

Методы машинного обучения. Байесовские латентные модели и тематическое моделирование

Воронцов Константин Вячеславович

www.MachineLearning.ru/wiki?title=User:Vokov

вопросы к лектору: voron@forecsys.ru

материалы курса:

github.com/MSU-ML-COURSE/ML-COURSE-22-23

орг.вопросы по курсу: ml.cmc@mail.ru

- 1 **Байесовские модели с латентными переменными**
 - Байесовская регуляризация
 - Байесовская теория EM-алгоритма
 - Задача разделения смеси распределений
- 2 **Вероятностное тематическое моделирование**
 - Постановка задачи BTM
 - Свойства тематических моделей
 - Регуляризованный EM-алгоритм для BTM
- 3 **Регуляризация тематических моделей**
 - Модели PLSA и LDA
 - Небайесовское обобщение LDA
 - Байесовская и классическая регуляризация

Напоминание. Вероятностные модели порождения данных

Дано:

$X = (x_1, \dots, x_\ell)$ — исходные данные, *наблюдаемые переменные*

Найти:

$p(X|\Omega)$ — модель порождения данных, с параметром Ω

$p(X|\Omega) = \prod_{i=1}^{\ell} p(x_i|\Omega)$ — в случае простой (i.i.d.) выборки

Критерии максимизации:

— *правдоподобия* (Maximum Likelihood Estimate, MLE):

$$\ln p(X|\Omega) = \sum_{i=1}^{\ell} \ln p(x_i|\Omega) \rightarrow \max_{\Omega}$$

— *апостериорной вероятности* (Maximum a Posteriori, MAP):

$$\ln p(X, \Omega) = \ln p(X|\Omega)p(\Omega|\gamma) = \sum_{i=1}^{\ell} \ln p(x_i|\Omega) + \ln p(\Omega|\gamma) \rightarrow \max_{\Omega}$$

где $R(\Omega) = \ln p(\Omega|\gamma)$ играет роль регуляризатора.

Порождающая модель с латентными переменными

Дано:

$X = (x_1, \dots, x_\ell)$ — исходные данные, *наблюдаемые переменные*

$Z = (z_1, \dots, z_m)$ — *латентные (скрытые) переменные*

Найти:

$p(X, Z|\Omega)$ — модель совместного порождения наблюдаемых данных и скрытых переменных, с параметром Ω

Критерий:

максимум правдоподобия (Maximum Likelihood Estimate, MLE):

$$\ln p(X|\Omega) = \ln \int_Z p(X, Z|\Omega) dZ \rightarrow \max_{\Omega}$$

Для дискретных переменных Z вместо интеграла \int_Z сумма \sum_Z

Договоримся далее dZ опускать

Bishop C. M. Pattern recognition and machine learning. Springer, 2006.

Пример. Задача разделения смеси распределений

Порождающая модель смеси k вероятностных распределений:

$$p(x) = \sum_{j=1}^k p(x, j) = \sum_{j=1}^k p(j)p(x|j) = \sum_{j=1}^k w_j \varphi(x, \theta_j)$$

$X = (x_1, \dots, x_\ell)$ — исходные данные, *наблюдаемые переменные*

$Z = (j_1, \dots, j_\ell)$ — *компоненты смеси j_i , порождающие объекты x_i*

$\Omega = (w_j, \theta_j)_{j=1}^k$ — параметры смеси, $w_j \geq 0$, $\sum_{j=1}^k w_j = 1$

Порождающая модель наблюдаемых и скрытых i.i.d. данных:

$$p(X, Z|\Omega) = \prod_{i=1}^{\ell} p(x_i, j_i) = \prod_{i=1}^{\ell} p(j_i)p(x_i|\theta_{j_i}) = \prod_{i=1}^{\ell} w_{j_i} \varphi(x_i, \theta_{j_i})$$

Задача разделения смеси — максимизация \log правдоподобия:

$$\ln p(X|\Omega) = \ln \prod_{i=1}^{\ell} p(x_i) = \sum_{i=1}^{\ell} \ln \sum_{j=1}^k w_j \varphi(x_i, \theta_j) \rightarrow \max_{\Omega}$$

Байесовская порождающая модель с латентными переменными

$X = (x_1, \dots, x_\ell)$ — исходные данные, наблюдаемые переменные

$Z = (z_1, \dots, z_m)$ — латентные (скрытые) переменные

$p(X, Z|\Omega)$ — модель наблюдаемых и скрытых переменных

$p(\Omega|\gamma)$ — априорное распределение с гиперпараметрами γ

Задача: по X найти Ω .

