максимизации функции на единичных симплексах

Перед выведением решения обратной задачи выпишем лемму, позволяющую это решение найти.

Введём операцию нормировки вектора:

$$p_i = \left( x_i \right) = \frac{\max x_i, 0}{\sum_{k \in I} \max x_k, 0} \quad (2)$$

#### Лемма о максимизации функции на единичных симплексах:

Пусть функция  $f(\Omega)$  непрерывно дифференцируема по набору векторов  $\Omega = (w_i)_{j \in J}$ ,  $w_j = (w_{ij})_{i \in I_j}$  различных размерностей  $|I_j|$ . Тогда векторы  $w_j$  локального экстремума задачи

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

при условии  $1^0$  :  $(\exists i \in I_j) w_{ij} \frac{\partial f}{\partial w_{ij}} > 0$  удовлетворяют уравнениям

$$w_{ij} = \text{norm}_{i \in I_j} \left( w_{ij} \frac{\partial f}{\partial w_{ij}} \right), \quad i \in I_j; \quad (3)$$

при условии  $2^0$  :  $(\forall i \in I_j) w_{ij} \frac{\partial f}{\partial w_{ij}} \leq 0$  и  $(\exists i \in I_j) w_{ij} \frac{\partial f}{\partial w_{ij}} < 0$  удовлетворяют уравнениям

$$w_{ij} = \text{norm}_{i \in I_j} \left( -w_{ij} \frac{\partial f}{\partial w_{ij}} \right), \quad i \in I_j; \quad (4)$$

в противном случае (условие  $3^0$ ) — однородным уравнениям

$$w_{ij} \frac{\partial f}{\partial w_{ij}} = 0, \quad i \in I_j. \quad (5)$$

Данная лемма служит для оптимизации любых моделей, параметрами которых являются неотрицательные нормированные векторы.

#### 1.4.2 Сведение обратной задачи к задаче максимизации функционала

Чтобы вычислить величину  $p(w|d)$  воспользуемся принципом максимума правдоподобия, согласно которому будут подобраны параметры  $\Phi$ ,  $\Theta$  такие,

что  $p(w|d)$  примет наибольшее значение.

$$\prod_{i=1}^n p(d_i, w_i) = \prod_{d \in D} \prod_{w \in d} p(d, w)^{n_{dw}} \quad (6)$$

Прологарифмировав правдоподобие, перейдём к задаче максимизации логарифма правдоподобия.

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln p(w|d) \xrightarrow{const} n_{dw} \rightarrow \max \quad (7)$$

Данная задача эквивалентна задаче максимизации функционала

$$L(\Phi, \Theta) = \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} \rightarrow \max_{\Phi, \Theta} \quad (8)$$

при ограничениях неотрицательности и нормировки

$$\phi_{wt} \geq 0; \quad \sum_{w \in W} \phi_{wt} = 1; \quad \theta_{td} \geq 0; \quad \sum_{t \in T} \theta_{td} = 1 \quad (9)$$

#### 1.4.3 Аддитивная регуляризация тематических моделей

Задача не соответствует определению корректно поставленной задачи по Адамару, так как она в общем случае имеет бесконечно много решений, следовательно задачу нужно доопределить.

Для доопределения некорректно поставленных задач используют регуляризацию: к основному критерию добавляют дополнительный критерий — регуляризатор, соответствующий решаемой задаче.

ARTM: аддитивная регуляризация тематических моделей основана на максимизации линейной комбинации логарифма правдоподобия и регуляризаторов  $R_i(\Phi, \Theta)$  с неотрицательными коэффициентами регуляризации  $t\tau_i$ ,  $i = 1, \dots, k$ .

Преобразуем задачу к ARTM виду:

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}; \quad R(\Phi, \Theta) = \sum_{i=1}^k \tau_i R_i(\Phi, \Theta) \quad (10)$$

при ограничениях неотрицательности и нормировки.

#### 1.4.4 Е-М алгоритм

Из ограничений видно, что столбцы матриц можно принять за неотрицательные единичные векторы, а, следовательно, задача является задачей максимизации функции на единичных симплексах.

