

# Лекция 11

Методы кластерного анализа  
проектирование данных на плоскость, метод главных  
компонент

*Лектор – Сенько Олег Валентинович*

Курс «Математические основы теории прогнозирования»  
4-й курс, III поток

- 1 Кластерный анализ
- 2 Визуализация многомерных данных
- 3 Метод главных компонент

Важное прикладное значение имеют методы анализа данных, связанные с теорией распознавания. К их числу относятся методы кластерного анализа и методы визуализации многомерных данных. Целью методов кластерного анализа является разбиение выборок многомерных данных на группы объектов близких в смысле некоторой заданной меры сходства. Такие компактные группы называются кластерами, классами или таксонами. Методы кластерного анализа называют также методами обучения без учителя, автоматической группировки или таксономии. Методы кластерного анализа могут использоваться в качестве вспомогательных инструментов при решении задач прогнозирования или распознавания. Так с помощью кластеризации могут отбираться эталонные объекты. Однако нередко кластеризация может иметь самостоятельное значение. Можно выделить задачи кластерного анализа, для которых число кластеров задано, а также задачи, в которых число кластеров следует определить в ходе решения кластеризации.

Большинство известных алгоритмов кластеризации предполагает задание неотрицательной функцией близости  $\rho(\mathbf{x}, \mathbf{y})$  между произвольными векторами  $\mathbf{x}$  и  $\mathbf{y}$ . В качестве функций близости могут выступать евклидова метрика или метрика Хэмминга. Одним из наиболее известных методов кластеризации является **алгоритм  $k$  внутригрупповых средних**. Предположим, что у нас задана выборка многомерных векторов-объектов  $\tilde{S}_{ini} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ . Алгоритм находит такие кластеры, для объектов которых центр «своего кластера» будет ближе центра любого «чужого кластера». На начальном этапе произвольным образом выбирается начальная кластеризация  $\tilde{G}_0 = \{G_1^0, \dots, G_k^0\}$  с содержанием объектов  $(m_1^0, \dots, m_k^0)$  соответственно.

Предположим, что на шаге  $(l - 1)$  получены группы  $\{G_1^{l-1}, \dots, G_k^{l-1}\}$  с содержанием объектов  $(m_1^{l-1}, \dots, m_k^{l-1})$ . На шаге  $l$  для каждой из групп  $G_i^{l-1}$  вычисляется центр

$$\bar{\mathbf{x}}_i^{l-1} = \frac{1}{m_i^{l-1}} \sum_{\mathbf{x}_j \in G_i^{l-1}} \mathbf{x}_j$$

при  $i = 1, \dots, k$ .

Произвольный объект  $x_{j'}$  из выборки  $\tilde{S}_{ini}$  переносится в группу  $G_{i''}^{l-1}$ , если при произвольном  $i'$  из множества  $\{1, \dots, k\} \setminus i''$  выполняется неравенство

$$\rho(x_{j'}, \bar{x}_{i''}^{l-1}) < \rho(x_{j'}, \bar{x}_{i'}^{l-1}).$$

. В результате мы получаем группы  $\{G_1^l, \dots, G_k^l\}$  и переходим к шагу  $(l + 1)$ . Процесс останавливается, если на каком-то шаге оказывается, что  $\bar{x}_i^l = \bar{x}_i^{l+1}$  при  $i = 1, \dots, k$ .

Другим методом кластеризации, основанным на итерационной процедуре является алгоритм Форель, основанный на движении гипершаров фиксированного радиуса в сторону мест «сгущения» объектов. Пусть фиксировано некоторое положительное число  $R$ . Выбирается случайный вектор  $x_{j'} \in \tilde{S}_{ini}$  и гипершар  $\mathbf{R}_1$  радиуса  $R$  с центром в  $z_1 = x_{j'}$ . То есть  $\mathbf{R}_1 = \{x \mid \rho(x, z_1) < R\}$ . Полагаем  $G_1 = \tilde{S}_{ini} \cap \mathbf{R}_1$  и вычисляем центр новой сферы по формуле