Апостериорное распределение, по формуле Байеса:

$$p(\Omega|X, \gamma) \propto p(X|\Omega) p(\Omega|\gamma) = \int_Z p(X, Z|\Omega) p(\Omega|\gamma)$$

Принцип максимума апостериорной вероятности:

$$\ln p(\Omega|X, \gamma) = \ln \int_Z p(X, Z|\Omega) + \underbrace{\ln p(\Omega|\gamma)}_{R(\Omega)} \rightarrow \max_{\Omega}$$

$R(\Omega) = \ln p(\Omega|\gamma)$ — байесовский регуляризатор, хотя

$R(\Omega)$ может и не иметь вероятностной интерпретации.

Общий ЕМ-алгоритм для задачи со скрытыми переменными

Теорема. Точка Ω локального максимума регуляризованного маргинализованного правдоподобия (Marginal log-Likelihood)

$$\ln \int_Z p(X, Z|\Omega) + R(\Omega) \rightarrow \max_{\Omega} \quad (\text{RML})$$

удовлетворяет системе уравнений, решение которой методом простых итераций сводится к чередованию двух шагов:

$$\text{Е-шаг: } q(Z) = p(Z|X, \Omega);$$

$$\text{М-шаг: } \int_Z q(Z) \ln p(X, Z|\Omega) + R(\Omega) \rightarrow \max_{\Omega}.$$

Общий ЕМ-алгоритм используется не только для разделения смесей, но и в анализе сигналов, изображений, текстов и др.

A.P.Dempster, N.M.Laird, D.B.Rubin. Maximum likelihood from incomplete data via the EM algorithm. 1977.

Доказательство теоремы

Необходимые условия локального экстремума:

$$\frac{\partial}{\partial \Omega} \left(\ln \int_Z p(X, Z|\Omega) + R(\Omega) \right) = \frac{1}{p(X|\Omega)} \int_Z \frac{\partial p(X, Z|\Omega)}{\partial \Omega} + \frac{\partial R(\Omega)}{\partial \Omega} = 0$$

По формуле условной вероятности $p(X|\Omega) = \frac{p(X, Z|\Omega)}{p(Z|X, \Omega)}$, подставляем:

$$\begin{aligned} \int_Z \frac{p(Z|X, \Omega)}{p(X, Z|\Omega)} \frac{\partial p(X, Z|\Omega)}{\partial \Omega} + \frac{\partial R(\Omega)}{\partial \Omega} &= 0 \\ \int_Z \underbrace{p(Z|X, \Omega)}_{q(Z)} \frac{\partial}{\partial \Omega} \ln p(X, Z|\Omega) + \frac{\partial R(\Omega)}{\partial \Omega} &= 0 \end{aligned}$$

Это уравнение совпадает с необходимым условием локального экстремума для задачи M-шага, при этом $q(Z)$ рассматривается как переменная, не зависящая от Ω .



Ещё более общий ЕМ-алгоритм и его сходимость

Теорема. Значение маргинализованного правдоподобия

$$\ln \int_Z p(X, Z|\Omega) + R(\Omega) \rightarrow \max_{\Omega} \quad (\text{RML})$$

не убывает на каждом шаге итерационного процесса

$$\text{Е-шаг: } \text{KL}(q(Z) \parallel p(Z|X, \Omega)) \rightarrow \min_q;$$

$$\text{М-шаг: } \int_Z q(Z) \ln p(X, Z|\Omega) + R(\Omega) \rightarrow \max_{\Omega}.$$

$q(Z) = p(Z|X, \Omega)$ является точным решением задачи Е-шага.

Минимизация KL-дивергенции на Е-шаге используется в случаях, когда не удаётся вычислить $p(Z|X, \Omega)$ в явном виде.

Сходимость *в слабом смысле*: глобальный max не гарантируется.

Доказательство теоремы

По формуле условной вероятности $p(X|\Omega) = \frac{p(X, Z|\Omega)}{p(Z|X, \Omega)}$.

