
Занятие № 11

Feature Selection



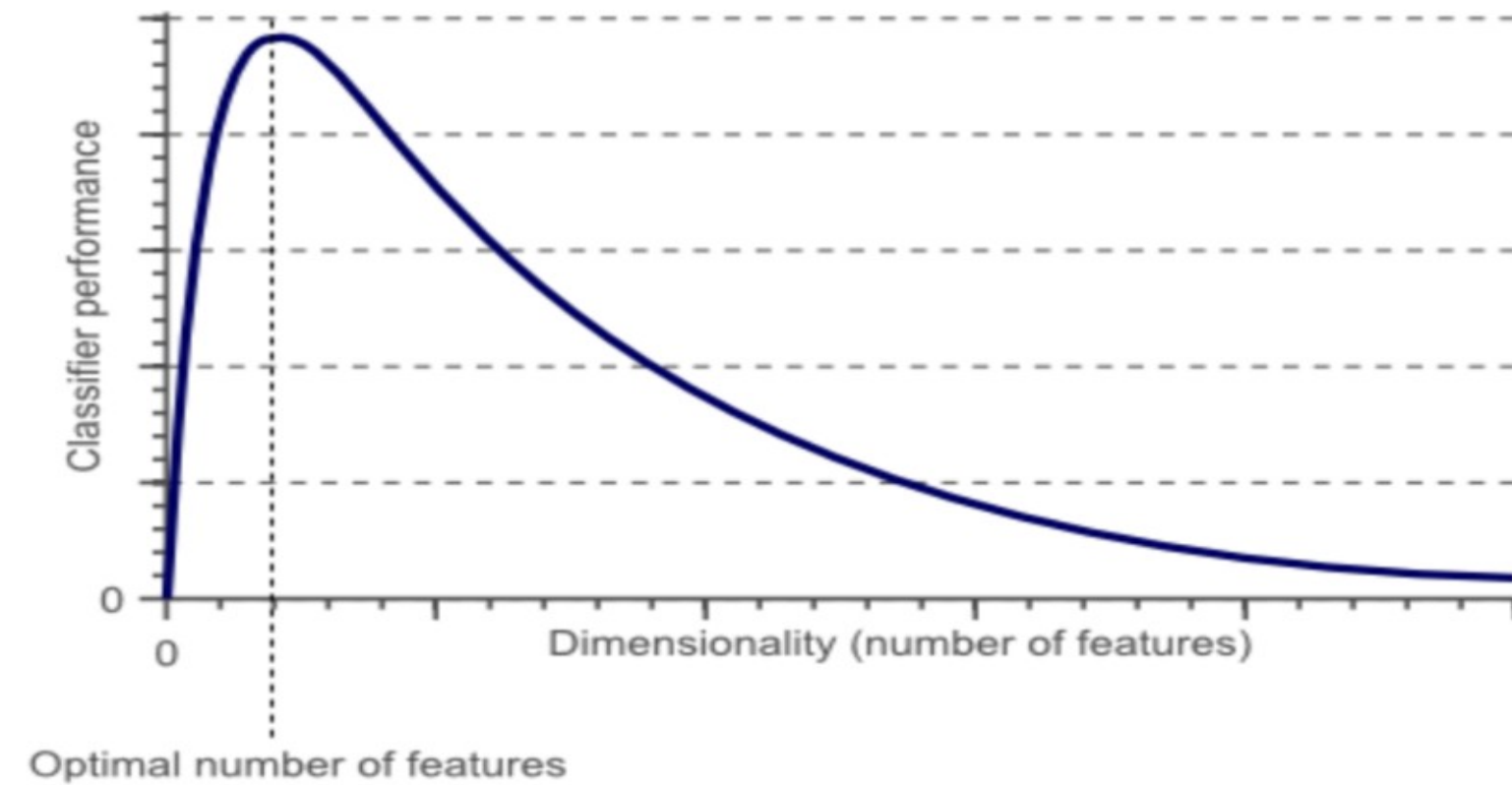
Содержание

- 1 Введение. Зачем всё это?
- 2 Статистика в отборе признаков
- 3 Декомпозиция данных
- 4 Практика.



Введение. Зачем всё это?

Проклятие размерности



x x x x x

Одно измерение - 5 точек

x x x x x
x x x x x
x x x x x
x x x x x
x x x x x

Два измерения - 25 точек

x x x x x
x x x x x
x x x x x
x x x x x
x x x x x
x x x x x
x x x x x
x x x x x
x x x x x

