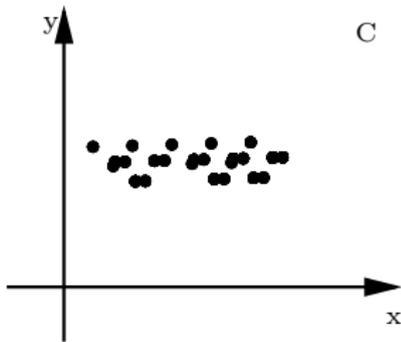
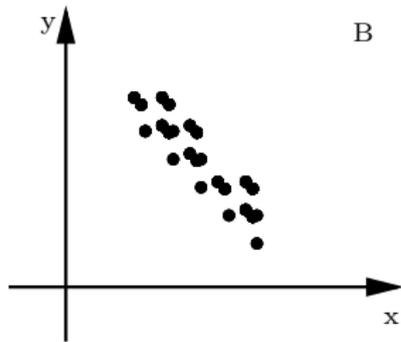
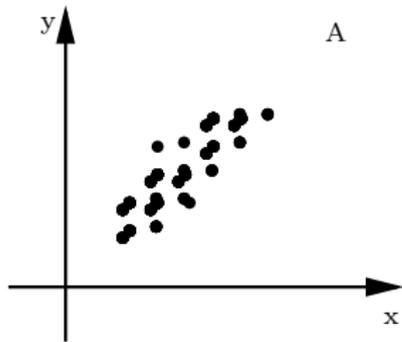
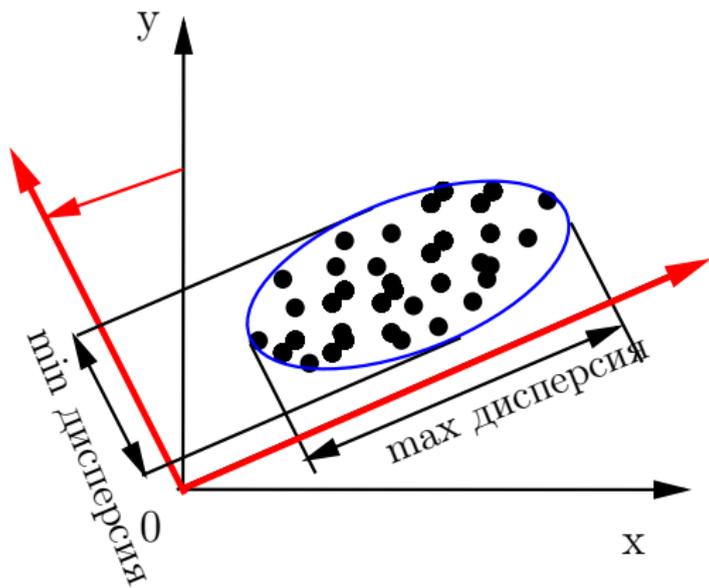
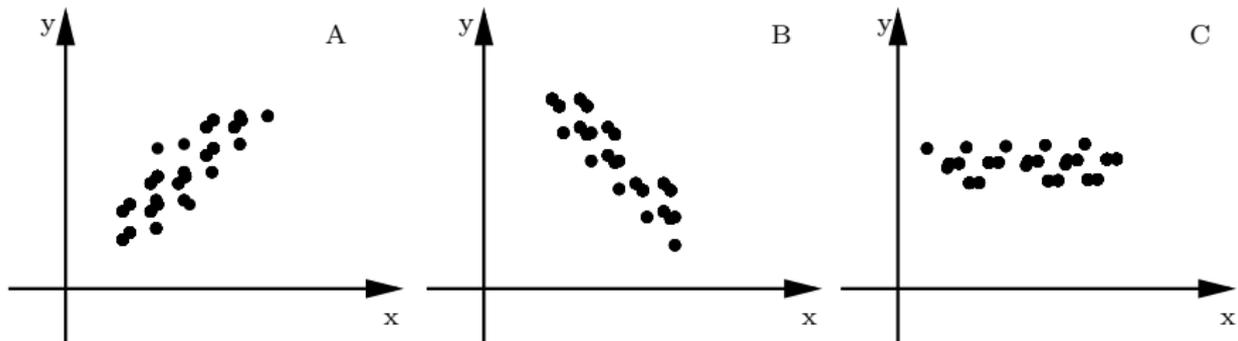