Для произвольного распределения $q(Z)$

$$\begin{aligned} \ln p(X|\Omega) &= \int_Z q(Z) \ln p(X|\Omega) = \int_Z q(Z) \ln \frac{p(X, Z|\Omega)}{p(Z|X, \Omega)} = \\ &= \underbrace{\int_Z q(Z) \ln \frac{p(X, Z|\Omega)}{q(Z)}}_{L(q, \Omega)} + \underbrace{\int_Z q(Z) \ln \frac{q(Z)}{p(Z|X, \Omega)}}_{\text{KL}(q(Z) \| p(Z|X, \Omega)) \geq 0} \end{aligned}$$

Максимизируем достижимую нижнюю оценку RML то по q , то по Ω :

$$\text{Е-шаг: } L(q, \Omega) + \cancel{R(\Omega)} \rightarrow \max_q \Leftrightarrow \text{KL}(q(Z) \| p(Z|X, \Omega)) \rightarrow \min_q$$

$$\text{М-шаг: } L(q, \Omega) + R(\Omega) \rightarrow \max_{\Omega} \Leftrightarrow \int_Z q(Z) \ln p(X, Z|\Omega) + R(\Omega) \rightarrow \max_{\Omega}$$

На каждом шаге значение функционала может только возрастать, откуда и следует сходимость в слабом смысле. ■

ЕМ-алгоритм для разделения смеси распределений

$X = (x_1, \dots, x_\ell)$ — исходные данные, *наблюдаемые переменные*

$Z = (j_1, \dots, j_\ell)$ — *компоненты смеси j_i , порождающие объекты x_i*

$$p(X, Z|\Omega) = \prod_{i=1}^{\ell} p(x_i, j_i|\Omega) = \prod_{i=1}^{\ell} p(j_i)p(x_i|\theta_{j_i}) = \prod_{i=1}^{\ell} w_{j_i} \varphi(x_i, \theta_{j_i})$$

Е-шаг: в силу независимости элементов выборки и формулы Байеса

$$q(Z) = p(Z|X, \Omega) = \prod_{i=1}^{\ell} p(j_i|x_i, \Omega); \quad p(j|x, \Omega) = \frac{p(x, j|\Omega)}{p(x|\Omega)} = \frac{w_j \varphi(x, \theta_j)}{\sum_t w_t \varphi(x, \theta_t)}$$

М-шаг: подставим $q(Z)$ и $p(X, Z|\Omega)$ в общую формулу М-шага:

$$\begin{aligned} & \sum_Z q(Z) \ln p(X, Z|\Omega) + R(\Omega) = \\ & = \sum_{j_1=1}^k \cdots \sum_{j_\ell=1}^k \prod_{t=1}^{\ell} p(j_t|x_t, \Omega) \sum_{i=1}^{\ell} \ln p(x_i, j_i|\Omega) + R(\Omega) = \\ & = \sum_{i=1}^{\ell} \sum_{j=1}^k \underbrace{p(j|x_i, \Omega)}_{g_{ij}} \underbrace{\ln p(x_i, j|\Omega)}_{w_j \varphi(x_i, \theta_j)} + R(\Omega) \rightarrow \max_{\Omega}, \quad \Omega = (w_j, \theta_j)_{j=1}^k \end{aligned}$$

ЕМ-алгоритм для разделения смеси распределений

М-шаг распадается на $2k$ подзадач по $w_j, \theta_j, j = 1, \dots, k$:

$$\sum_{j=1}^k \sum_{i=1}^{\ell} \left(g_{ij} \ln w_j + g_{ij} \ln \varphi(x_i, \theta_j) \right) + R(\Omega) \rightarrow \max_{\Omega}$$

Необходимые условия локального экстремума:

$$\begin{cases} \frac{\partial}{\partial \theta_j} \left(\sum_{i=1}^{\ell} g_{ij} \ln \varphi(x_i, \theta_j) + R(\Omega) \right) = 0; \\ \frac{\partial}{\partial w_j} \left(\sum_{i=1}^{\ell} g_{ij} \ln w_j + R(\Omega) + \left(1 - \sum_{j=1}^k w_j \right) \lambda_j - w_j \mu_j \right) = 0; \\ w_j \geq 0, \quad \sum_{j=1}^k w_j = 1, \quad w_j \mu_j = 0, \quad \mu_j \geq 0 \end{cases}$$

Относительно w_j решение аналитическое, из условий ККТ:

$$w_j = \text{norm}_j \left(\sum_{i=1}^{\ell} g_{ij} + w_j \frac{\partial R}{\partial w_j} \right), \quad j = 1, \dots, k.$$

ЕМ-алгоритм для разделения смеси распределений

Задача разделения смеси: $(X, Z) = (x_i, j_i)_{i=1}^{\ell}$, $\Omega = (w_j, \theta_j)_{j=1}^k$