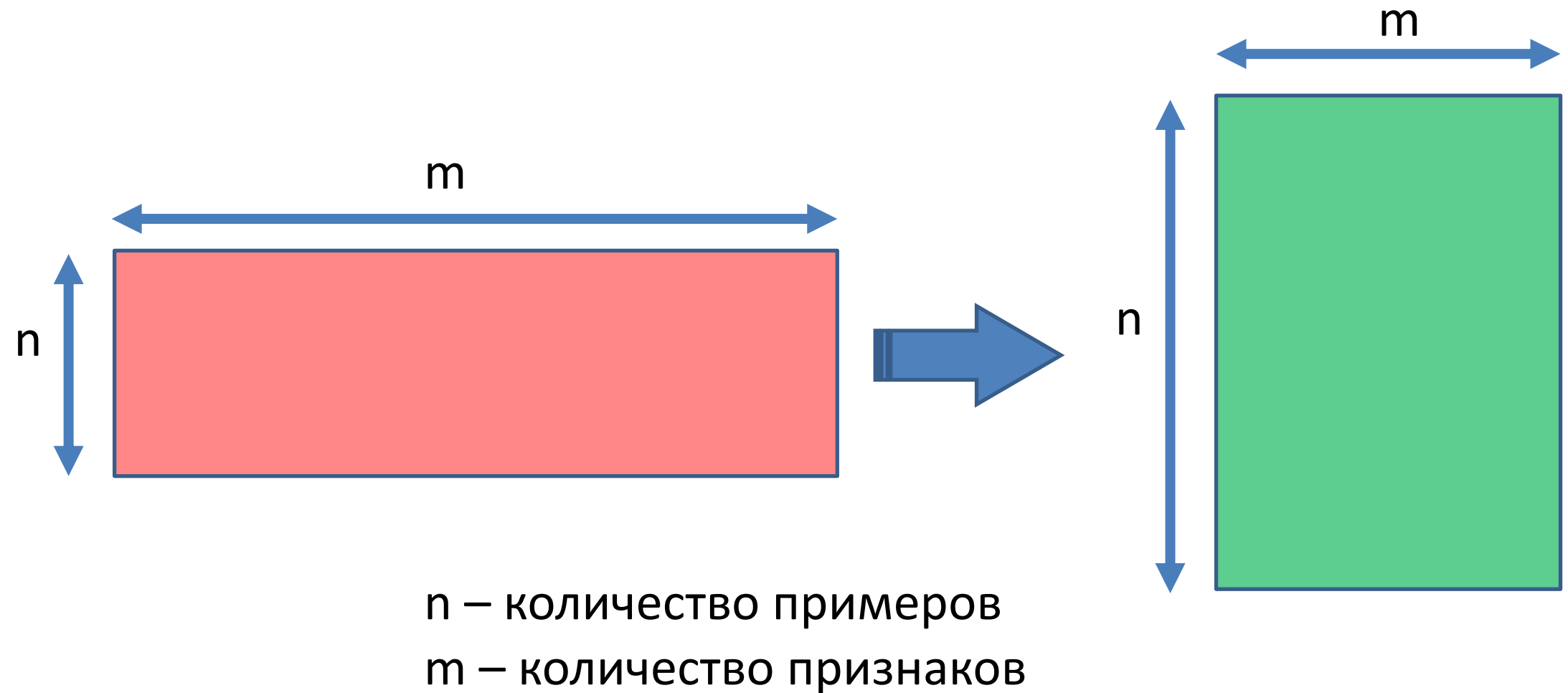
Три измерения - 125 точек



Методы отбора признаков

Позволит получить:

- упрощение моделей для того, чтобы сделать их проще для интерпретации исследователями или пользователями
- более короткое время тренировки
- уменьшения влияния проклятия размерности
- улучшение обобщения путём сокращения переобучения
- фильтрацию шумных признаков



Что можно сделать?

- Отобрать признаки
- Преобразовать признаки



Методы отбора признаков

Методы отбора

Задача – найти подмножество признаков на котором выбранная модель покажет лучшее качество

Фильтры

основаны на некоторых показателях, которые не зависят от метода классификации (коэффициент корреляции, взаимная информация, WOE, IG)

Обертки

опираются на информацию о важности признаков полученную от других методов или моделей ML (последовательный отбор и последовательное исключение признаков и др.)

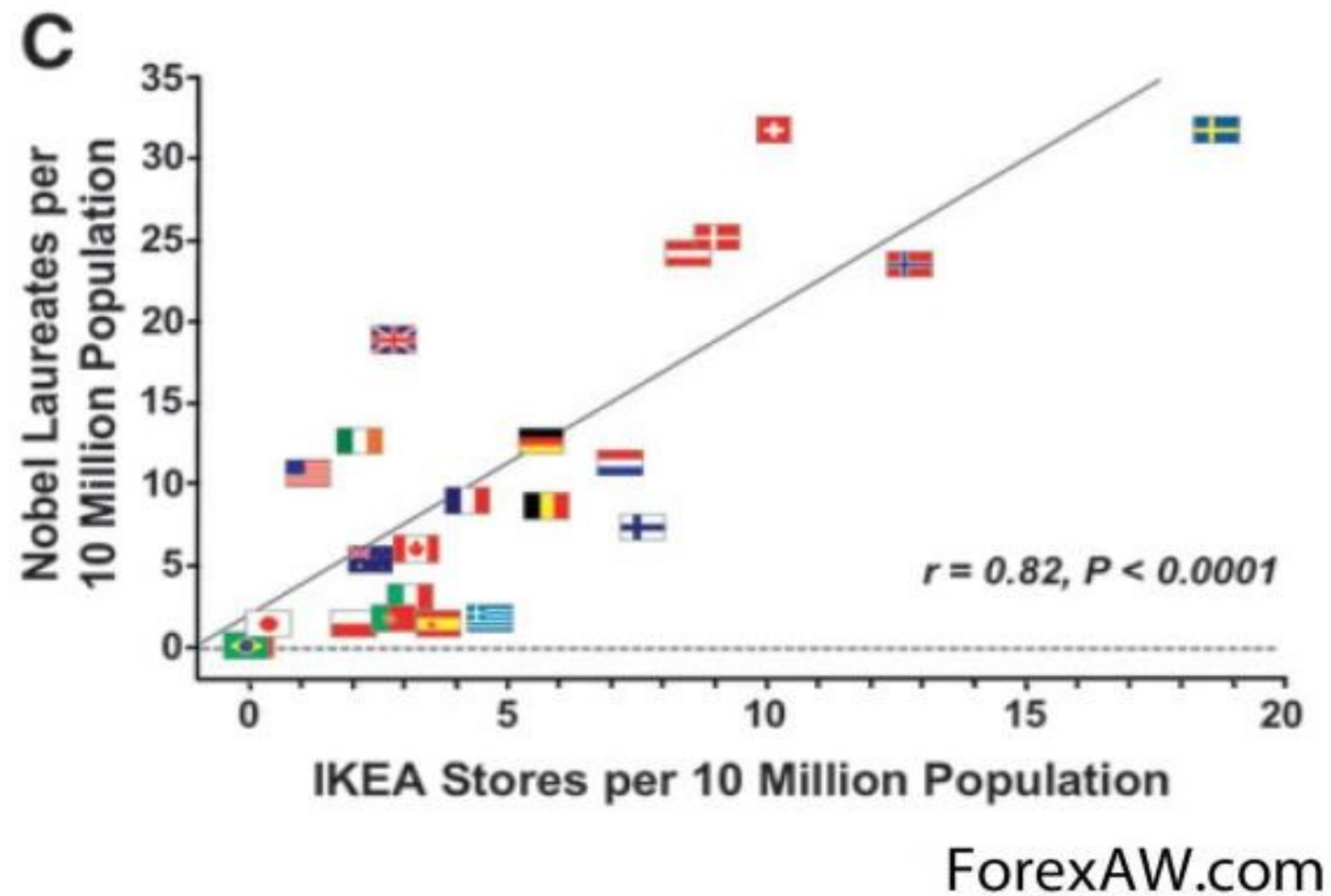
Встроенные в алгоритмы

выполняют отбор признаков во время процедуры обучения классификатора, и именно они явно оптимизируют набор используемых признаков для достижения лучшей точности (регрессия с L1-регуляризацией, Random Forest, SHAP)



Корреляция

Корреляция — статистическая взаимосвязь двух или более случайных величин. При этом изменения значений одной или нескольких из этих величин сопутствуют систематическому изменению значений другой или других величин.

