

Мультипликативный метод неотрицательного матричного разложения с АБ-дивергенцией и его сходимостью*

Е. А. Рябенко

riabenko.e@gmail.com

Москва, Вычислительный центр им. А. А. Дородницына РАН

Мультипликативный метод неотрицательного матричного разложения для случая, когда точность приближения модели измеряется с помощью АБ-дивергенции, вблизи границы неотрицательной области может сходиться к нестационарной точке. Предлагается модифицированный мультипликативный метод, в котором за счет отделения элементов матриц от нуля константой ε удается показать не только монотонность невозрастания функции потерь, но и тот факт, что любая предельная точка этого метода является стационарной точкой отделенной от нуля задачи. Разреживание получаемых таким методом матриц дает решение, являющееся стационарной точкой исходной задачи с точностью до $\mathcal{O}(\varepsilon)$. Для частного случая, соответствующего норме Фробениуса, показано, что метод всегда сходится.

Ключевые слова: неотрицательное матричное разложение, мультипликативные обновления, АБ-дивергенция, сходимость.

Multiplicative Method for Nonnegative Matrix Factorization with AB-Divergence and its Convergence*

E. A. Riabenko

Moscow, Dorodnicyn Computing Centre of RAS

Multiplicative method for nonnegative matrix factorization with AB-divergence could converge to nonstationary point when the elements of matrices approach zero. A modified multiplicative method that bounds the matrices from zero by constant ε is proposed, and there proved not only monotonic descent, but the fact that its every limiting point is a stationary point of the modified bounded problem. Setting particular elements of the resulting matrix to zero yields the solution that is $\mathcal{O}(\varepsilon)$ -close to the stationary point of the original problem. For a special case of Frobenius norm, it is proved that the method always converges.

Keywords: nonnegative matrix factorization, multiplicative updates, AB-divergence, convergence.

Введение

Матричные разложения используются для решения задач сжатия данных, восстановления сигналов, заполнения пропусков, для выявления структурных особенностей коллекций данных. Приложения, связанные с получением и анализом матричных разложений, различаются ограничениями, накладываемыми на факторы. В задаче неотрицательного матричного разложения, рассматриваемой в данной работе, ключевую роль играют ограничения на знак элементов факторных матриц. Впервые она была рассмотрена в ра-

*Работа выполнена при финансовой поддержке РФФИ (проект № 11-07-00480) и программы ОМН РАН «Алгебраические и комбинаторные методы математической кибернетики и информационные системы нового поколения».

боте [1] в приложении к задаче византийских генералов из теории отказоустойчивости, однако основной интерес к этой теме возник после работ [2, 3], авторы которых обобщили постановку задачи и предложили простой алгоритм ее решения. Неотрицательные матричные разложения используются при анализе изображений, текстов, в вычислительной биологии, медицине и других прикладных областях. Обзор применений можно найти в [4].

Задача формулируется следующим образом. Дана матрица $P \in \mathbb{R}_+^{m \times n}$ с неотрицательными элементами (подразумевается, что в ней отсутствуют нулевые строки и столбцы), а также некоторое натуральное число $r < \min(m, n)$. Требуется найти матрицы $A^* \in \mathbb{R}_+^{m \times r}$ и $X^* \in \mathbb{R}_+^{r \times n}$ с неотрицательными элементами, такие, что их произведение $Q^* = A^*X^*$ максимально близко к исходной матрице P в смысле некоторой функции потерь $D(P, Q)$:

$$(A^*, X^*) = \arg \min_{A \geq 0, X \geq 0} D(P, AX). \quad (1)$$

Выбор функции потерь оказывает существенное влияние как на саму задачу, так и на получаемое решение [5].

Наибольший интерес представляют аддитивные функции потерь:

$$D(P, Q) = \sum_{i,j} d(p_{ij}, q_{ij}),$$

где $d(p, q) \geq 0$, причем $d(p, q) = 0$ тогда и только тогда, когда $p = q$.

Чаще всего в качестве функции потерь используется норма Фробениуса:

$$D_F(P, Q) = \sum_{i,j} (p_{ij} - q_{ij})^2.$$

Одна из причин ее популярности является оптимальность получаемых в результате ее минимизации оценок для моделей с аддитивным гауссовским шумом — такие оценки совпадают с оценками максимального правдоподобия. Однако для других видов шума, а также в присутствии выбросов, оценки, доставляющие минимум норме Фробениуса, могут оказываться несостоятельными. В случае аддитивного пуассоновского шума задача максимизации правдоподобия эквивалентна задаче минимизации дивергенции Кульбака–Лейблера между данными и моделью [6]:

$$D_{KL}(P, Q) = \sum_{i,j} \left(p_{ij} \ln \frac{p_{ij}}{q_{ij}} - p_{ij} + q_{ij} \right).$$

Кроме того, если столбцы матриц P , A и X нормированы и рассматриваются как вероятностные распределения, минимизация дивергенции Кульбака–Лейблера также соответствует максимизации правдоподобия модели. Первые алгоритмы неотрицательного матричного разложения, предложенные в работах [2, 3], были построены именно для этих мер качества.

В данной работе мы рассмотрим более широкий класс функций потерь — АБ-дивергенции, параметрическое семейство, содержащее как приведенные функции, так и многие другие [7].

Для решения оптимизационной задачи (1) наиболее популярными являются итерационные методы с мультипликативным шагом, в которых обновления переменных A и X имеют вид умножения на положительное число. Преимущество мультипликативных обновлений заключается в том, что неотрицательность решения сохраняется естественным образом без дополнительных вычислительных затрат.

В данной работе предлагается мультипликативный метод получения неотрицательного матричного разложения с АБ-дивергенцией в качестве функции потерь и доказывается несколько фактов о его сходимости, в частности, монотонность невозрастания ошибки и близость получаемого решения к стационарной точке задачи.

АБ-дивергенция

Для измерения качества моделей в процессе настройки ее параметров существуют разнообразные функции потерь, как возникшие из контекста прикладных задач, так и мотивированные соображениями теории информации, выпуклого анализа или информационной геометрии. Представление различных функций потерь в виде параметрических семейств позволяет унифицировать работу с ними, увидеть между ними связь и лучше понять, как выбор меры качества влияет на получаемое решение. Среди таких семейств широко используются альфа- [8] и бета-дивергенции [9], задающие непрерывные по параметрам множества сепарабельных функций потерь, включающих в том числе норму Фробениуса, дивергенции Кульбака-Лейблера, Итакура-Саито и другие. Класс АБ-дивергенций, предложенный в работе [7], обобщает эти семейства и задается в виде двухпараметрического семейства функций следующего вида:

$$D_{AB}^{(\alpha, \beta)}(P, Q) = \sum_{i,j} d_{AB}^{(\alpha, \beta)}(p_{ij}, q_{ij}),$$

$$d_{AB}^{(\alpha, \beta)}(p, q) = \begin{cases} \frac{1}{\alpha\beta} \left(\frac{\alpha}{\alpha+\beta} p^{\alpha+\beta} + \frac{\beta}{\alpha+\beta} q^{\alpha+\beta} - p^\alpha q^\beta \right), & \text{если } \alpha, \beta, \alpha + \beta \neq 0; \\ \frac{1}{\alpha^2} \left(p^\alpha \ln \frac{p^\alpha}{q^\alpha} - p^\alpha + q^\alpha \right), & \text{если } \alpha \neq 0, \beta = 0; \\ \frac{1}{\alpha^2} \left(\ln \frac{q^\alpha}{p^\alpha} + \left(\frac{q^\alpha}{p^\alpha} \right)^{-1} - 1 \right), & \text{если } \alpha = -\beta \neq 0; \\ \frac{1}{\beta^2} \left(q^\beta \ln \frac{q^\beta}{p^\beta} - q^\beta + p^\beta \right), & \text{если } \alpha = 0, \beta \neq 0; \\ \frac{1}{2} (\ln p - \ln q)^2, & \text{если } \alpha = \beta = 0. \end{cases}$$

Частные случаи (2)–(5) здесь являются предельными для случая (1) и могут быть получены по правилу Лопиталю.