Теорема. Точка Ω локального максимума (RML)

$$\ln \sum_Z p(X, Z | \Omega) + R(\Omega) = \sum_{i=1}^{\ell} \ln \sum_{j=1}^k w_j \varphi(x_i, \theta_j) + R(\Omega)$$

удовлетворяет системе уравнений, решение которой методом простых итераций сводится к чередованию двух шагов:

Е-шаг: $g_{ij} \equiv p(j | x_i) = \text{norm}_j(w_j \varphi(x_i, \theta_j)), \quad i = 1..{\ell}, \quad j = 1..k;$

М-шаг: $\theta_j = \arg \max_{\theta} \left(\sum_{i=1}^{\ell} g_{ij} \ln \varphi(x_i, \theta) + R(\Omega) \right), \quad j = 1..k;$

$$w_j = \text{norm}_j \left(\sum_{i=1}^{\ell} g_{ij} + w_j \frac{\partial R}{\partial w_j} \right), \quad j = 1..k.$$

Регуляризация в ЕМ-алгоритме для разделения смеси

- При $R(\Omega) = 0$ это уже знакомый нам ЕМ-алгоритм для разделения смеси вероятностных распределений
- При разделении смеси n -мерных гауссиан $\mathcal{N}(x; \mu_j, \Sigma_j)$ регуляризация ковариационных матриц: $\Sigma_j + \tau I_n$
- $R(w) = \tau \sum_j p_j \ln w_j$ — регуляризатор сглаживания, приближает веса w_j к априорному распределению p_j :

$$w_j = \text{norm}_j \left(\sum_{i=1}^{\ell} g_{ij} + \tau p_j \right)$$

- $R(w) = -\tau \sum_j \ln w_j$ — регуляризатор разреживания, удаляет компоненты с весами $w_j < \tau$:

$$w_j = \text{norm}_j \left(\sum_{i=1}^{\ell} g_{ij} - \tau \right)$$

Преимущества общего ЕМ-алгоритма

- Широкий класс задач связан с выявлением латентных структур, порождающих наблюдаемые данные:
 - разделение смеси вероятностных распределений;
 - «мягкая» или «жёсткая» кластеризация;
 - вероятностное тематическое моделирование;
 - сегментация временных рядов, сигналов, изображений;
 - восстановление пропущенных данных;
- Готовый рецепт для вывода вычислительных формул Е и М шагов по известной порождающей модели $p(X, Z|\Omega)$
- Можно добавлять какие угодно регуляризаторы $R(\Omega)$, причём не обязательно вероятностные вида $\ln p(\Omega|\gamma)$
- Гарантируется сходимость в слабом смысле
- Это обучение без учителя, не требующее разметки

Постановка задачи тематического моделирования

Дано: коллекция текстовых документов

- W — конечное множество (словарь) термов (слов)
- D — конечное множество документов
- T — конечное множество тем
- $X = (d_i, w_i)_{i=1}^n$ — данные, *наблюдаемые переменные*
- $Z = (t_i)_{i=1}^n$ — *латентные (скрытые) переменные*, $t_i \in T$
- n_{dw} — частота терма $w \in W$ в документе $d \in D$
- n_d — длина документа $d \in D$
- n — суммарная длина всех документов коллекции

Найти: *вероятностную тематическую модель (ВТМ)*

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

где $\phi_{wt} = p(w|t)$, $\theta_{td} = p(t|d)$ — параметры модели

Принцип максимума правдоподобия

Правдоподобие — плотность распределения выборки $(d_i, w_i)_{i=1}^n$:

$$\prod_{i=1}^n p(d_i, w_i) = \prod_{d \in D} \prod_{w \in d} p(d, w)^{n_{dw}}$$

Максимизация логарифма правдоподобия

$$\sum_{d \in D} \sum_{w \in d} n_{dw} \ln p(w|d) \xrightarrow{\text{const}} \max_{\Phi, \Theta}$$

эквивалентна максимизации функционала

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

при ограничениях неотрицательности и нормировки

$$\phi_{wt} \geq 0; \quad \sum_{w \in W} \phi_{wt} = 1; \quad \theta_{td} \geq 0; \quad \sum_{t \in T} \theta_{td} = 1$$