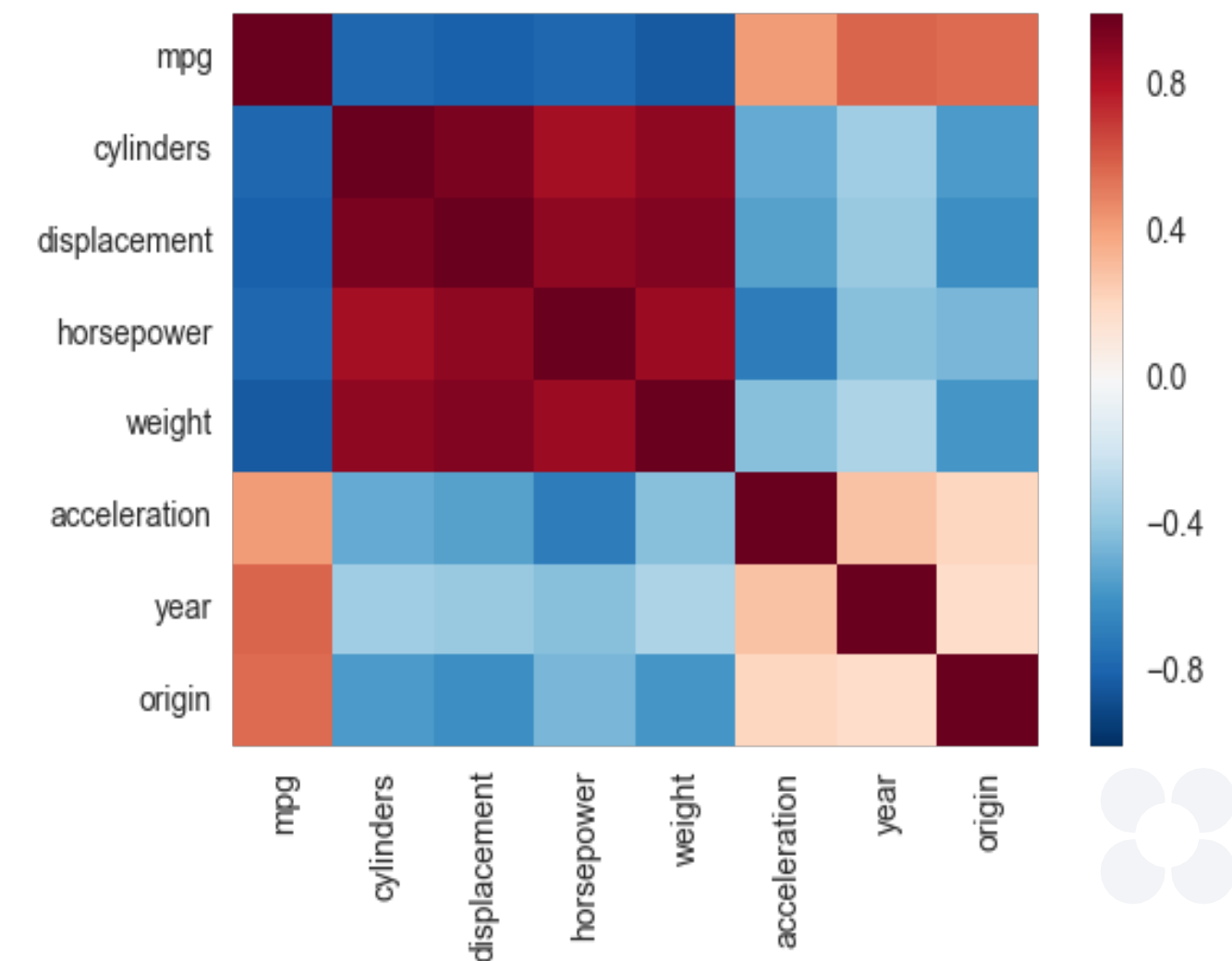


Ковариация

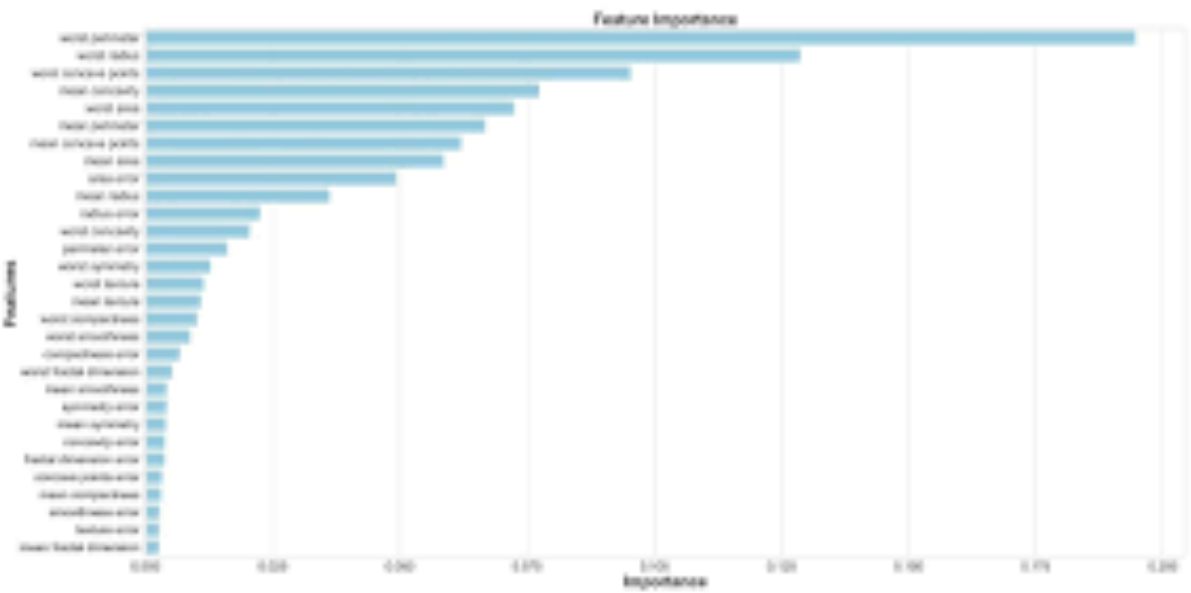
$$\text{cov}_{XY} = \mathbf{M}[(X - \mathbf{M}(X))(Y - \mathbf{M}(Y))] = \mathbf{M}(XY) - \mathbf{M}(X)\mathbf{M}(Y)$$

