



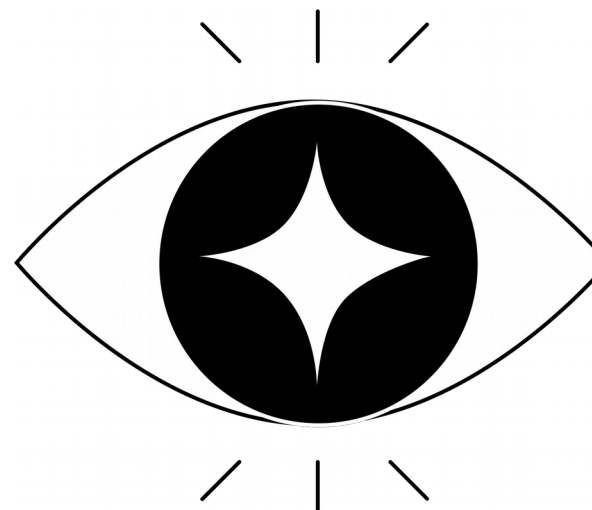
Тематическое моделирование ARTM



Начало работы

На этом уроке

- ARTM теория;
- Практика BigARTM.



Операционная задача

Функционал правдоподобия:

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

Ограничения стохастичности на матрицы:

$$\left. \begin{aligned} \sum_{w \in W} \phi_{wt} &= 1, & \phi_{wt} &\geq 0 \\ \sum_{t \in T} \theta_{td} &= 1, & \theta_{td} &\geq 0 \end{aligned} \right\} \quad (2)$$

Можно решить с помощью EM-алгоритма:



Регуляризация модели PLSA

PLSA позволяет получить из коллекции D при заданном числе тем $|T|$ матрицы параметров Φ и Θ , но задача матричного разложения $F \approx \Phi \Theta$ имеет бесконечное множество решений!

Пусть $S \in \mathbb{R}^{|T| \times |T|}$ — произвольная невырожденная квадратная матрица. Тогда верно:

$$F \approx \Phi \Theta = (\Phi S)(S^{-1} \Theta) = \Phi' \Theta'$$

Задача PLSA является некорректно поставленной. Попробуем наложить на неё больше ограничений — это позволит получить Φ и Θ , которые будут обладать дополнительными полезными свойствами.



Additive Regularization of Topic Models

Стандартный приём доопределения некорректно поставленной задачи — регуляризация.

ARTM (аддитивная регуляризация тематических моделей) расширяет задачу PLSA слагаемыми-регуляризаторами в функционале правдоподобия:

$$\sum_{d \in D} \sum_{w \in W} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta} \quad (3)$$
$$R(\Phi, \Theta) = \sum_i \tau_i R_i(\Phi, \Theta)$$

Здесь R_i — выраженное в виде функционала ограничение, τ_i — весовой гиперпараметр.



ЕМ-алгоритм для ARTM

Теорема: пусть функция $R(\Phi, \Theta)$ непрерывно дифференцируема. Тогда точка (Φ, Θ) локального экстремума задачи (3) с ограничениями (2) удовлетворяет системе уравнений со вспомогательными переменными, если из решения исключить нулевые столбцы Φ и Θ :

$$\begin{array}{l} \text{Е-шаг} \left\{ \begin{array}{l} p_{tdw} = \text{norm}_{t \in T}(\phi_{wt} \theta_{td}) \end{array} \right. \\ \text{М-шаг} \left\{ \begin{array}{l} \phi_{wt} = \text{norm}_{w \in W} \left(n_{wt} + \phi_{wt} \frac{\partial R}{\partial \phi_{wt}} \right), \quad n_{wt} = \sum_{d \in D} n_{dw} p_{tdw} \\ \theta_{td} = \text{norm}_{t \in T} \left(n_{td} + \theta_{td} \frac{\partial R}{\partial \theta_{td}} \right), \quad n_{td} = \sum_{w \in W} n_{dw} p_{tdw} \end{array} \right. \end{array}$$

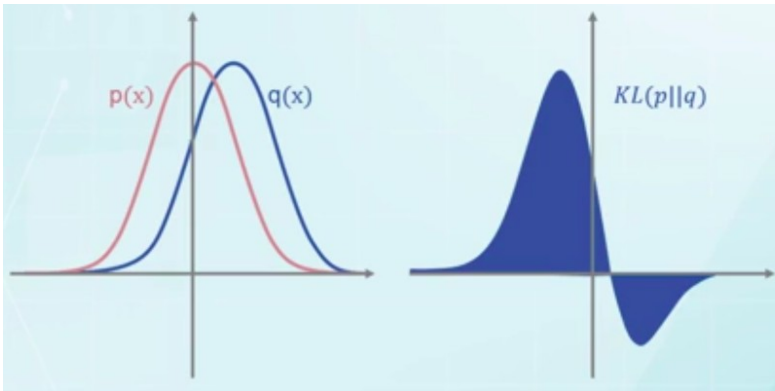
где $\text{norm}_{i \in I}(x_i) = \frac{\max\{x_i, 0\}}{\sum_{j \in I} \max\{x_j, 0\}}$



KL-дивергенция

KL-дивергенция определяется для пары непрерывных распределений p и q без нулевых элементов на одинаковом носителе:

$$KL(p||q) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$



KL-дивергенция

KL-дивергенция определяется для пары непрерывных распределений p и q без нулевых элементов на одинаковом носителе:

$$KL(p||q) = \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

и показывает степень вложенности одного распределения в другое. Это несимметричная мера различия распределений.