$$z_2 = \frac{1}{|G_1|} \sum_{x_j \in G_1} x_j.$$

Формируем группу  $G_2$  из объектов  $\tilde{S}_{ini}$ , попавших в сферу  $R_2 = \{x \mid \rho(x, z_2) < R\}$ . Процесс заканчивается на некотором шаге  $l^*$  при выполнении условия  $G_{l^*+1} = G_{l^*}$ . Полученное множество объектов объявляется первым кластером  $G_1^f$ . Оно исключается из  $\tilde{S}_{ini}$ , а вышеописанная процедура повторяется относительно оставшейся части выборки, в результате чего формируется не пересекающаяся с  $G_1^f$  выборка  $G_2^f$ . Процесс кластеризации заканчивается на итерации  $k^*$ , на которой достигается условие

$$\tilde{S}_{ini} \setminus \bigcup_{i=1}^{k^*} G_i^f = \emptyset.$$

Таким образом формируется набор кластеров

$$G_1^f, \dots, G_{k^*}^f.$$

Полученное число кластеров зависит от выбора радиуса  $R$ , который является параметром алгоритма.

Метод иерархической группировки позволяет не только осуществить кластеризацию с заранее выбранным числом классов и выявить иерархию кластеров. На начальном этапе в качестве кластеров рассматриваются отдельные объекты выборки  $\tilde{S}_{ini}$ . Дальнейшая кластеризация производится с использованием функции близости между кластерами, которая задаётся на основе функции близости между векторными описаниями объектов. На практике используется несколько типов функций близости между кластерами  $G_{i'}$  и  $G_{i''}$ :

- минимальное расстояние между объектами из двух кластеров

$$P_{min}(G_{i'}, G_{i''}) = \min_{\mathbf{x}_\mu \in G_{i'}, \mathbf{x}_\nu \in G_{i''}} \rho(\mathbf{x}_\mu, \mathbf{x}_\nu);$$

- максимальное расстояние между объектами из двух кластеров

$$P_{max}(G_{i'}, G_{i''}) = \max_{\mathbf{x}_\mu \in G_{i'}, \mathbf{x}_\nu \in G_{i''}} \rho(\mathbf{x}_\mu, \mathbf{x}_\nu);$$

- расстояние между центрами двух кластеров

$$P_c(G_{i'}, G_{i''}) = \rho(\bar{\mathbf{x}}_{i'}, \bar{\mathbf{x}}_{i''});$$

- среднее расстояние между объектами двух классов

$$P_{av}(G_{i'}, G_{i''}) = \frac{1}{|G_{i'}|} \frac{1}{|G_{i''}|} \sum_{\mathbf{x}_\mu \in G_{i'}} \sum_{\mathbf{x}_\nu \in G_{i''}} \rho(\mathbf{x}_\mu, \mathbf{x}_\nu).$$

На втором шаге два ближайших кластера объединяются в один. Процесс объединения повторяется до выделения заранее фиксированного числа кластеров. Для остановки процесса объединения кластеров могут быть использованы дополнительные условия, задаваемые экспертом, и связанные со спецификой конкретной задачи. В этом случае число кластеров устанавливается в ходе решения.

Используются также методы кластеризации, основанные на поиске разбиений  $\tilde{S}_{ini}$ , для которых достигают максимума специальные функционалы качества.



Так качество разбиения на набор кластеров  $\tilde{G} = \{G_1, \dots, G_k\}$  может быть описано с помощью функционала внутренних дисперсий  $F_{VS}(\tilde{G})$  представляющего собой взвешенную сумму средних отклонений от центра внутри каждой из групп

$$F_{VS}(\tilde{G}) = \sum_{i=1}^k |G_i| \sum_{\mathbf{x}_\mu \in G_i} \frac{\rho(\mathbf{x}_\mu, \bar{\mathbf{x}}_i)}{|G_i|} = \sum_{i=1}^k \sum_{\mathbf{x}_\mu \in G_i} \rho(\mathbf{x}_\mu, \bar{\mathbf{x}}_i).$$

Нетрудно видеть, что “вес” каждой из групп пропорционален числу объектов в ней. Поскольку число всевозможных разбиений  $\tilde{S}_{ini}$  на  $k$  групп оценивается как  $\frac{k^m}{k!}$  полный перебор разбиений здесь заведомо исключен. Поэтому обычно применяют методы частичного перебора с использованием случайного выбора начальных разбиений и последующей локальной оптимизацией