Коэффициент корреляции Пирсона

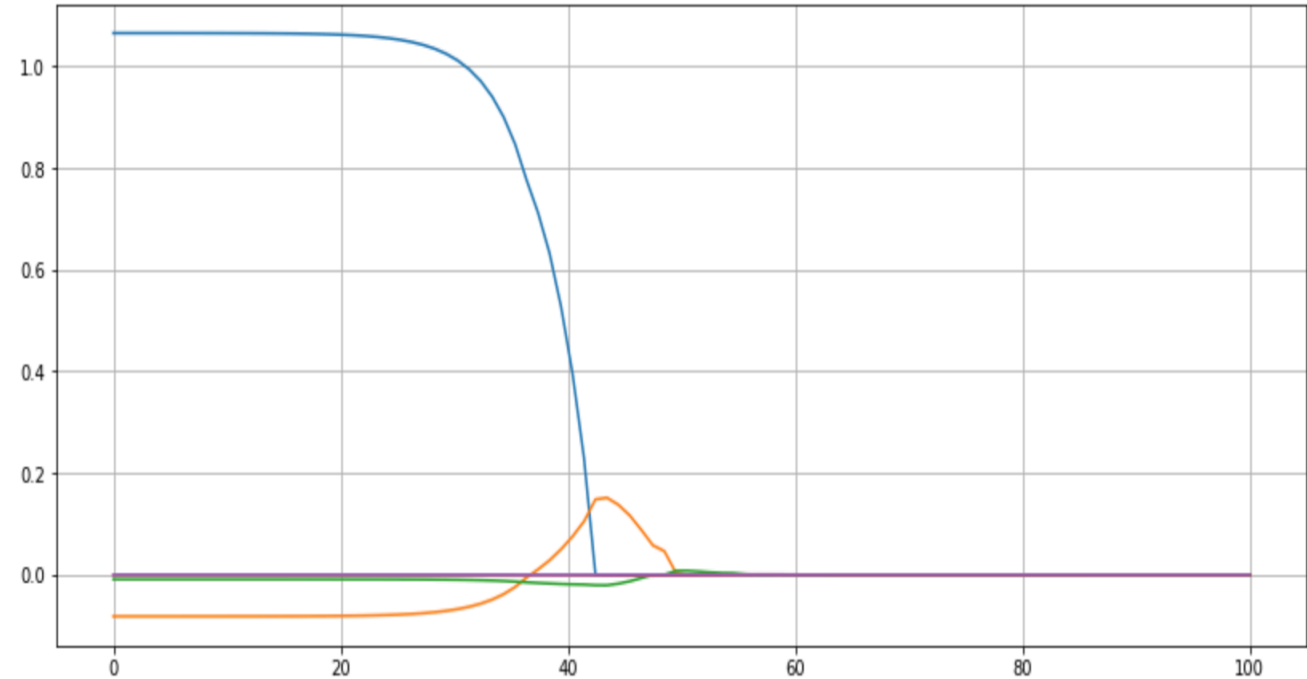
$$\mathbf{r}_{XY} = \frac{\text{cov}_{XY}}{\sigma_X \sigma_Y} = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{\sqrt{\sum (X - \bar{X})^2 \sum (Y - \bar{Y})^2}}$$



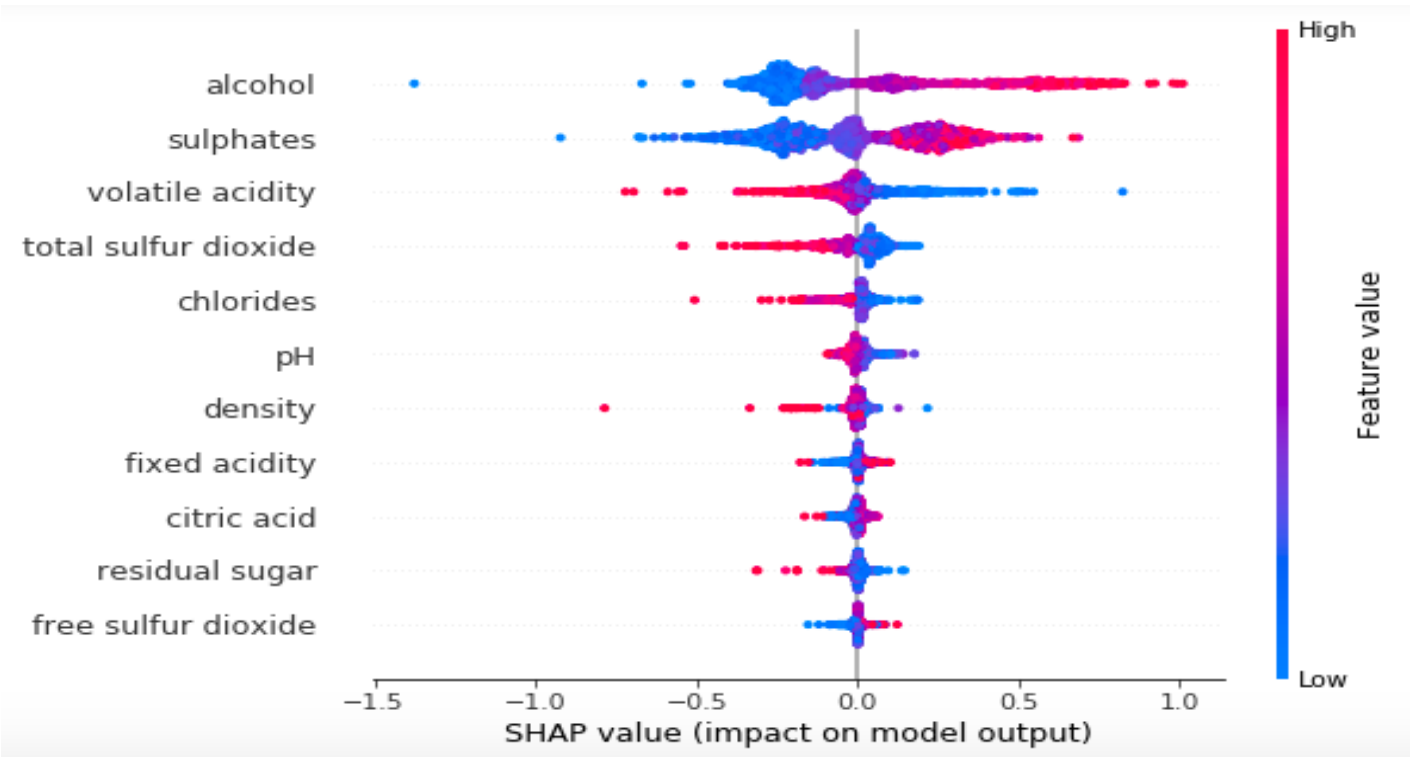
Random Forest



L1 - регуляризация



SHAP (SHapley Additive exPlanations)