Несколько интерпретаций постановки задачи ВТМ

1. **Вероятностная языковая модель** $p(w|d)$
2. **Мягкая би-кластеризация** по кластерам-темам $t \in T$
как документов: $p(t|d)$, так и термов: $p(t|w) = p(w|t) \frac{p(t)}{p(w)}$
3. **Векторные представления (эмбединги)** — вероятностные, интерпретируемые, разреженные: $p(t|d), p(t|w), p(t|d, w), \dots$
4. **Автокодировщик** документов в тематические эмбединги:
кодировщик: $f_{\Phi}: (\hat{p}(w|d) = \frac{n_{dw}}{n_d}) \rightarrow (p(t|d) = \theta_d)$
декодировщик: $g_{\Phi}: \theta_d \rightarrow \Phi \theta_d$
задача реконструкции: $\sum_d n_d \text{KL}(\hat{p}(w|d) \parallel \langle \phi_w, \theta_d \rangle) \rightarrow \min_{\Phi, \Theta}$
5. **Матричное разложение** — низкоранговое,
неотрицательное (стохастическое), приближённое: $(\frac{n_{dw}}{n_d}) \approx \Phi \Theta$

Свойство интерпретируемости тематических моделей

Тематические векторные представления (эмбединги) текста:

- $p(t|d) = \theta_{td}$ для документа d
- $p(t|w) = \phi_{wt} \frac{p(t)}{p(w)}$ для термина w
- $p(t|d, w)$ для локального контекста (d, w)
- $p(t|x)$ для нетекстового объекта x

Интерпретируемость тематических векторов:

- каждая тема t описывается *семантическим ядром* — частотным словарём слов $\{w: p(w|t) > \gamma p(w)\}$, встречающихся в данной теме в γ раз чаще обычного
- любой объект x с вектором $p(t|x)$ описывается частотным словарём слов $\left\{w: p(w|x) = \sum_{t \in T} p(w|t)p(t|x) > \gamma p(w)\right\}$

Цели и не-цели тематического моделирования

Цели:

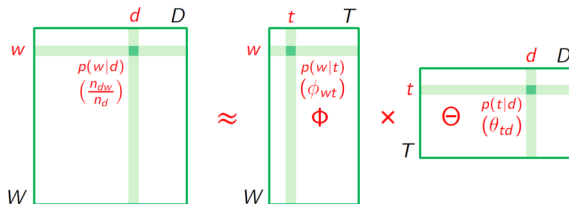
- Выяснить тематическую кластерную структуру текстовой коллекции, сколько в ней тем и какие они
- Получать интерпретируемые тематические векторные представления (эмбединги) документов, фрагментов, слов $p(t|d)$, $p(t|w)$, $p(t|d, w)$ и нетекстовых объектов $p(t|x)$
- Решать задачи поиска, категоризации, сегментации, суммаризации с помощью тематических эмбедингов

Не-цели:

- Угадывать следующие слова (ТМ — слабые модели языка)
- Генерировать связный текст
- Понимать смысл текста

Некорректно поставленная задача матричного разложения

Низкоранговое стохастическое матричное разложение:



Если Φ, Θ — решение, то стохастические Φ', Θ' — тоже решения

- $\Phi' \Theta' = (\Phi S)(S^{-1} \Theta)$, $\text{rank } S = |T|$
- $L(\Phi', \Theta') = L(\Phi, \Theta)$ — линейно не зависимые решения
- $L(\Phi', \Theta') \geq L(\Phi, \Theta) - \varepsilon$ — приближённые решения

Регуляризация необходима для доопределения решения
 Аддитивная регуляризация — сумма регуляризаторов

BTM как порождающая модель с латентными переменными

$X = (d_i, w_i)_{i=1}^n$ — исходные данные, *наблюдаемые переменные*

$Z = (t_i)_{i=1}^n$ — латентные (скрытые) переменные

$\Omega = (\Phi, \Theta)$ — параметры порождающей модели $p(X|\Omega)$

$R(\Omega)$ — регуляризатор; не обязательно байесовский $\ln p(\Omega|\gamma)$

Задача максимизации правдоподобия — по X найти Ω :

$$\ln p(X|\Omega) + R(\Omega) = \sum_{d,w} n_{dw} \ln \sum_{t \in T} \phi_{wt} \theta_{td} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

Порождающая модель наблюдаемых и скрытых i.i.d. данных:

$$\begin{aligned} p(X, Z|\Omega) &= \prod_{i=1}^n p(d_i, w_i, t_i) = \\ &= \prod_{i=1}^n p(w_i|t_i)p(t_i|d_i)p(d_i) = \prod_{i=1}^n \phi_{w_i t_i} \theta_{t_i d_i} p(d_i) \end{aligned}$$