В методах локальной оптимизации (для определенности, минимизации) строится последовательность разбиений

$$\tilde{G}_1, \tilde{G}_2, \dots, \tilde{G}_l, \dots$$

, для которых

$$F_{VS}(\tilde{G}_1) > F_{VS}(\tilde{G}_2), \dots, F_{VS}(\tilde{G}_l) > F_{VS}(\tilde{G}_{l+1}), \dots$$

а разбиение  $\tilde{G}_{l+1}$  вычисляется непосредственно по предшествующему разбиению  $\tilde{G}_l = \{G_1^l, \dots, G_k^l\}$  путем его «локального» изменения – переноса некоторого объектов из одного кластера в другой. Ищется такой объект  $x_{J(l)}$ , при переносе которого из кластера  $G_\mu^l$ , содержащего  $x_{J(l)}$  в разбиении  $\tilde{G}_l$ , в некоторый кластер  $G_\nu^l$  уменьшение функционала  $F_{VS}$  максимально среди всевозможных переносов такого рода. В результате разбиение отличается от только составом кластеров с номерами  $G_\mu^l$  и  $G_\nu^l$ . Процесс завершается, когда никакой последующий перенос не уменьшает функционал или достигнуто указанное пользователем максимальное число итераций

Пусть в результате применения разнообразных методов кластеризации получено множество различных решений для одних и тех же данных. При отсутствии внешнего критерия, выбор одного решения из данного множества кластеризаций может быть не ясен. Поэтому представляет интерес применение методов обработки полученных множеств кластеризаций с целью построения коллективных решений, более предпочтительных и обоснованных, чем полученные отдельными алгоритмами кластеризации. Кластеризацию выборки  $\tilde{S}_{ini}$ , включающую кластеры  $\{G_1, \dots, G_k\}$  можно описать с помощью информационной матрицы  $\|\alpha_{ji}\|_{m \times k}$ , где  $\alpha_{ji} = 1$ , если  $x_j \in G_i$ , и  $\alpha_{ji} = 0$  в противном случае. Наличие нескольких единиц в одной строке соответствует принадлежности объекта сразу нескольким кластерам. Нулевая строка означает отказ от кластеризации соответствующего объекта.

**Определение 1.** Информационные матрицы  $I = \|\alpha_{ji}\|_{m \times k}$  и  $I = \|\alpha'_{ji}\|_{m \times k}$  называются эквивалентными, если они равны с точностью до перестановки столбцов.

Таким образом произвольная информационная матрица  $I = \|\alpha_{ji}\|_{m \times k}$  определяет класс всех эквивалентных ей матриц  $\tilde{K}(I)$ .

**Определение 2.** Алгоритмом кластеризации  $A^c$  называется алгоритм, переводящий выборку в класс эквивалентности  $\tilde{K}(I)$  некоторой информационной матрицы  $I$ .

Иными словами  $A^c = \tilde{K}(\|\alpha_{ji}\|_{m \times k})$ . Данное определение отражает возможности произвола в обозначении полученных алгоритмом кластеров. Пусть существует  $r$  кластеризаций выборки  $\tilde{S}_{ini}$  алгоритмами  $A_1^c, \dots, A_r^c$  на  $k$  кластеров. Задача построения оптимальной коллективной кластеризации состоит в вычислении по множеству из  $r$  исходных кластеризаций, задающих классы эквивалентности

$$\tilde{K}(\|\alpha_{ji}^1\|_{m \times k}), \dots, \tilde{K}(\|\alpha_{ji}^r\|_{m \times k})$$

некоторого нового коллективного решения  $\tilde{K}(\|\hat{\alpha}_{ji}\|_{m \times k})$ , где  $\hat{\alpha}_{ji} \in [0, 1]$ .

Оператор  $\mathbf{B}(I_1, \dots, I_r) = \|\beta_{ji}\|_{m \times k}$ , где  $\beta_{ji} \in \{1, \dots, \}$ , называется сумматором, если

$$\beta_{ji} = \sum_{t=1}^r \alpha_{ji}^t.$$

Матрицу, полученную в результате применения сумматора к некоторому набору информационных матриц, будем называть матрицей оценок. Оператор  $\mathbf{C}$  называется решающим правилом, если

$$\mathbf{C}(\|\beta_{ji}\|_{m \times k}) = \|\alpha_{ji}^s\|_{m \times k},$$

где при произвольном  $j \in \{1, \dots, m\}$   $\alpha_{ji}^s = 1$ , если  $\beta_{ji} > \beta_{jt}$  при  $t = \{1, \dots, k\} \setminus \{i\}$ , и  $\alpha_{ji}^s = 0$  в противном случае.

