

Исследование алгоритмов уменьшения размерности данных для задачи классификации

Руслан Сайфутдинов, 344 гр.
Научный руководитель: Константин Невоструев
Кафедра системного программирования

Введение

- Задача: распознавать что-нибудь на изображениях и не только на изображениях
- Изображения слишком большие, чтобы запустить обучения прямо на них
- Можно просто уменьшить разрешение
- А можно попытаться сохранить больше важной информации

Уменьшение размерности









Постановка задачи

- Реализовать несколько интересных алгоритмов уменьшения размерности данных
- Сравнить их работу

Данные

- MNIST

Тренировочная выборка: 60,000

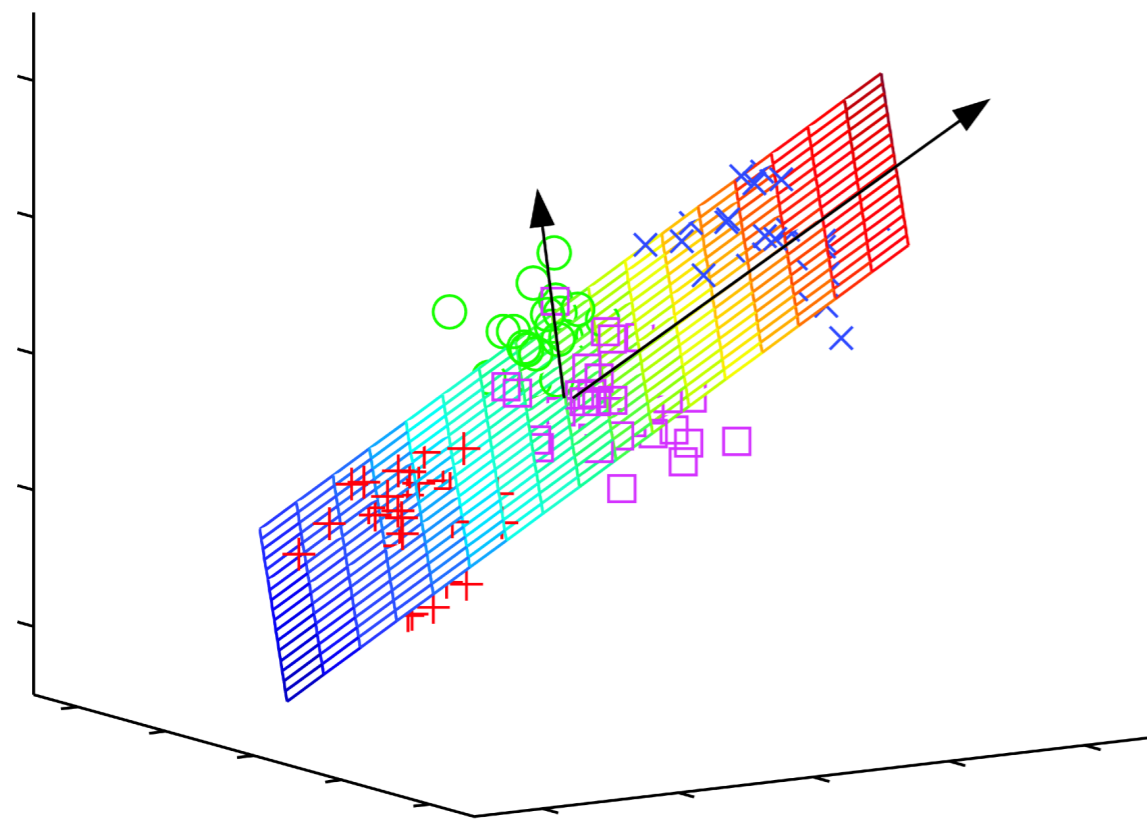
Тестовая: 10,000

784 признака

Алгоритмы

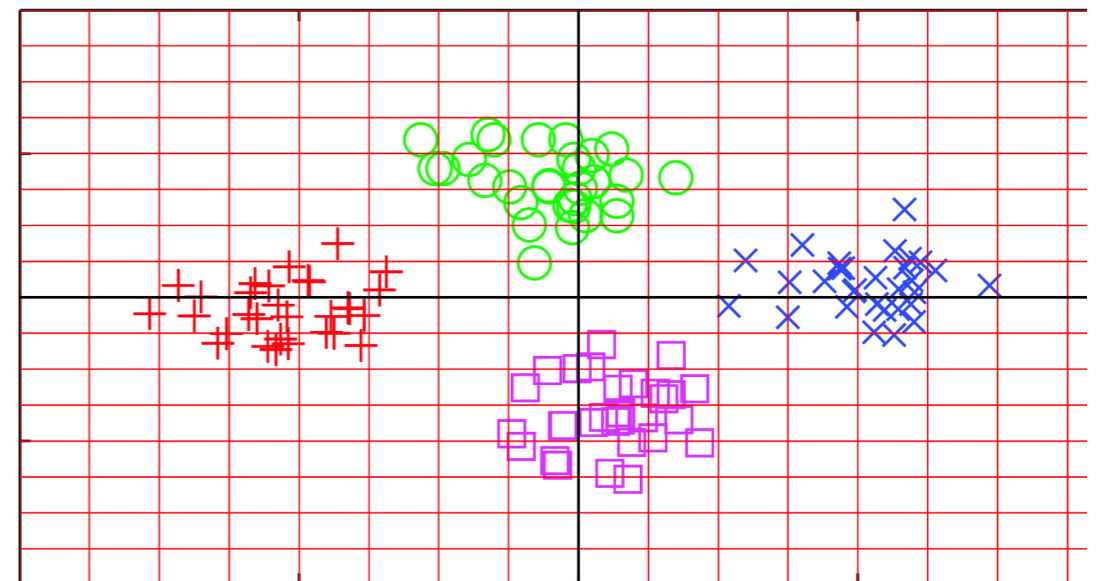
- PCA (Principal Component Analysis)
- ICA (Independent Component Analysis)
- SVD (Singular Value Decomposition)

PCA



исходное пространство данных

PCA
→



преобразованное пространство

ICA

- Пытается разделить данные на максимально независимые компоненты
- Критерий независимости — негауссовость
- Мера негауссовости — коэффициент эксцесса:

$$\text{kurt}(v) = \mathbb{E}[v^4] - 3(\mathbb{E}[v^2])^2$$

SVD

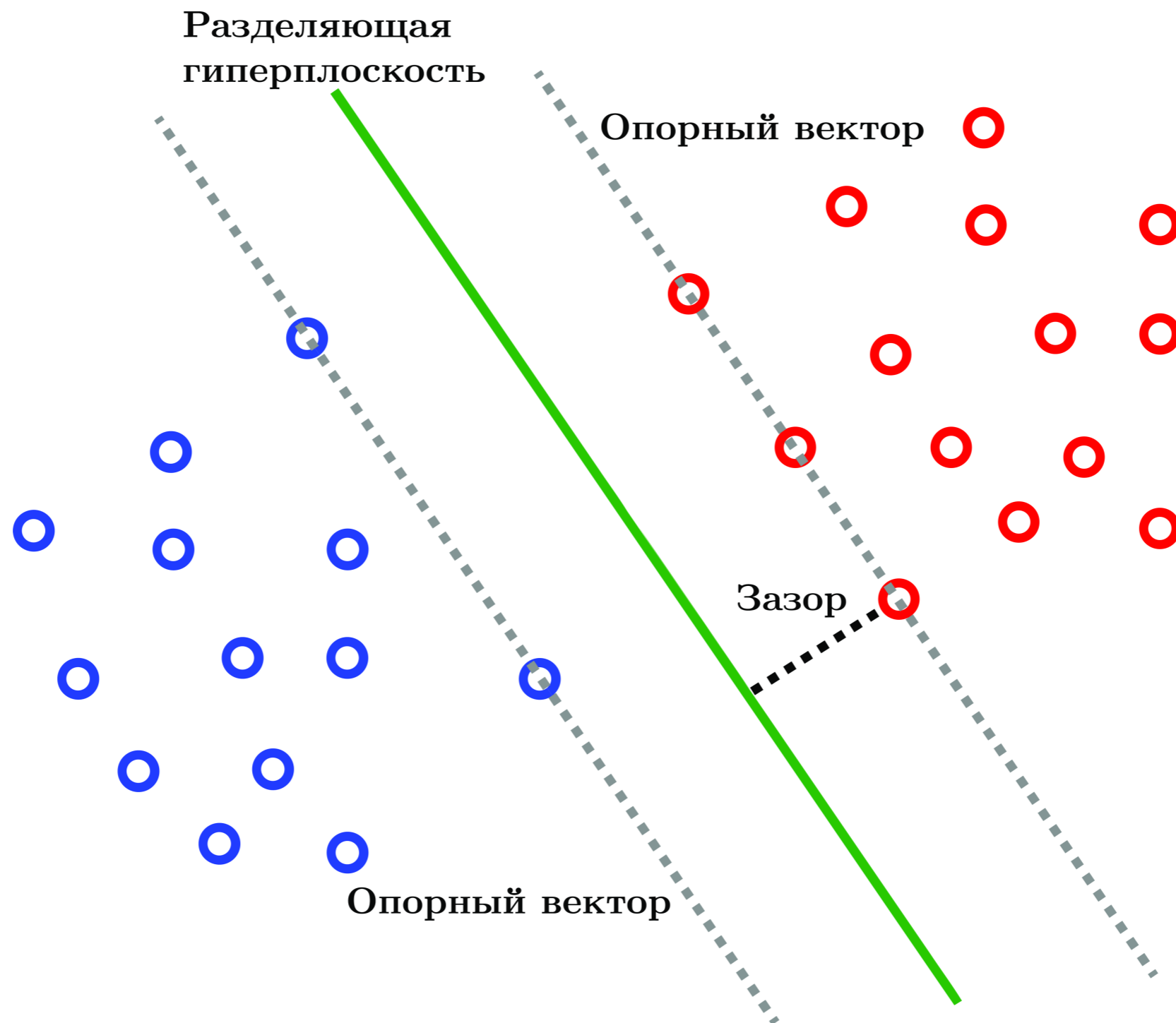
- Раскладываем матрицу X :