Для дискретных распределений p и q на одном носителе мощности k определение аналогичное:

$$KL(p||q) = \sum_{i=1}^k p_i \log \frac{p_i}{q_i}$$



Пример регуляризатора: сглаживание

Базовый вариант регуляризации — сглаживание параметров заданными вероятностными распределениями. Потребуем близости по KL-дивергенции столбцов ϕ_t и θ_d к распределениям

$\beta = (\beta_w)_{w \in W}$ и $\alpha = (\alpha_t)_{t \in T}$:

$$\sum_{t \in T} KL(\beta || \phi_t) \rightarrow \min_{\Phi} \quad \sum_{d \in D} KL(\alpha || \theta_d) \rightarrow \min_{\Theta}$$

Сложим две суммы с коэффициентами τ_1 и τ_2 и удалим константные слагаемые.

$$R(\Phi, \Theta) = \tau_1 \sum_{t \in T} \sum_{w \in W} \beta_w \log \phi_{wt} + \tau_2 \sum_{d \in D} \sum_{t \in T} \alpha_t \log \theta_{td} \rightarrow \max_{\Phi, \Theta}$$



Пример регуляризатора: сглаживание

Базовый вариант регуляризации — сглаживание параметров заданными вероятностными распределениями.

$$R(\Phi, \Theta) = \tau_1 \sum_{t \in T} \sum_{w \in W} \beta_w \log \phi_{wt} + \tau_2 \sum_{d \in D} \sum_{t \in T} \alpha_t \log \theta_{td} \rightarrow \max_{\Phi, \Theta}$$

Применим общие формулы EM-алгоритмы из теоремы. Получим итоговые выражения для M-шага:

$$\phi_{wt} = \text{norm}_{w \in W}(n_{wt} + \tau_1 \beta_w)$$

$$\theta_{td} = \text{norm}_{t \in T}(n_{td} + \tau_2 \alpha_t)$$



Пример регуляризатора: разреживание

Противоположная стратегия регуляризации — разреживание. В ряде случаев оно приводит к получению более интерпретируемых тем и полезно с точки зрения оптимизации ресурсов при обучении.

Потребуем, чтобы столбцы φ_t и Θ_d были далеки от заданных распределений KL-дивергенции.

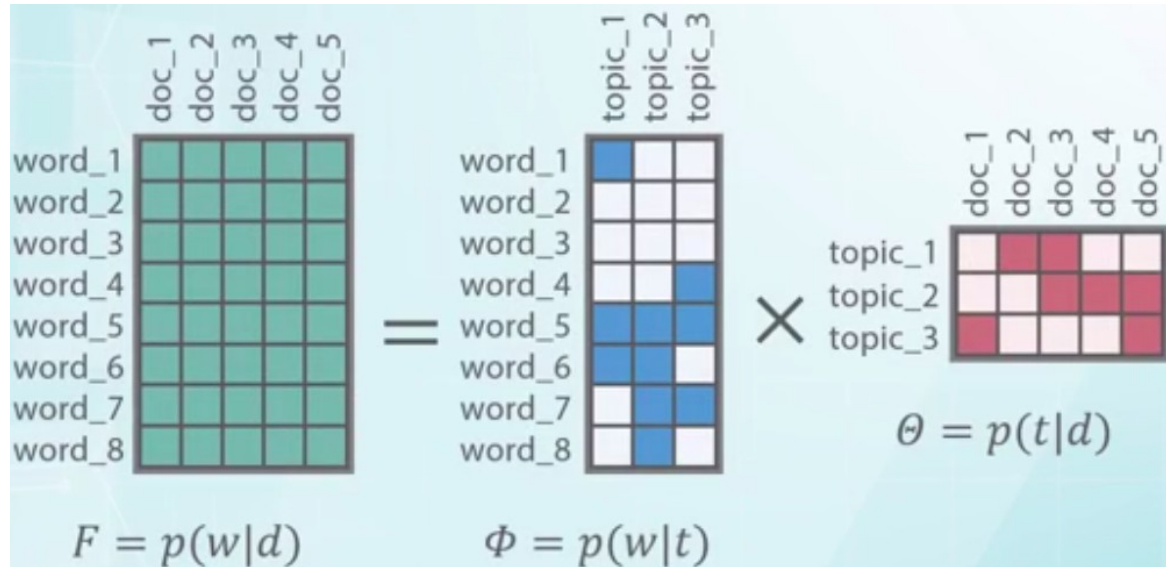
Итоговый регуляризатор и формулы M-шага получатся такими же, как и при сглаживании, но с противоположным знаком.

Разреживание и сглаживание можно успешно комбинировать в разных частях одной модели.



Пример регуляризатора: разреживание

На практике разрежённые модели полезны, однако ещё лучше, если модель разрежена так, чтобы темы получались как можно более различными.



Пример регуляризатора: декорреляция

Формализуем требование различности тем с помощью такого регуляризатора:

$$R(\Phi, \Theta) = -\frac{\tau_3}{2} \sum_{t \in T} \sum_{s \in T \setminus t} \sum_{w \in W} \phi_{wt} \phi_{ws} \rightarrow \max_{\Phi, \Theta}$$

Применим теорему и получим формулы М-шага:

$$\phi_{wt} = \text{norm}_{w \in W} \left(n_{wt} - \tau_3 \phi_{wt} \sum_{s \in T \setminus t} \phi_{ws} \right)$$

Получился разреживающий регуляризатор. Напрямую он работает только с матрицей Φ .



ОСНОВНЫЕ ВЫВОДЫ

- Задача тематического моделирования имеет бесконечно много решений;
- Регуляция позволяет наложить дополнительные требования на модель;
- Подход аддитивной регуляризации ARTM позволяет обучать модели с любым набором дополнительных требований;
- Простейшие регуляризаторы позволяют сглаживать или разреживать модель;
- Регуляризатор декорреляции даёт возможность обучать модели с различными темами.

