

# МАТЕМАТИЧНЕ МОДЕЛЮВАННЯ В ПРИРОДНИЧИХ НАУКАХ ТА ІНФОРМАЦІЙНІ ТЕХНОЛОГІЇ



А.А. ШУМЕЙКО, д.т.н., профессор

М.Ю. БРАТУГА, магистр

Днепропетровский государственный технический университет, г. Каменское

## Использование метода главных компонент для кластеризации изображений лиц

В работе предложена кластеризация изображений лиц методом  $k$ -средних, использующим в качестве центра кластера главную компоненту элементов, входящих в кластер, а в качестве критерия близости, значения соответствующих коэффициентов. Рассмотрено приложение полученной кластеризации к решению задачи поиска лиц, наиболее похожих на тестовое изображение.

In the paper, clustering faces images by the  $k$ -means method is suggested. The center of the cluster using the main component of the Principal Component Analysis clusters elements. As a criterion of proximity, the values of the corresponding coefficients are used. The application of the resulting clusterization to the solution of the problem of finding the persons most similar to the test image is considered.

### Введение

Круг задач, рассматриваемых в обработке изображений (Image Processing), достаточно большой, и среди них важную роль играют задачи, связанные с особенностями восприятия человеческих лиц [1—2]. Это нахождение на изображении человеческого лица, идентификация этого человека или анализ полученного изображения для оценки агрессивности данной личности, пола, возраста и других характеристик. Одной из задач является нахождение в имеющейся выборке человека похожего да данного. Что есть «похожесть» человеческих лиц. Это вопрос не только и не столько математический, сколько психологический, пример тому — тест Леопольда Сонди. Но, тем не менее, задача «похожести» всегда будет интересна человеку, вечная проблема — на кого похож ребенок. С другой стороны, известны многолетние наблюдения, которые говорят о том, что достаточно часто в счастливых семьях, супруги похожи друг на друга. Что первично — счастливы, потому, что похожи или похожи, потому что счастливы? Да и не важно. Важно то, что сама по себе задача похожести лиц имеет смысл и интересна.

### Постановка задачи

Пусть дано множество изображений человеческих лиц  $\mathfrak{S} = \{I_i\}_{i=1}^n$  и изображение  $I^0$ , для которого из  $\mathfrak{S}$  нужно выбрать  $m$  наиболее похожих и, упорядочить их по рейтингу (степени похожести). Прежде всего, нужно заметить, что так как изображения представляют собой проекции трехмерного тела на плоскость в условиях разной освещенности, то первой задачей является

выравнивание и масштабирование изображения (рис. 1), то есть использование аффинного преобразования, позволяющего получить изображения анфас в одном и том же масштабе. Сама по себе эта задача нетривиальная и не является предметом нашего исследования. Будем считать, что это уже проведено.

Для того, чтобы нивелировать искажения, внесенные освещением, будем использовать не полноцветную палитру, а люминисцентную составляющую [3].

$$Y = \frac{1}{256}(77\text{Red} + 150\text{Green} + 29\text{Blue} + 128).$$

### Вычисление степени похожести лиц

Самым простым подходом к определению рейтинга похожести лиц является использование метода наименьших квадратов

$$\left( I^0 - \sum_{i=1}^n \alpha_i I_i \right)^2 \xrightarrow{\alpha_i} \min.$$

Решение этой задачи сводится к нахождению корней системы линейных уравнений

$$\sum_{i=1}^n \alpha_i \langle I_i, I_j \rangle = \langle I^0, I_j \rangle (j = 1, 2, \dots, n).$$

Большее значение  $\alpha_i$  будет соответствовать степени «большой похожести», то есть более высокому рейтингу данного изображения.

В случае малого размера выборки этот подход дает приемлемые результаты, но увеличение выборки достаточно быстро приводит как к увеличению времени работы алгоритма, так и к неустойчивости решения.