$$X = U\Sigma V^*$$

- Приближаем матрицей меньшего ранга:

$$X_k = U_k \Sigma_k V_k^*$$

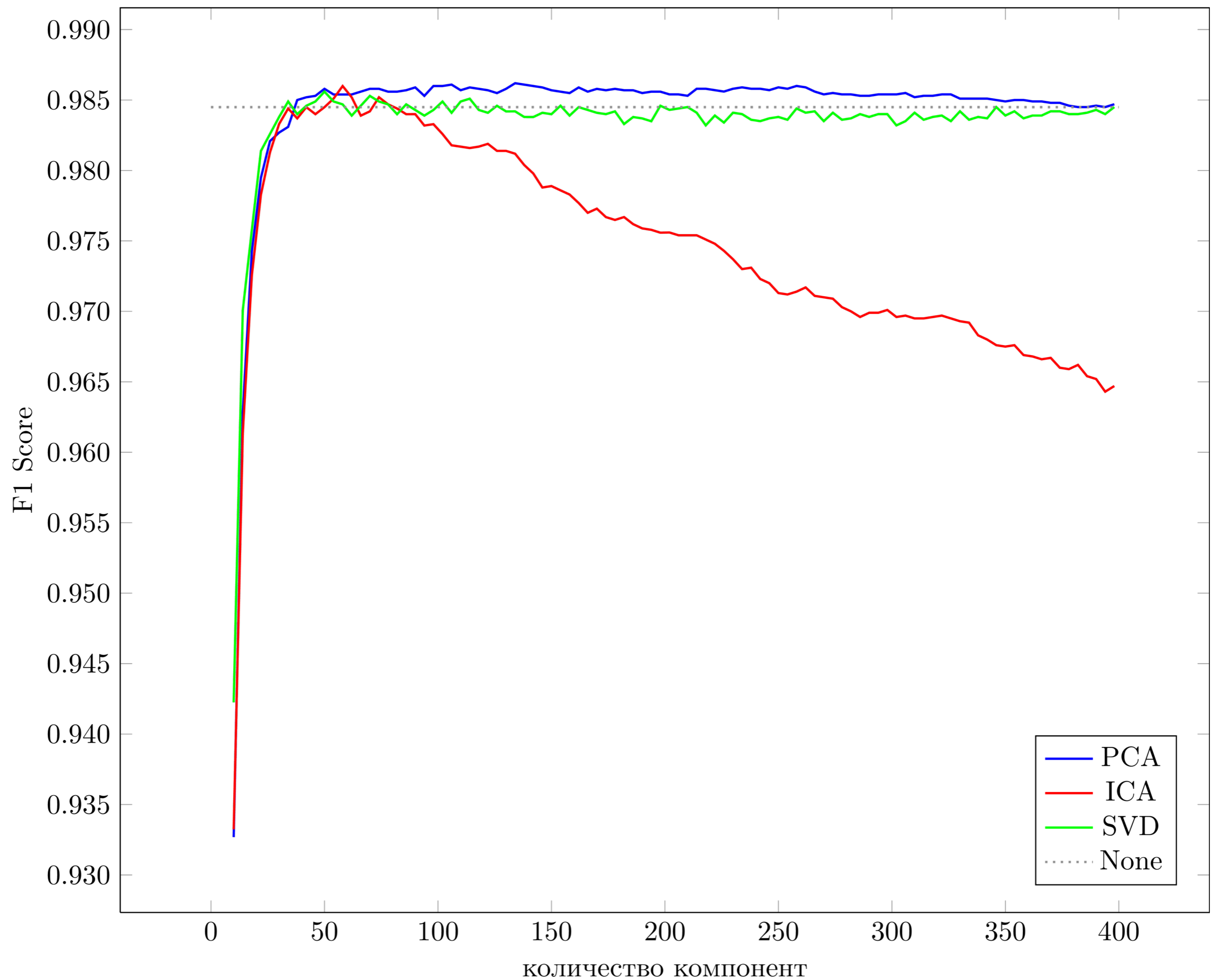
SVM



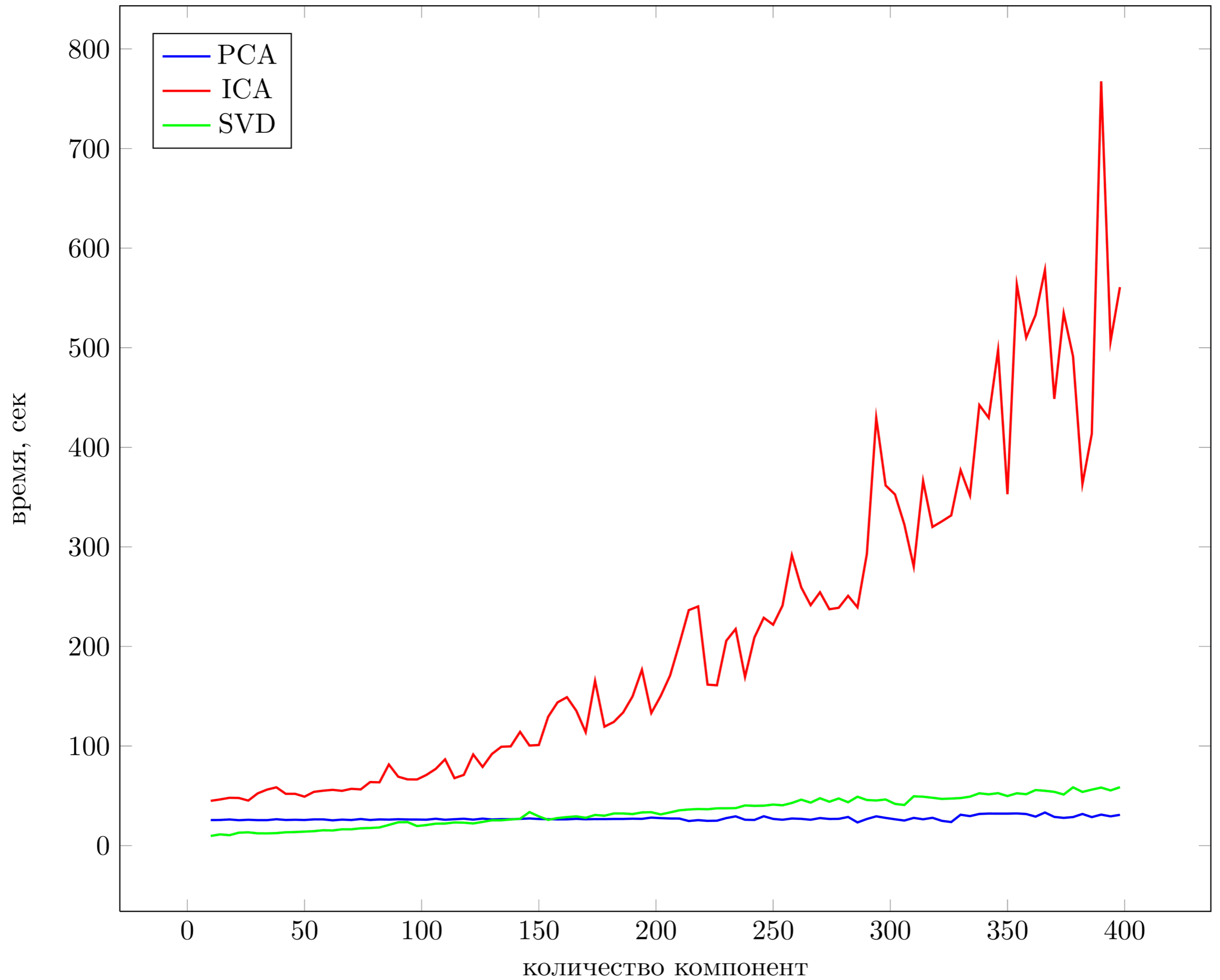
Методика сравнения

- Python 2.7, numpy, scipy, libsvm
- Intel Core i7 2.9 ГГц, 8 ГБ RAM
- $F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$, больше — лучше

