

# Оценивание способов борьбы с несбалансированностью тем в тематическом моделировании

Рогозина Анна

Московский Физико-технический институт

Факультет управления и прикладной математики

Кафедра интеллектуальных систем

Научный руководитель: д. ф.-м. н. Воронцов Константин Вячеславович

2019

# План

- 1 Постановка задачи
  - Тематическое моделирование
  - Проблема несбалансированности
  - Оценивание качества моделей
  - Предложенные способы решения
- 2 Описание эксперимента
  - Общая постановка
  - Генерация коллекций
  - Квантильная регрессия
- 3 Результаты
  - Проверка качества предобученной модели

# Постановка задачи

## Задача тематического моделирования

### Дано:

- Множество токенов  $W$ , коллекция текстовых документов  $D$ , множество тем  $T$
- $n_{wd}$  - частоты токенов в документах
- $D \times W \times T$  - дискретное вероятностное пространство

### Предположение:

- Гипотеза условной независимости:  $p(w | d, t) = p(w | t)$

**Найти:**  $p(w | d) = \sum_{t \in T} \varphi_{wt} \theta_{td}$

- $\varphi_{wt} = p(w | t)$  - вероятность токенов  $w$  в теме  $t$
- $\theta_{td} = p(t | d)$  - вероятность тем  $t$  в документе  $d$

# Принцип максимума правдоподобия

Оптимизируемый функционал:

$$\mathcal{L}(\Phi, \Theta) = \sum_{d \in D} \sum_{w \in d} n_{dw} \ln \left( \sum_{t \in T} \varphi_{wt} \theta_{td} \right) \rightarrow \max_{\Phi, \Theta}$$

- Мощность темы  $\hat{p}(t) = \frac{n_t}{n}$
- Гипотеза:

Отношение максимальной к минимальной мощности  
 $\frac{\hat{p}_{max}(t)}{\hat{p}_{min}(t)} \leq 3 - 4$  раза

## Задача

Научиться выявлять «несбалансированность» тем в построенной модели. Сравнить способы борьбы с появлением «несбалансированности» тем.

# Кластерная структура распределений слов

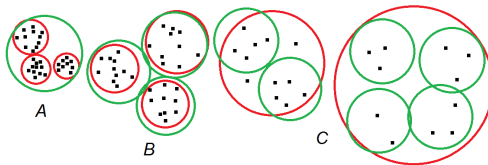


Рис.: Иллюстрация работы Е-М алгоритма для несбалансированных тем

- Точки на графике - распределения слов в документах  $p(w | t, d)$
- Центры кластеров - распределение слов в теме  $p(w | t)$



## Обозначения

- $S_{dt}$  - значение дивергенции  $CR_\lambda$  для документа  $d$  и темы  $t$
- Радиус семантической однородности  $R_t^\alpha(n_{td})$  темы  $t$  с уровнем значимости  $\alpha$

## Степень семантической неоднородности

$$\text{SemHeterogeneity}(t) = \frac{\sum_{d \in D} p(t|d) [S_{dt} > R_t^\alpha(n_{td})]}{\sum_{d \in D} p(t|d)}$$

## Степень семантической загрязненности

$$\text{SemImpurity}(t) = \frac{\sum_{d \in D} p(t|d) [S_{dt} < R_t^\alpha(n_{td})] [S_{dt'} < R_{t'}^\alpha(n_{td})]}{\sum_{d \in D} p(t|d)}$$

# Итеративная балансировка тем

## Локальный экстремум задачи тематического моделирования:

$$p_{tdw} = \text{norm}_{t \in T}(\varphi_{wt} \theta_{td}); \quad (1)$$

$$\varphi_{wt} = \text{norm}_{w \in W} \left( n_{wt} + \varphi_{wt} \frac{\partial R}{\partial \varphi_{wt}} \right); \quad n_{wt} = \sum_{d \in D} n_{dw} p_{tdw}; \quad (2)$$

$$\theta_{td} = \text{norm}_{t \in T} \left( n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad n_{td} = \sum_{w \in d} n_{dw} p_{tdw}. \quad (3)$$

$n_{tdw} = n_{dw} p_{tdw}$  - оценка числа употреблений слова  $w$  в документе  $d$  по теме  $t$ .

### Идея:

Домножим  $n_{tdw}$  на величину, обратно пропорциональную  $n_t$ .



## Изменение вероятностей:

$$p'_{tdw} = \frac{n'_{tdw}}{\sum_s n'_{sdw}} = \frac{\frac{1}{n_t} n_{tdw}}{\sum_s \frac{1}{n_s} n_{sdw}} = \frac{1}{n_t} p_{tdw} \frac{\sum_s n_{sdw}}{\sum_s \frac{1}{n_s} n_{sdw}} = \text{norm}_{t \in T} \left( \frac{1}{n_t} p_{tdw} \right).$$

## Изменение формул Е-М алгоритма:

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

$$Z_{dw} = \sum_{t \in T} \frac{p_{tdw}}{n_t}.$$

$$\varphi_{wt} = \text{norm}_{w \in W} \left( n_{wt} + \varphi_{wt} \frac{\partial R}{\partial \varphi_{wt}} \right); \quad n_{wt} = \sum_{d \in D} \frac{n_{dw}}{n_t Z_{dw}} p_{tdw};$$

$$\theta_{td} = \text{norm}_{t \in T} \left( n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right); \quad n_{td} = \sum_{w \in d} \frac{n_{dw}}{n_t Z_{dw}} p_{tdw}.$$