Рис. 1. Выравненная и масштабированная база лиц

### Кластеризация изображений

Одним из подходов решения этой проблемы является кластеризация исходного множества данных, то есть получение подвыборок, в которых собраны изображения людей, которые в данном кластере более похожи между собой, чем с любым другим из другого кластера [4].

Пусть  $\mathfrak{X} = \{x_i\}_{i=1}^n$  множество объектов, представленное набором атрибутов  $x_i = \{t_1^i, t_2^i, \dots, t_m^i\}$ , где  $t_v^i$  принимает значения из заданного множества  $T_v^i$ . Задача кластеризации состоит в построении множества  $C = \{c_v\}_{v=1}^k$  и отображения  $F: \mathfrak{X} \rightarrow C$  заданного множества объектов на множество кластеров. Кластер содержит объекты из  $\mathfrak{X}$  похожие (по заданному критерию) друг на друга

$$x_i \in c_v, x_j \in c_v \Rightarrow d(x_i, x_j) < \varepsilon,$$

где  $d(\bullet, \circ)$  — мера близости между объектами (расстояние), а  $\varepsilon$  — максимальное значение порога, формирующего один кластер.

Задача кластеризации данных является достаточно популярной, существует множество разных методов и подходов, среди которых наиболее популярным является метод  $k$ -средних. В методах кластеризации краеугольной характеристикой является мера близости  $d(\bullet, \circ)$ , так для  $k$ -средних это евклидово расстояние. Приведем этот метод.

Пусть  $C = \{c_i\}_{i=1}^k$  множество кластеров с центрами

$$\mu_i = \frac{\sum \{x_j \mid x_j \in c_i\}}{\sum \{1 \mid x_j \in c_i\}} = \frac{\sum_{j=1}^n u_j^i x_j}{\sum_{j=1}^n u_j^i},$$

где  $u_j^i$  индикаторная функция, то есть

$$u_j^i = \begin{cases} 1, & \text{если } x_j \in c_i, \\ 0, & \text{в противном случае.} \end{cases}$$

Целевая функция

$$S(C, \mathfrak{X}) = \sum_{i=1}^k \sum_{j=1}^n u_j^i d(x_j, \mu_i),$$

и условия

$$\sum_{i=1}^k u_j^i = 1, 0 < \sum_{j=1}^k u_j^i \leq n,$$

то есть, каждый элемент может быть только в одном кластере, и, кластер не может быть пустым или содержать элементов больше, чем их исходное количество.

Условие остановки выполнения алгоритма после  $v$ -го шага будет иметь вид

$$|S^v(C, \mathfrak{X}) - S^{v-1}(C, \mathfrak{X})| < \varepsilon,$$

где  $\varepsilon$  выбранный порог.

Непосредственное использование метода  $k$ -средних для кластеризации изображений лиц весьма проблематично, во-первых, что есть центром тяжести данного кластера, а во-вторых, что есть расстояние между двумя изображениями. Решению этой задачи посвящена данная работа.

Для решения этой задачи будем использовать метод главных компонент, то есть нахождение минимума  $\varepsilon(e_1, \dots, e_k, \alpha_{1,1}, \dots, \alpha_{n,k})$

$$\varepsilon(\underbrace{e_1, \dots, e_k, \alpha_{1,1}, \dots, \alpha_{n,k}}_{\text{unknowns}}) = \sum_{j=1}^n \varepsilon_j^2 = \sum_{j=1}^n \left\| x_j - \sum_{i=1}^k \alpha_{j,i} e_i \right\|_2^2. \quad (1)$$

по всем  $e_1, \dots, e_k$  и  $\alpha_{1,1}, \dots, \alpha_{n,k}$ .

Точное решение этой задачи сводится к нахождению собственных чисел и векторов ковариационной матрицы, размер которой совпадает с размером выборки. Ясно, что решение этой задачи неустойчиво и для больших выборок весьма проблематично. Поэтому будем использовать итерационный подход, позволяющий, хотя и не точно, но зато гарантировано получить решение этой задачи [4].