**Определение 3.** Комитетным синтезом информационной матрицы  $\|\alpha_{ji}^s\|_{m \times k}$  по множеству исходных кластеризаций, задаваемых набором информационных матриц  $\tilde{I} = \{\tilde{I}_1, \dots, \tilde{I}_r\}$ , называется последовательное применение к  $\tilde{I}$  сумматора  $\mathbf{B}$  и решающего правила  $\mathbf{C}$ .

Для оценивания коллективного решения вводится понятие контрастных матриц оценок, соответствующих случаям, когда все исходные решения задач классификации оказались одинаковыми. Очевидно, что для произвольной контрастной матрицы  $\|\beta_{ji}^c\|_{m \times k}$  выполняется условие  $\beta_{ji}^c \in \{0, r\}$ . Пусть  $\tilde{\mathbf{B}}^c = \|\beta_{ji}^c\|_{m \times k}$  - множество всевозможных контрастных матриц. Качеством произвольной матрицы  $B = \|\beta_{ji}\|_{m \times k}$ , вычисленной сумматором, определяется как минимальное расстояние до матриц из множества  $\tilde{\mathbf{B}}^c$ .

$$\Phi(B) = \min_{\tilde{\mathbf{B}}^c} \sum_{j=1}^m \sum_{i=1}^k |\beta_{ji} - \beta_{ji}^c|$$

Набор информационных матриц  $\{I'_1, \dots, I'_r\}$  назовём эквивалентным  $\{I_1, \dots, I_r\}$ , если

$$I'_1 = \tilde{K}(I_1), \dots, I'_r = \tilde{K}(I_r).$$

Задача оптимального коллективного синтеза сводится к поиску эквивалентного  $\{I_1, \dots, I_r\}$  набора  $\{I_1^m, \dots, I_r^m\}$ , для которого в результате применения сумматора получается матрица  $\mathbf{B}_m$  с минимальным значением  $\Phi$  среди всевозможных матриц, вычисляемых сумматором по наборам, эквивалентным  $\{I_1, \dots, I_r\}$ .

При решении задач распознавания, классификации и анализа данных важное значение имеет наличие средств визуализации многомерных данных, позволяющих наглядно получать представление о конфигурации классов, кластеров и расположении отдельных объектов. Предполагаем опять, что у нас задана выборка  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , состоящая из элементов пространства  $\mathbb{R}^n$ . Требуется найти отображение этого набора точек на плоскость  $\mathbb{R}^2$  так, чтобы метрические соотношения между образами точек на плоскости максимально соответствовали бы метрическим соотношениям между ними в исходном признаковом пространстве.

- "Близкие" точки в исходном пространстве должны быть по-возможности "близкими" на плоскости.
- Соответственно "удалённые" точки в исходном пространстве должны быть по-возможности "удалёнными" на плоскости.



Пусть точки  $\{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  являются образами точек  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  на плоскости  $\mathbb{R}^2$ . Пусть  $\delta_{ij}$  - расстояние между векторами  $\mathbf{x}_i$  и  $\mathbf{x}_j$ ,  $d_{ij}$  - расстояние между векторами  $\mathbf{y}_i$  и  $\mathbf{y}_j$ . Ищется такое отображение, для которого сумма различий расстояний между точками будет минимальна

$$\mathcal{J}(\tilde{\mathbf{y}}) = \sum_{i=1}^m \sum_{j=1}^m (\delta_{ij} - d_{ij})^2 \rightarrow \min,$$

где  $\tilde{\mathbf{y}} = \{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  является вектором размерности  $2m$ , содержащим последовательно координаты двумерных образов многомерных объектов.

Минимизация функционала  $\mathcal{J}(\tilde{\mathbf{y}})$  проводится с помощью стандартной процедуры градиентного спуска.