Метод главных компонент

М.И.Зуев, Ю.Л.Калиновский, О.И.Стрельцова

IT SCHOOL JINR
7 - 11 октября 2024





Целями PCA являются:

1. извлечь наиболее важную информацию из таблицы данных;
2. сократить размер набора данных, сохранив только эту важную информацию;
3. упростить описание набора данных;
4. Проанализировать структуру наблюдений и переменных.
5. Сжать данные, уменьшив число измерений, без значительной потери информации.

Математические основы

Обсудим статистику, которая рассматривает измерения распределения, как распределяются данные, а также матричную алгебру, вычисляя собственные векторы и собственные значения.

Дисперсия (variance)

- X - случайная величина
- $D[X] = M[X - M[X]]^2$
- $\text{Var}(X) = E[(X - E(X))^2]$
- Для вещественных значений: $D[X] = M[X^2] - (M[X])^2$
- Среднеквадратическое отклонение: $\sigma_X = \sqrt{D[X]}$

Ковариация (covariance)

- X, Y - случайные величины

$$\text{cov}(X, Y) = M[X - M[X]]M[Y - M[Y]]$$

$$\text{cov}(X, Y) = E[X - E[X]]E[Y - E[Y]]$$

- Для выборки $X_{(n)}, Y_{(n)}$

$$\text{cov}(X_{(n)}, Y_{(n)}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}),$$

где \bar{x}, \bar{y} - средние значения выборок

Ковариационная матрица

- Пусть X, Y - случайные векторы размерности m и n соответственно

$$C = \text{cov}(X, Y) = E[(X - EX)(Y - EY)^T]$$

$$C = [c_{ij}]$$

$$c_{ij} = \text{cov}(X_i, Y_j) = E[(X_i - EX_i)(Y_j - EY_j)],$$

$i = 1, 2, \dots, n, j = 1, 2, \dots, m$

- Если $X = Y$
 C - матрица ковариации вектора X
- Для выборок

$$c_{ij} = \frac{1}{m-1} \sum_{k=1}^m (x_{ki} - \bar{X}_i)(x_{kj} - \bar{X}_j),$$

где \bar{X}_i и \bar{X}_j - средние соответствующих компонентов векторов

Корреляция (correlation)

- Для выборки $X_{(n)}, Y_{(n)}$
- Коэффициент корреляции Пирсона

$$r_{XY} = \frac{(X, Y)}{\sigma_X \sigma_Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$r_{XY} \in [-1, +1]$$

- Мера линейной зависимости
 - $\|r_{XY}\| = 1, \implies x, y$ - линейно зависимы
 - $r_{XY} = 0, \implies x, y$ - линейно независимы

Корреляция и коэффициент корреляции

Неравенство Коши-Буняковского

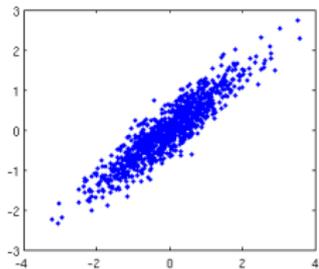
$$\text{cov}^2(X, Y) \leq D(X)D(Y)$$

Линейный коэффициент корреляции

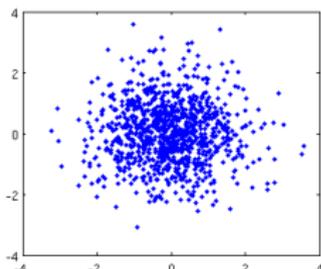
$$r_{XY} = \frac{\text{cov}_{XY}}{\sigma_x \sigma_y} = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{\sqrt{\sum(X - \bar{X})^2 \sum(Y - \bar{Y})^2}}$$

