

# Метод главных компонент и его применение

Работу выполнил  
студент 4 курса 7 группы  
Стец Вадим

# Содержание метода главных компонент

Метод главных компонент в настоящее время представляет эффективный аппарат комплексного анализа геоданных, и его программное обеспечение является составной частью многих компьютерных технологий по обработке геофизической информации.

Математической моделью метода служит, как и для корреляционно-регрессионного анализа, система случайных величин  $x_1, x_2, \dots, x_N$ . При этом каждая случайная величина обычно содержит  $n$  – наблюдений, т.е. исходный массив геоданных представлен матрицей  $X$  размерностью  $n * N$ :

$$X = \begin{vmatrix} x_{11} & x_{21} & \dots & x_{N1} \\ x_{12} & x_{22} & \dots & x_{N2} \\ \dots & \dots & \dots & \dots \\ x_{1n} & x_{2n} & \dots & x_{Nn} \end{vmatrix}$$

В качестве числа  $N$  может быть число профилей съемки, и тогда с помощью главных компонент решается задача оценки регионального тренда (региональной составляющей). Если число  $N$  представлено совокупностью различных методов и атрибутов, то путем метода главных компонент решается задача комплексного анализа по разделению исследуемой территории на классы.

Значение  $n$  определяет число точек наблюдений по отдельным профилям съемки. Суть метода главных компонент состоит в переходе от системы случайных величин  $x_1, \dots, x_N$  к новой системе случайных величин  $y_1, \dots, y_N$ , ориентируясь на поведение дисперсий  $y_i$ . При этом главная компонента определяется как линейная комбинация исходных случайных величин  $x_i$ ,

$$y_j = \sum_{i=1}^N a_{ij} x_i; \quad j = 1, \dots, N \quad (1.1)$$

причем первая главная компонента  $y_1$  :  $y_1 = \sum_{i=1}^N a_{i1} x_i$ , обладает максимальной дисперсией среди всех

возможных линейных комбинаций вида (1.1).

Величины  $a_{ij}$  являются коэффициентами перехода от одной системы случайных величин  $x_i$  к другой системе случайных величин  $y_i$ . Дисперсии линейных комбинаций  $y_i$  располагаются в убывающем порядке, т.е.:  $\sigma^2(y_1) > \sigma^2(y_2) > \dots > \sigma^2(y_N)$

Переход от системы величин  $x_i$  к системе  $y_j$  сопровождается нормировкой коэффициентов  $a_{ij}$  в виде:

$$\sum_{i=1}^N a_{ij}^2 = 1$$

Математически метод главных компонент сводится либо к вычислению ковариационной матрицы  $B$  системы случайных величин  $x_i$ , если все  $x_i$  измерены в одних и тех же единицах, либо к вычислению корреляционной матрицы  $R$ , если случайные величины  $x_i$  измерены в разных физических единицах.

**Первый случай** соответствует измерениям одного и того же поля по  $N$ -профилям съемки, **второй случай** соответствует измерениям  $N$ -разных полей и (или) их атрибутов. Далее для матрицы  $B$  или матрицы  $R$  находят их собственные значения  $\lambda_i$  и соответствующие этим собственным значениям собственные векторы, которыми являются коэффициенты перехода  $a_{ij}$ .

Обычно ограничиваются вычислением первых двух-трех главных компонент, поскольку в этих компонентах сосредоточена основная энергия исходных данных.

Физическое истолкование главных компонент является весьма неоднозначным. Однако первая главная компонента практически всегда имеет однозначное истолкование, поскольку ее дисперсия отражает основную энергию поля при обработке данных по площади или энергию нескольких полей и (или) их атрибутов при комплексном анализе данных.

# Метод главных компонент при оценке региональной составляющей поля

Выделение регионального тренда является распространенной процедурой обработки практически для всех методов геофизики. Однако при решении этой задачи приходится задавать те или иные параметры. Так, при осреднении поля в скользящем окне надо задать размеры окна, при пересчете поля на высоту надо задать высоту пересчета, при оценке региональной составляющей путем регрессии надо задать степень полинома регрессии и т.д.

Метод главных компонент не требует задания той или иной априорной информации. Единственное предположение, при котором происходит применение метода главных компонент, состоит в том, что региональная составляющая обладает наибольшей дисперсией по сравнению с локальными составляющими, что обычно на практике выполняется.

Алгоритм оценки региональной составляющей на основе метода главных компонент сводится к реализации следующих процедур:

1. Вычисление средних значений поля по каждому профилю  $\bar{x}_i = \frac{1}{n} \sum_{k=1}^n x_{ki}$ , где  $n$  – число точек

наблюдений,  $i=1, \dots, N$ ,  $N$  – число профилей, и ковариаций данных различных пар профилей :

$$b_{ij} = \frac{1}{n} \sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j), \quad i, j=1, \dots, N.$$

Поскольку система  $x_1, \dots, x_N$  представлена измерениями одного и того же поля, нет необходимости рассчитывать коэффициенты корреляции, которые используются при комплексном анализе данных, измеренных в разных физических единицах.

