

Метод главных компонент (principal component analysis) - Метод понижения размерности выборки, с наименьшей потерей информации. Идея метода коррелирует с LS решением (linear regression), поскольку главная компонента (ось) совпадает с линейной аппроксимацией выборки.

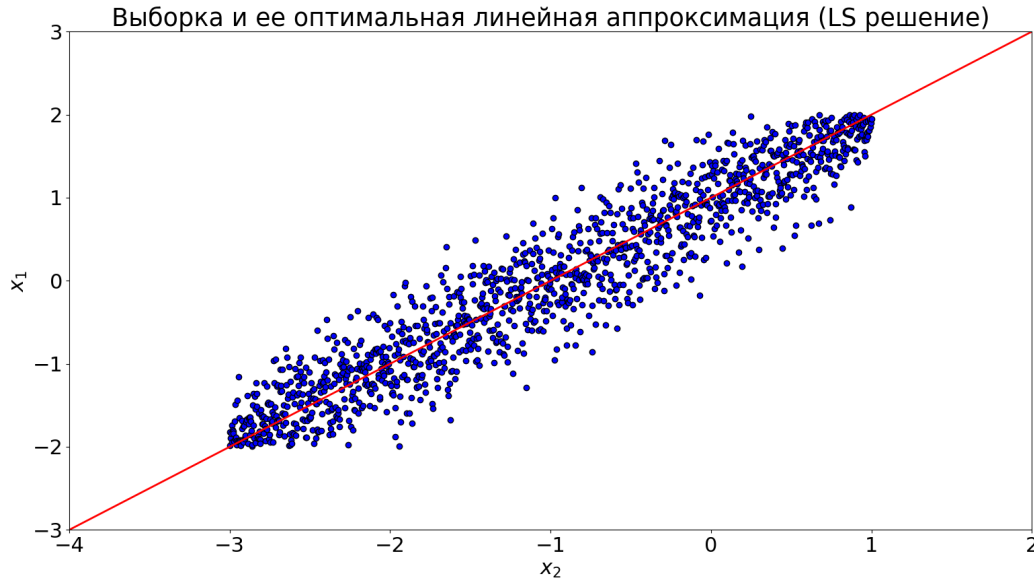


Figure 1: Двумерная выборка и линейная регрессионная модель

Цель этого метода заключается в поиске таких осей, на которые можно спроектировать многомерную выборку и сохранить при этом наибольшее количество информации. Сравним два варианта. Проекция выборки на оптимальную линейную поверхность аппроксимирующую выборку и на произвольную прямую:

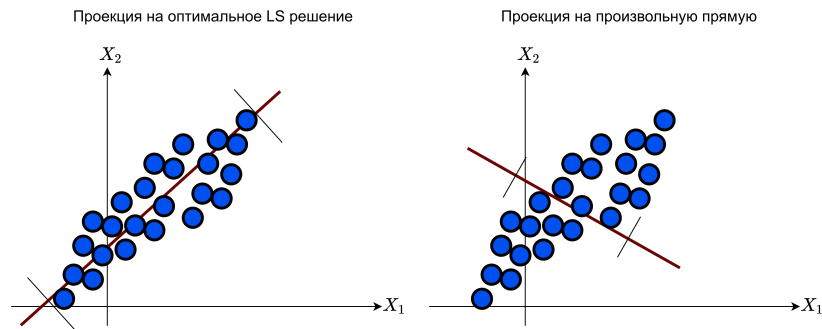


Figure 2: Сравнительное представление двух вариантов проектирования выборки на разные прямые

С точки зрения теории информации - степень неопределенности характеризует информативность данных. Известно, что дисперсия данных определяет их квадрат отклонения от среднего или же с точки зрения теории информации это понятие коррелирует с неопределенностью. Таким образом, чтобы уменьшить размерность выборки, сохраняя ее информативность, необходимо найти такое аффинное множество, проектируя на которое данные, их дисперсия будет максимальна. Сформулируем критерий.

