

# Вероятностное тематическое моделирование несбалансированных текстовых коллекций

Панкратов Виктор Владимирович

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

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

Москва,  
2021г

# Постановка задачи: вероятностная модель

Заданы три множества:

- $D$  — множество документов
- $W$  — множество слов
- $T$  — множество тем

$n_{wd}$  — число вхождений слова  $w \in W$  в документ  $d \in D$

Требуется найти матрицы  $\Phi, \Theta$ :  $\phi_{wt} = p(w|t)$ ,  $\theta_{td} = p(t|d)$

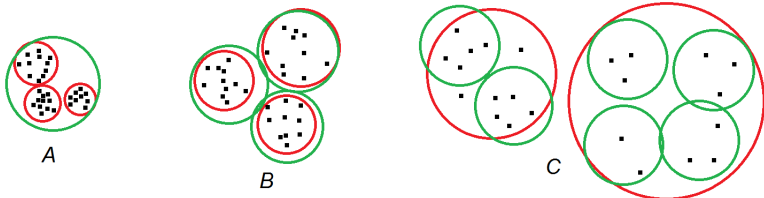
Критерий — максимизация регуляризованного правдоподобия:

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + \sum_i \tau_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}, \quad (1)$$

где  $R_i$  — регуляризаторы,  $\tau_i$  — коэффициенты регуляризации

Решение (1) осуществляется с помощью EM алгоритма.

# Проблема несбалансированности

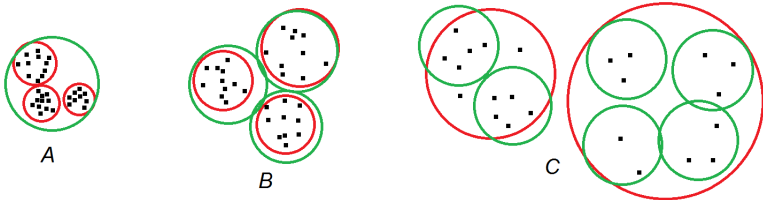


В случае тематически несбалансированной коллекции матричное разложение путём максимизации правдоподобия может приводить к дроблению крупных тем (А) и слиянию мелких (С).

## Цель работы:

Предложить и экспериментально проверить решение проблемы несбалансированности с помощью регуляризатора

# Семантическая неоднородность



Гипотеза условной независимости:

$$p(w|t) = p(w|d, t); p(w, d|t) = p(w|t)p(d|t)$$

Проверка - статистика семантической неоднородности:

$$S_t = \text{KL}(p(w, d|t) || p(w|t)p(d|t))$$

Тема — кластер размерности  $|W|$ , центр которого -  $p(w|t)$ .

$S_t$  — удаленность  $p(w|d, t)$  от центра кластера.

# Регуляризатор семантической неоднородности

Статистика семантической неоднородности

$$S_t = \text{KL}(\hat{p}(w, d|t) || p(w|t)p(d|t)) = \sum_{d \in D} \sum_{w \in d} \hat{p}(w, d|t) \ln \frac{\hat{p}(w, d|t)}{p(w|t)p(d|t)}$$

Здесь  $\hat{p}$  — частотные оценки вероятности  
Преобразовывая и суммируя по всем темам:

$$\sum_{t \in T} S_t = \sum_{d \in D} \sum_{w \in d} \left( \sum_{t \in T} \hat{p}(w, d|t) \right) \ln \frac{\hat{p}(w|d)}{p(w|d)}$$

В результате преобразований получаем регуляризатор:

$$R = \sum_{t \in T} S_t = \sum_{d \in D} \sum_{w \in d} \beta_{dw} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td}, \quad \beta_{dw} = \sum_{t \in T} \frac{p(t|d, w)}{p(t)}$$

# Эксперимент: генерация коллекции

Для эксперимента будем генерировать синтетическую коллекцию:

- Генерируются столбцы  $\Phi_0, \Theta_0$ , из симметричных распределений Дирихле
- Генерируется одна фоновая тема из распределения Ципфа
- Фоновая тема — столбец в  $\Phi_0$
- Фоновая тема — константная строка в  $\Theta_0$
- Из полученных матриц  $\Phi_0, \Theta_0$  генерируются документы

$$t_i \sim \text{Dir}(t|d_i), w_i \sim \text{Dir}(w|t_i), i = 1 \dots \quad (2)$$

