

Нелинейные методы снижения размерности

Кантонистова Елена

elena.kantonistova@yandex.ru

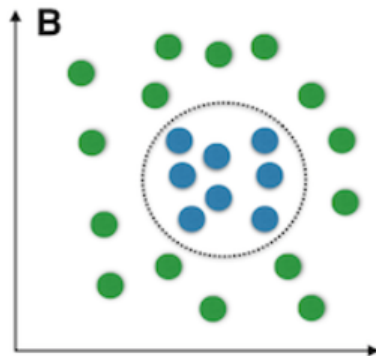
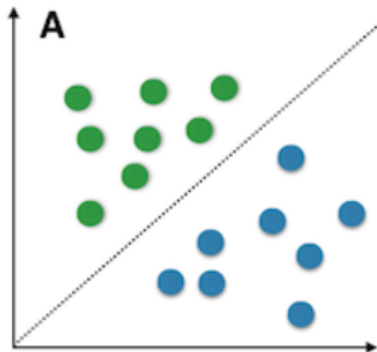
1 ноября 2017

1 Ядровой PCA

2 t-SNE

Метод главных компонент (РСА) преобразует базовые признаки в новые, каждый из которых является линейной комбинацией изначальных таким образом, что разброс данных (то есть среднеквадратичное отклонение от среднего значения) вдоль них максимален. Метод применяется для визуализации данных и для уменьшения размерности данных (сжатия).

Linear vs. nonlinear problems

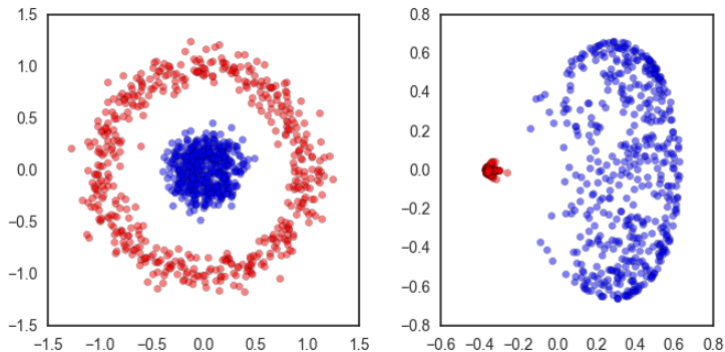


РСА с ядром (Kernel PCA) основан на том же принципе, что и обычный РСА, однако он позволяет создавать новые признаки не только линейной комбинацией исходных, но и более сложными преобразованиями.

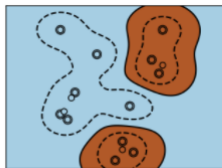
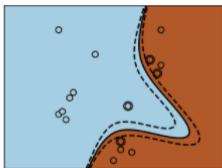
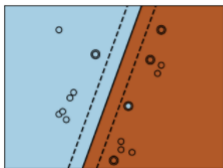
Это происходит благодаря замене обычного скалярного произведения (x, y) , возникающего при решении задачи РСА, на некоторую функцию $K(x, y)$, обладающую всеми свойствами скалярного произведения.

Зачем нужны ядра

Иногда данные в задаче не являются линейно разделимыми. Однако, после некоторого нелинейного преобразования признаков данные можно разделить прямой. Введение ядер - это как раз и есть то самое нелинейное преобразование признаков, которое помогает нам впоследствии разделить данные прямой.



- Линейное
- Полиномиальное
- RBF $K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$
-



t-SNE = t-distributed stochastic neighbor embedding

Данный метод нужен для визуализации многомерных данных в двумерном или трехмерном пространстве.

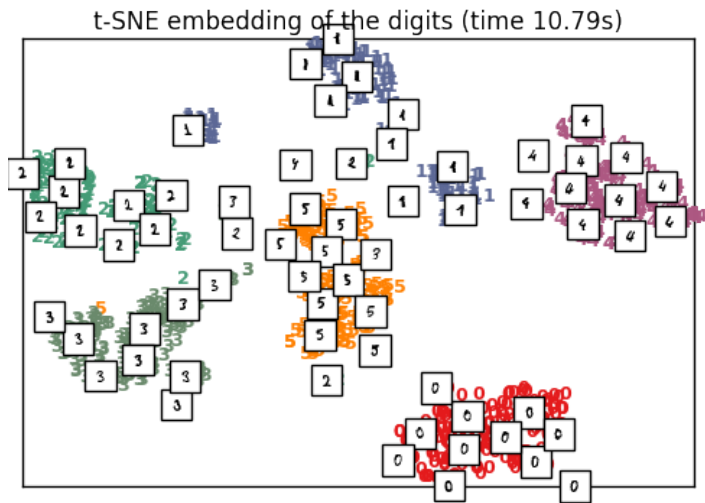
У нас есть набор данных с точками, описываемыми многомерной переменной с размерностью пространства существенно больше трех. Необходимо получить новую переменную, существующую в двумерном или трехмерном пространстве, которая бы в максимальной степени сохраняла структуру и закономерности в исходных данных. Другими словами, точки, которые в исходном пространстве находятся далеко друг от друга, в новом 2-х или 3-х мерном тоже должны оказаться далеко друг от друга (и наоборот, близкие - близко).

Пример использования: MNIST

A selection from the 64-dimensional digits dataset

0	1	2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	5
5	5	0	4	1	3	5	1	0	0	2	2	2	0	1	2	3	3	3	3
4	4	1	5	0	5	2	2	0	0	1	3	2	1	4	3	1	3	1	4
3	1	4	0	5	3	1	5	4	4	2	2	2	5	5	4	4	0	0	1
2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5
0	4	1	3	5	1	0	0	2	2	1	0	1	2	3	3	3	3	4	4
1	5	0	5	2	2	0	0	1	3	2	1	3	1	3	4	4	3	1	4
0	5	3	4	5	4	4	1	2	1	5	5	4	4	0	0	1	2	3	4
5	0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1
3	5	1	0	0	2	2	2	0	1	2	3	3	3	3	4	4	1	5	0
5	2	2	0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5
3	1	5	4	4	2	2	2	5	5	4	4	0	3	0	1	2	3	4	5
0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1	3
5	1	0	0	1	2	2	0	1	1	3	3	3	3	4	4	1	5	0	5
1	2	0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5	3
1	5	4	4	2	2	2	5	5	4	4	0	0	1	2	3	4	5	0	1
1	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1	3	5	1
0	0	1	2	2	0	1	1	3	3	3	3	4	4	1	5	0	5	1	2
0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5	3	1	5
4	4	2	2	1	5	5	4	4	0	0	1	2	3	4	5	0	1	2	3

Проекция в двумерное пространство с помощью t-SNE



Плюсы и минусы t-SNE

- Основной плюс: с помощью t-SNE можно получить очень хорошее наглядное представление о том, как выглядят многомерные данные
- Основной минус: метод работает очень долго, даже реализация из sklearn