Воспользуемся леммой о максимизации функции на единичных симплексах и перепишем задачу.

Пусть функция  $R(\Phi, \Theta)$  непрерывно дифференцируема. Тогда точка  $(\Phi, \Theta)$  локального экстремума задачи с ограничениями, удовлетворяет системе уравнений ссc вспомогательными переменными  $p_{tdw} = p(t|d, w)$ , если из решения исключить нулевые столбцы матриц  $\Phi$  и  $\Theta$ :

$$\begin{cases} p_{tdw} = \underset{t \in T}{\text{norm}}(\phi_{wt}\theta_{td}) \\ \phi_{wt} = \underset{w \in W}{\text{norm}}\left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}}\right); \\ \theta_{td} = \underset{t \in T}{\text{norm}}\left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}}\right) \end{cases} \quad (11)$$

Полученная модель соответствует Е-М алгоритму, где первая строка системы уравнений соответствует Е шагу, а вторая и третья строки — М шагу.

Решив полученную систему уравнений, методом простых итерации получим искомые матрицы  $\Phi$  и  $\Theta$ .

### 1.5 Регуляризаторы в тематическом моделировании

В этом разделе будут рассмотрены некоторые возможные варианты регуляризаторов.

#### 1.5.1 Дивергенция Кульбака-Лейблера

Чтобы оценить близость тем можно воспользоваться дивергенцией Кульбака-Лейблера (KL или KL-дивергенция). KL-дивергенция позволяет оценить степень вложенности одного распределения в другое, в случае тематического моделирования будет оцениваться вложенность матриц.

Определим KL-дивергенцию:

Пусть  $P = (p_i)_{i=1}^n$  и  $Q = (q_i)_{i=1}^n$  некоторые распределения. Тогда дивергенция Кульбака-Лейблера имеет следующий вид:

$$KL(P||Q) = KL_i(p_i||q_i) = \sum_{i=1}^n p_i \ln \frac{p_i}{q_i}. \quad (12)$$



Свойства KL-дивергенции:

1.  $KL(P||Q) \geq 0$ ;
2.  $KL(P||Q) = 0 \Leftrightarrow P = Q$ ;
3. Минимизация KL эквивалентна максимизации правдоподобия:

$$KL(P||Q(\alpha)) = \sum_{i=1}^n p_i \ln \frac{p_i}{q_i(\alpha)} \rightarrow \min_{\alpha} \Leftrightarrow \sum_{i=1}^n p_i \ln q_i(\alpha) \rightarrow \max_{\alpha};$$

4. Если  $KL(P||Q) < KL(Q||P)$ , то  $P$  сильнее вложено в  $Q$ , чем  $Q$  в  $P$ .

### 1.5.2 Регуляризатор сглаживания

Сглаживание предполагает сближение тем, это может быть полезно в следующих случаях:

1. Темы могут быть похожи между собой по терминологии, например, основы теории вероятностей и линейной алгебры обладают рядом одинаковых терминов;
2. При выделении фоновых тем важно максимально вобрать в них слова, следовательно, сглаживание поможет решить эту задачу.

Определим регуляризатор сглаживания:

Пусть распределения  $\phi_{wt}$  близки к заданному распределению  $\beta_w$  и пусть распределения  $\theta_{td}$  близки к заданному распределению  $\alpha_t$ . Тогда в форме KL-дивергенции выразим задачу сглаживания:

$$\sum_{t \in T} KL(\beta_w || \phi_{wt}) \rightarrow \min_{\Phi}; \quad \sum_{d \in D} KL(\alpha_t || \theta_{td}) \rightarrow \min_{\Theta}. \quad (13)$$

Согласно свойству 3 KL-дивергенции перейдём к задаче максимизации правдоподобия:

$$R(\Phi, \Theta) = \beta_o \sum_{t \in T} \sum_{w \in W} \beta_w \ln \phi_{wt} + \alpha_0 \sum_{d \in D} \sum_{t \in T} \alpha_t \ln \theta_{td} \rightarrow \max. \quad (14)$$