# Проверка качества борьбы с несбалансированностью

## Алгоритм

- Генерируем список коллекций с разной степенью сбалансированности тем  $\{D_i\}_{i=1}^n$
- Для каждой коллекции  $D_i$  обучаем модель ARTM  $\{(\Phi_i, \Theta_i)\}_{i=1}^n$
- Для каждой модели оцениваем качество всех тем с помощью  $\text{SemHeterogeneity}(t)$  и  $\text{SemImpurity}(t)$

# Генерация коллекций

## Дано:

- Коллекция документов из «ПостНауки» (3404 документа)
- Предобученная на этой коллекции матрица  $\Phi$  на 20 тем

## Алгоритм генерации коллекции длины $N$ :

- Из предобученной модели берем распределение  $p(w | t)$
- Задаем дискретное распределение тем в коллекции  $p(t)$
- Для каждого документа  $\{d_i\}_{i=1}^N$ :
  - Из списка длин реальных документов «ПостНауки» равновероятно выбираем длину генерируемого документа  $n_i$
  - Генерируем распределение  $p(t | d) \in \text{Dir}(p(t))$
  - Генерируем слова  $\{w_j^i\}_{j=1}^{n_i}$ :
    - Выбираем тему  $t_j \in p(t | d)$
    - Добавляем в документ слово  $w_j \in p(w | t_j)$

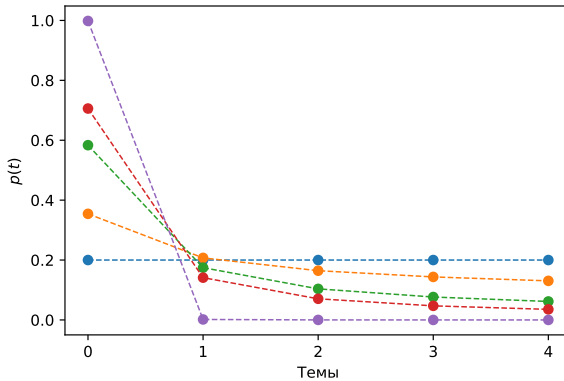


Рис.: Пример градации  $p(t)$  для 5 тем

# Вычисление $R_t^\alpha(n_{td})$

## Проблема:

Распределение статистики Кресси-Рида  $CR_\lambda \in \chi^2(|W|)$ , но только если  $\forall w, d : n_{wd} \geq 5$

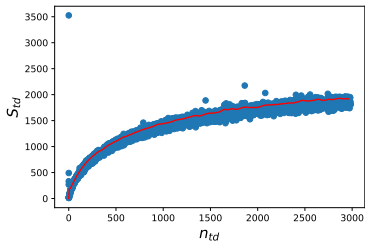
## На самом деле

Распределение статистики зависит от  $t$  и  $n_{td}$

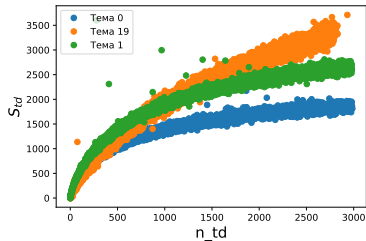
## Предложение

Для каждой модели  $(\Phi_i, \Theta_i)$ , для каждой темы  $t_j^i$

- Генерируем коллекцию  $D_{t_j^i}$ , состоящую из документов  $d_{t_i}$  разной длины, содержащих только тему  $t_j^i$
- Вычисляем значение статистики  $S_{t_j^i d}$
- По полученной эмпирической зависимости  $S_{t_j^i d}(n_{td})$  считаем непарметрическую квантильную регрессию  $R_t^\alpha(n_{td})$ .



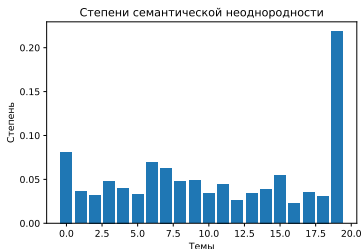
а) Квантильная регрессия  
для темы 0



б) Распределение  $S_{td}$  для  
разных тем

# Проверка качества предобученной модели

Подсчитаны степени семантической однородности и  
загрязненности



## Выводы

- Видим, что степень загрязненности  $\sim 1$  для каждой темы
- Значит, либо каждый документ содержит больше двух тем
- Либо критерий в данном виде слишком маломощный и вообще ничего не отвергает (Сейчас статистика считается для всех слов  $w : p(w | t) > \frac{1}{|W|}$ )

## Возможные модификации

- Из всех 20 тем отобрать те, что не пересекаются в смысле степеней загрязненности (если это возможно)
- Начать считать статистику  $S_{td}$  только для слов, входящих в документ  $d$



## Заключение

Что сделано:

- Поставлена задача
- Предложен способ оценивания качества моделей
- Разработан дизайн эксперимента
- Проведена оценка качества предобученной модели

Нужно:

- Увеличить мощность используемого критерия  $CR_\lambda$
- Провести эксперимент для итеративной балансировки тем
- Реализовать оптимизацию гиперпараметров сглаживания
- Провести эксперимент и оптимизации гиперпараметров