Регуляризованный ЕМ-алгоритм для тематической модели

Наблюдаемые $X = (d_i, w_i)_{i=1}^n$, латентные $Z = (t_i)_{i=1}^n$

Теорема. Точка $\Omega = (\Phi, \Theta)$ локального максимума RML (регуляризованного маргинализованного log-правдоподобия)

$$\ln \sum_Z p(X, Z | \Omega) + R(\Omega) = \sum_{d, w} n_{dw} \ln \sum_t \phi_{wt} \theta_{td} + R(\Phi, \Theta)$$

удовлетворяет системе уравнений, решение которой методом простых итераций сводится к чередованию двух шагов:

Е-шаг: $p(t | d, w) = \text{norm}_{t \in T}(\phi_{wt} \theta_{td}), \quad \forall (d \in D, w \in d, t \in T)$

М-шаг: $\sum_{d, w, t} n_{dw} p(t | d, w) \ln(\phi_{wt} \theta_{td}) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$

Доказательство леммы

Е-шаг: в силу независимости элементов выборки и формулы Байеса

$$q(Z) = p(Z|X, \Omega) = \prod_{i=1}^n p(t_i|d_i, w_i) = \prod_{i=1}^n \text{norm}_{t_i \in T}(\phi_{w_i t_i} \theta_{t_i d_i})$$

М-шаг: подставим $q(Z)$ и $p(X, Z|\Omega)$ в общую формулу М-шага:

$$\begin{aligned} & \sum_{Z \in T^n} q(Z) \ln p(X, Z|\Omega) + R(\Omega) = \\ & \sum_{t_1 \in T} \cdots \sum_{t_n \in T} \prod_{k=1}^n p(t_k|d_k, w_k) \sum_{i=1}^n \ln p(d_i, w_i, t_i|\Omega) + R(\Omega) = \\ & \sum_{i=1}^n \sum_{t_1 \in T} \cdots \sum_{t_n \in T} \prod_{k=1}^n p(t_k|d_k, w_k) \ln p(d_i, w_i, t_i|\Omega) + R(\Omega) = \\ & \sum_{i=1}^n \sum_{t \in T} p(t|d_i, w_i) \ln p(d_i, w_i, t|\Omega) + R(\Omega) = \\ & \sum_{d \in D} \sum_{w \in W} \sum_{t \in T} n_{dw} \underbrace{p(t|d, w)}_{p_{tdw}} \ln(\phi_{wt} \theta_{td}) + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta} \end{aligned}$$

ARTM: аддитивная регуляризация тематических моделей

Задача М-шага декомпозируется на независимые подзадачи:

$$\sum_{w,t} \ln \phi_{wt} \sum_d n_{dw} p_{tdw} + \sum_{d,t} \ln \theta_{td} \sum_w n_{dw} p_{tdw} + R(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

ЕМ-алгоритм: метод простой итерации для системы уравнений

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

Воронцов К. В. Аддитивная регуляризация тематических моделей коллекций текстовых документов. Доклады РАН, 2014.

Два наиболее известных частных случая: модели PLSA и LDA

PLSA: probabilistic latent semantic analysis [Hofmann, 1999]
(вероятностный латентный семантический анализ):

$$R(\Phi, \Theta) = 0.$$

М-шаг — частотные оценки условных вероятностей:

$$\phi_{wt} = \text{norm}_w(n_{wt}), \quad \theta_{td} = \text{norm}_t(n_{td}).$$

LDA: latent Dirichlet allocation (латентное размещение Дирихле):

$$R(\Phi, \Theta) = \sum_{t,w} (\beta_w - 1) \ln \phi_{wt} + \sum_{d,t} (\alpha_t - 1) \ln \theta_{td}.$$

М-шаг — частотные оценки с поправками $\beta_w > 0$, $\alpha_t > 0$:

$$\phi_{wt} = \text{norm}_w(n_{wt} + \beta_w - 1), \quad \theta_{td} = \text{norm}_t(n_{td} + \alpha_t - 1).$$

Hofmann T. Probabilistic latent semantic indexing. SIGIR 1999.

Blei D., Ng A., Jordan M. Latent Dirichlet allocation. 2003.