$$-1 \leq r_{XY} \leq 1$$

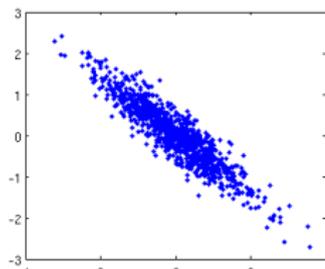
**Положительная
корреляция**



**Отсутствие
корреляции**



**Отрицательная
корреляция**



Матричная алгебра

Множество \mathbb{R}^2 имеет геометрическую интерпретацию как евклидова плоскость, в которой вектор $\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}$ в \mathbb{R}^2 представляет собой точку с координатами (a_1, a_2) на плоскости (рис. 1).

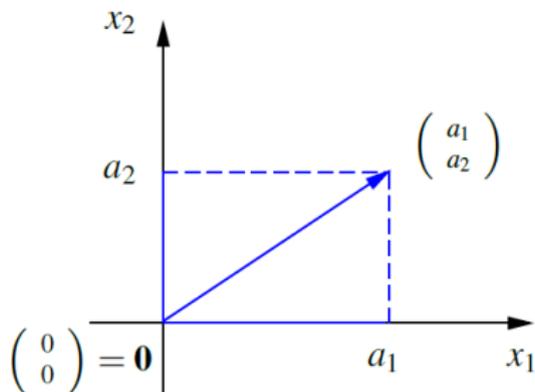


Figure 1: Векторы на плоскости.

Точно также отождествим \mathbb{R}^3 с трехмерным пространством, запишем точку с координатами (a_1, a_2, a_3) как вектор $\begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix}$ в \mathbb{R}^3 , выходящий из начала координат к точке (рис. 2).

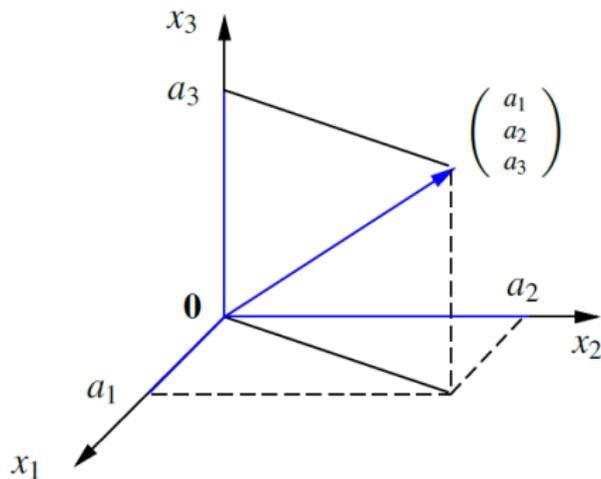


Figure 2: Векторы в пространстве.

Для описания преобразования $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ необходимо задать вектор $T(x)$ в \mathbb{R}^m для каждого x в \mathbb{R}^n . Это называется определением T , или определением действия T .

Говоря, что действие определяет преобразование, мы имеем в виду, что два преобразования $S : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ являются равными, если их действие одинаково.

Более формально

$$S = T, \text{ если, и только если } S(x) = T(x) \text{ для всех } x \in \mathbb{R}^n.$$

Умножение матриц является важным способом определения преобразований.

Если A - матрица $m \times n$, то умножение на A и дает преобразование

$$T_A : \mathbb{R}^n \rightarrow \mathbb{R}^m,$$

$$\text{определяемое как } T_A(x) = Ax \text{ для каждого } x \in \mathbb{R}^n.$$

T_A называется преобразованием, **индуцированным** A .