АБ-дивергенции — одно из наиболее обширных известных параметрических семейств функций потерь. Оно включает многие широко применяемые меры близости (табл. 1). Оценки, получаемые при минимизации функций этого класса, являются оценками максимального правдоподобия при различных распределениях шума как аддитивного и мультипликативного, так и смешанного, состоящего из обеих этих компонент. Выбирая параметры α и β , мы тем самым определяем такие свойства получаемой модели, как разреженность и устойчивость, а также неявно задаем вид шума, подходящий для исследуемой задачи. Пример такого использования АБ-дивергенций можно найти в работе [10].

Будем решать задачу неотрицательного матричного разложения, используя АБ-дивергенцию в качестве функции потерь:

$$(A^*, X^*) = \arg \min_{A \geq 0, X \geq 0} D_{AB}^{(\alpha, \beta)}(P, AX). \quad (2)$$

Мультипликативные алгоритмы минимизации функции потерь

Поскольку используемые на практике функции потерь не выпуклы по совокупности аргументов A, X , задача (1) решается поочередной минимизацией $D(P, AX)$ по A и X с помощью алгоритмов следующего вида:

Таблица 1: Некоторые функции из семейства АБ-дивергенций

(α, β)	Функция потерь $d_{AB}^{(\alpha, \beta)}(p, q)$	Название
(1, 0)	$d_{KL}(p, q) = p \ln \frac{p}{q} - p + q$	Дивергенция Кульбака-Лейблера
(1, -1)	$d_{IS}(p, q) = \ln \frac{q}{p} + \frac{p}{q} - 1$	Дивергенция Итакура-Саито
(1, 1)	$\frac{1}{2}d_E(p, q) = \frac{1}{2}(p - q)^2$	Евклидово расстояние
(0,5, 0,5)	$2d_H(p, q) = 2(\sqrt{p} - \sqrt{q})^2$	Расстояние Хелингера
(2, -1)	$\frac{1}{2}d_P(p, q) = \frac{1}{2}\frac{(p - q)^2}{q}$	χ^2 Пирсона
(-1, 2)	$\frac{1}{2}d_N(p, q) = \frac{1}{2}d_H(q, p) = \frac{1}{2}\frac{(p - q)^2}{p}$	χ^2 Неймана
(0, 0)	$\frac{1}{2}d_E(\ln p, \ln q) = \frac{1}{2}(\ln(p) - \ln(q))^2$	Лог-евклидово расстояние

Вход: $A^0 \geq 0, X^0 \geq 0$

цикл // $t = 0, 1, 2, \dots$

зафиксировав A^t , найдем такое X^{t+1} , что $D(P, A^t X^{t+1}) \leq D(P, A^t X^t)$;

зафиксировав X^{t+1} , найдем такое A^{t+1} , что $D(P, A^{t+1} X^{t+1}) \leq D(P, A^t X^{t+1})$;

Существуют различные подходы к построению таких обновлений. Можно на каждом шаге находить значение матрицы-множителя, минимизирующее $D(P, AX)$, когда вторая матрица-множитель фиксирована. Алгоритм с такими обновлениями будет реализовывать метод блочно-покоординатного спуска. В случае двух блоков любая предельная точка такого алгоритма является стационарной [11]. Для некоторых функций потерь точные минимумы по каждой из матриц найти легко: например, когда точность приближения измеряется при помощи нормы Фробениуса, покомпонентные минимумы находятся аналитически (это свойство лежит в основе метода HALS [12]). При использовании других функций потерь, в том числе рассматриваемых здесь АБ-дивергенций, покомпонентные минимумы можно найти только численно и приближенно. Это может потребовать значительных вычислительных затрат, притом что предельная точка такого алгоритма может и не являться стационарной. Вместо поиска точного покомпонентного минимума функции $D(P, AX)$ достаточно на каждой итерации делать только один шаг в направлении ее уменьшения. Для этого может использоваться обычный градиентный метод:

Вход: $A^0 \geq 0, X^0 \geq 0$

цикл // $t = 0, 1, 2, \dots$

$$x_{kj}^{t+1} = x_{kj}^t - \nu \frac{\partial D(P, AX)}{\partial x_{kj}}, \quad k = 1, \dots, r, \quad j = 1, \dots, n;$$

$$a_{ik}^{t+1} = a_{ik}^t - \eta \frac{\partial D(P, AX)}{\partial a_{ik}}, \quad i = 1, \dots, m, \quad k = 1, \dots, r;$$

где ν, η — положительные константы, задающие длины шага в направлении градиентного спуска. В общем случае длина шага в направлении спуска может выбираться индивидуаль-

но для каждого элемента матриц, а также меняться в зависимости от номера итерации t :

$$\begin{aligned}x_{kj}^{t+1} &= x_{kj}^t - \nu_{kj}^t \frac{\partial D(P, AX)}{\partial x_{kj}}; \\a_{ik}^{t+1} &= a_{ik}^t - \eta_{ik}^t \frac{\partial D(P, AX)}{\partial a_{ik}}.\end{aligned}$$

Кроме того, градиентный спуск может вестись в трансформированном пространстве параметров:

$$\begin{aligned}x_{kj}^{t+1} &= \varphi^{-1} \left(\varphi(x_{kj}^t) - \nu_{kj}^t \frac{\partial D(P, AX)}{\partial \varphi(x_{kj})} \right); \\a_{ik}^{t+1} &= \varphi^{-1} \left(\varphi(a_{ik}^t) - \eta_{ik}^t \frac{\partial D(P, AX)}{\partial \varphi(a_{ik})} \right),\end{aligned}$$

где $\varphi(z)$ — некоторое биективное отображение. Переход в трансформированное пространство позволяет увеличить скорость сходимости, когда покомпонентные гессианы минимизируемого функционала в нем имеют меньшее число обусловленности (известно, что скорость сходимости градиентных методов ограничивается разностью минимального и максимального собственных значений гессиана [13]).

Поскольку градиентный спуск не гарантирует сохранения неотрицательности компонент, на каждом шаге необходимо проецировать обновленные значения A и X на неотрицательную область. В таком случае для того, чтобы обеспечить монотонное убывание функции потерь, необходимо дополнительно оптимизировать параметры ν_{kj}^t, η_{ik}^t вдоль направления спуска, что может потребовать значительных вычислительных затрат. Вместо этого можно выбрать ν_{kj}^t, η_{ik}^t так, чтобы обновления были мультипликативными, т. е. чтобы x_{kj}, a_{ik} на каждом шаге умножались на некоторое положительное число. В этом случае неотрицательность решения будет сохраняться без дополнительных операций. Наиболее распространенные алгоритмы получения неотрицательного матричного разложения основаны именно на этом подходе.

Простейший способ выбора ν_{kj}^t, η_{ik}^t , приводящий к мультипликативным обновлениям в исходном пространстве параметров, заключается в следующем. Пусть ∇_A — градиент $D(P, AX)$ по A , $\nabla_A^+ = \max(\nabla_A, 0)$, $\nabla_A^- = \max(-\nabla_A, 0)$, где максимумы берутся покомпонентно, т. е. $\nabla_A = \partial D / \partial A = \nabla_A^+ - \nabla_A^-$. Выбирая $\eta_{ik}^t = a_{ik}^t / [\nabla_A^+]_{ik}$, получим мультипликативные обновления вида

$$a_{ik}^{t+1} = a_{ik}^t \frac{[\nabla_A^-]_{ik}}{[\nabla_A^+]_{ik}}.$$

Аналогично для X , взяв $\nu_{kj}^t = x_{kj}^t / [\nabla_X^+]_{kj}$, получаем:

$$x_{kj}^{t+1} = x_{kj}^t \frac{[\nabla_X^-]_{kj}}{[\nabla_X^+]_{kj}}.$$

Алгоритм, решающий задачу (2) путем поочередных мультипликативных шагов по A и X , был предложен в работе [7]. Градиентный спуск в нем ведется в пространстве, трансформированном функцией деформированного логарифма:

$$\varphi(z) = \ln_{1-\alpha}(z) = \begin{cases} \frac{z^\alpha - 1}{\alpha}, & \text{если } \alpha \neq 0; \\ \ln z, & \text{если } \alpha = 0, \end{cases}$$

обратной функцией к которой является деформированная экспонента:

$$\varphi^{-1}(z) = \exp_{1-\alpha}(z) = \begin{cases} \exp(z), & \text{если } \alpha = 0; \\ (1 + \alpha z)^{1/\alpha}, & \text{если } \alpha \neq 0, 1 + \alpha z \geq 0; \\ 0, & \text{если } \alpha \neq 0, 1 + \alpha z < 0. \end{cases}$$