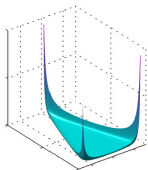
Распределение Дирихле

Гипотеза. Вектор-столбцы $\phi_t = (\phi_{wt})$ и $\theta_d = (\theta_{td})$ порождаются распределениями Дирихле, $\alpha \in \mathbb{R}^{|T|}$, $\beta \in \mathbb{R}^{|W|}$:

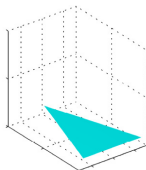
$$\text{Dir}(\phi_t | \beta) = \frac{\Gamma(\beta_0)}{\prod_w \Gamma(\beta_w)} \prod_w \phi_{wt}^{\beta_w - 1}, \quad \phi_{wt} > 0; \quad \beta_0 = \sum_w \beta_w, \quad \beta_t > 0;$$

$$\text{Dir}(\theta_d | \alpha) = \frac{\Gamma(\alpha_0)}{\prod_t \Gamma(\alpha_t)} \prod_t \theta_{td}^{\alpha_t - 1}, \quad \theta_{td} > 0; \quad \alpha_0 = \sum_t \alpha_t, \quad \alpha_t > 0;$$

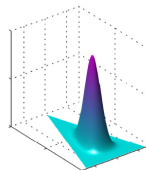
Пример. Распределение $\text{Dir}(\theta | \alpha)$ при $|T| = 3$, $\theta, \alpha \in \mathbb{R}^3$



$$\alpha_1 = \alpha_2 = \alpha_3 = 0.1$$

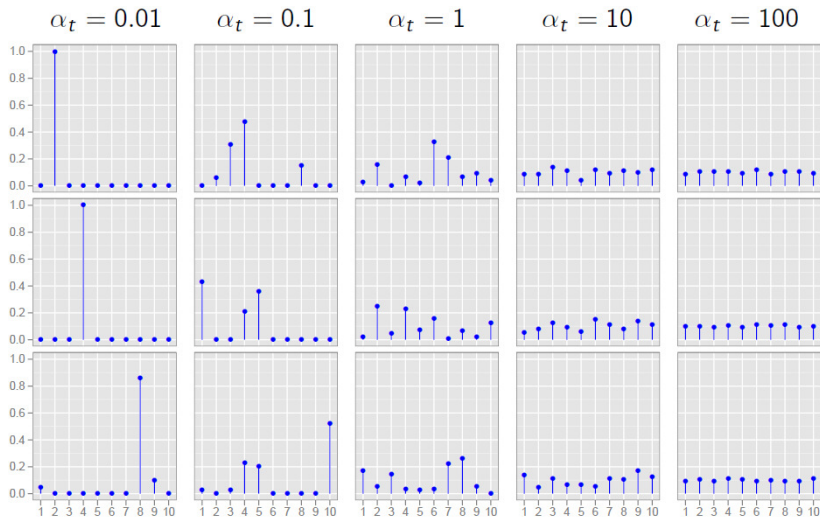


$$\alpha_1 = \alpha_2 = \alpha_3 = 1$$



$$\alpha_1 = \alpha_2 = \alpha_3 = 10$$

Пример. Выборки из трёх 10-мерных векторов $\theta \sim \text{Dir}(\theta|\alpha)$



Максимизация апостериорной вероятности для модели LDA

Совместное правдоподобие данных и модели:

$$\ln \prod_{d \in D} \prod_{w \in d} p(w, d | \Phi, \Theta)^{n_{dw}} \prod_{t \in T} \text{Dir}(\phi_t | \beta) \prod_{d \in D} \text{Dir}(\theta_d | \alpha) \rightarrow \max_{\Phi, \Theta}$$

Регуляризатор — логарифм априорного распределения:

$$R(\Phi, \Theta) = \sum_{t, w} (\beta_w - 1) \ln \phi_{wt} + \sum_{d, t} (\alpha_t - 1) \ln \theta_{td}$$

М-шаг — сглаженные или разреженные частотные оценки:

$$\phi_{wt} = \text{norm}_w(n_{wt} + \beta_w - 1), \quad \theta_{td} = \text{norm}_t(n_{td} + \alpha_t - 1).$$

при $\beta_w > 1$, $\alpha_t > 1$ — сглаживание,

при $0 < \beta_w < 1$, $0 < \alpha_t < 1$ — слабое разреживание,

при $\beta_w = 1$, $\alpha_t = 1$ априорное распределение равномерно, PLSA.