Новое значение вектора координат двумерных образов  $\tilde{\mathbf{y}}^{l+1}$  на шаге  $l$  вычисляется по значению вектора координат двумерных образов  $\tilde{\mathbf{y}}^l$ , вычисленному на предыдущем шаге, по формуле

$$\tilde{\mathbf{y}}^{l+1} = \tilde{\mathbf{y}}^l + \kappa \star \mathbf{grad}[\mathbb{J}(\tilde{\mathbf{y}}^l)]$$

где  $\mathbf{grad}[\mathbb{J}(\tilde{\mathbf{y}})]$  - градиент функционала  $\mathbb{J}(\tilde{\mathbf{y}})$  в точке  $\tilde{\mathbf{y}}$ ,  $\kappa$  - шаг градиентного спуска. В качестве начальной конфигурации может использоваться проекция точек  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  на  $n$  плоскость, соответствующую некоторой паре признаков. Пример проекции на плоскость из пространства размерности 26, полученной описанным методом, приведён на рисунке 1.

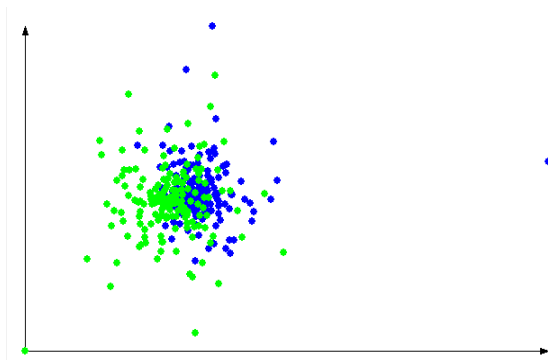


Рис.1. Точкам зелёного и синего цвета соответствуют описания двух классов объектов. Таким образом видно, что точки двух классов являются достаточно разделёнными.

## Методы преобразования признакового пространства. Метод главных компонент

Описанный метод многомерной визуализации фактически является методом нелинейного преобразования исходного признакового пространства. Вместе с тем существует эффективный метод линейной трансформации признакового пространства, позволяющий получить существенную информацию о структуре данных, а также получить новые признаки, удобные и эффективные при решении задач прогнозирования или распознавания. Данный метод называется Методом главных компонент (Principal component analysis), а также преобразованием Карунена-Лоэв (Karhunen–Loeve transform). Метод главных компонент основан на переходе от исходного множества вообще говоря коррелированных переменных  $X_1, \dots, X_n$  к новому набору переменных  $Z_1, \dots, Z_n$  с нулевыми коэффициентами ковариации между ними. То есть  $cov(Z_i, Z_j) = 0$  при  $i \neq j$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, n$ . Переход к некоррелированным переменным может быть осуществлён с помощью линейного преобразования

Данное преобразование задаётся матрицей вещественных коэффициентов  $W = \|w_{ij}\|_{n \times n}$ . Предположим, что у нас имеется исходная выборка  $\tilde{S}_{ini} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , которая может быть представлена в

виде  $\mathbf{X} = \begin{pmatrix} x_{11} & \dots & x_{1n} \\ \dots & \dots & \dots \\ x_{j1} & \dots & x_{jn} \\ \dots & \dots & \dots \\ x_{m1} & \dots & x_{mn} \end{pmatrix}$ . Далее будем предполагать, что

признаки в матрице являются центрированными, то есть  $\sum_{j=1}^m x_{ji} = 0$  при  $i = 1, \dots, n$ . Переход к центрированной выборке может быть всегда легко осуществлён с помощью простого линейного преобразования.

Отметим, что  $\mathbf{X}^t \mathbf{X} = (m - 1) \Sigma$ , где  $\Sigma = \|\hat{\sigma}_{ij}\|_{n \times n}$  - выборочная ковариационная матрица, элементы которой вычисляются по формуле

$$\hat{\sigma}_{ij} = \frac{1}{m - 1} \sum_{k=1}^m x_{ki} x_{kj}.$$

Подобная запись для коэффициентов ковариации возможна из за центрированности данных.

Предположим, что  $\mathbf{Z} = \begin{pmatrix} z_{11} & \dots & z_{1n} \\ \dots & \dots & \dots \\ z_{j1} & \dots & z_{jn} \\ \dots & \dots & \dots \\ z_{m1} & \dots & z_{mn} \end{pmatrix}$  . - матрица значения

признаков  $Z_1, \dots, Z_n$ ,

полученных с помощью линейного преобразования. Очевидно, что  $\mathbf{Z} = \mathbf{XW}$  . Отметим, что матрица  $\mathbf{Z}$  также является центрированной, поскольку линейное преобразование не приводит к утрате свойства центрированности. Вследствие требования отсутствия корреляции между переменными  $Z_1, \dots, Z_n$  матрица  $\mathbf{Z}^t\mathbf{Z}$  , являющаяся ковариационной матрицей для  $Z$ -переменных, является диагональной.