Алгоритм выглядит следующим образом:

$$\begin{aligned} x_{kj} &\leftarrow x_{kj} \left(\exp_{1-\alpha} \left(\frac{\sum_{i=1}^m a_{ik} q_{ij}^{\alpha+\beta-1} \ln_{1-\alpha}(p_{it}/q_{it})}{\sum_{i=1}^m a_{ik} q_{ij}^{\alpha+\beta-1}} \right) \right)^{\omega(\alpha,\beta)} ; \\ a_{ik} &\leftarrow a_{ik} \left(\exp_{1-\alpha} \left(\frac{\sum_{j=1}^n x_{kj} q_{ij}^{\alpha+\beta-1} \ln_{1-\alpha}(p_{it}/q_{it})}{\sum_{j=1}^n x_{kj} q_{ij}^{\alpha+\beta-1}} \right) \right)^{\omega(\alpha,\beta)}, \end{aligned} \quad (3)$$

где

$$\omega(\alpha, \beta) = \begin{cases} 1, & \text{если } \alpha = 0, \beta = 1; \\ 0, & \text{если } \alpha = 0, \beta \neq 1; \\ \frac{\alpha}{1-\beta}, & \text{если } \alpha \neq 0, \frac{\beta}{\alpha} < \frac{1}{\alpha} - 1; \\ 1, & \text{если } \alpha \neq 0, \frac{\beta}{\alpha} \in \left[\frac{1}{\alpha} - 1, \frac{1}{\alpha} \right]; \\ \frac{\alpha}{\alpha+\beta-1}, & \text{если } \alpha \neq 0, \frac{\beta}{\alpha} > \frac{1}{\alpha}. \end{cases}$$

При $\alpha \neq 0$ вид обновлений упрощается:

$$\begin{aligned} x_{kj} &\leftarrow x_{kj} \left(\frac{\sum_{i=1}^m a_{ik} p_{ij}^{\alpha} q_{ij}^{\beta-1}}{\sum_{i=1}^m a_{ik} q_{ij}^{\alpha+\beta-1}} \right)^{\omega'(\alpha,\beta)} ; \\ a_{ik} &\leftarrow a_{ik} \left(\frac{\sum_{j=1}^n x_{kj} p_{ij}^{\alpha} q_{ij}^{\beta-1}}{\sum_{j=1}^n x_{kj} q_{ij}^{\alpha+\beta-1}} \right)^{\omega'(\alpha,\beta)} ; \end{aligned}$$

$$\omega'(\alpha, \beta) = \begin{cases} \frac{1}{1-\beta}, & \text{если } \frac{\beta}{\alpha} < \frac{1}{\alpha} - 1; \\ \frac{1}{\alpha}, & \text{если } \frac{\beta}{\alpha} \in \left[\frac{1}{\alpha} - 1, \frac{1}{\alpha} \right]; \\ \frac{1}{\alpha + \beta - 1}, & \text{если } \frac{\beta}{\alpha} > \frac{1}{\alpha}. \end{cases}$$

В матричном виде их можно записать следующим образом:

$$\begin{aligned} X &\leftarrow X \otimes ((A^T Z) \otimes (A^T Q^{[\alpha+\beta-1]}))^{\omega'(\alpha,\beta)} ; \\ A &\leftarrow A \otimes ((Z X^T) \otimes (Q^{[\alpha+\beta-1]} X^T))^{\omega'(\alpha,\beta)}. \end{aligned} \quad (4)$$

Здесь $Z = P^{[\alpha]} \otimes Q^{[\beta-1]}$, символом \otimes обозначается операция поэлементного (Адамарова) произведения матриц, символом \oslash — поэлементного деления, $[\cdot]$ — поэлементного возведения в степень.

Видно, что при $\alpha = 0$, $\beta \neq 1$ алгоритм (3) не модифицирует A и X ; авторы [7] предлагают ограничить $\omega(\alpha, \beta)$ снизу небольшой положительной константой, чтобы алгоритм можно было использовать при α , близких к нулю, и β , не равных единице. В данной работе мы не будем подробнее останавливаться на этом частном случае.

Отметим, что в случаях $\alpha = 1$, $\beta = 1$, когда АБ-дивергенция превращается в норму Фробениуса, и $\alpha = 1$, $\beta = 0$, когда получается дивергенция Кульбака-Лейблера, алгоритм (3) совпадает с мультипликативными алгоритмами для соответствующих функций потерь из работ [2, 3].

Сходимость

Рассмотрим определения и понятия сходимости, используемые для неотрицательного матричного разложения [4].

Лагранжиан оптимизационной задачи (1) имеет вид $L(A, X, \mu, \nu) = D(P, AX) - \mu \otimes A - \nu \otimes X$, где μ и ν — матрицы множителей Лагранжа размеров $m \times r$ и $r \times n$ соответственно. Согласно теореме Каруша-Куна-Таккера [13], если (A, X) — локальный минимум, то найдутся такие μ и ν с неотрицательными элементами, что выполняются условия

$$\begin{aligned} A &\geq 0, & X &\geq 0; \\ \nabla_A L &= 0, & \nabla_X L &= 0; \\ \mu \otimes A &= 0, & \nu \otimes X &= 0. \end{aligned}$$

Упрощая вторую пару условий, получаем

$$\mu = -\nabla_A D, \quad \nu = -\nabla_X D.$$

С учетом этого, а также требований $\mu_{ik} \geq 0$, $\nu_{kj} \geq 0$, условия Каруша-Куна-Таккера можно записать в следующем виде:

$$A \geq 0, \quad X \geq 0; \tag{5-а}$$

$$\nabla_A D \geq 0, \quad \nabla_X D \geq 0; \tag{5-б}$$

$$A \otimes \nabla_A D = 0, \quad X \otimes \nabla_X D = 0. \tag{5-в}$$

Определение 1. Будем называть (A, X) стационарной точкой задачи неотрицательного матричного разложения с функцией потерь $D(P, AX)$ тогда и только тогда, когда для A и X выполняются условия Каруша-Куна-Таккера (5-а), (5-б), (5-в).

Поскольку рассматриваемые функции потерь $D(P, AX)$ не являются выпуклыми по совокупности аргументов (A, X) , выполнение (5-а)–(5-в) является только необходимым условием локального минимума. Однако, как правило, в задачах невыпуклой оптимизации, к которым относится и рассматриваемая задача, для большинства алгоритмов удается показать только сходимость к стационарной точке, а более сильных утверждений о сходимости доказать не удается [14]. На практике при минимизации невыпуклых функций именно такой вид сходимости используется как доказательство применимости алгоритма.

Мультипликативные алгоритмы широко используются в неотрицательном матричном разложении несмотря на то, что их сходимость к стационарной точке показать достаточно

сложно — для многих алгоритмов она не доказана вообще. Вместо этого как некоторый аналог сходимости зачастую используется тот факт, что на каждой итерации алгоритма значение функции потерь не возрастает. Поскольку $D(P, A^t, X^t)$ ограничена снизу (нулем), ее невозрастание позволяет надеяться, что алгоритм сойдется. При этом невозрастание функции потерь не является ни необходимым, ни достаточным условием сходимости к стационарной точке: во-первых, не гарантируется, что предельная точка алгоритма с невозрастающей функцией потерь является стационарной, во-вторых, нет гарантии, что алгоритм сойдется к своей предельной точке [15].