1) Вычисление средних значений поля по каждому профилю  
число точек

$$\bar{x}_i = \frac{1}{n} \sum_{k=1}^n x_{ki}, \text{ где } n -$$

наблюдений,  $i=1, \dots, N$ ,  $N$  – число профилей, и ковариаций данных различных пар профилей :

$$b_{ij} = \frac{1}{n} \sum_{k=1}^n (x_{ki} - \bar{x}_i)(x_{kj} - \bar{x}_j), \quad i, j=1, \dots, N.$$

Поскольку система  $x_1, \dots, x_N$  представлена измерениями одного и того же поля, нет необходимости рассчитывать коэффициенты корреляции, которые используются при комплексном анализе данных, измеренных в разных физических единицах.

2) Составление ковариационной матрицы исходных данных по их коэффициентам ковариаций  $b_{ij}$  :

$$B = \begin{vmatrix} b_{11} & b_{12} & \dots & b_{1N} \\ b_{21} & b_{22} & \dots & b_{2N} \\ \dots & \dots & \dots & \dots \\ b_{N1} & b_{N2} & \dots & b_{NN} \end{vmatrix}$$

Матрица  $B$  симметрична относительно главной диагонали, т.е.  $b_{ij} = b_{ji}$  , а по диагонали расположены дисперсии значений поля каждого профиля.

3) Нахождение максимального собственного значения  $\lambda_{\max}$  из уравнения:

$$|B - \lambda_{\max} I| = 0 \quad \text{или} \quad |B - \lambda_{\max} I| = \begin{vmatrix} b_{11} - \lambda_{\max} & b_{12} & \dots & b_{1N} \\ b_{21} & b_{22} - \lambda_{\max} & \dots & b_{2N} \\ \dots & \dots & \dots & \dots \\ b_{N1} & b_{N2} & \dots & b_{NN} - \lambda_{\max} \end{vmatrix} = 0$$

т.е. после раскрытия определителя из этого уравнения достаточно найти его корень с максимальным значением  $\lambda_{\max}$ .

- 4) Вычисление значений собственного вектора матрицы  $(B - \lambda_{\max} * I)$ , соответствующего максимальному собственному значению  $\lambda_{\max}$  из системы линейных уравнений

$$(B - \lambda_{\max} I) \bar{a}_1 = 0 \text{ или } \begin{bmatrix} b_{11} - \lambda_{\max} & b_{12} & \dots & b_{1N} \\ b_{21} & b_{22} - \lambda_{\max} & \dots & b_{2N} \\ \dots & \dots & \dots & \dots \\ b_{N1} & b_{N2} & \dots & b_{NN} - \lambda_{\max} \end{bmatrix} \begin{bmatrix} a_{11} \\ a_{21} \\ \dots \\ a_{N1} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

Значения собственного вектора  $\bar{a}_1$  определяются с учетом нормировки  $\sum a_{i1}^2 = 1$

Физический смысл этой нормировки состоит в том, чтобы преобразованные данные, т.е. значения региональной составляющей, не отличались бы по масштабу от исходных значений поля, а физический смысл значений  $a_{i1}$  заключается в определении весовых коэффициентов для каждого профиля.



5) Нахождение значений первой главной компоненты  $y_1 = \sum_{i=1}^N a_{i1} x_i$ , то есть

$$y_{1k} = (a_{11}, a_{21}, \dots, a_{N1}) \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{N1} & x_{N2} & \dots & x_{Nn} \end{bmatrix} = \begin{bmatrix} y_{11} \\ y_{21} \\ \cdot \\ \cdot \\ y_{1n} \end{bmatrix}$$