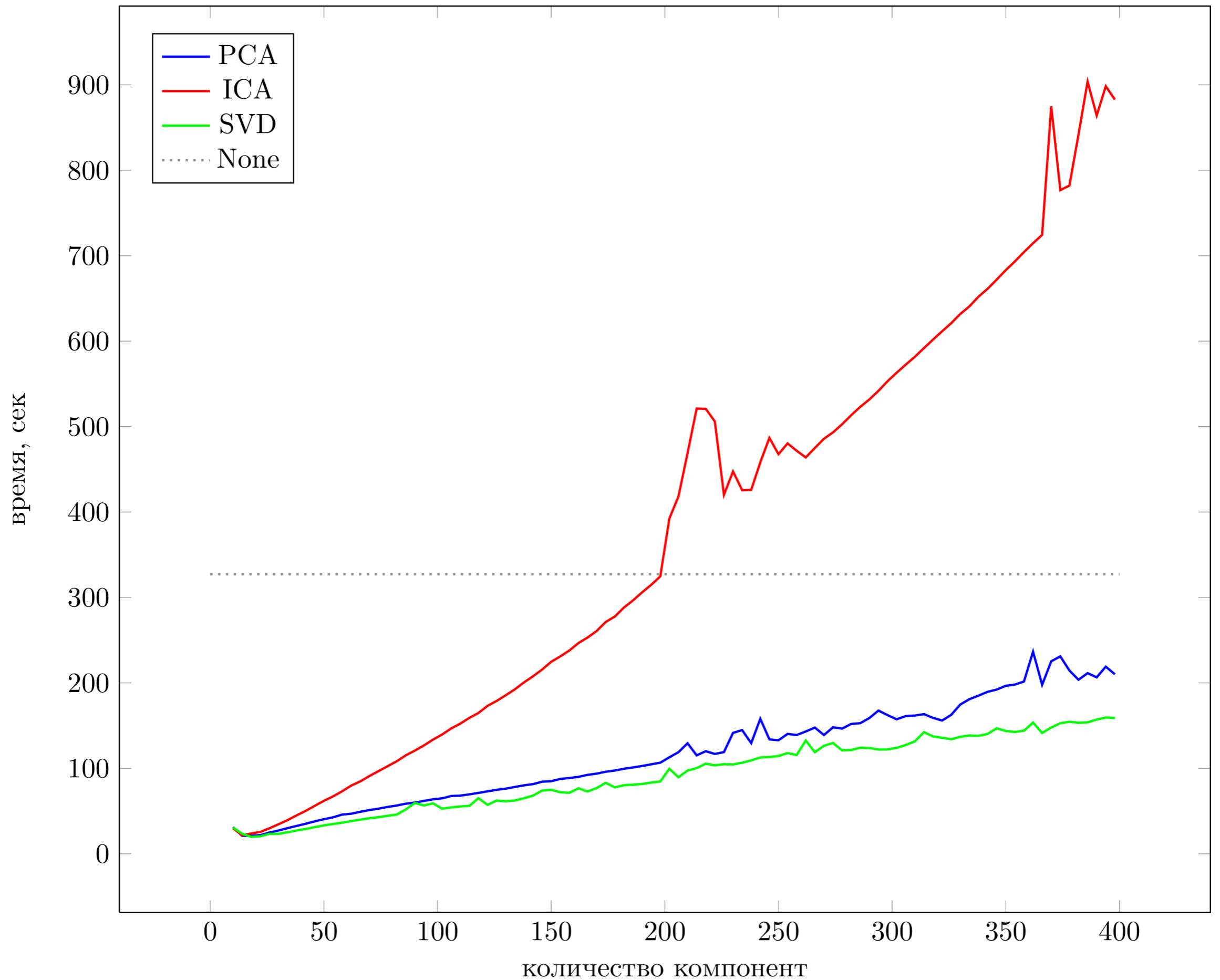
Точность



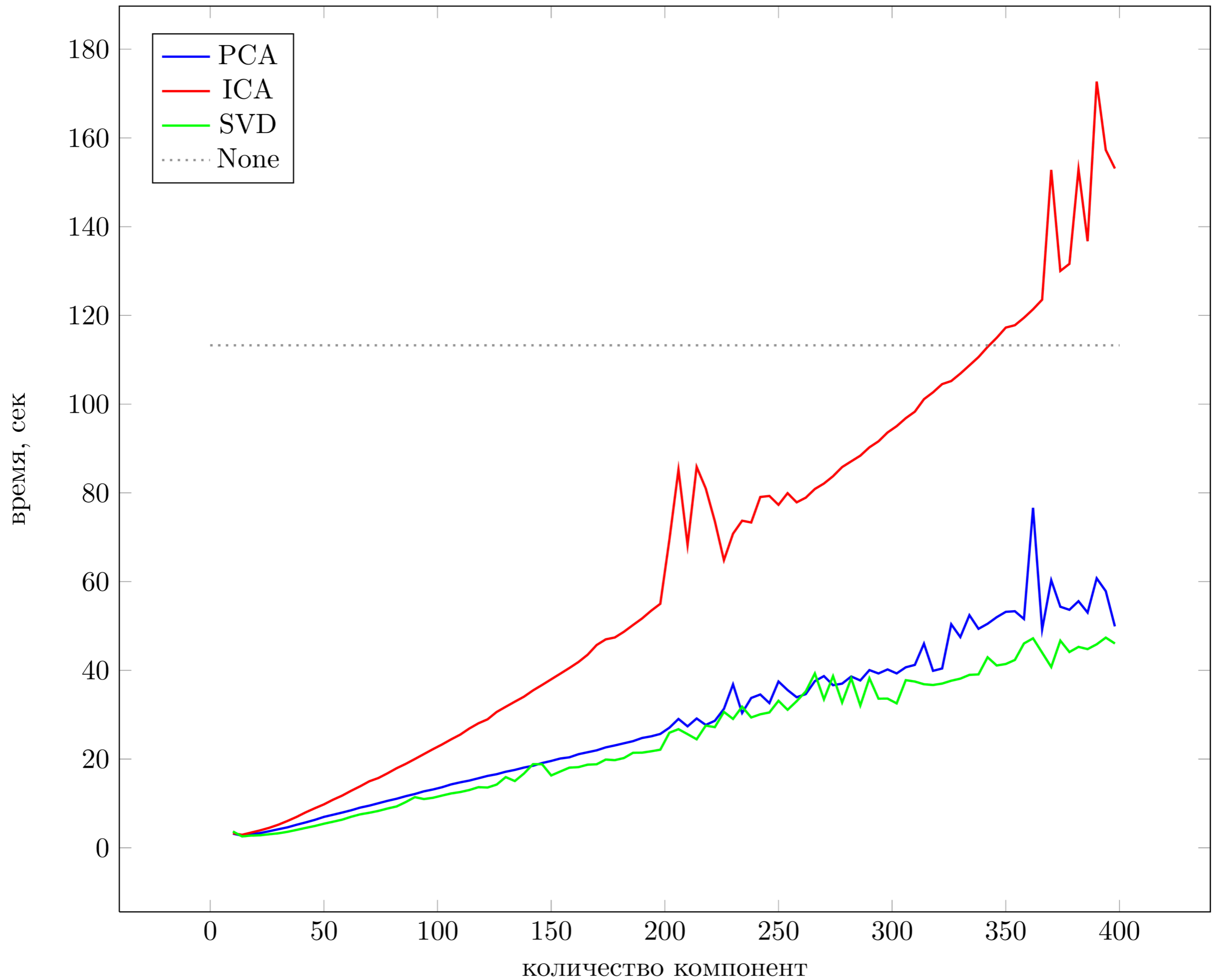
Время работы