Для доказательства невозрастания функции потерь чаще всего используются дополнительные функции. Пусть при фиксированной матрице A минимизируется $D(X) = D(P, AX)$; функция $G(X, Y)$ называется дополнительной для $D(X)$, если

$$G(X, Y) \geq D(X), \quad G(X, X) = D(X) \quad \forall X, Y.$$

Определим

$$X^{t+1} = \arg \min_X G(X, X^t). \quad (6)$$

Тогда по построению

$$D(X^t) = G(X^t, X^t) \geq G(X^{t+1}, X^t) \geq G(X^{t+1}, X^{t+1}) = D(X^{t+1}).$$

В силу симметричности задачи факторизации функция $G(X, Y)$, дополнительная для $D(X)$, будет дополнительной и для $D(A^T) = D(P^T, X^T A^T)$:

$$G(A^T, B^T) \geq D(A^T), \quad G(A^T, A^T) = D(A^T) \quad \forall A, B.$$

Таким образом, для доказательства невозрастания функции потерь алгоритма достаточно привести такую функцию $G(X, Y)$, что обновления X можно представить как решение задачи (6), а обновления A — как решение аналогичной задачи:

$$A^{t+1} = \arg \min_A G(A, A^t). \quad (7)$$

Мультипликативный алгоритм минимизации нормы Фробениуса

Рассмотрим сначала наиболее исследованный частный случай $\alpha = 1$, $\beta = 1$, при котором АБ-дивергенция превращается в норму Фробениуса:

$$(A^*, X^*) = \arg \min_{A \geq 0, X \geq 0} D_F(P, AX) = \arg \min_{A \geq 0, X \geq 0} \|P - AX\|_F^2. \quad (8)$$

Именно для этого случая в одной из первых работ по неотрицательному матричному разложению [2] был впервые предложен мультипликативный алгоритм следующего вида:

$$\begin{aligned} X &\leftarrow X \otimes ((A^T P) \oslash (A^T A X)); \\ A &\leftarrow A \otimes ((P X^T) \oslash (A X X^T)), \end{aligned} \quad (9)$$

а в [3] показано монотонное невозрастание функции потерь при его использовании.

У этого алгоритма существует несколько проблем, возникающих на границе области неотрицательности элементов A и X . Чтобы выполнялись условия Каруша-Куна-Таккера (5-б), (5-в), необходимо, чтобы в предельной точке алгоритма не было таких k, j , что

$\nabla_{x_{jk}} D < 0$, $x_{jk} = 0$. Если элементы матрицы начального приближения X^0 строго положительны, то это не может произойти, так как, если на шаге с номером t выполняется

$$\nabla_{x_{kj}} D = \left[(A^t)^T A^t X^t \right]_{kj} - \left[(A^t)^T P \right]_{kj} < 0,$$

то на следующем шаге

$$x_{kj}^{t+1} = x_{kj}^t \frac{\left[(A^t)^T P \right]_{kj}}{\left[(A^t)^T A^t X^t \right]_{kj}} > x_{kj}^t > 0.$$

Аналогичные рассуждения справедливы и для A . Следовательно, чтобы гарантировать выполнение условий стационарности, достаточно, чтобы все элементы искомым матриц были инициализированы строго положительными, а в ходе мультипликативных обновлений они автоматически будут оставаться такими.

Однако это избавляет нас не от всех проблем: сходимость алгоритма оказывается под вопросом не только тогда, когда какие-то из элементов матриц в точности равны нулю, но и тогда, когда они к нулю стремятся. В этом случае алгоритм также может останавливаться в нестационарных точках. Чтобы показать это, рассмотрим обновление для столбца \mathbf{x}_j матрицы X , обозначая как \mathbf{p}_j соответствующий столбец матрицы P , и перепишем мультипликативное обновление в аддитивном виде:

$$\begin{aligned} \mathbf{x}_j^{t+1} &= \mathbf{x}_j^t \otimes \left((A^T \mathbf{p}_j) \otimes (A^T A \mathbf{x}_j^t) \right) = \\ &= \mathbf{x}_j^t \otimes \left((A^T \mathbf{p}_j + A^T A \mathbf{x}_j^t - A^T A \mathbf{x}_j^t) \otimes (A^T A \mathbf{x}_j^t) \right) = \\ &= \mathbf{x}_j^t \otimes \left(\mathbf{1} - (A^T A \mathbf{x}_j^t - A^T \mathbf{p}_j) \otimes (A^T A \mathbf{x}_j^t) \right) = \\ &= \mathbf{x}_j^t - \mathbf{x}_j^t \otimes (A^T A \mathbf{x}_j^t) \otimes (A^T A \mathbf{x}_j^t - A^T \mathbf{p}_j) = \\ &= \mathbf{x}_j^t - d_{\mathbf{x}_j}^T \nabla_{\mathbf{x}_j} D(\mathbf{p}_j, A \mathbf{x}_j). \end{aligned}$$

Здесь $d_{\mathbf{x}_j} \in \mathbb{R}_+^r$ — вектор-столбец с элементами $d_k = \frac{x_{kj}^t}{[A^T A \mathbf{x}_j^t]_k}$. Чтобы алгоритм обеспечивал достаточный спуск¹, необходимо, чтобы все d_k были отделены от нуля, что для данного вектора в общем случае неверно; поэтому обновления X могут прекращаться в точке, которая не является стационарной.

Чтобы избежать описанных проблем со сходимостью, достаточно явным образом отделить решения от нуля положительной константой. Это делает следующая модификация алгоритма (9):

$$\begin{aligned} X &\leftarrow \max(\varepsilon, X \otimes ((A^T P) \otimes (A^T A X))); \\ A &\leftarrow \max(\varepsilon, A \otimes ((P X^T) \otimes (A X X^T))), \end{aligned} \tag{10}$$

где $\varepsilon > 0$ — произвольная константа, а максимум берется покомпонентно.

В работе [16] доказаны следующие свойства такого алгоритма.

Утверждение 1. Для любых $\varepsilon > 0$, $P \geq 0$, $(A^0, X^0) \geq \varepsilon$ функция потерь $\|P - AX\|_F^2$ монотонно не возрастает в процессе обновлений (10).

¹В задаче минимизации функции $f(x)$, $x \in \mathbb{R}^n$, направление $d \in \mathbb{R}^n$ является направлением спуска, если $d^T \nabla f(x) < 0$, и направлением достаточного спуска, если существует такая константа C , что $d^T \nabla f(x) \leq -C \|\nabla f(x)\|^2$.

Утверждение 2. Любая предельная точка алгоритма (10) является стационарной точкой оптимизационной задачи

$$(A_\varepsilon^*, X_\varepsilon^*) = \arg \min_{A \geq \varepsilon, X \geq \varepsilon} \|P - AX\|_F^2.$$

Утверждение 3. Пусть $(A_\varepsilon, X_\varepsilon)$ — предельная точка алгоритма (10). Проведем разреживание, обнуляя элементы полученных матриц, в точности равные ε . Определим $A = A_\varepsilon \otimes [A_\varepsilon > \varepsilon]$, $X = X_\varepsilon \otimes [X_\varepsilon > \varepsilon]$. Для всех i, j, k выполняется следующее:

$$\begin{cases} \begin{cases} a_{ik} = 0, & \nabla_{a_{ik}} D \geq -\mathcal{O}(\varepsilon); \\ a_{ik} > 0, & |\nabla_{a_{ik}} D| \leq \mathcal{O}(\varepsilon); \end{cases} \\ \begin{cases} x_{kj} = 0, & \nabla_{x_{kj}} D \geq -\mathcal{O}(\varepsilon); \\ x_{kj} > 0, & |\nabla_{x_{kj}} D| \leq \mathcal{O}(\varepsilon). \end{cases} \end{cases}$$

То есть, хотя полученное таким образом решение (A, X) и не является в точности стационарной точкой исходной задачи (8), условия Каруша–Куна–Таккера (5-а)–(5-в) выполняются с точностью до $\mathcal{O}(\varepsilon)$.

Таким образом, для алгоритма (10) открытым остается только вопрос существования в генерируемой им последовательности (A^t, X^t) предельной точки. Используем для доказательства ее существования следующее свойство, установленное в работе [17].

Утверждение 4. Любая последовательность (A^t, X^t) , генерируемая алгоритмом (10), принадлежит компакту C при условии, что $(A^0, X^0) \in C$.

Отсюда, с учетом теоремы Больцано–Вейерштрасса, получаем следствие.

Следствие 1. Последовательность (A^t, X^t) , порождаемая алгоритмом (10), имеет хотя бы одну предельную точку.

Таким образом, для модифицированного мультипликативного алгоритма (10) неотрицательного матричного разложения с нормой Фробениуса показано, что он всегда сходится, а его предельная точка лежит сколь угодно близко к стационарной точке исходной задачи (8). То есть небольшими изменениями в алгоритме (9) удастся обеспечить его вычислительную устойчивость (отсутствие деления на ноль) и гарантировать сходимость в окрестность стационарной точки, причем параметром, определяющим размер этой окрестности, можно управлять, обеспечивая баланс между вычислительной устойчивостью алгоритма и точностью получаемого решения.

