

# Поиск полного набора тем с помощью обучения нескольких тематических моделей

Алексеев Василий Антонович

Научный руководитель:

д.ф.-м.н. Воронцов К. В.

МФТИ, группа 874

17 июня 2020

# Задача тематического моделирования

Вероятностная тематическая модель:

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

Максимизация log-правдоподобия по параметрам  $\Phi, \Theta$ :

$$\mathcal{L}(\Phi, \Theta | D) = \sum_{d \in D} \sum_{w \in W_d} n_{wd} \log \sum_{t \in T} \varphi_{wt}\theta_{td} \rightarrow \max_{\Phi, \Theta}$$

Аддитивная регуляризация<sup>1</sup>:

$$\begin{cases} \mathcal{L} + \sum_{i=1}^n \tau_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta} \\ \tau_i \geq 0, i = 1, \dots, n \end{cases}$$

---

<sup>1</sup>Воронцов К. В., Потапенко А. А. *Аддитивная регуляризация тематических моделей* // Доклады Академии наук. – 2014. – Т. 456. – №. 3. – С. 268-271.

- 1 Постановка задачи
- 2 Построение банка тем
- 3 Использование банка тем для оценки качества моделей

- Тематические модели *неполны и неустойчивы*.
- Получение хорошей тематической модели, как правило, требует больших затрат времени.
- Не существует идеального автоматического способа оценивания качества тематических моделей.

## Решение

Банк тем — инструмент для сохранения интерпретируемых тем, построенных при многократных запусках, с целью последующего их использования для оценки качества моделей.

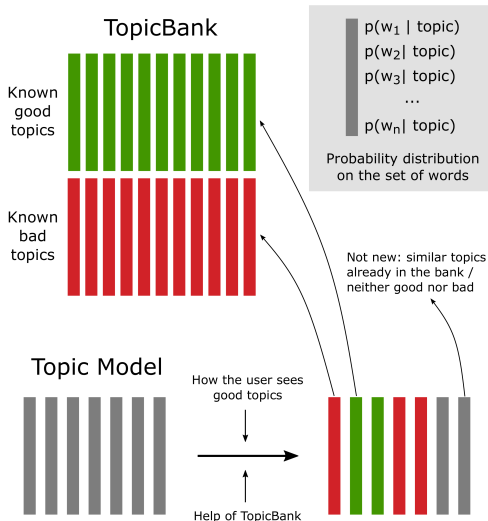
## Цели

Реализовать метод построения банка тем и оценивания качества тематических моделей с помощью банка тем.

# Банк тем: сохранение интерпретируемых тем

Банк тем — модель  
полного набора тем:  
таких тем, которые

- 1) интерпретируемы,
- 2) существенно  
различны,
- 3) обеспечивают  
высокое  
правдоподобие  
модели  
 $p(\Phi, \Theta | D)$ .

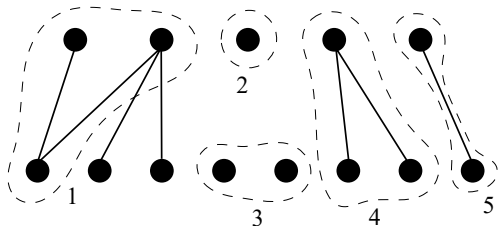


# Построение банка тем

Аналогично построению двухуровневой иерархической тематической модели:

$$\underbrace{p(w | t)}_{\varphi_{wt}^{parent}} = \sum_{s \in S} \underbrace{p(w | s)}_{\varphi_{ws}^{child}} \underbrace{p(s | t)}_{\psi_{st}} \quad \text{Hierarchy}$$

$$\underbrace{p(w | t)}_{\varphi_{wt}^{bank}} = \sum_{s \in S} \underbrace{p(w | s)}_{\varphi_{ws}^{new}} \underbrace{p(s | t)}_{\psi_{st}} \quad \text{TopicBank}$$



№	Hierarchy	TopicBank
1	ok	no
2	ok	ok
3	no	ok
4	ok	maybe
5	ok	maybe

- 1 Постановка задачи
- 2 Построение банка тем
- 3 Использование банка тем для оценки качества моделей

# Построение банка тем

- Многократное обучение тематических моделей
- Бинарная классификация тем на интерпретируемые и нет на основании одного признака – когерентности<sup>1</sup>

# B: topic bank

# N: number of model trainings

for i = 1 to N:

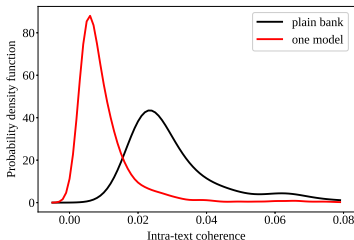
  model ← train\_model(i)

  new\_topics ← get\_new(model, B)

  good\_topics ← get\_good(model)

  B ← update(B,

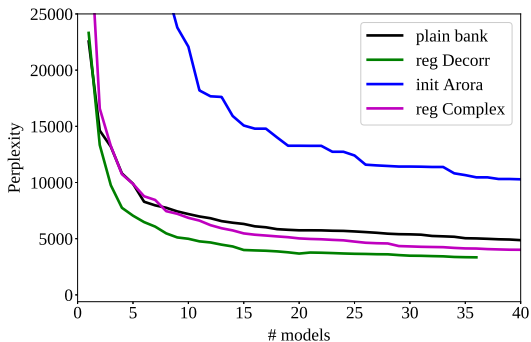
    new\_topics ∩ good\_topics)



<sup>1</sup>Alekseev V., Bulatov V., and Vorontsov K. *Intra-text coherence as a measure of topic models' interpretability*. Dialogue, 2018



# Зависимость банка тем от числа обучаемых моделей



- Наилучшая перплексия – при декоррелировании.
- При добавлении моделей число тем в банке постепенно увеличивается; скорость же пополнения банка снижается.

Vorontsov K. et al. *Additive regularization of topic models*, 2015.

Arora S. et al. *Learning topic models—going beyond SVD*, 2012.

Hofmann T. *Probabilistic latent semantic indexing*, 1999.

- 1 Постановка задачи
- 2 Построение банка тем
- 3 Использование банка тем для оценки качества моделей**

Использование банка тем для оценивания моделей:

- $B$  – множество тем в банке тем
- $\mathcal{D}$  – датасеты,  $|\mathcal{D}| < \infty$
- $\mathcal{M}$  – множество моделей,  $|\mathcal{M}| < \infty$

Модель  $m$  по набору данных даёт множество тем

$$\begin{cases} m : d \mapsto T \\ m \in \mathcal{M}, d \in \mathcal{D} \end{cases}$$

Качество модели можно оценить с помощью банка тем по тому, насколько темы модели  $T$  похожи на темы банка  $B$ .

## Полнота модели $m$ относительно банка $B$

$$\text{recall@bank} = \frac{|t \in B \mid \exists \tau \in T : \rho(t, \tau) < h|}{|B|}$$

$B$  – множество тем в банке тем,  $T$  – множество тем во вновь обученной модели  $m$ , а  $h \in \mathbb{R}_{>0}$  – порог

Где  $\rho(t_1, t_2)$  – расстояние между темами (по мере Жаккара):

$$\rho(t_1, t_2) = 1 - \frac{\sum_{w \in \text{Ker}_{12}} \min_{i \in \{1,2\}} (p(w \mid t_i))}{\left( \sum_{i=1}^2 \sum_{w \in \text{Ker}_i \setminus \text{Ker}_{12}} p(w \mid t_i) + \sum_{w \in \text{Ker}_{12}} \max_{i \in \{1,2\}} (p(w \mid t_i)) \right)}$$

Где  $\text{Ker}_i \equiv \text{Ker}(t_i)$ ,  $\text{Ker}_{12} \equiv \text{Ker}(t_1) \cap \text{Ker}(t_2)$

и  $\text{Ker}(t) = \{w \in t : p(w \mid t) > 1/|W|\}$  – ядро темы  $t$

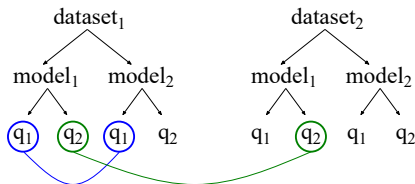
Введённая функция качества зависит не только от модели:

$$\text{recall@bank} = f(d, m, h), \quad d \in \mathcal{D}, \quad m \in \mathcal{M}, \quad h \in H \subseteq \text{Im}(\rho), \quad |H| < \infty$$

Необходимо провести усреднения по  $\mathcal{D}$  и  $H$ :