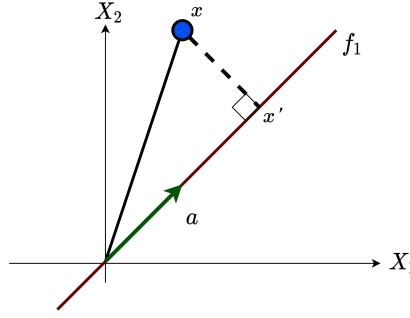


Figure 3: Геометрическая предпосылка к формулированию критерия

Положим, что $\mathbf{a} = (a_1, a_2, \dots, a_n)^T$ - конечномерный базис линейного множества, на которое планируется спроектировать выборку, где $\|\mathbf{a}\| = 1$. Естественным образом будет думать, что длина спроектированного вектора $\|x'_i\|$ это взвешенный базис $\mathbf{a} = (a_1, a_2, \dots, a_n)^T$.

$$\|x'_i\| = \frac{\bar{a}}{\|\mathbf{a}\|} \bar{x}_i = \|x_i\| \cos \theta = \text{proj}_{\mathbf{a}} x_i = \mathbf{a}^T x_i$$

Поскольку цель - найти прямую с максимальной дисперсией проекций - это задача условной оптимизации:

$$\begin{cases} \text{var}(\mathbf{a}^T x_i - \bar{x}) \rightarrow \max_a \\ \|\mathbf{a}\| - 1 = 0 \end{cases}$$

Положим, что данные смещены к центру на координату МО, поэтому задача упрощается до:

$$\begin{cases} \frac{1}{n} \sum_i (\mathbf{a}^T x_i)^2 \rightarrow \max_a \\ \|\mathbf{a}\| - 1 = 0 \end{cases}$$

Раскладывая $\frac{1}{n} [(a_1 x_1)^2 + (a_2 x_2)^2 + \dots + (a_n x_n)^2] = \frac{1}{n} \left\langle \mathbf{a}^T X \quad (\mathbf{a}^T x)^T \right\rangle = \frac{1}{n} \mathbf{a}^T X (\mathbf{a}^T X)^T = \mathbf{a}^T \frac{XX^T}{n} \mathbf{a} = \mathbf{a}^T S \mathbf{a}$ и по свойству: $\|\mathbf{a}\| = \mathbf{a}^T \mathbf{a}$ задача трансформируется к:

$$\begin{cases} \mathbf{a}^T S \mathbf{a} \rightarrow \max_a \\ \mathbf{a}^T \mathbf{a} - 1 = 0 \end{cases} \quad (1)$$

Задачи условной оптимизации решаются методом множителей Лагранжа:(ищем $\mathbf{a} \Rightarrow$ дифференцируем по \mathbf{a})

$$L(\mathbf{a}, S, \lambda) = \mathbf{a}^T S \mathbf{a} - \lambda (\mathbf{a}^T \mathbf{a} - 1)$$

$$\frac{dL(\mathbf{a}, S, \lambda)}{d\mathbf{a}} = \mathbf{a}^T S \mathbf{a} - \lambda (\mathbf{a}^T \mathbf{a} - 1)$$

$$\frac{dL(\mathbf{a}, S)}{d\mathbf{a}} = \langle \mathbf{a}, S \mathbf{a} \rangle - \lambda \langle \mathbf{a}, \mathbf{a} \rangle + 1$$

$$\frac{dL(\mathbf{a}, S, \lambda)}{d\mathbf{a}} = \langle \mathbf{a}, S \mathbf{a} \rangle - \lambda \langle \mathbf{a}, \partial \mathbf{a} \rangle + 1$$

$$dL(\mathbf{a}, S, \lambda) = d \langle \mathbf{a}, S \mathbf{a} \rangle - \lambda d \langle \mathbf{a}, \mathbf{a} \rangle$$

$$dL(\mathbf{a}, S, \lambda) = \langle d\mathbf{a}, S \mathbf{a} \rangle + \langle \mathbf{a}, dS \mathbf{a} \rangle - \lambda \langle d\mathbf{a}, \mathbf{a} \rangle - \lambda \langle \mathbf{a}, d\mathbf{a} \rangle$$

$$dL(\mathbf{a}, S, \lambda) = \langle S \mathbf{a}, d\mathbf{a} \rangle + \langle S^T \mathbf{a}, d\mathbf{a} \rangle - \lambda \langle \mathbf{a}, d\mathbf{a} \rangle - \lambda \langle \mathbf{a}, d\mathbf{a} \rangle$$

$$dL(\mathbf{a}, S, \lambda) = S\mathbf{a} + S^T\mathbf{a} - 2\lambda\mathbf{a}$$