Мультипликативный алгоритм минимизации АБ-дивергенции

Вернемся к рассмотрению задачи (2) неотрицательного матричного разложения с АБ-дивергенцией. Относительно алгоритма (3) в работе [7] доказано следующее утверждение.

Утверждение 5. Функция $D_{AB}^{(\alpha, \beta)}(P, AX)$ монотонно не возрастает при обновлениях (4) для любого начального приближения $(A^0, X^0) \geq 0$.

Для доказательства используется следующая дополнительная функция:

$$G(Q^{t+1}, Q^t, P) = \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^r \frac{a_{ik}^t x_{kj}^t}{q_{ij}^t} \bar{d}_{AB}^{(\alpha, \beta)}(\hat{q}_{ij}, q_{ij}^t, p_{ij});$$

$$\bar{d}_{AB}^{(\alpha, \beta)}(\hat{q}_{ij}, q_{ij}^t, p_{ij}) = \begin{cases} \frac{p_{ij}^{\alpha+\beta}}{\beta(\alpha+\beta)} + \frac{(q_{ij}^t)^{\alpha+\beta}}{\alpha(\alpha+\beta)} - \frac{p_{ij}^\alpha \hat{q}_{ij}^\beta}{\alpha\beta} + \frac{(q_{ij}^t)^{\alpha+\beta-1}}{\alpha} (\hat{q}_{ij} - q_{ij}^t), \\ \text{если } \frac{\beta}{\alpha} < \frac{1}{\alpha} - 1; \\ \frac{p_{ij}^{\alpha+\beta}}{\beta(\alpha+\beta)} + \frac{\hat{q}_{ij}^{\alpha+\beta}}{\alpha(\alpha+\beta)} - \frac{p_{ij}^\alpha \hat{q}_{ij}^\beta}{\alpha\beta}, \\ \text{если } \frac{\beta}{\alpha} \in \left[\frac{1}{\alpha} - 1, \frac{1}{\alpha} \right]; \\ \frac{p_{ij}^{\alpha+\beta}}{\beta(\alpha+\beta)} + \frac{\hat{q}_{ij}^{\alpha+\beta}}{\alpha(\alpha+\beta)} - \frac{p_{ij}^\alpha (q_{ij}^t)^\beta}{\alpha\beta} - \frac{p_{ij}^\alpha (q_{ij}^t)^{\beta-1}}{\alpha} (\hat{q}_{ij} - q_{ij}^t); \\ \text{если } \frac{\beta}{\alpha} > \frac{1}{\alpha}, \end{cases} \quad (11)$$

$$\hat{q}_{ij}^t = q_{ij}^t \frac{a_{ik}^{t+1} x_{kj}^{t+1}}{a_{ik}^t x_{kj}^t}.$$

При $\beta = 0$ и $\alpha = -\beta$ функция доопределяется при помощи правила Лопиталья.

Заметим, что знаменатели в (4) могут принимать нулевые значения. Чтобы избежать деления на ноль, в работе [7] было предложено прибавлять небольшую положительную константу ε к знаменателям выражений в (4), однако вопрос влияния такой модификации на сходимость алгоритма и свойства получаемого решения не рассматривался. Вместо добавления ε к знаменателям множителей мультипликативного алгоритма модифицируем (4) аналогично (10), отделив принимаемые элементами матриц значения от нуля:

$$\begin{aligned} X &\leftarrow \max \left(\varepsilon, X \otimes ((A^T Z) \otimes (A^T Q^{[\alpha+\beta-1]}))^{[\omega'(\alpha, \beta)]} \right); \\ A &\leftarrow \max \left(\varepsilon, A \otimes ((Z X^T) \otimes (Q^{[\alpha+\beta-1]} X^T))^{[\omega'(\alpha, \beta)]} \right). \end{aligned} \quad (12)$$

Рассмотрим задачу минимизации функции потерь в отделенной от нуля области:

$$(A_\varepsilon^*, X_\varepsilon^*) = \arg \min_{A \geq \varepsilon, X \geq \varepsilon} D_{AB}^{(\alpha, \beta)}(P, AX). \quad (13)$$

Теорема 1. Для любого $\varepsilon > 0$ функция $D_{AB}^{(\alpha, \beta)}(P, AX)$ монотонно не возрастает при обновлениях (12) для любого начального приближения $(A^0, X^0) \geq \varepsilon$.

Доказательство. Пусть $A^t \geq \varepsilon$. Обозначим за $X^{t+1} = \arg \min_{X \geq \varepsilon} G(A^t X, A^t X^t, P)$, где G — дополнительная функция (11). По определению дополнительной функции имеем:

$$\begin{aligned} D_{AB}^{(\alpha, \beta)}(P, A^t X^t) &= G(A^t X^t, A^t X^t, P) \geq \min_{X \geq \varepsilon} G(A^t X, A^t X^t, P) = \\ &= G(A^t X^{t+1}, A^t X^t, P) \geq G(A^t X^{t+1}, A^t X^{t+1}, P) = D_{AB}^{(\alpha, \beta)}(P, A^t X^{t+1}). \end{aligned}$$

Покажем, что обновления (12) совпадают с решениями задач (6), (7) минимизации дополнительной функции. Рассмотрим функцию

$$G(A^t X, A^t X^t, P) = \sum_{i,j,k} \frac{a_{ik}^t x_{kj}^t}{q_{ij}^t} \bar{d}_{AB}^{(\alpha,\beta)} \left(q_{ij}^t \frac{x_{kj}}{x_{kj}^t}, q_{ij}^t, p_{ij} \right).$$

Обозначим

$$g_{kj}(x_{kj}, x_{kj}^t, A^t, P) = x_{kj}^t \sum_{i=1}^m \frac{a_{ik}^t}{q_{ij}^t} \bar{d}_{AB}^{(\alpha,\beta)} \left(q_{ij}^t \frac{x_{kj}}{x_{kj}^t}, q_{ij}^t, p_{ij} \right).$$

Поскольку функция $G(A^t X, A^t X^t, P)$ состоит из независимых слагаемых по x_{kj} ,

$$\begin{aligned} X^{t+1} &= \arg \min_{X \geq \varepsilon} G(A^t X, A^t X^t, P) \iff \\ x_{kj}^{t+1} &= \arg \min_{x_{kj} \geq \varepsilon} g_{kj}(x_{kj}, x_{kj}^t, A^t, P) \quad \forall k, j. \end{aligned}$$

Далее, поскольку функция

$$\bar{d}_{AB}^{(\alpha,\beta)} \left(q_{ij}^t \frac{x_{kj}}{x_{kj}^t}, q_{ij}^t, p_{ij} \right) = \begin{cases} \frac{p_{ij}^{\alpha+\beta}}{\beta(\alpha+\beta)} + \frac{(q_{ij}^t)^{\alpha+\beta}}{\alpha(\alpha+\beta)} - \frac{p_{ij}^\alpha (q_{ij}^t)^\beta x_{kj}^\beta}{\alpha\beta (x_{kj}^t)^\beta} + \frac{(q_{ij}^t)^{\alpha+\beta}}{\alpha} \left(\frac{x_{kj}}{x_{kj}^t} - 1 \right), & \text{если } \frac{\beta}{\alpha} < \frac{1}{\alpha} - 1; \\ \frac{p_{ij}^{\alpha+\beta}}{\beta(\alpha+\beta)} + \frac{(q_{ij}^t)^{\alpha+\beta} x_{kj}^{\alpha+\beta}}{\alpha(\alpha+\beta) (x_{kj}^t)^{\alpha+\beta}} - \frac{p_{ij}^\alpha \hat{q}_{ij}^\beta}{\alpha\beta}, & \text{если } \frac{\beta}{\alpha} \in \left[\frac{1}{\alpha} - 1, \frac{1}{\alpha} \right]; \\ \frac{p_{ij}^{\alpha+\beta}}{\beta(\alpha+\beta)} + \frac{(q_{ij}^t)^{\alpha+\beta} x_{kj}^{\alpha+\beta}}{\alpha(\alpha+\beta) (x_{kj}^t)^{\alpha+\beta}} - \frac{p_{ij}^\alpha (q_{ij}^t)^\beta}{\alpha\beta} - \frac{p_{ij}^\alpha (q_{ij}^t)^\beta}{\alpha} \left(\frac{x_{kj}}{x_{kj}^t} - 1 \right), & \text{если } \frac{\beta}{\alpha} > \frac{1}{\alpha} \end{cases}$$