## Эксперимент: оценивание качества решения

- $\Phi_0$  — исходная матрица вероятностей  $p(w|t)$
- $\Phi$  — матрица вероятностей  $p(w|t)$ , найденная алгоритмом

Для всех пар  $i, j$  будем проверять равенства:

$$\arg \min_k (\text{dist}(\Phi[i], \Phi_0[k])) = j \quad (3)$$

$$\arg \min_k (\text{dist}(\Phi[k], \Phi_0[j])) = i \quad (4)$$

$\text{dist}$  — заданная функция расстояния между столбцами  $\Phi, \Phi_0$

Критерии качества восстановления матриц — число тем:

- взаимно близких: (3), (4) выполнены для некоторых  $i, j$
- невосстановленных:  $\Phi_0[j]$ , если (3) не выполнено для всех  $i$
- ложных:  $\Phi[i]$ , если (4) не выполнено для всех  $j$

# Эксперимент: параметры

Параметры генерации синтетической коллекции:

- Число тем  $|T| = 100$
- Число документов  $|D| = 2000$
- Число слов в каждом документе  $|w \in (d \in D)| = 1000$
- Коэффициент регуляризации  $\tau = 0.5$

Степень несбалансированности коллекции:

- Отношение минимальной и максимальной сумм вероятностей:

$$S_1 = \frac{\max_{t \in T} \sum_d \theta_{td}}{\min_{t \in T} \sum_d \theta_{td}}$$

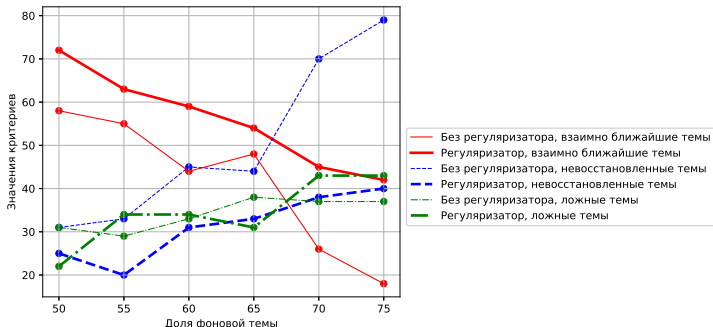
- Отклонение столбцов  $\Theta$  от равномерного распределения:

$$S_2 = \left| \Theta - \frac{1}{|T|} \right|$$



# Эксперимент: сбалансированные коллекции

$S_1 \approx 1.2$ ,  $S_2 \approx 20$ , изменяем доли фоновых тем

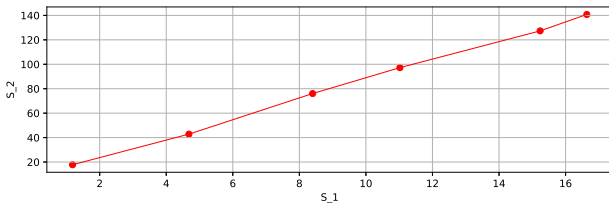


- Добавление регуляризатора улучшает восстановление тем
- С ростом доли фоновой темы темы восстанавливаются хуже

# Эксперимент: несбалансированные коллекции

Генерируется шесть коллекций с одной темой, мощность которой во много раз превышает мощность остальных.

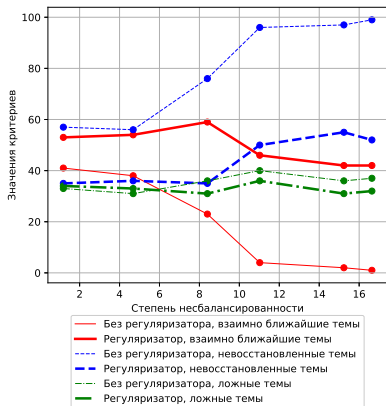
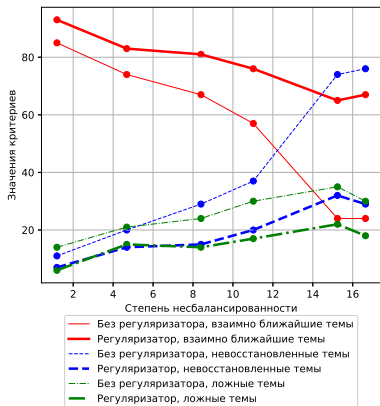
Степени несбалансированности шести синтетических коллекций:



Для сбалансированных коллекций:  $S_1 \approx 1.2$ ,  $S_2 \approx 20$

# Эксперимент: несбалансированные коллекции

Доля фоновой темы: 0.5 — первый эксперимент, 0.7 — второй



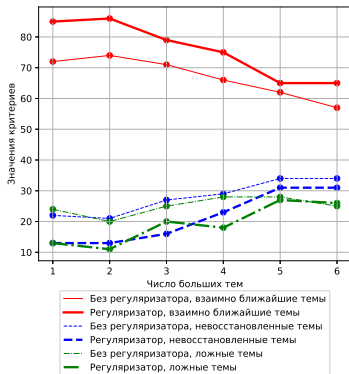
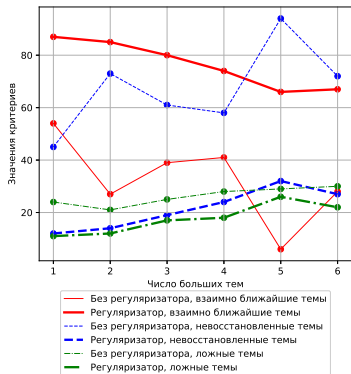
- Модель с регуляризатором — темы лучше восстанавливаются
- Темы восстанавливаются хуже при увеличении  $S_1, S_2$
- Темы восстанавливаются хуже с ростом доли фоновой темы

# Эксперимент: несбалансированные коллекции

Изменяем число больших тем в коллекции

$S_1 \approx 3.5$ ,  $S_2 \approx 25 \cdot (\text{число больших тем} + 1)$

Расстояния: косинусные (слева), Йенсена-Шеннона (справа)

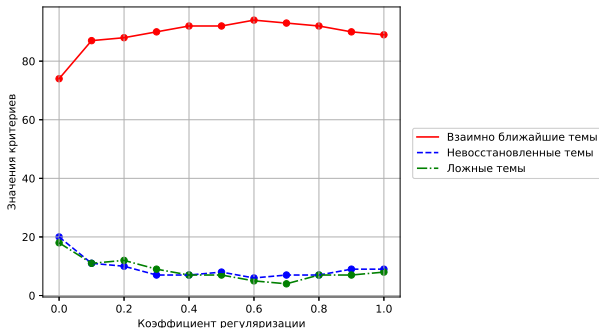


- Расстояние Йенсена-Шеннона даёт более устойчивые результаты
- Качество восстановления тем падает с ростом  $S_2$

# Эксперимент: несбалансированные коллекции

Коллекция: половина крупных, половина мелких тем.

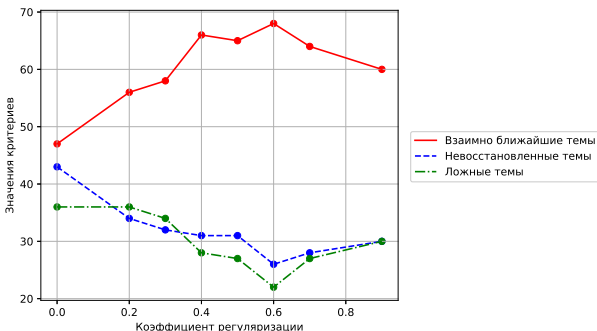
$S_1 = 1.3$ ,  $S_2 = 107$ , изменяем коэффициент регуляризации



- Большая  $S_2$  — не всегда несбалансированная коллекция
- Оптимум в окрестности  $\tau \approx 0.6$

# Эксперимент: несбалансированные коллекции

Коллекция: по  $\frac{1}{20}$  крупных и мелких тем, остальные равномощны.  
 $S_1 = 1.65$ ,  $S_2 = 48$ , изменяем коэффициент регуляризации



- Малое значение  $S_2$  — не всегда сбалансированная коллекция
- Оптимум в окрестности  $\tau \approx 0.6$

# Результаты, выносимые на защиту

- На синтетических данных показано, что тематическая несбалансированность коллекции приводит к дроблению крупных тем и слиянию мелких
- Предложен алгоритм устранения проблемы несбалансированности путем добавления регуляризатора на основе семантической неоднородности тем
- В экспериментах показано, что предложенный регуляризатор отчасти устраняет проблему несбалансированности