Метод главных компонент (РСА)

- Один из основных практических способов уменьшить размерность данных
- Дана матрица $X_{m \times n}$ - матрица «объекты - признаки»
- Реализация метода
 - вычисление собственных векторов и собственных значений ковариационной матрицы исходных данных
 - сингулярное разложение центрированной матрицы исходных данных

Изменение базиса

Представим набор исходным данных: каждый столбец представляет собой единичный образец (или момент времени) нашего набора данных (т.е. X).

Пусть Y - другая $m \times n$ матрица, связанная линейным преобразованием P .

X - это исходный записанный набор данных. Y - представляет собой новое представление этого набора данных.

$$PX = Y$$

r_i - строки P ,

x_i - столбцы X ,

y_i - столбцы Y .

$$PX = \begin{bmatrix} p_1 \\ \vdots \\ p_m \end{bmatrix} \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}$$

$$Y = \begin{bmatrix} p_1 x_1 & \dots & p_1 x_n \\ \vdots & \ddots & \vdots \\ p_m x_1 & \dots & p_m x_n \end{bmatrix}$$

$$y_i = \begin{bmatrix} p_1 x_i \\ \vdots \\ p_m x_i \end{bmatrix}$$

Матрица ковариантности

$$C_X = \frac{1}{n}XX^T$$

C_X является квадратной симметричной матрицей

Диагональные элементы C_X являются дисперсиями определенных типов измерений.

Недиагональные термины C_X являются ковариациями между типами измерений.

- Выбрать нормированное направление в m - мерное пространство, вдоль которого дисперсия в X максимальна. Сохранить этот вектор как p_1 .
- Найти другое направление, вдоль которого дисперсия максимальна, однако, из-за условия ортонормальности, ограничить поиск всеми направлениями, ортогональными всем ранее выбранным направлениям. Сохраните этот вектор как p_i
- Повторить эту процедуру до тех пор, пока не выбраны все m векторов.

Полученный упорядоченный набор представляет собой **основные компоненты**.

Найдем некоторую ортонормальную матрицу P в $Y = PX$ такую, что $C_Y = \frac{1}{n}YY^T$ - диагональная матрица.
Строки в P являются основными компонентами для X .

$$\begin{aligned}C_Y &= \frac{1}{n}YY^T \\ &= \frac{1}{n}(PX)(PX)^T \\ &= \frac{1}{n}PXX^TP^T \\ &= \frac{1}{n}PC_XP^T \\ C_Y &= PC_XP^T\end{aligned}$$

Метод SVD

Пусть X произвольная матрица $n \times m$ и $X^T X$ будет квадратной симметричной матрицей ранга r .

Определим все интересующие нас величины.

- $\{\hat{v}_1, \hat{v}_2, \dots, \hat{v}_r\}$ - ортонормированный набор $m \times 1$ собственных векторов с соответствующими собственными значениями $\{\lambda_1, \lambda_2, \dots, \lambda_r\}$.

Для симметричной матрицы $X^T X$:

$$(X^T X)\hat{v}_i = \lambda_i \hat{v}_i.$$

- $\sigma_i = \sqrt{\lambda_i}$ являются положительными действительными и называются сингулярными числами.
- $\{\hat{u}_1, \hat{u}_2, \dots, \hat{u}_r\}$ - набор $n \times 1$ векторов $\hat{u}_i = \frac{1}{\sigma_i} X \hat{v}_i$.

•

$$\hat{u}_i \cdot \hat{u}_j = 1, i = j$$

$$\hat{u}_i \cdot \hat{u}_j = 0, i \neq j$$

- $\|X \hat{v}_i\| = \sigma_i$.

Отсюда следует, что мы можем построить матрицу Σ , у которой на главной диагонали будут стоять собственные значения, а все остальные элементы будут равны нулю .

$$XV = U\Sigma \implies X = U\Sigma V^T.$$