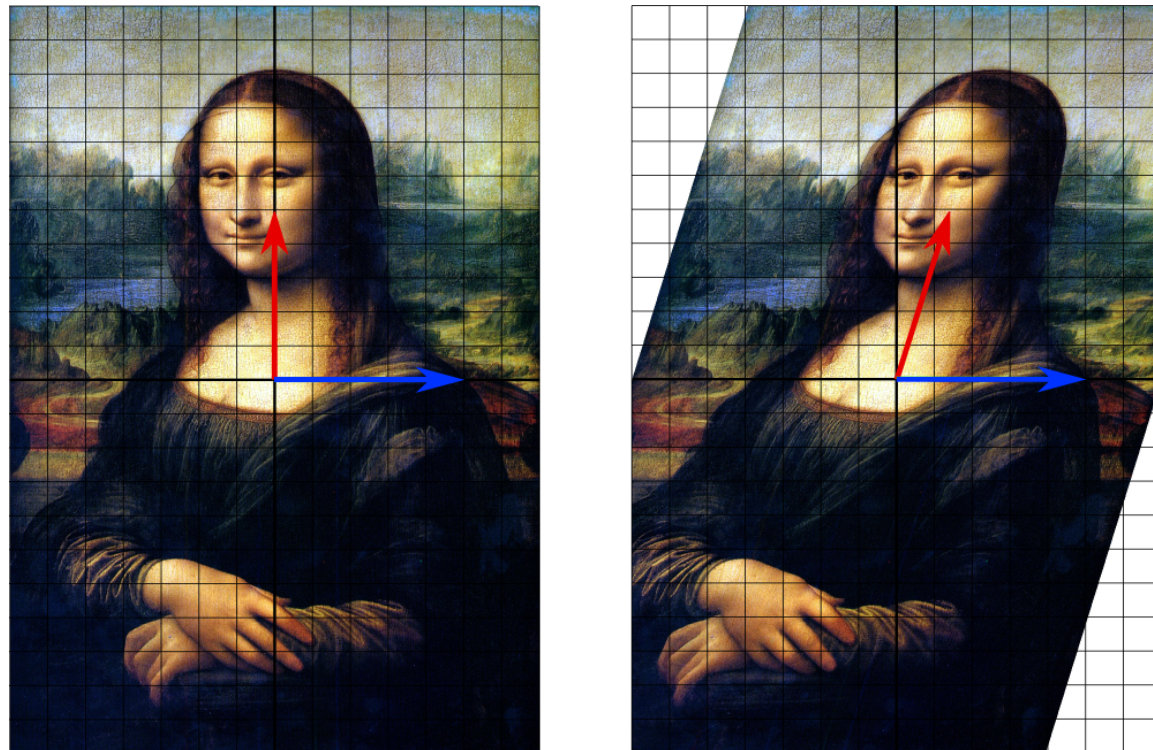


Преобразование признаков

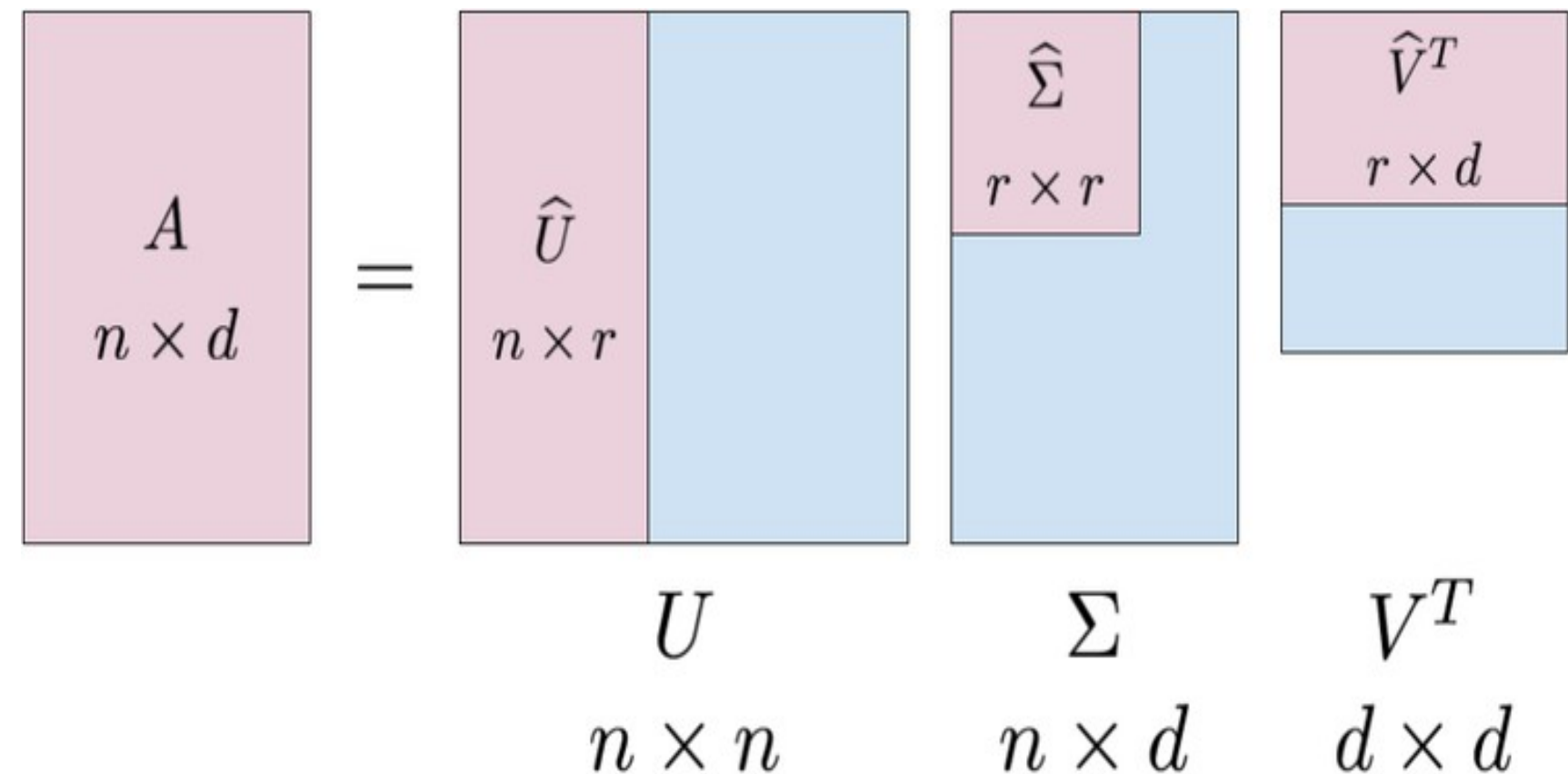
Метод главных компонент (principal component analysis, PCA):
позволяет уменьшить размерность данных с помощью преобразования
на основе линейной алгебры