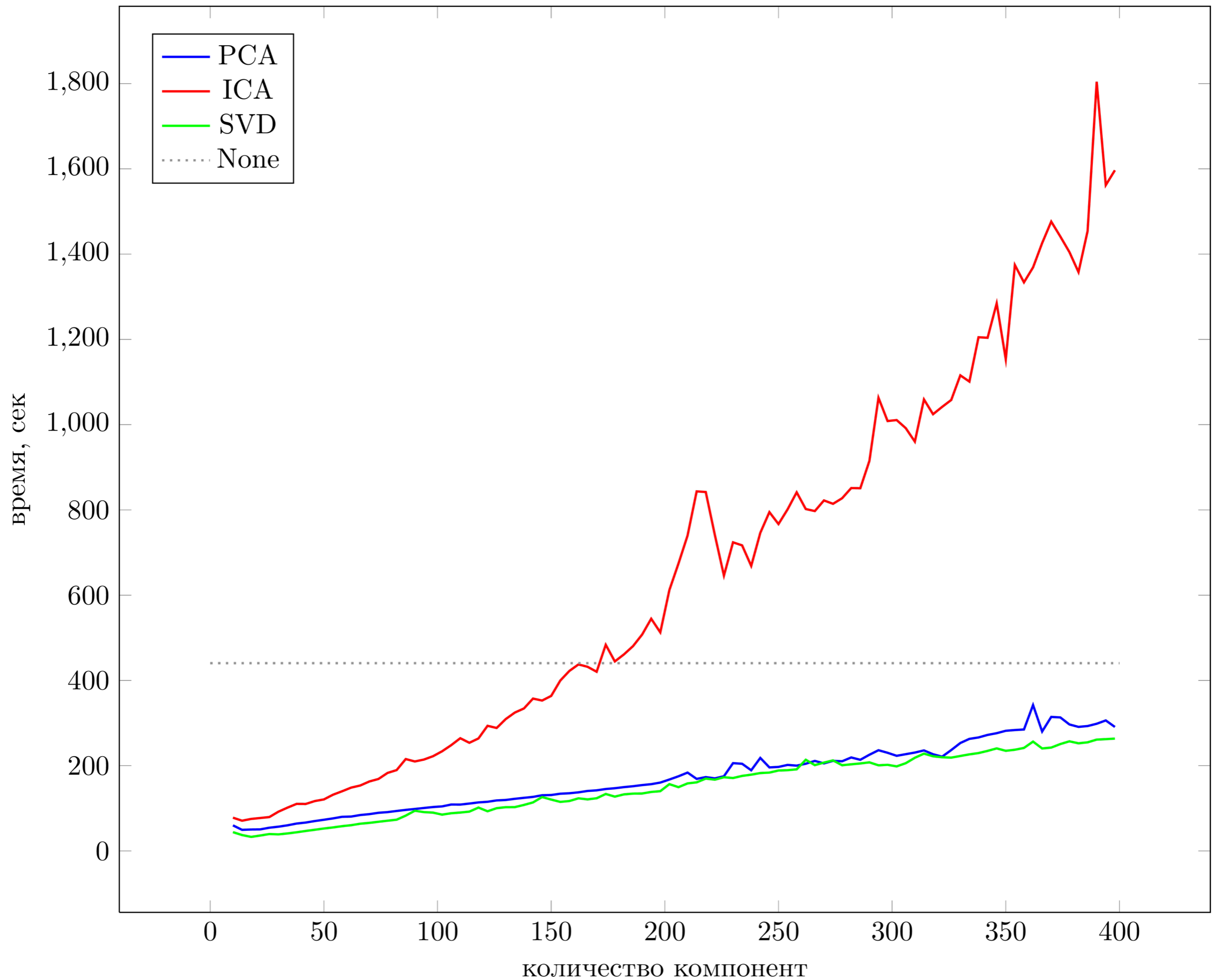
Время обучения



Время предсказания



Суммарное время работы



Полученные результаты

- Реализованы алгоритмы PCA, ICA, SVD
- Проведено сравнение их работы в зависимости от количества компонент