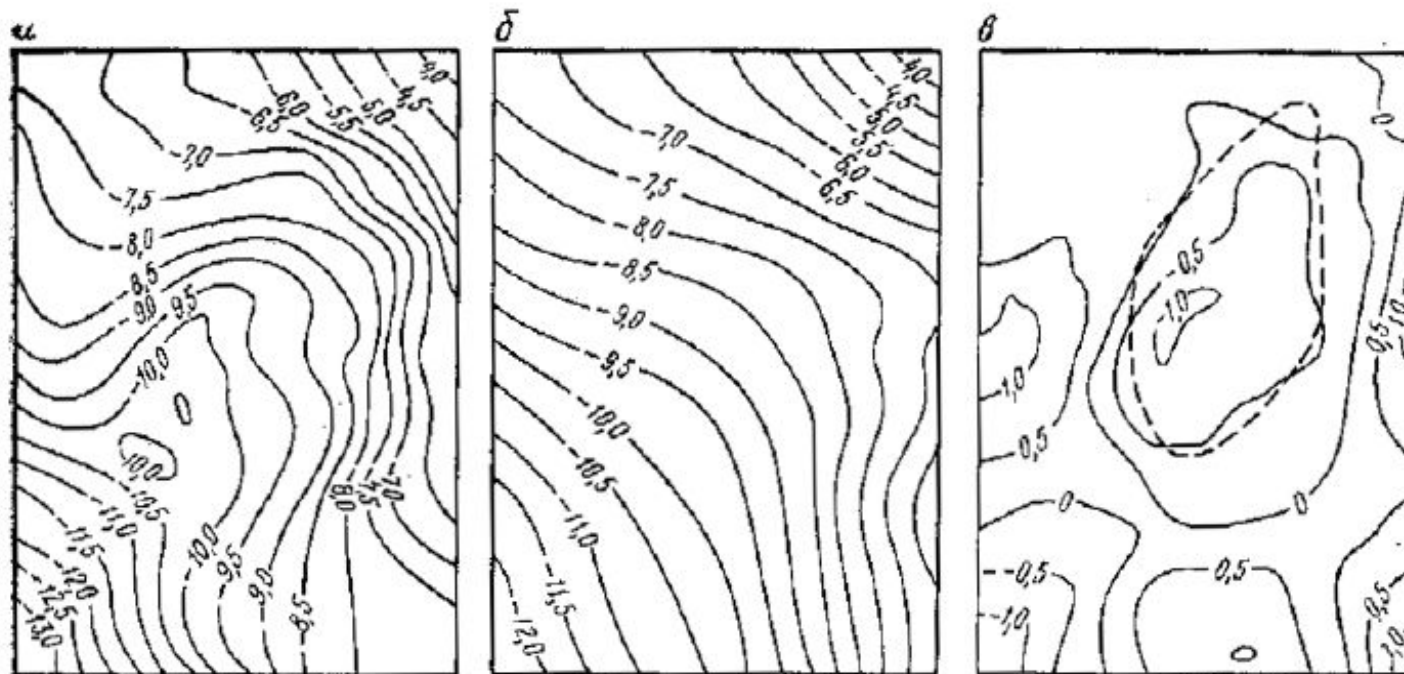
Физический смысл значений первой главной компоненты  $y_{1k} (k_i = 1, \dots, n)$  состоит в том, что они определяют весовые коэффициенты для каждого пикета исходных данных, аналогично тому, как значения  $a_{i1} (i = 1, \dots, N)$  определяют весовые коэффициенты для каждого профиля сьс<sub>СМННН</sub>.

б) Оценка региональной составляющей исходного поля характеризующейся  $y_{1i}$  наибольшей дисперсией. Эта оценка равна произведению вектора-столбца на вектор-строку  $a_{i1}$  с добавлением к каждому элементу образующейся матрицы среднего значения поля по профилю  $\bar{x}_i$ , то есть


$$x_{ki}^{pez} = \begin{bmatrix} y_{11} \\ y_{21} \\ \cdot \\ \cdot \\ y_{1n} \end{bmatrix} (a_{11}, a_{21}, \dots, a_{N1}) + \bar{x}_1 = \begin{bmatrix} y_{11}a_{11} + \bar{x}_1 & y_{11}a_{21} + \bar{x}_2 & \dots & y_{11}a_{N1} + \bar{x}_N \\ y_{21}a_{11} + \bar{x}_1 & y_{12}a_{21} + \bar{x}_2 & \dots & y_{12}a_{N1} + \bar{x}_N \\ \dots & \dots & \dots & \dots \\ y_{N1}a_{11} + \bar{x}_1 & y_{1n}a_{21} + \bar{x}_2 & \dots & y_{1n}a_{N1} + \bar{x}_N \end{bmatrix}$$

Поскольку значения  $x_{ki}^{pez}$  представляют оценку региональной составляющей, то разность  $x_{ki}^{ocm} = x_{ki} - x_{ki}^{pez}$  оценивает поле локальных составляющих.

Эффективность метода главных компонент иллюстрируется на (рис.3.1) , на котором приведены исходное поле силы тяжести (а), оценка региональной составляющей (б) и локальная составляющая (в). Пунктиром на рисунке показана область рудного объекта.



**Рис.3.1. Исходное поле силы тяжести (а), оценка региональной составляющей (б) и локальной составляющей (в) с использованием метода главных компонент.**



Метод главных компонент эффективен при обработке данных на достаточно ограниченных площадях, поскольку не учитывается изменение корреляционных свойств, т.е. структуры корреляционных матриц, по площади.

Следует отметить эффективность применения метода главных компонент при решении задач интерполяции. При этом задача интерполяции физического поля, представленного в виде функции двух переменных  $x$  и  $y$ , сводится к интерполяции функций, зависящих от одного аргумента.

Однако чаще всего метод главных компонент используется при решении задач комплексного анализа данных.

Ну всё, чё