При этом на диагонали в строке  $i$  находится величина  $(m - 1)\delta_i$ , где

$$\delta_i = \frac{1}{m-1} \sum_{j=1}^m z_{ji}^2.$$

Однако

$$\mathbf{z}^t \mathbf{Z} = \mathbf{w}^t \mathbf{X}^t \mathbf{X} \mathbf{W} = (m - 1) \mathbf{w}^t \mathbf{\Sigma} \mathbf{W}.$$

Таким образом  $\mathbf{W}^t \mathbf{\Sigma} \mathbf{W}$  является диагональной матрицей. Однако из теории матриц известно, что диагонализация симметрической вещественной матрицы  $\mathbf{\Sigma}$  может быть осуществлена с помощью квадратной матрицы  $\mathbf{E}$ , столбцами которой являются ортонормированные собственные вектора  $\mathbf{\Sigma}$ . Иными словами справедливо равенство  $\mathbf{E}^t \mathbf{\Sigma} \mathbf{E} = \mathbf{V}$ , где  $\mathbf{V}$  - диагональная матрица, на диагонали которой лежат собственные значения  $\mathbf{\Sigma}$ . Отсюда можно сделать вывод, что использование в качестве  $\mathbf{W}$  матрицы  $\mathbf{E}$  позволяет осуществлять переход к некоррелированным переменным  $Z_1, \dots, Z_n$ . При этом дисперсиям переменных  $Z_1, \dots, Z_n$  будут соответствовать собственные значения матрицы  $\mathbf{\Sigma}$ .

Столбцы матрицы  $\mathbf{E}$  могут быть интерпретированы как новый ортонормированный базис в пространстве исходных переменных. При этом переменные  $Z_1, \dots, Z_n$  являются проекциями на оси нового базиса. Отметим, что ось, соответствующая максимальному собственному значению, является одновременно тем направлением в исходном пространстве, для которого дисперсия проекций на него векторов обучающей выборки максимальна. Следует также отметить, что преобразование, задаваемое матрицей  $\mathbf{E}$  является унитарным преобразованием, не изменяющим длины векторов. Вследствие этого полная выборочная дисперсия остаётся после преобразования  $\mathbf{E}$  неизменной и равной сумме собственных значений ковариационной матрицы  $\Sigma$ . То есть для полной дисперсии  $D(\tilde{S}_{ini})$  выборки  $\tilde{S}_{ini}$  справедливо разложение

$$D(\tilde{S}_{ini}) = \frac{1}{m-1} \sum_{j=1}^m (\mathbf{x}_j - \bar{\mathbf{x}})^2 = \frac{1}{m-1} \sum_{j=1}^m (\mathbf{z}_j - \bar{\mathbf{z}})^2 \quad (1)$$



Однако из некоррелированности переменных  $Z_1, \dots, Z_n$  следует, что

$$\frac{1}{m-1} \sum_{j=1}^m (z_j - \bar{z})^2 = \sum_{i=1}^n \Lambda_i$$

где  $\Lambda_1, \dots, \Lambda_n$  являются неотрицательными собственными значениями  $\Sigma$ . В формуле (1) используются обозначения  $\bar{\mathbf{x}} = \frac{1}{m} \sum_{j=1}^m \mathbf{x}_j$ ,  $\bar{z} = \frac{1}{m} \sum_{j=1}^m z_j$ . Принимая во внимание условие центрированности переменных  $X$ , а значит и переменных  $Z$ , получаем

$$\frac{1}{m-1} \sum_{j=1}^m (z_j)^2 = \sum_{i=1}^n \Lambda_i$$

Полученные в результате преобразования  $\mathbf{E}$  переменные называют главными компонентами. Главные компоненты ранжируются в зависимости от величин соответствующих собственных значений.

Переменная, соответствующая максимальному собственному значению  $\Lambda_1$  и задаваемая соответствующим собственному вектором  $e_1$  называется первой главной компонентой. Она обладает максимальной дисперсией, равной  $\Lambda_1$ . Переменная, соответствующая второму по величине собственному значению  $\Lambda_2$  и задаваемая соответствующим собственному вектором  $e_2$  называется второй главной компонентой и т.д.