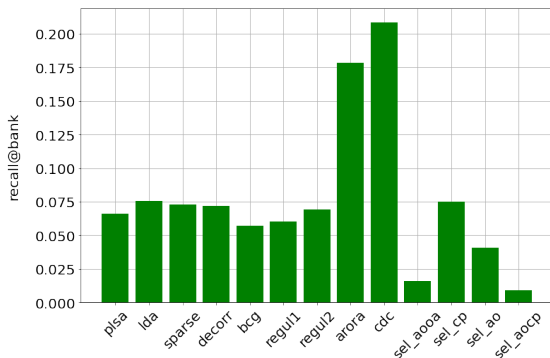
$$\langle \text{recall@bank}(m, d, h) \rangle_h = \frac{\sum_{\eta \in H} w(\eta) \cdot \text{recall@bank}(m, d, \eta)}{\sum_{\eta \in H} w(\eta)}, \quad w(\eta) \geq 0$$

$$\langle \text{recall@bank}(m, d, h) \rangle_{d,h} = \frac{1}{|\mathcal{D}|} \sum_{\delta \in \mathcal{D}} \frac{\langle \text{recall@bank}(m, \delta, h) \rangle_h}{\sum_{\mu \in \mathcal{M}} \langle \text{recall@bank}(\mu, \delta, h) \rangle_h}$$



# Результат, усреднённый по датасетам

**Цель:** по множеству датасетов  $\mathcal{D}$  получить оценки качества моделей  $\text{recall@bank}(m)$ ,  $m \in \mathcal{M}$ .



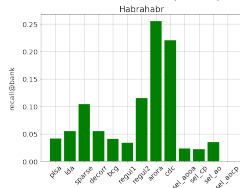
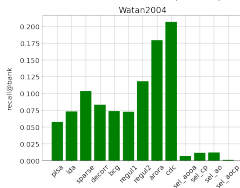
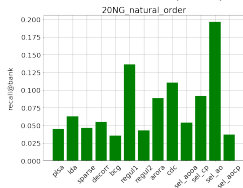
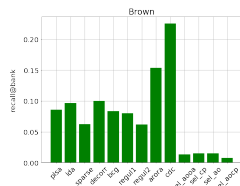
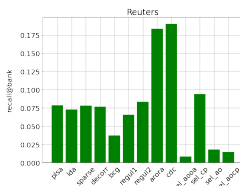
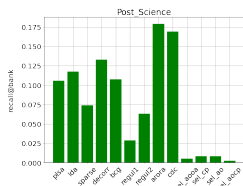
**Вывод:** банк тем помог из фиксированного ряда моделей  $\mathcal{M}$  найти те модели, которые лучше подстраиваются под данные.

---

Dobrynin V. et al. *Contextual document clustering*, 2004.

Blei D. et al. *Latent dirichlet allocation*, 2003.

# Результаты по разным датасетам



**Вывод:** под каждый набор данных нужна своя модель, но в большинстве случаев модели с неслучайной инициализацией превосходят по качеству остальные модели из  $M$ .

[habr.com](http://habr.com), [postnauka.ru](http://postnauka.ru), [nltk.org/book/ch02.html](http://nltk.org/book/ch02.html)

[sites.google.com/site/mouradabbas9/corpora](http://sites.google.com/site/mouradabbas9/corpora)

[scikit-learn.org/0.19/datasets/twenty\\_newsgroups.html](http://scikit-learn.org/0.19/datasets/twenty_newsgroups.html)

- Предложен алгоритм создания банка тем с использованием многократного обучения моделей
- Предложена методика оценивания качества тематических моделей с помощью банка тем
- Реализована система для использования банка тем<sup>1</sup>

## Публикации

- Alekseev V. et al. *TopicNet: Making Additive Regularisation for Topic Modelling Accessible*. LREC, 2020.<sup>2</sup>
- Alekseev V. et al. *Topic Modelling for Extracting Behavioral Patterns from Transactions Data*. IEEE, 2019.<sup>3</sup>
- Alekseev V. et al. *Intra-Text Coherence as a Measure of Topic Models' Interpretability*. Dialogue, 2018.<sup>4</sup>

---

<sup>1</sup>[github.com/machine-intelligence-laboratory/OptimalNumberOfTopics](https://github.com/machine-intelligence-laboratory/OptimalNumberOfTopics)

<sup>2</sup>[aclweb.org/anthology/2020.lrec-1.833](https://aclweb.org/anthology/2020.lrec-1.833)

<sup>3</sup>[ieeexplore.ieee.org/abstract/document/9007329](https://ieeexplore.ieee.org/abstract/document/9007329)

<sup>4</sup>[dialog-21.ru/media/4281/alekseevva.pdf](https://dialog-21.ru/media/4281/alekseevva.pdf)