является выпуклой по x_{kj} по построению [7], то выпукла по x_{kj} и $g_{kj}(x_{kj}, x_{kj}^t, A^t, P)$; следовательно, ее безусловный минимум можно получить из условия

$$\frac{\partial g_{kj}}{\partial x_{kj}} = x_{kj}^t \sum_{i=1}^m \frac{a_{ik}^t}{q_{ij}^t} \frac{\partial \bar{d}_{AB}}{\partial x_{kj}} = 0.$$

Распишем его:

$$\frac{\partial \bar{d}_{AB}}{\partial x_{kj}} = \begin{cases} \frac{(q_{ij}^t)^\beta}{\alpha} \frac{1}{x_{kj}^t} \left((q_{ij}^t)^\alpha - p_{ij}^\alpha \left(\frac{x_{kj}}{x_{kj}^t} \right)^{\beta-1} \right), & \text{если } \frac{\beta}{\alpha} < \frac{1}{\alpha} - 1; \\ \frac{(q_{ij}^t)^\beta}{\alpha} \frac{x_{kj}^{\beta-1}}{(x_{kj}^t)^\beta} \left((q_{ij}^t)^\alpha \left(\frac{x_{kj}}{x_{kj}^t} \right)^\alpha - p_{ij}^\alpha \right), & \text{если } \frac{\beta}{\alpha} \in \left[\frac{1}{\alpha} - 1, \frac{1}{\alpha} \right]; \\ \frac{(q_{ij}^t)^\beta}{\alpha} \frac{1}{x_{kj}^t} \left((q_{ij}^t)^\alpha \left(\frac{x_{kj}}{x_{kj}^t} \right)^{\alpha+\beta-1} - p_{ij}^\alpha \right), & \text{если } \frac{\beta}{\alpha} > \frac{1}{\alpha}. \end{cases}$$

$$\frac{\partial g_{kj}}{\partial x_{kj}} = \begin{cases} \frac{1}{\alpha} \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\beta-1} \left((q_{ij}^t)^\alpha - p_{ij}^\alpha \left(\frac{x_{kj}}{x_{kj}^t} \right)^{\beta-1} \right), & \text{если } \frac{\beta}{\alpha} < \frac{1}{\alpha} - 1; \\ \frac{1}{\alpha} \left(\frac{x_{kj}}{x_{kj}^t} \right)^\beta \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\beta-1} \left((q_{ij}^t)^\alpha \left(\frac{x_{kj}}{x_{kj}^t} \right)^\alpha - p_{ij}^\alpha \right), & \text{если } \frac{\beta}{\alpha} \in \left[\frac{1}{\alpha} - 1, \frac{1}{\alpha} \right]; \\ \frac{1}{\alpha} \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\beta-1} \left((q_{ij}^t)^\alpha \left(\frac{x_{kj}}{x_{kj}^t} \right)^{\alpha+\beta-1} - p_{ij}^\alpha \right), & \text{если } \frac{\beta}{\alpha} > \frac{1}{\alpha}; \end{cases}$$

$$\frac{\partial g_{kj}}{\partial x_{kj}} = 0 \Leftrightarrow \begin{cases} \left(\frac{x_{kj}}{x_{kj}^t} \right)^{1-\beta} \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\alpha+\beta-1} = \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\beta-1} p_{ij}^\alpha, & \text{если } \frac{\beta}{\alpha} < \frac{1}{\alpha} - 1; \\ \left(\frac{x_{kj}}{x_{kj}^t} \right)^\alpha \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\alpha+\beta-1} = \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\beta-1} p_{ij}^\alpha, & \text{если } \frac{\beta}{\alpha} \in \left[\frac{1}{\alpha} - 1, \frac{1}{\alpha} \right]; \\ \left(\frac{x_{kj}}{x_{kj}^t} \right)^{\alpha+\beta-1} \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\alpha+\beta-1} = \sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\beta-1} p_{ij}^\alpha, & \text{если } \frac{\beta}{\alpha} > \frac{1}{\alpha}. \end{cases}$$

Отсюда получаем:

$$x_{kj}^{t+1} = \arg \min_{x_{kj} \geq \varepsilon} g_{kj}(x_{kj}, x_{kj}^t, A^t, P) = \max \left(\varepsilon, x_{kj}^t \left(\frac{\sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\beta-1} p_{ij}^\alpha}{\sum_{i=1}^m a_{ik}^t (q_{ij}^t)^{\alpha+\beta-1}} \right)^{\omega'(\alpha, \beta)} \right),$$

что в точности совпадает с обновлениями алгоритма (12); следовательно, обновления по X из (12) приводят к монотонному невозрастанию функции потерь.

Для обновлений по A доказательство строится аналогично. \blacksquare

Рассмотрим теперь вопрос качества решения, получаемого модифицированным алгоритмом (12). Для задачи (13) условия Каруша–Куна–Таккера записываются следующим образом:

$$\begin{aligned} A &\geq \varepsilon, & X &\geq \varepsilon, \\ \nabla_A D_{AB}^{(\alpha, \beta)}(P, AX) &\geq 0, & \nabla_X D_{AB}^{(\alpha, \beta)}(P, AX) &\geq 0, \\ (A - \varepsilon) \otimes \nabla_A D_{AB}^{(\alpha, \beta)}(P, AX) &= 0, & (X - \varepsilon) \otimes \nabla_X D_{AB}^{(\alpha, \beta)}(P, AX) &= 0. \end{aligned}$$

Следующая теорема утверждает, что в предельной точке последовательности, получаемой при помощи алгоритма (12), эти условия выполняются.

Теорема 2. Любая предельная точка последовательности, порождаемой алгоритмом (12) для любого начального приближения $(A^0, X^0) \geq \varepsilon$, является стационарной точкой отделенной от нуля задачи (13).

Доказательство. Пусть (\bar{A}, \bar{X}) — предельная точка последовательности (A^t, X^t) , порождаемой алгоритмом (12), а $\bar{Q} = \bar{A}\bar{X}$. Поскольку функция потерь ограничена снизу (нулем), из ее монотонного невозрастания следует сходимость $D_{AB}^{(\alpha, \beta)}(P, A^t X^t)$ к $D_{AB}^{(\alpha, \beta)}(P, \bar{A}\bar{X})$. Кроме того,

$$\bar{x}_{kj} = \max \left(\varepsilon, c_{kj}^{\omega'(\alpha, \beta)} \bar{x}_{kj} \right), \quad (15)$$

где

$$c_{kj} \equiv \frac{\sum_{i=1}^m \bar{a}_{ik} \bar{q}_{ij}^{\beta-1} p_{ij}^{\alpha}}{\sum_{i=1}^m \bar{a}_{ik} \bar{q}_{ij}^{\alpha+\beta-1}},$$

причем это выражение всегда определено, так как $(A^0, X^0) \geq \varepsilon$ и выражение в знаменателе не может быть равным нулю. Пользуясь выражением для градиента АБ-дивергенции:

$$\left[\nabla_X D_{AB}^{(\alpha, \beta)}(P, AX) \right]_{kj} = \frac{1}{\alpha} \sum_{i=1}^m q_{ij}^{\beta-1} a_{ik} (q_{ij}^{\alpha} - p_{ij}^{\alpha}) = \frac{1}{\alpha} \left(1 - \frac{\sum_{i=1}^m a_{ik} q_{ij}^{\beta-1} p_{ij}^{\alpha}}{\sum_{i=1}^m a_{ik} q_{ij}^{\alpha+\beta-1}} \right),$$

необходимые и достаточные условия стационарности для компоненты \bar{x}_{kj} решения задачи (13) можно записать следующим образом:

$$\begin{aligned} \bar{x}_{kj} &\geq \varepsilon; \\ c_{kj} &\leq 1; \\ (\bar{x}_{kj} - \varepsilon)(c_{kj} - 1) &= 0. \end{aligned}$$

В то же время, из выражения (15) следует, что, поскольку \bar{x}_{kj} — предельная точка, либо $c_{kj}^{\omega'(\alpha, \beta)} \leq 1$ (а раз $\omega'(\alpha, \beta) \leq 1$, то и $c_{kj} \leq 1$) и $\bar{x}_{kj} = \varepsilon$, либо $\bar{x}_{kj} > \varepsilon$ и $c_{kj} = 1$; следовательно, условия стационарности для \bar{x}_{kj} выполняются. Это справедливо для всех k, j , т. е. для всей матрицы \bar{X} .