Перепишем ЕМ-флгоритм в соответствии с полученной формулой:

$$\begin{cases} p_{tdw} = \underset{t \in T}{\text{norm}}(\phi_{wt}\theta_{td}) \\ \phi_{wt} = \underset{w \in W}{\text{norm}}(n_{wt} + \beta_0\beta_w); \\ \theta_{td} = \underset{t \in T}{\text{norm}}(n_{td} + \alpha_0\alpha_t) \end{cases} \quad (15)$$

### 1.5.3 Регуляризатор разреживания

Разреживание предполагает разделение тем и документов, исключение из них общих слов. Данный тип регуляризации отталкивается от того, что в большинстве своём темы и документы специфичны и описываются относительно небольшим набором терминов, не встречающихся в других темах.

Определим регуляризатор разреживания:

Пусть распределения  $\phi_{wt}$  далеки от заданного распределения  $\beta_w$  и пусть распределения  $\theta_{td}$  далеки от заданного распределения  $\alpha_t$ . Тогда в форме KL-дивергенции выразим задачу сглаживания:

$$\sum_{t \in T} KL(\beta_w || \phi_{wt}) \rightarrow \max_{\Phi}; \quad \sum_{d \in D} KL(\alpha_t || \theta_{td}) \rightarrow \max_{\Theta}. \quad (16)$$

Согласно свойству 3 KL-дивергенции перейдём к задаче максимизации правдоподобия:

$$R(\Phi, \Theta) = -\beta_o \sum_{t \in T} \sum_{w \in W} \beta_w \ln \phi_{wt} - \alpha_0 \sum_{d \in D} \sum_{t \in T} \alpha_t \ln \theta_{td} \rightarrow \max. \quad (17)$$

Перепишем ЕМ-флгоритм в соответствии с полученной формулой:

$$\begin{cases} p_{tdw} = \underset{t \in T}{\text{norm}}(\phi_{wt}\theta_{td}) \\ \phi_{wt} = \underset{w \in W}{\text{norm}}(n_{wt} - \beta_0\beta_w); \\ \theta_{td} = \underset{t \in T}{\text{norm}}(n_{td} - \alpha_0\alpha_t) \end{cases} \quad (18)$$

### 1.5.4 Регуляризатор декоррелирования тем

Декоррелятор тем — это частный случай разреживания, призванный выделить для каждой темы лексическое ядро — набор термов, отличающий её от других тем:

Определим регуляризатор декоррелирования:

Минимизируем ковариации между вектор-столбцами  $\phi_t$ :

$$R(\Phi) = -\frac{\tau}{2} \sum_{t \in T} \sum_{s \in T \setminus t} \sum_{w \in W} \phi_{wt} \phi_{ws} \rightarrow \max. \quad (19)$$

Перепишем ЕМ-алгоритм в соответствии с полученной формулой:

$$\begin{cases} p_{tdw} = \underset{t \in T}{\text{norm}}(\phi_{wt} \theta_{td}) \\ \phi_{wt} = \underset{w \in W}{\text{norm}} \left( n_{wt} - \tau \phi_{wt} \sum_{t \in T \setminus t} \phi_{ws} \right); \\ \theta_{td} = \underset{t \in T}{\text{norm}} \left( n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right) \end{cases} \quad (20)$$

## 1.6 Оценка качества моделей тематического моделирования

После обучения модели, очевидно, нужно оценить её качество.

Перечислим основные критерии оценки качества тематических моделей:

1. Внешние критерии (оценка производится экспертами):
  - а) Полнота и точность тематического поиска;
  - б) Качество ранжирования при тематическом поиске;
  - в) Качество классификации / категоризации документов;
  - г) Качество суммаризации / сегментации документов;
  - д) Экспертные оценки качества тем.
2. Внутренние критерии (оценка производится программно):
  - а) Правдоподобие и перплексия;
  - б) Средняя когерентность (согласованность тем);
  - в) Разреженность матриц  $\Phi$  и  $\Theta$ ;
  - г) Различность тем;
  - д) Статистический тест условной независимости.

