

Автоматическая настройка параметров BigARTM под широкий класс задач

Гришанов А. В.

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра «Интеллектуальные системы»

Задачу поставил к.ф.-м.н., н.с. ВЦ РАН К. В. Воронцов
Консультант Мурат Апишев

Москва,
2019 г.

Проблема

BigARTM — продвинутая библиотека для тематического моделирования. В ней реализовано много регуляризаторов, что повышает гибкость, но при этом усложняет настройку параметров. В результате на практике популярнее более простые методы, такие как LDA.

Цель работы

Проверить гипотезу о существовании конфигураций, хорошо работающих на широком классе задач.

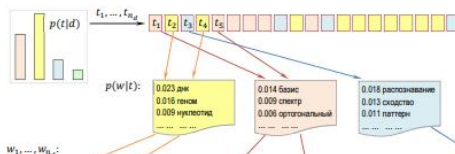
- ❶ Логистическая регрессия с L_1 регуляризацией
- ❷ PLSA
David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. Journal of machine Learning research, 3(Jan):993–1022, 2003.
- ❸ LDA Thomas Hofmann. Probabilistic latent semantic analysis. In Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence, UAI'99, pages 289–296, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.

- D — коллекция текстовых документов, состоящая из документов d
- W — словарь, состоящий из терминов w ;
- T — множество тем, состоящее из тем t .

Согласно формуле полной вероятности и гипотезе условной независимости, распределение термов в документе $p(w|d)$ описывается вероятностной смесью распределений термов в темах $\varphi_{wt} = p(w|t)$ с весами $\theta_{td} = p(t|d)$

Распределение термов в документах

$$p(w|d) = \sum_{t \in T} p(w|t)p(t|d) = \sum_{t \in T} \phi_{wt}\theta_{td}. \quad (1)$$



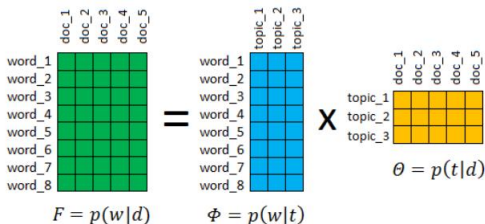
Разработан спектрально-аналитический подход к выявлению размытых протяженных повторов в **геномных последовательностях**. Метод основан на разномасштабном оценивании **сходства нуклеотидных последовательностей** в пространстве коэффициентов разложения фрагментов кривых GC- и GA-содержания по классическим ортогональным базисам. Найлены условия оптимальной аппроксимации, обеспечивающие автоматическое распознавание повторов различных видов (прямых и инвертированных, а также **тандемных**) на спектральной матрице сходства. Метод одинаково хорошо работает на разных масштабах данных. Он позволяет выявлять следы **сегментных дупликаций** и **мегасателлитные участки** в **геноме**, районы **синтении** при сравнении пары **геномов**. Его можно использовать для детального изучения фрагментов **хромосом** (поиска размытых участков с умеренной длиной повторяющегося **паттерна**).

Распределение термов в документах

$$p(w|d) = \sum_{t \in T} p(w|t)p(t|d) = \sum_{t \in T} \phi_{wt}\theta_{td}. \quad (2)$$

Постановка задачи тематического моделирования

Ставится задача разложения матрицы F в произведение двух матриц Φ и Θ меньшего размера



Поставленная задача ($F \approx \Phi \Theta$) эквивалентна поиску матриц Φ и Θ , максимизирующих следующий функционал:

Задача

$$L(\Phi, \Theta) = \sum_{d \in D} \sum_{w \in d} n_{dw} \sum_{t \in T} \phi_{wt} \theta_{td} \rightarrow \max_{\Phi, \Theta} \quad (3)$$

Проблема

Разложение матрицы F в произведение матриц Φ и Θ не единственно. В частности, для любой невырожденной матрицы S размера $T \times T$ верно, что $F = (\Phi S)(S^{-1}\Theta)$.

Таким образом, из-за сложившейся неопределённости, невозможен поиск произвольного матричного разложения, нужно уточнять модель. Можно наложить дополнительные ограничения, что приведёт к сокращению произвольности выбора или же сделать некоторые предположения о вероятностном распределении коллекции. Рассматриваются следующие два подхода к решению проблемы:

Латентное размещение Дирихле

Тематическая модель латентного размещения Дирихле (latent Dirichlet allocation, LDA) основана на дополнительном предположении, что векторы документов $\theta_d = (\theta_{td}) \in \mathbb{R}^{|T|}$ и векторы тем $\phi_t = (\phi_{wt}) \in \mathbb{R}^{|W|}$ порождаются распределениями Дирихле с параметрами $\alpha \in \mathbb{R}^{|T|}$ и $\beta \in \mathbb{R}^{|W|}$.

Аддитивная регуляризация

Тематическая модель аддитивной регуляризации (additive regularization of topic models, ARTM) получается при наложении на модель дополнительных требований (регуляризаторов).

$$L(\Phi, \Theta) + \sum_{i=1}^n \tau_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta} \quad (4)$$

- Рассмотрим набор датасетов $\{\mathcal{D}_{ex}, \mathcal{D}_{in}\}$, где \mathcal{D}_{ex} имеют внешний критерий качества, а \mathcal{D}_{in} — только внутренние.
- Необходимо проверить гипотезу о том, что существуют коэффициенты регуляризации $\tau_{general}$, которые можно считать «универсальными», т.е. для которых метрики качества отличаются от лучших на том же датасете не более чем на 5%.
- Для каждого из первых найдём лучшие параметры, затем будем искать общие.
- В конце проверим выполнение гипотезы на всех данных, для неразмеченных будем сравнивать внутренние критерии качества.

Цель работы — построить модель, которая **не хуже** чем PLSA и **лучше** PLSA по нескольким критериям.

① 20news groups

Best f1_score: 0.9155

General params f1_score: 0.9148

② Victorian Era

Best f1_score: 0.9777

General params f1_score: 0.9777

③ Toxic comments

Best f1_score: 0.9539

General params f1_score: 0.9582