Обобщение LDA: регуляризатор сглаживания и разреживания

Общий вид регуляризатора сглаживания и разреживания:

$$R(\Phi, \Theta) = \beta_0 \sum_{t \in T} \sum_{w \in W} \beta_{wt} \ln \phi_{wt} + \alpha_0 \sum_{d \in D} \sum_{t \in T} \alpha_{td} \ln \theta_{td} \rightarrow \max,$$

где $\beta_0 > 0$, $\alpha_0 > 0$ — коэффициенты регуляризации,

β_{wt} , α_{td} — параметры, задаваемые пользователем:

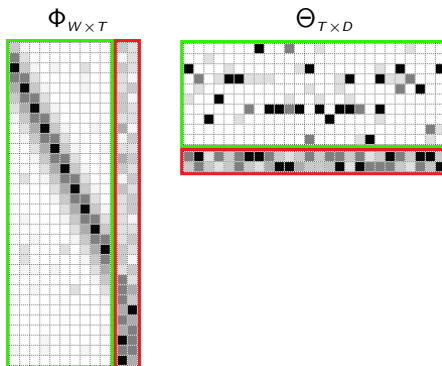
- $\beta_{wt} > 0$, $\alpha_{td} > 0$ — сглаживание
- $\beta_{wt} < 0$, $\alpha_{td} < 0$ — разреживание

Возможные применения сглаживания и разреживания:

- задать фоновые темы с общей лексикой языка
- задать шумовую тему для нетематичных термов
- задать псевдо-документ с ключевыми термами темы
- скорректировать состав термов и документов темы

Разделение тем на предметные и фоновые

Предметные темы S содержат термины предметной области,
 $p(w|t)$, $p(t|d)$, $t \in S$ — разреженные, существенно различные
Фоновые темы B содержат слова общей лексики,
 $p(w|t)$, $p(t|d)$, $t \in B$ — существенно отличные от нуля



Регуляризатор декоррелирования тем

Цель: усилить различность тем; выделить в каждой теме лексическое ядро, отличающее её от других тем; вывести слова общей лексики из предметных тем в фоновые.

Минимизируем ковариации между вектор-столбцами ϕ_t :

$$R(\Phi) = -\frac{\tau}{2} \sum_{t \in T} \sum_{s \in T \setminus t} \sum_{w \in W} \phi_{wt} \phi_{ws} \rightarrow \max.$$

Подставляем в формулы М-шага, получаем ещё один вариант разреживания — контрастирование строк матрицы Φ (малые вероятности ϕ_{wt} в строке становятся ещё меньше):

$$\phi_{wt} = \text{norm}_w \left(n_{wt} - \tau \phi_{wt} \sum_{s \in T \setminus t} \phi_{ws} \right).$$

Байесовская и классическая регуляризация

Байесовский вывод апостериорного распределения $p(\Omega|X)$ (громоздкий, приближённый) ради точечной оценки Ω :

$$\text{Posterior}(\Omega|X, \gamma) \propto p(X|\Omega) \text{Prior}(\Omega|\gamma)$$
$$\Omega := \arg \max_{\Omega} \text{Posterior}(\Omega|X, \gamma)$$

Максимизация апостериорной вероятности (MAP) даёт точечную оценку Ω напрямую, без вывода Posterior:

$$\Omega := \arg \max_{\Omega} (\ln p(X|\Omega) + \ln \text{Prior}(\Omega|\gamma))$$

Многокритериальная аддитивная регуляризация (ARTM) обобщает MAP на любые регуляризаторы и их комбинации:

$$\Omega := \arg \max_{\Omega} (\ln p(X|\Omega) + \sum_{i=1} \tau_i R_i(\Omega))$$

- Тематическое моделирование — «мягкая кластеризация», автокодировщик или стохастическое матричное разложение
- Стандартные методы — PLSA и LDA
- Нестандартные — огромное разнообразие регуляризаторов
- Аддитивная регуляризация позволяет комбинировать модели
- Обычно в ТМ используется байесовское обучение.

Почему оно не нужно в ТМ: на практике используются не апостериорные распределения $p(\Omega|X, \gamma)$, а их точечные оценки Ω по максимуму правдоподобия

- В ARTM те же модели выводятся намного проще — с помощью Леммы о максимизации на симплексах,
- она применима для оптимизации любых моделей, параметры которых — неотрицательные нормированные векторы (дискретные вероятностные распределения)

Воронцов К.В. Вероятностное тематическое моделирование: теория, модели, алгоритмы и проект BigARTM. 2017–2023.

<http://www.machinelearning.ru/wiki/images/d/d5/Voron17survey-artm.pdf>