

# Отбор признаков и линейные методы снижения размерности

Кантонистова Елена

[elena.kantonistova@yandex.ru](mailto:elena.kantonistova@yandex.ru)

27 октября 2017

## 1 Методы отбора признаков

- VarianceThreshold
- Отбор по корреляции с целевой переменной
- Отбор признаков по различным статистическим тестам

## 2 Линейные методы снижения размерности

- Метод главных компонент
- Линейный дискриминантный анализ

Можем удалить признаки, которые имеют очень маленькую дисперсию, т.е. практически константы.

# Отбор по корреляции с целевой переменной

Для каждого признака вычислим его корреляцию с целевой переменной. Будем выкидывать признаки, имеющие маленькую корреляцию.

В sklearn есть сразу несколько методов, использующих отбор по статистическим критериям. Среди них выделим следующие:

- SelectKBest - оставляет  $k$  признаков с наибольшим значением выбранной статистики
- SelectPercentile - оставляет признаки со значениями выбранной статистики, попавшими в заданную пользователем квантиль
- и другие (см.sklearn)

- mutual information: для векторов  $X$  и  $Y$  статистика вычисляется по формуле

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \left( \frac{p(x, y)}{p(x) p(y)} \right)$$

- хи-квадрат:

$$\chi^2(X, Y) = \sum_{i=1}^n \frac{(Y_i - X_i)^2}{X_i}$$

- f-regression - тест, основанный на корреляции линейного регрессора с целевой переменной

Предыдущие методы отбирали из исходных признаков некоторое подмножество признаков. Теперь мы хотим придумать новые признаки, каким-то образом выражающиеся через старые, причем новых признаков хочется меньше, чем старых. Сегодня будем рассматривать только случай, когда новые признаки линейно выражаются через старые.

Постановка задачи:

$f_1(x), \dots, f_n(x)$  — исходные числовые признаки;

$g_1(x), \dots, g_m(x)$  — новые числовые признаки,  $m \leq n$ ;

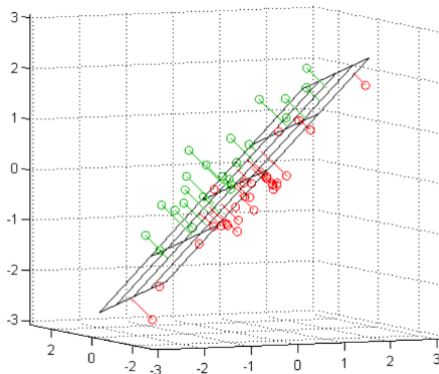
Мы хотим, чтобы новые числовые признаки  $g_i(x)$  линейно выражались через исходные признаки  $f_j(x)$ , при этом чтобы исходные признаки также линейно восстанавливались по новым признакам. При этом мы хотим, чтобы при переходе к новым признакам было потеряно наименьшее количество исходной информации.



Метод главных компонент работает только с признаками. Для него не важна целевая переменная (если она есть). Таким образом, метод главных компонент - это обучение без учителя.

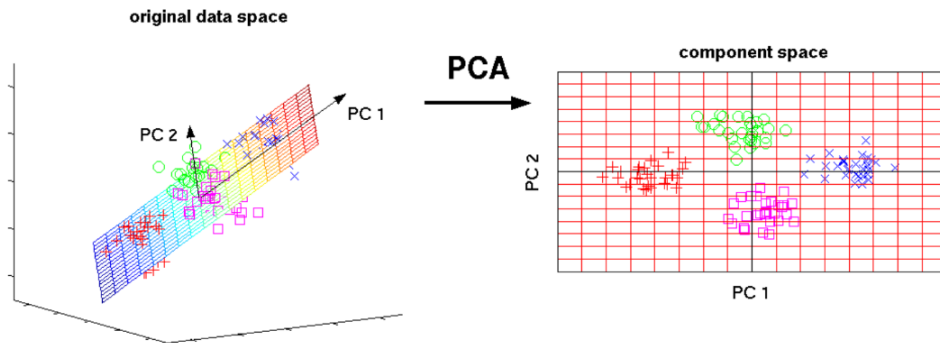
# Геометрическая интерпретация PCA

Геометрически метод главных компонент ищет гиперплоскость заданной размерности, при проекции на которую сумма квадратов расстояний от исходных точек будет минимальной.



# Визуализация проекции на гиперплоскость

Точки, плохо разделимые в исходном пространстве, могут быть лучше разделимы при проекции на некоторую гиперплоскость.



# Faces dataset



# Faces dataset (main components)

Первые главные компоненты после применения PCA



LDA - это обучение с учителем. При помощи метода линейного дискриминантного анализа выбирается проекция исходного пространства признаков на новое пространство признаков таким образом, чтобы минимизировать внутриклассовый разброс точек и максимизировать межклассовое расстояние в пространстве признаков.

- Классификация между  $\omega_1$  и  $\omega_2$ .
- Пусть  $C_1 = \{i : x_i \in \omega_1\}$ ,  $C_2 = \{i : x_i \in \omega_2\}$  и

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n, \quad m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

$$\mu_1 = w^T m_1, \quad \mu_2 = w^T m_2$$

- Определим дисперсии спроецированных на подпространство  $w$  классов:

$$s_1 = \sum_{n \in C_1} (w^T x_n - w^T m_1)^2, \quad s_2 = \sum_{n \in C_2} (w^T x_n - w^T m_2)^2$$

- Критерий LDA Фишера:  $\frac{(\mu_1 - \mu_2)^2}{s_1^2 + s_2^2} \rightarrow \max_w$