Собственный вектор

$$M\vec{x} = \lambda\vec{x}$$

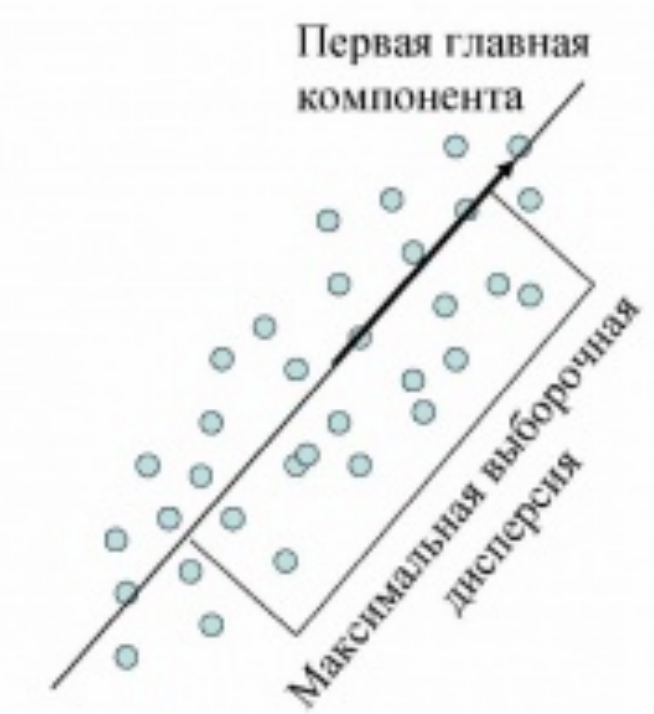


Сингулярное разложение (SVD)

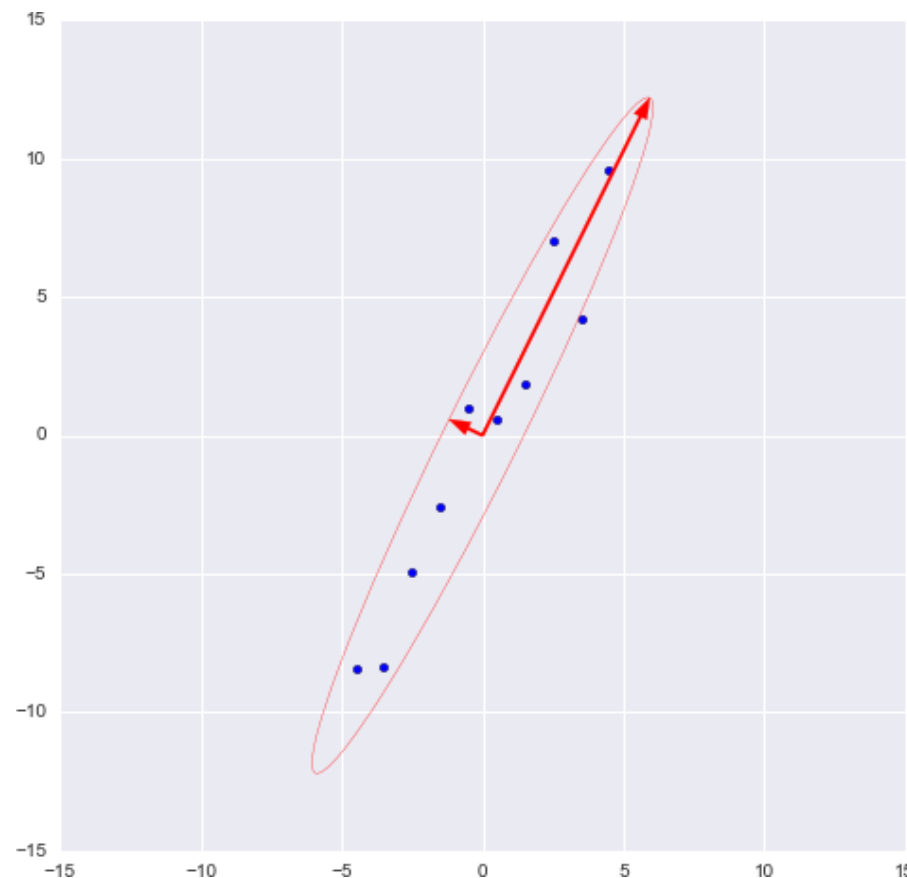


РСА

Зачем он нужен? Он уменьшает размерность с минимумом потери информации



- перевести данные в пространство меньшей размерности
- найти такое преобразование при котором разброс данных и дисперсия в ортогональных проекциях максимален
- корреляция между отдельными координатами обратится в ноль.



$$\text{Cov}(X_i, X_j) = E\left[(X_i - E(X_i)) \cdot (X_j - E(X_j))\right] = E(X_i X_j) - E(X_i) \cdot E(X_j)$$

$$\begin{aligned} \text{Var}(X^*) &= \Sigma^* = E(X^* \cdot X^{*T}) = E\left((\vec{v}^T X) \cdot (\vec{v}^T X)^T\right) = \\ &= E(\vec{v}^T X \cdot X^T \vec{v}) = \vec{v}^T E(X \cdot X^T) \vec{v} = \vec{v}^T \Sigma \vec{v} \end{aligned}$$



LDA

Линейный дискриминантный анализ

Метод уменьшения размерности, используемый в качестве этапа предварительной обработки в приложениях машинного обучения и классификации.

Первый шаг - вычислить разделимость между разными классами (то есть расстояние между средними значениями разных классов), также называемое межклассовой дисперсией.

$$S_b = \sum_{i=1}^g N_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})^T$$

Второй шаг - вычислить расстояние между средним значением и выборкой каждого класса, которое называется внутриклассовой дисперсией.

$$S_w = \sum_{i=1}^g (N_i - 1) S_i = \sum_{i=1}^g \sum_{j=1}^{N_i} (x_{i,j} - \bar{x}_i)(x_{i,j} - \bar{x}_i)^T$$

Третий шаг - построить пространство более низкой размерности, которое максимизирует дисперсию между классами и минимизирует дисперсию внутри класса.