Для случая  $i=1$  задача (1) сводится к определению одной компоненты  $e_1$ , которая наилучшим образом восстанавливает все исходные данные  $\{x_1, \dots, x_n\}$

$$\varepsilon(e_1, \alpha_{1,1}, \dots, \alpha_{n,1}) = \sum_{j=1}^n \left\| x_j - \alpha_{j,1} e_1 \right\|_2^2 \rightarrow \min \quad (2)$$

по всем  $e_1$  и  $\{\alpha_{i,1}\}_{i=1}^n$  при условии  $\sum_{i=1}^n \alpha_{i,1}^2 = 1$ .

Если  $\{\tilde{\alpha}_{i,1}\}_{i=1}^n$  и  $\tilde{e}_1$  есть решение этой задачи и  $\Delta x_j = x_j - \tilde{\alpha}_{j,1} \tilde{e}_1$  — ошибка восстановления данных одной первой главной компонентой, то решая задачу

$$\sum_{j=1}^n \left\| \Delta x_j - \alpha_{j,2} e_2 \right\|_2^2 \rightarrow \min$$

по всем  $e_2$  и  $\{\alpha_{i,2}\}_{i=1}^n$  при условии  $\sum_{i=1}^n \alpha_{i,2}^2 = 1$ , получаем вторую главную компоненту  $\tilde{e}_2$  и соответствующий вектор  $\{\tilde{\alpha}_{i,2}\}_{i=1}^n$  и т.д.

При фиксированных  $\{\alpha_{i,1}\}_{i=1}^n$  задача (2) решается методом наименьших квадратов. В силу того, что функция цели представляет собой квадратичный функционал, необходимое и достаточное условия экстремума совпадают. Таким образом, решение задачи сводится к поиску решения уравнения

$$\begin{aligned} \frac{\partial}{\partial e_1} \varepsilon(e_1, \alpha_{1,1}, \dots, \alpha_{n,1}) &= -2 \sum_{j=1}^n (x_j - \alpha_{j,1} e_1) \alpha_{j,1} = \\ &= -2 \left( \sum_{j=1}^n x_j \alpha_{j,1} - \sum_{j=1}^n \alpha_{j,1}^2 e_1 \right). \end{aligned}$$

Отсюда получаем

$$e_1 = \frac{\sum_{j=1}^n x_j \alpha_{j,1}}{\sum_{j=1}^n \alpha_{j,1}^2},$$

учитывая условие нормирования единицей, то есть

$$\sum_{i=1}^n \alpha_{i,1}^2 = 1, \text{ имеем}$$

$$e_1 = \sum_{j=1}^n x_j \alpha_{j,1}.$$

Следующий шаг будем делать исходя из предположения, что в задаче (2) нам известна компонента  $e_1$  и требуется найти экстремум по  $\{\alpha_{i,1}\}_{i=1}^n$

$$\begin{aligned} \frac{\partial}{\partial \alpha_{v,1}} \varepsilon(e_1, \alpha_{1,1}, \dots, \alpha_{n,1}) &= -2(x_v - \alpha_{v,1} e_1) e_1 = \\ &= -2(\langle x_v, e_1 \rangle - \alpha_{v,1} \langle e_1, e_1 \rangle) = 0, \end{aligned}$$

то есть

$$\alpha_{v,1} = \frac{\langle x_v, e_1 \rangle}{\langle e_1, e_1 \rangle},$$

где, как обычно,  $\langle x, y \rangle$  — скалярное произведение векторов  $x$  и  $y$ .

Далее, считая найденные  $\{\alpha_{i,1}\}_{i=1}^n$  известными, повторяем весь процесс, пока не произойдет стабилизация ошибки. Полученные  $e_1$  будем считать первой главной компонентой  $\tilde{e}_1$ . Тогда  $\Delta x_j = x_j - \tilde{\alpha}_{j,1} \tilde{e}_1$  — ошибка восстановления данных одной первой главной компонентой.