Поскольку матрица квадратичной формы симметрична: $S = S^T$, то приравняв дифференциал к 0, получим, что:

$$2S\mathbf{a}^T - 2\lambda\mathbf{a} = 0$$

$$(S - \lambda I)\mathbf{a} = 0 \quad (2)$$

Отсюда видно, что \mathbf{a} - собственный вектор и λ - собственные числа. Решая систему уравнений, и извлекая максимальное собственное значение (определяет длину вектора) мы можем найти собственный вектор согласно критерию. Таким образом, найден вектор, на который можно спроектировать данные и сформировать сжатие с максимальным сохранением информации. Другим вопросом является то, что получен единственный вектор. Пусть этот вектор будет обозначаться как $\mathbf{a} = \mathbf{a}_1$. Если мы планируем спроектировать данные на плоскость, нам необходимо решить другую задачу оптимизации:

$$\begin{cases} \mathbf{a}_2^T S \mathbf{a}_2 \rightarrow \max_{\mathbf{a}_2} \\ \mathbf{a}_1^T \mathbf{a}_2 = 0 \\ \|\mathbf{a}_2\| - 1 = 0 \end{cases}$$

Известно, что результат этого решения: \mathbf{a}_2 , (который, кстати соответствует второму по величине собственному вектору из задачи (1)) ставля новую задачу оптимизации с новыми условиями, включающими все предыдущие результаты, мы можем произвольно уменьшать/увеличивать размерность выборки. Более простым пояснением является то, что собственные векторы симметрической матрицы являются ортогональными. Поэтому решить задачу можно и в одну итерацию.

Поиск решения разложением SVD

Для симметрических матриц существует взаимосвязь между собственными значениями и сингулярными значениями SVD разложения. Известно, что SVD раскладывает матрицу по закону:

$$X = U\Sigma V^H \quad (3)$$

Где

U - самосопряженный оператор вращения ($UU^T = I$)

Σ - оператор растяжения (Диагональная матрица)

$V^H = V^T$ - самосопряженный оператор вращения для $X \in R^{n \times p}$ (Эрмитово сопряжение к V) и ($VV^T = I$)

Согласно (3), имеем $X^T = (U\Sigma V^T)^T = V(\Sigma U^T)^T = V\Sigma^T U^T = V\Sigma U^T$

$S = X^T X = V\Sigma U^T U \Sigma V^T = V\Sigma I \Sigma V^T = V\Sigma^2 V^T$

$X^T X V = S V = V\Sigma^2$ (по свойству самосопряженности: $VV^T = I$)

Сравнивая (2) $S\mathbf{a} = \lambda\mathbf{a}$ и $S V = V\Sigma^2$, видим, что $\mathbf{a} = \Sigma^2$, то есть:

$$\begin{pmatrix} \lambda_1 & \cdots & \cdots & 0 \\ \vdots & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_N \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \cdots & \cdots & 0 \\ \vdots & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_N^2 \end{pmatrix} \quad (4)$$

Таким образом, формулировка алгоритма метода главных компонент следующая:

Algorithm 1 Метод главных компонент

- Нормировать данные по среднему
 - Рассчитать корреляционную матрицу для данных: $S = \frac{1}{N} X^T X$
 - Определить собственные значения и собственные вектора корреляционной матрицы (можно разложить по собственным векторам или использовать SVD)
 - Отсортировать собственные вектора в порядке убывания, выберите требуемое количество и спроектируйте данные на новые вектора $\mathbf{a}^T x_i \forall i = \overline{1, n}$
-

Вывод:

Был рассмотрен метод главных компонент. Его целью является уменьшить размерность выборки, с максимальным сохранением информации. Это позволяет визуализировать многомерные данные, энтропийно кодировать их и так далее. Было показано то, что наилучшими плоскостями, на которые проектируется выборка - те плоскости, на которых данные будут иметь максимальную дисперсию. Поставлена и решена задача оптимизации, а также показана связь между собственными числами симметрической матрицы и оператором растяжения SVD разложения.