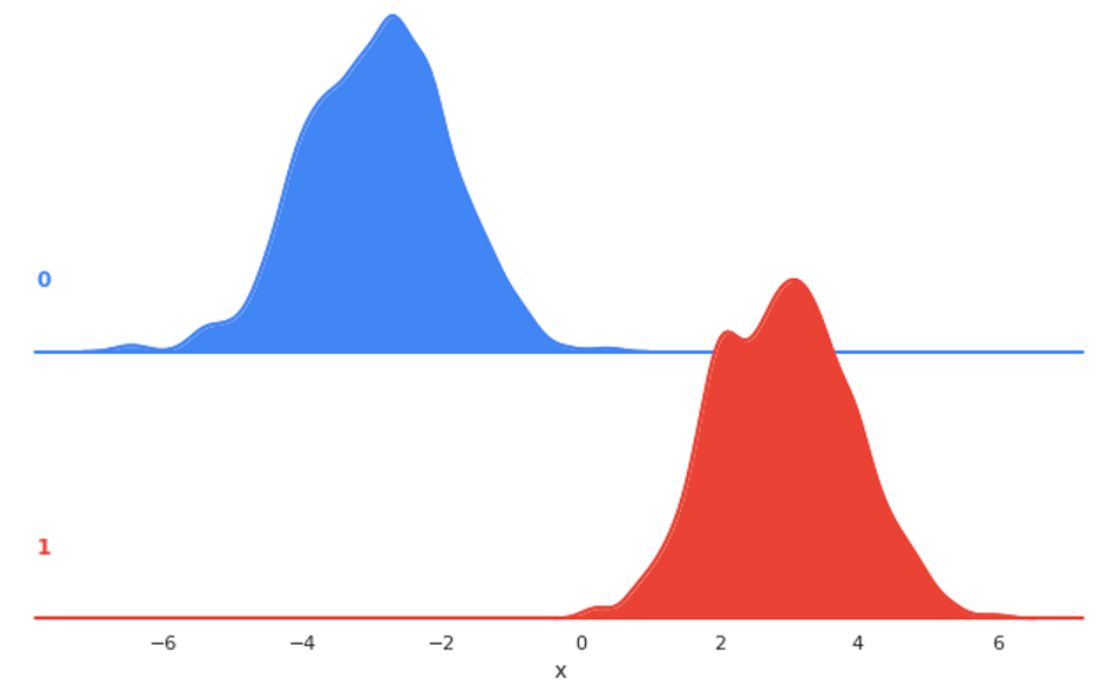
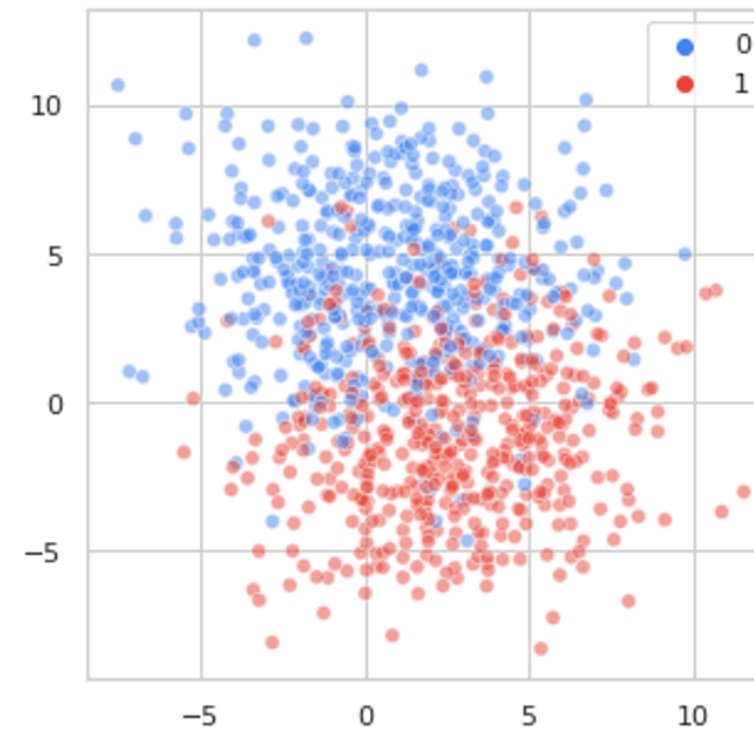
P - проекция пространства нижней размерности, которая называется критерием Фишера

$$P_{lda} = \arg \max_P \frac{|P^T S_b P|}{|P^T S_w P|}$$

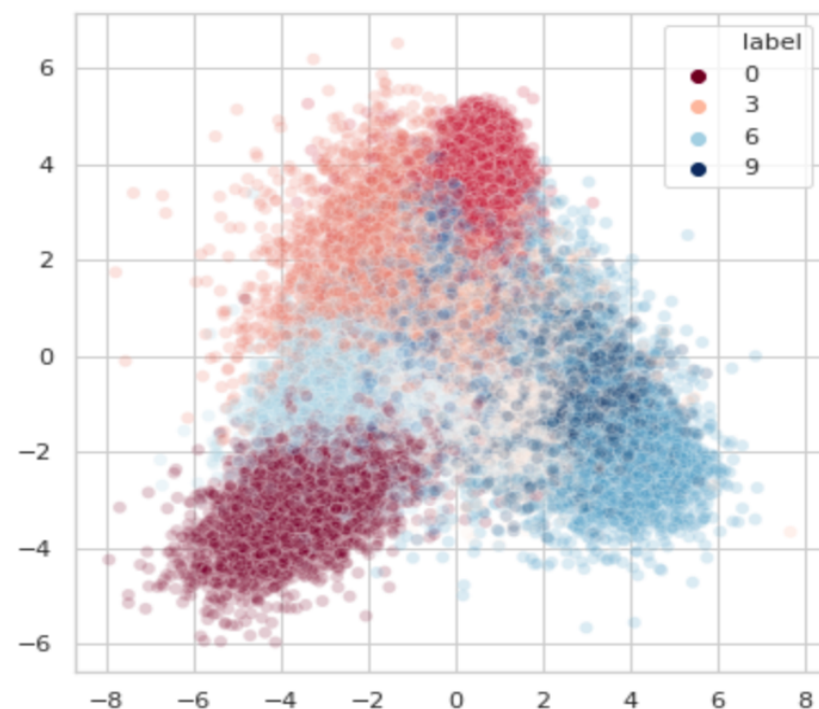


LDA

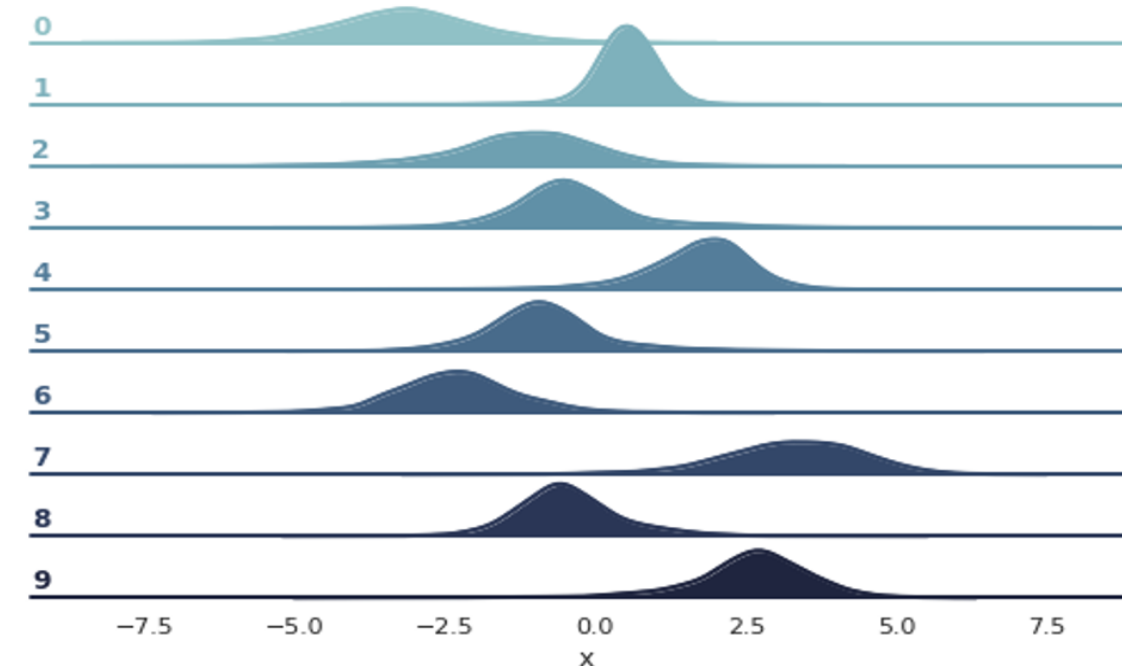
Отображение распределение
в 1- мерное пространство



Two-Dimensional Representation



One-Dimensional Representation



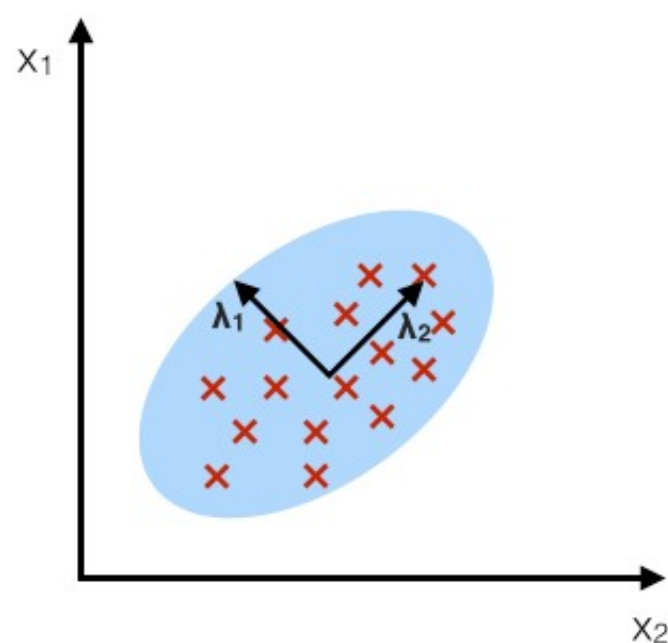
Отображение картинок MNIST
в 2- и 1- мерное пространство



Сравнение LDA и PCA

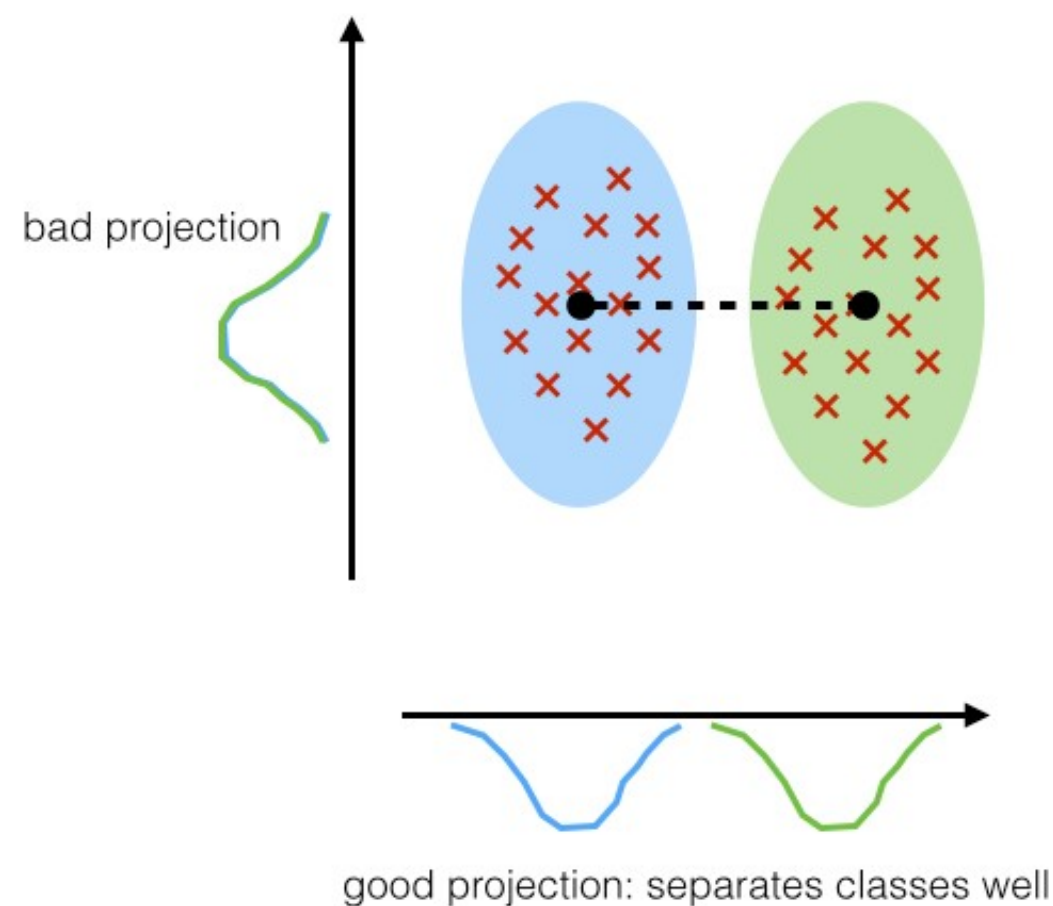
PCA:

component axes that maximize the variance



LDA:

maximizing the component axes for class-separation



ПРАКТИКА



Спасибо
за
Внимание!