Для \bar{A} доказательство строится аналогично. ■

Пусть $(A_{\varepsilon}, X_{\varepsilon})$ — предельная точка алгоритма (12). Определим следующие матрицы:

$$\begin{aligned} A_0 &= A_{\varepsilon} \otimes [A_{\varepsilon} > \varepsilon]; \\ X_0 &= X_{\varepsilon} \otimes [X_{\varepsilon} > \varepsilon], \end{aligned}$$

т. е. обнулим в $A_{\varepsilon}, X_{\varepsilon}$ все элементы, равные ε . Покажем, что (A_0, X_0) аппроксимирует стационарную точку исходной задачи (2) с точностью порядка ε .

Теорема 3. Для матриц (A_0, X_0) , полученных из $(A_{\varepsilon}, X_{\varepsilon})$ обнулением элементов, равных ε , верно следующее:

$$\left\{ \begin{array}{l} \left[\begin{array}{l} a_{0ik} = 0, \quad \left(\nabla_A D_{AB}^{(\alpha, \beta)}(P, A_0 X_0) \right)_{ik} \geq -\mathcal{O}(\varepsilon); \\ a_{0ik} > 0, \quad \left| \nabla_A D_{AB}^{(\alpha, \beta)}(P, A_0 X_0) \right|_{ik} \leq \mathcal{O}(\varepsilon); \end{array} \right. \\ \left[\begin{array}{l} x_{0kj} = 0, \quad \left(\nabla_X D_{AB}^{(\alpha, \beta)}(P, A_0 X_0) \right)_{kj} \geq -\mathcal{O}(\varepsilon); \\ x_{0kj} > 0, \quad \left| \nabla_X D_{AB}^{(\alpha, \beta)}(P, A_0 X_0) \right|_{kj} \leq \mathcal{O}(\varepsilon), \end{array} \right. \end{array} \right.$$

т. е. в точке (A_0, X_0) условия Каруша-Куна-Таккера для исходной задачи (2) практически выполняются.

Доказательство. Для стационарной точки $(A_\varepsilon, X_\varepsilon)$ задачи (13) выполняются условия Каруша–Куна–Таккера:

$$\left\{ \begin{array}{l} \left[\begin{array}{l} a_{\varepsilon ik} = \varepsilon, \quad \left(\nabla_A D_{AB}^{(\alpha, \beta)} (P, A_\varepsilon X_\varepsilon) \right)_{ik} \geq 0; \\ a_{\varepsilon ik} > \varepsilon, \quad \left| \nabla_A D_{AB}^{(\alpha, \beta)} (P, A_\varepsilon X_\varepsilon) \right|_{ik} = 0; \end{array} \right. \\ \left. \left[\begin{array}{l} x_{\varepsilon kj} = \varepsilon, \quad \left(\nabla_X D_{AB}^{(\alpha, \beta)} (P, A_\varepsilon X_\varepsilon) \right)_{kj} \geq 0; \\ x_{\varepsilon kj} > \varepsilon, \quad \left| \nabla_X D_{AB}^{(\alpha, \beta)} (P, A_\varepsilon X_\varepsilon) \right|_{kj} = 0. \end{array} \right. \right. \end{array} \right. \quad (16)$$

Построим верхнюю оценку для расстояния между значениями градиента $\nabla_X D_{AB}^{(\alpha, \beta)}$ в точках (A_0, X_0) и $(A_\varepsilon, X_\varepsilon)$, обозначая $Q_0 = A_0 X_0$ и $Q_\varepsilon = A_\varepsilon X_\varepsilon$:

$$\begin{aligned} & \max_{k,j} \left| \nabla_X D_{AB}^{(\alpha, \beta)} (P, A_\varepsilon X_\varepsilon) - \nabla_X D_{AB}^{(\alpha, \beta)} (P, A_0 X_0) \right|_{k,j} = \\ & = \frac{1}{\alpha} \max_{k,j} \left| A_\varepsilon^\top (Q_\varepsilon^{[\alpha+\beta-1]} - P^{[\alpha]} \otimes Q_\varepsilon^{[\beta-1]}) - A_0^\top (Q_0^{[\alpha+\beta-1]} - P^{[\alpha]} \otimes Q_0^{[\beta-1]}) \right|_{k,j} \leq \\ & \leq \frac{1}{\alpha} \max_{k,j} \left| A_\varepsilon^\top Q_\varepsilon^{[\alpha+\beta-1]} - A_0^\top Q_0^{[\alpha+\beta-1]} \right|_{k,j} + \\ & \quad + \frac{1}{\alpha} \max_{k,j} \left| A_\varepsilon^\top (Q_\varepsilon^{[\beta-1]} \otimes P^{[\alpha]}) - A_0^\top (Q_0^{[\beta-1]} \otimes P^{[\alpha]}) \right|_{k,j} \leq \\ & \leq \frac{1}{\alpha} \max_{k,j} \left| \sum_{i=1}^m \left((a_{0ik} + \varepsilon) q_{\varepsilon ij}^{\alpha+\beta-1} - a_{0ik} q_{0ij}^{\alpha+\beta-1} \right) \right| + \\ & \quad + \frac{1}{\alpha} \max_{k,j} \left| \sum_{i=1}^m p_{ij}^\alpha \left((a_{0ik} + \varepsilon) q_{\varepsilon ij}^{\beta-1} - a_{0ik} q_{0ij}^{\beta-1} \right) \right|. \end{aligned} \quad (17)$$

Запишем выражения для элементов матриц Q_0 и Q_ε :

$$\begin{aligned} q_{0ij} &= \sum_{l=1}^r a_{0il} x_{0lj} = \sum_{l: \substack{a_{0il} \neq 0, \\ x_{0lj} \neq 0}} a_{0il} x_{0lj}; \\ q_{\varepsilon ij} &= \sum_{l=1}^r a_{\varepsilon il} x_{\varepsilon lj} = \sum_{l: \substack{a_{\varepsilon il} > \varepsilon, \\ x_{\varepsilon lj} > \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj} + \sum_{l: \substack{a_{\varepsilon il} > \varepsilon, \\ x_{\varepsilon lj} = \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj} + \sum_{l: \substack{a_{\varepsilon il} = \varepsilon, \\ x_{\varepsilon lj} > \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj} + \sum_{l: \substack{a_{\varepsilon il} = \varepsilon, \\ x_{\varepsilon lj} = \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj}. \end{aligned}$$

По построению матриц A_0, X_0 верно следующее:

$$\begin{aligned} & \sum_{l: \substack{a_{\varepsilon il} > \varepsilon, \\ x_{\varepsilon lj} > \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj} = \sum_{l: \substack{a_{0il} \neq 0, \\ x_{0lj} \neq 0}} a_{0il} x_{0lj} = q_{0ij}; \\ & \sum_{l: \substack{a_{\varepsilon il} > \varepsilon, \\ x_{\varepsilon lj} = \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj} = \varepsilon \sum_{l: \substack{a_{0il} \neq 0, \\ x_{0lj} = 0}} a_{0il}; \\ & \sum_{l: \substack{a_{\varepsilon il} = \varepsilon, \\ x_{\varepsilon lj} > \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj} = \varepsilon \sum_{l: \substack{a_{0il} = 0, \\ x_{0lj} \neq 0}} x_{0lj}; \\ & \sum_{l: \substack{a_{\varepsilon il} = \varepsilon, \\ x_{\varepsilon lj} = \varepsilon}} a_{\varepsilon il} x_{\varepsilon lj} = \varepsilon^2 \cdot \#\{l: a_{\varepsilon il} = \varepsilon, x_{\varepsilon lj} = \varepsilon\} = \varepsilon^2 \cdot \#\{l: a_{0il} = 0, x_{0lj} = 0\}, \end{aligned}$$