Применяя этот алгоритм к ошибке восстановления  $\Delta x_j$ , находим вторую главную компоненту  $e_2$  вместе с коэффициентами  $\alpha_{j,2}$ , и т.д.

Приведем алгоритмизацию этого алгоритма.

Вначале центрируем данные, вычитая из исходных данных среднее значение и в дальнейшем считаем, что данные в среднем равны нулю.

1. Положим номер итерации  $\nu = 1$ .

2. Выбираем стартовые значения  $\{\alpha_{i,1}^\nu\}_{i=1}^n$ , например, пусть все они между собой равны, то есть  $\alpha_{i,1}^\nu = \frac{1}{\sqrt{n}}, i=1, 2, \dots, n$ .

3. Вычисляем  $e_1^\nu = \sum_{j=1}^n x_j \alpha_{j,1}^\nu$ .

4. Далее находим  $\beta_i = \frac{\langle x_i, e_1^\nu \rangle}{\langle e_1^\nu, e_1^\nu \rangle}$ , и, нормируя

единицей, получаем

$$\alpha_{i,1}^{\nu+1} = \frac{\beta_i}{\sqrt{\sum_{j=1}^n \beta_j^2}}.$$

5. Полагаем  $\nu = \nu + 1$ .

6. Проводим проверку критерия останова, в качестве этого может быть либо стабилизация коэффициентов  $\{\alpha_{i,1}^\nu\}_{i=1}^n$ , либо стабилизация главной компоненты  $e_1^\nu$ , либо проверка на заранее заданное фиксированное число итераций. Если условие окончания итерационного процесса не выполнено, то переходим к пункту 3.

Применяя итерационный метод нахождения главных компонент в модифицированном методе  $k$ -средних, и используя в качестве центра кластера главную компоненту, получаем (рис. 2)



Рис. 2. Центры полученных кластеров

Заметим, что в качестве критерия близости, используется значение соответствующего масштабирующего коэффициента.

#### Выводы

Предложенный метод кластеризации изображений лиц людей показал достаточно высокую эффективность для решения задачи схожести лиц.

#### ЛИТЕРАТУРА

- Щеголева Н.Л. Применение алгоритмов двумерного анализа главных компонент для задач распознавания изображений лиц / Н.Л.Щеголева, Г.А.Кухарев // Бизнес-информатика .– 2011 .– №№4(18) .– С.31–38 .– Режим доступа: <https://cyberleninka.ru/article/n/primenenie-algoritmov-dvumernogo-analiza-glavnyh-komponent-dlya-zadach-raspoznavaniya-izobrazheniy-lits>
- Прокошев В.Г. Проблема автоматического распознавания лиц с одним эталонным изображением / В.Г.Прокошев, М.М.Рожков, П.Шамин // Научно-технические ведомости СПбГПУ .– 2010 .– №5 .– С. 13–18 .– Режим доступа: <https://cyberleninka.ru/article/n/problema-avtomaticheskogo-raspoznavaniya-lits-s-odnim-etalonnyim-izobrazheniem>
- Лигун А.О. Комп'ютерна графіка (Обработка та стиск зображень): навч. посіб./ А.О.Лигун, О.О.Шумейко.– Біла К.О., 2010.– 114 с .– Режим доступа: [http://pzs.dstu.dp.ua/Data/U\\_CompGraph.pdf](http://pzs.dstu.dp.ua/Data/U_CompGraph.pdf)
- Шумейко А.А. Интеллектуальный анализ данных (Введение в Data Mining) / А.А.Шумейко, С.Л.Сотник .– Днепропетровск: Белая Е.А., 2012 .– 212 с .– Режим доступа: <http://pzs.dstu.dp.ua/DataMining/bibl/DataMining.pdf>

пост. 15.11.2017