Сумму  $\sum_{i=1}^k \Lambda_i$  принято называть объяснённой дисперсией, а сумму  $\sum_{i=k+1}^n \Lambda_i$  - остаточной дисперсией для  $k$  первых главных компонент. Покажем справедливость равенства

$$\sum_{i=k+1}^n \Lambda_i = \frac{1}{m-1} \sum_{j=1}^m \left[ \mathbf{x}_j - \sum_{i=1}^k e_i (e_i, \mathbf{x}_j) \right]^2.$$

Действительно,  $\mathbf{x}_j = \sum_{i=1}^n \mathbf{e}_i(e_i, \mathbf{x}_j)$ . По определению переменных  $Z$  справедливо равенство  $\mathbf{x}_j = \sum_{i=1}^k \mathbf{e}_i z_{ji}$ . Очевидно

$$\mathbf{x}_j - \sum_{i=1}^k \mathbf{e}_i(e_i, \mathbf{x}_j) = \mathbf{x}_j - \sum_{i=1}^k \mathbf{e}_i z_{ji} = \sum_{i=k+1}^n \mathbf{e}_i z_{ji}$$

То есть

$$\begin{aligned} \frac{1}{m-1} \sum_{j=1}^m [\mathbf{x}_j - \sum_{i=1}^k \mathbf{e}_i(e_i, \mathbf{x}_j)]^2 &= \frac{1}{m-1} \sum_{j=1}^m \left( \sum_{i=k+1}^n \mathbf{e}_i z_{ji} \right)^2 = \\ &= \frac{1}{m-1} \sum_{j=1}^m \sum_{i=k+1}^n z_{ji}^2. \end{aligned}$$

Однако последняя сумма является суммой дисперсий переменных  $Z_{k+1}, \dots, Z_n$ , равной остаточной дисперсии  $\sum_{i=k+1}^n \Lambda_i$ .

Иными словами остаточная дисперсия равна среднему квадрату расстоянию векторов до линейного замыкания  $k$  собственных векторов, соответствующих  $k$  первым главным компонентам. Таким образом низкая величина остаточной дисперсии соответствует хорошей аппроксимации данных этим линейным замыканием. На рисунке 2 представлен пример использования метода главных компонент для задачи с исходным признаковым пространством размерности 2. Собственный вектор, соответствующей первой главной компоненте направлен вдоль прямой  $(c_1, c_2)$ . Собственный вектор, соответствующей второй главной компоненте, перпендикулярен  $(c_1, c_2)$ . Синими кружками отмечены точки, соответствующие данным. Из рисунка можно видеть значительную дисперсию проекций объектов выборки на направление первой главной компоненты. Следует отметить небольшую величину расстояний объектов обучающей выборки до прямой  $(c_1, c_2)$ .

Отметим, что изначально метод главных компонент рассматривался К.Пирсоном как метод поиска линейных многообразий в признаковом пространстве, наилучшим образом аппроксимирующих данные. Метод главных компонент имеет широкий спектр применений, включая формирование нового признакового пространства, снижение размерности, визуализацию данных. Главные компоненты могут быть эффективно использованы при решении задач распознавания, регрессионного анализа. Метод очень широко используется в различных областях, включая биоинформатику, экономику, гуманитарные науки.

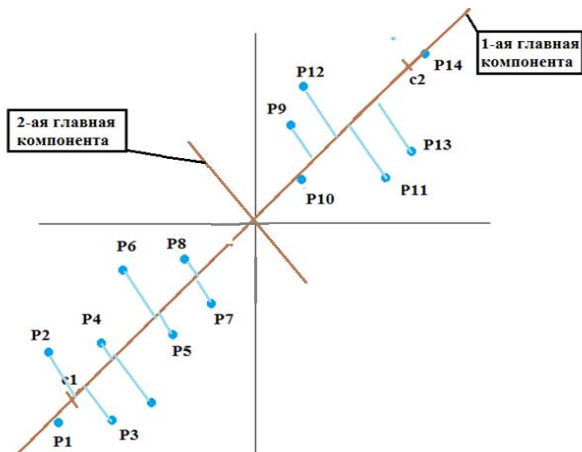


Рис. 2. пример использования метода главных компонент для задачи с исходным признаковым пространством размерности 2. Синими кружками отмечены точки, соответствующие данным.