откуда получаем:

$$q_{\varepsilon ij} = \sum_{\substack{l: a_{0il} \neq 0, \\ x_{0lj} \neq 0}} a_{0il} x_{0lj} + \varepsilon \left(\sum_{\substack{l: a_{0il} \neq 0, \\ x_{0lj} = 0}} a_{0il} + \sum_{\substack{l: a_{0il} = 0, \\ x_{0lj} \neq 0}} x_{0lj} \right) + \varepsilon^2 \cdot \# \{l: a_{0il} = 0, x_{0lj} = 0\} \equiv \\ \equiv q_{0ij} + \varepsilon d_{ij} + \varepsilon^2 f_{ij}.$$

Используя это представление, запишем первые несколько членов ряда Тейлора функции $q_{\varepsilon ij}^y$ при разложении по ε в нуле:

$$q_{\varepsilon ij}^y = (q_{0ij} + \varepsilon d_{ij} + \varepsilon^2 f_{ij})^y = q_{0ij}^y + y q_{0ij}^{y-1} d_{ij} \varepsilon + \mathcal{O}(\varepsilon^2).$$

Подставим это разложение при $y = \alpha + \beta - 1$ в выражение под знаком суммы в первом слагаемом правой части (17):

$$(a_{0ik} + \varepsilon) q_{\varepsilon ij}^{\alpha+\beta-1} - a_{0ik} q_{0ij}^{\alpha+\beta-1} = (a_{0ik} + \varepsilon) \left(q_{0ij}^{\alpha+\beta-1} + (\alpha + \beta - 1) q_{0ij}^{\alpha+\beta-2} d_{ij} \varepsilon + \mathcal{O}(\varepsilon^2) \right) - \\ - a_{0ik} q_{0ij}^{\alpha+\beta-1} = \varepsilon (\alpha + \beta - 1) q_{0ij}^{\alpha+\beta-2} d_{ij} (a_{0ik} + \varepsilon) + \mathcal{O}(\varepsilon^2) = \mathcal{O}(\varepsilon).$$

Аналогично для второго слагаемого при $y = \beta - 1$:

$$p_{ij}^\alpha \left((a_{0ik} + \varepsilon) q_{\varepsilon ij}^{\beta-1} - a_{0ik} q_{0ij}^{\beta-1} \right) = p_{ij} \left((a_{0ik} + \varepsilon) \left(q_{0ij}^{\beta-1} + (\beta - 1) q_{0ij}^{\beta-2} d_{ij} \varepsilon + \mathcal{O}(\varepsilon^2) \right) - \right. \\ \left. - a_{0ik} q_{0ij}^{\beta-1} \right) = p_{ij} \left(\varepsilon (\beta - 1) q_{0ij}^{\beta-2} d_{ij} (a_{0ik} + \varepsilon) + \mathcal{O}(\varepsilon^2) \right) = \mathcal{O}(\varepsilon).$$

Таким образом, для максимального расстояния между матрицами градиентов $\nabla_X D_{AB}^{(\alpha, \beta)}$ справедливо

$$\max_{k,j} \left| \nabla_X D_{AB}^{(\alpha, \beta)}(P, A_\varepsilon X_\varepsilon) - \nabla_X D_{AB}^{(\alpha, \beta)}(P, A_0 X_0) \right|_{k,j} \leq \mathcal{O}(\varepsilon).$$

Аналогичным образом легко показать, что

$$\max_{i,k} \left| \nabla_A D_{AB}^{(\alpha, \beta)}(P, A_\varepsilon X_\varepsilon) - \nabla_A D_{AB}^{(\alpha, \beta)}(P, A_0 X_0) \right|_{i,k} \leq \mathcal{O}(\varepsilon).$$

Подставляя полученное в условия Каруша–Куна–Таккера для точки $(A_\varepsilon, X_\varepsilon)$ (16), получаем утверждение теоремы. ■

Заключение

Функции потерь, которые могут быть применены в задаче неотрицательного матричного разложения, соответствуют различным вероятностным предположениям о характере шума в модели. Несмотря на то что выбор функции потерь существенным образом влияет на решение задачи, на практике исследователи зачастую выбирают норму Фробениуса или дивергенцию Кульбака–Лейблера, поскольку методы для них лучше всего известны и проработаны. Разработка методов, использующих такие параметрические семейства функций потерь, как АБ-дивергенции, позволяет унифицировать работу с различными функциями потерь, делая выбор меры качества более обоснованным. Имея в распоряжении такие методы, вопрос определения модели шума можно свести к выбору оптимальных значений параметров дивергенции.

Нулевые элементы матриц являются проблемой для многих методов неотрицательного матричного разложения. Во многих задачах естественно предположение о разреженности матриц-множителей, но добиться ее мультипликативными методами нельзя, а более сложные методы разработаны только для отдельных функций потерь. Предложенный в данной работе метод позволяет, с одной стороны, получить разреженные матрицы, а с другой — гарантировать, что получаемое решение близко к стационарной точке задачи.

Литература

- [1] *Paatero P., Tapper U.* Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values // *Environmetrics*, 1994. Vol. 5. Pp. 111–126.
- [2] *Lee D., Seung S.* Learning the parts of objects by non-negative matrix factorization // *Nature*, 1999. Vol. 401, no. 6755. Pp. 788–791.
- [3] *Lee D., Seung S.* Algorithms for non-negative matrix factorization // *Advances in neural information processing*, 2001. Vol. 13. Pp. 556–562.
- [4] *Wang Y., Zhang Y.* Non-negative matrix factorization: A comprehensive review // *IEEE Transactions on Knowledge and Data Engineering*, 2012. Vol. 25, no. 6. Pp. 1336–1353.
- [5] *Zhang Z.-Y.* Divergence functions of non negative matrix factorization: A comparison study // *Communications in Statistics — Simulation and Computation*, 2011. Vol. 40, no. 10. Pp. 1594–1612.
- [6] *Févotte C., Cemgil A.* Nonnegative matrix factorizations as probabilistic inference in composite models // *17th European Signal Processing Conference*, 2009. Glasgow, Scotland. Pp. 1913–1917.
- [7] *Cichocki A., Cruces S., Amari S.* Generalized alpha-beta divergences and their application to robust nonnegative matrix factorization // *Entropy*, 2011. Vol. 13, no. 1. Pp. 134–170.
- [8] *Cichocki A., Lee H., Kim Y., Choi S.* Nonnegative matrix factorization with α -divergence // *Pattern Recognition Letters*, 2008. Vol. 29, no. 9. Pp. 1433–1440.
- [9] *Kompass R.* A generalized divergence measure for nonnegative matrix factorization // *Neural Computation*, 2007. Vol. 19, no. 3. Pp. 780–791.
- [10] *Рябенко Е. А.* Настройка нелинейной модели данных экспериментов с экспрессионными ДНК-микрочипами // *Математическая биология и биоинформатика*, 2012. Т. 7. № 2. С. 554–566.
- [11] *Grippo L., Sciandrone M.* On the convergence of the block nonlinear Gauss–Seidel method under convex constraints // *Operations Research Letters*, 2000. Vol. 26, no. 3. Pp. 127–136.
- [12] *Cichocki A., Anh-Huy P.* Fast local algorithms for large scale nonnegative matrix and tensor factorizations // *IEICE Transactions on Fundamentals of Electronics Communications and Computer Sciences*, 2009. Vol. E92-A, no. 3. Pp. 708–721.
- [13] *Bertsekas D. P.* *Nonlinear programming*. Massachusetts: Athena Scientific Belmont, 1999. 780 p.
- [14] *Horst R., Pardalos P. M., Nguyen V. T.* *Introduction to global optimization*. 2nd ed. New York: Springer, 2001. 358 p.
- [15] *Badeau R., Bertin N., Vincent E.* Stability analysis of multiplicative update algorithms and application to nonnegative matrix factorization // *IEEE Transactions on Neural Networks*, 2010. Vol. 21, no. 12. Pp. 1869–1881.
- [16] *Gillis N.* Nonnegative matrix factorization: Complexity, algorithms and applications. PhD Thesis. Departement d’Ingenierie Mathematique, Universite catholique de Louvain, 2011.
- [17] *Hibi R., Takahashi N.* A modified multiplicative update algorithm for euclidean distance-based nonnegative matrix factorization and its global convergence // *Neural Information Processing*, 2011. Vol. 7063. Pp. 655–662.