Так как оценка по внешним критериям не представляется возможной в рамках данной работы, то рассмотрим внутренние критерии оценки, так как их можно вычислять автоматически.

### 1.6.1 Правдоподобия и перплексия

Перплексия основывается на логарифме правдоподобия и является его некоторой модификацией.

$$P(D) = \exp \left( -\frac{1}{n} \sum_{d \in D} \sum_{w \in d} n_{dw} \ln p(w|d) \right), \quad n = \sum_{d \in D} \sum_{w \in d} n_{dw} \quad (21)$$

Не трудно заметить, что, при равномерном распределении слов в тексте  $p(w|d) = \frac{1}{|W|}$ , значение перплексии равно мощности словаря  $P = |W|$ . Тогда можно сделать вывод, что перплексия — это мера различности и неопределённости слов в тексте, то есть, чем меньше перплексия, тем различнее вероятности появления слов в тексте.

Таким образом, чем меньше перплексия, тем больше слов с большей вероятностью  $p(w|d)$ , которые модель умеет лучше предсказывать, следовательно, чем меньше перплексия, тем лучше.

### 1.6.2 Когерентность

Когерентность является мерой, коррелирующей с экспертной оценкой интерпретируемости тем, которую можно вычислять автоматически.

Когерентность (согласованность) темы  $t$  по  $k$  топовым словам:

$$PNI_t = \frac{2}{k(k-1)} \sum_{i=1}^{k-1} \sum_{j=i+1}^k PMI(w_i, w_j), \quad (22)$$

где  $w_i$  —  $i$ -ое слово в порядке убывания  $\phi_{wt}$ ,  $PMI(u, v) = \ln \frac{|D|N_{uv}}{N_u N_v}$  — потоковая взаимная информация,  $N_{uv}$  — число документов, в которых слова  $u, v$  хотя бы один раз встречаются рядом (расстояние определяется отдельно),  $N_u$  — число документов, в которых  $u$  встретился хотя бы один раз.

Гипотезу когерентности можно выразить так: когда человек говорит о какой-либо теме, то часто употребляет достаточно ограниченный набор слов, относящийся к этой теме, следовательно, чем чаще будут встречаться вместе слова этой темы, тем лучше её можно будет интерпретировать.

Сама когерентность берёт самые часто встречающиеся слова из тем, и вычисляет для каждой пары из них насколько они часто встречаются, соответственно, чем выше будет значение взаимовстречаемости, тем лучше.

### 1.6.3 Разреженность и различность

Разреженность — доля нулевых элементов в матрицах  $\Phi$  и  $\Theta$ . Разреженность служит для выявления различности тем, так как каждая тема состоит из небольшого набора слов, то и остальные слова в ней должны встречаться нечасто, что соответствует нулевым элементам в матрицах. Разреженность должна быть в рамках оптимальных значений, высокой, но не слишком, тогда темы будут хорошо различимы, в противном случае, они либо не будут различаться (разреженность слишком низкая), либо будут содержать слишком мало слов (разреженность слишком высокая).

- Чистота темы:  $\sum_{w \in W_t} p(w|t)$ , где  $W_t$  — ядро темы:  $W_t = \{w : p(w|t) > \alpha\}$ , где  $\alpha$  подбирается по разному, например  $\alpha = 0.25$  или  $\alpha = \frac{1}{|W|}$ . Данная характеристика показывает как вероятно относится ядро темы к фоновым словам темы, следовательно, чем больше вероятность ядра, тем лучше;
- Контрастность темы:  $\frac{1}{|W_T|} \sum_{w \in W_t} p(t|w)$ . Данная характеристика показывает насколько часто слова из ядра темы встречаются в других темах, очевидно, что чем меньше ядро будет встречаться в других темах, тем лучше.

## **2 Тематическое моделирование новостей**

### **ЗАКЛЮЧЕНИЕ**

### **СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ**