

**МИНОБРНАУКИ РОССИИ**  
**САНКТ-ПЕТЕРБУРГСКИЙ ГОСУДАРСТВЕННЫЙ**  
**ЭЛЕКТРОТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ**  
**«ЛЭТИ» ИМ. В.И. УЛЬЯНОВА (ЛЕНИНА)**  
**Кафедра МОЭВМ**

**ОТЧЕТ**  
**по лабораторной работе №8**  
**по дисциплине «Машинное обучение»**  
**Тема: Классификация (линейный дискриминантный анализ, метод**  
**опорных векторов)**

Студент гр. 6304

Ковынев М.В.

Преподаватель

Жангиров Т.Р.

Санкт-Петербург

2020

## Цель работы

Ознакомиться с методами классификации модуля *Sklearn*.

## Ход работы

### Загрузка данных

Датасет загружен в датафрейм. Вид данных представлен на рис. 1.

	0	1	2	3	4
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa
...	...	...	...	...	...
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

150 rows × 5 columns

Рисунок 1 – Исходные данные

Выделены данные и их метки, тексты меток преобразованы к числам. Выборка разбита на обучающую и тестовую *train\_test\_split*.

## Линейный дискриминантный анализ

1. Проведена классификация наблюдений с помощью *LinearDiscriminantAnalysis*. Выявлено 3 неправильно классифицированных наблюдения. Параметры классификатора представлены в табл. 1. Атрибуты классификатора представлены в табл.2.

Таблица 1 – Параметры *LinearDiscriminantAnalysis*

Параметр	Описание
solver	<ul style="list-style-type: none"><li>• «svd»: Разложение по сингулярным числам. Не вычисляет ковариационную матрицу, поэтому рекомендуется для данных с большим количеством признаков.</li><li>• «lsqr»: Решение наименьших квадратов, можно комбинировать с параметром <i>shrinkage</i>.</li><li>• «eigen»: Разложение на собственные значения, можно комбинировать с параметром <i>shrinkage</i>.</li></ul>
shrinkage	<ul style="list-style-type: none"><li>• «auto»: Автоматическое сжатие по лемме Ледуа-Вольфа.</li><li>• float from [0, 1]</li></ul>
priors	Класс априорных вероятностей. По умолчанию пропорции классов выводятся из данных обучения.
n_components	Количество компонентов ( $\leq \min(n\_classes - 1, n\_features)$ ) для уменьшения размерности. Если None, будет установлено значение $\min(n\_classes - 1, n\_features)$ . Этот параметр влияет только на метод преобразования <i>transform</i> .
store_covariance	Если True, явно вычислить взвешенную ковариационную матрицу внутри класса, когда решатель – «svd». Матрица всегда вычисляется и сохраняется для других решателей.

tol	Абсолютный порог для того, чтобы единичное значение X считалось значимым, используется для оценки ранга X. Измерения, единичные значения которых не значимы, отбрасываются. Используется только если решатель - «svd».
-----	--

Таблица 2 – Атрибуты *LinearDiscriminantAnalysis*

Атрибут	Описание
coef_	Весовые вектора.
intercept_	Массив прерывания.
covariance_	Взвешенная внутриклассовая ковариационная матрица.
explained_variance_ratio_	Процент дисперсии, объясняемой каждым из выбранных компонентов. Если n_components не задано, то все компоненты сохраняются, а сумма объясненных дисперсий равна 1,0. Доступно только при использовании собственного решателя или «svd».
means_	Средние в классах.
priors_	Вероятности классов.
scalings_	Масштабирование объектов в пространстве, охватываемом центроидами классов. Доступно только для решателей «svd» и «eigen».
xbar_	Общее среднее. Присутствует, только если решатель - «svd».
classes_	Уникальные метки классов.

2. Точность классификации получена с помощью функции score() и составляет 98%.

3. Построен график зависимости неправильно классифицированных наблюдений и точности классификации от размера тестовой выборки. График представлен на рис. 2.

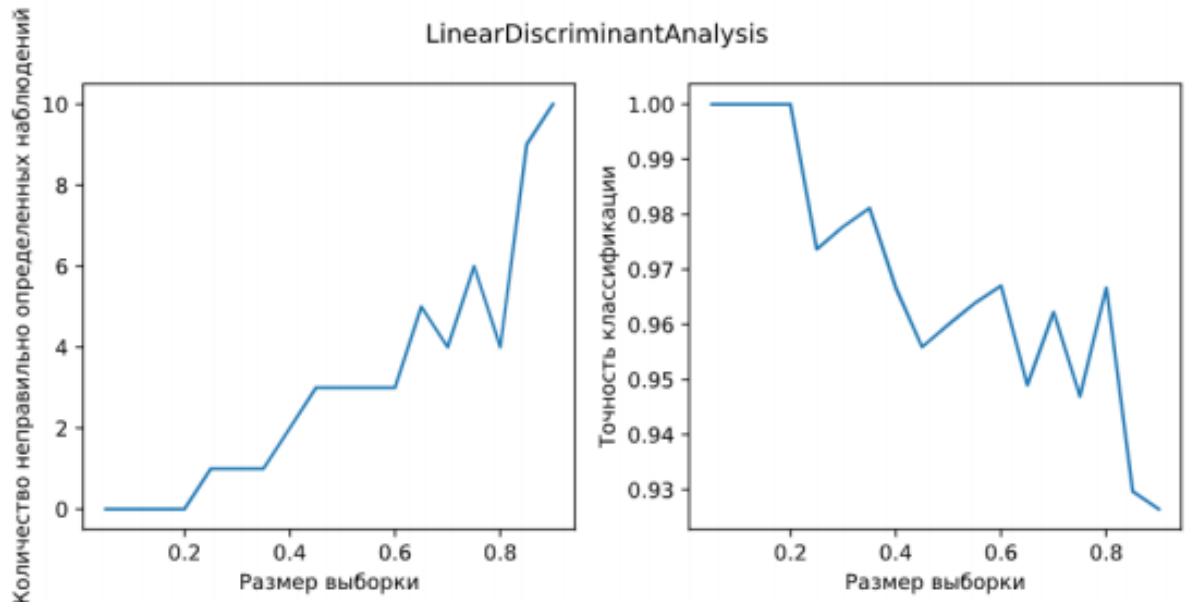


Рисунок 2 – Классификация *LinearDiscriminantAnalysis*

4. Функция `transform` проецирует данные для максимизации разбиения классов. LDA пытается определить атрибуты, на которые приходится наибольшая разница между классами. В частности, LDA, в отличие от PCA, является контролируемым методом, использующим известные метки классов.

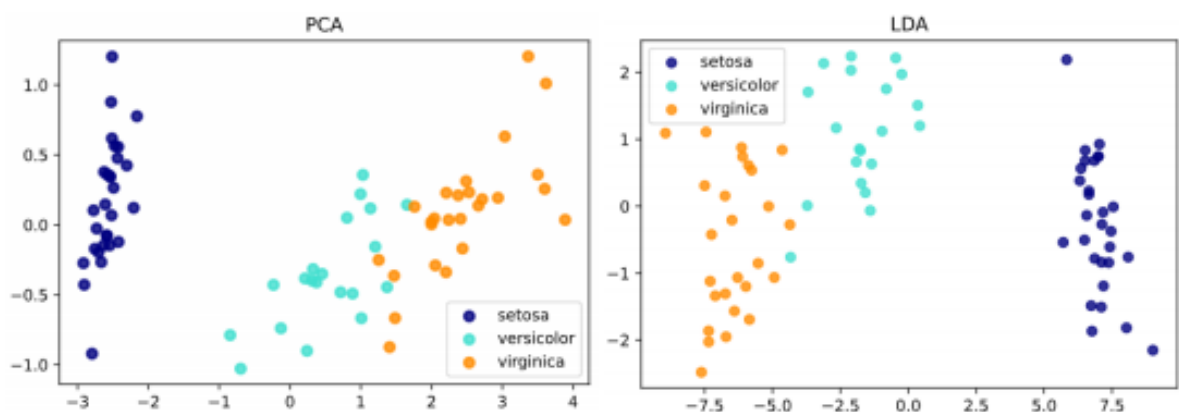


Рисунок 3 – Сравнение PCA и LDA

5. Работа классификатора исследована при различных параметрах solver, shrinkage. Результаты представлены на рис. 4-7.

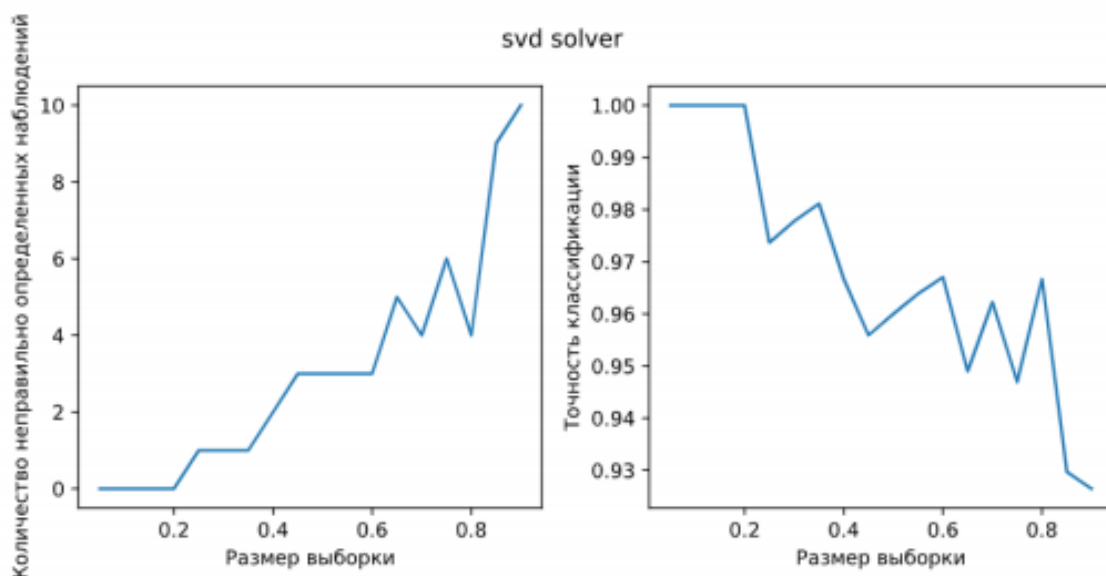


Рисунок 4 – Классификация *LinearDiscriminantAnalysis* с svd solver

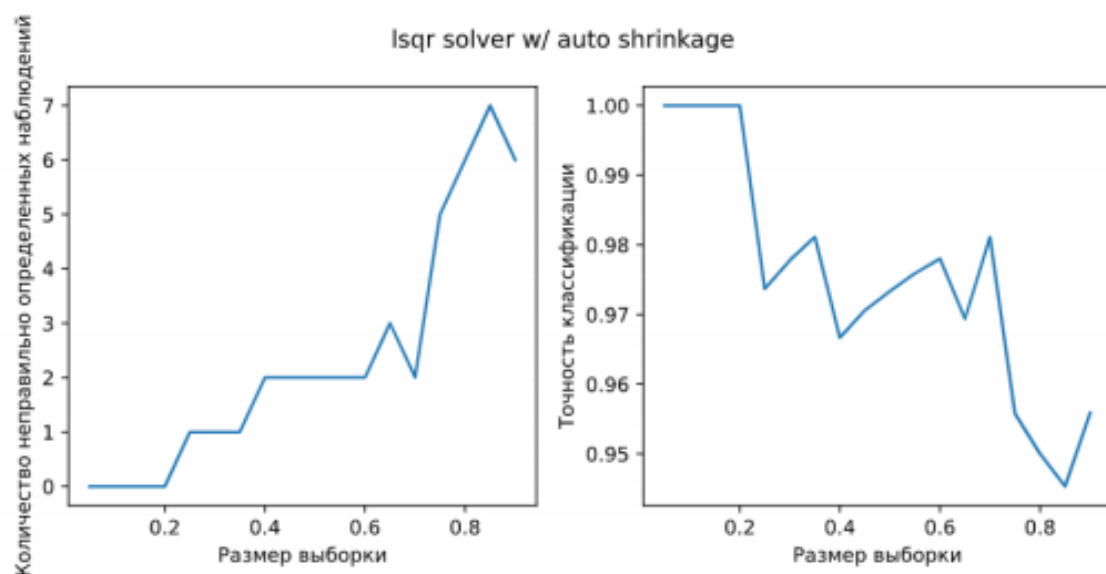


Рисунок 5 – Классификация *LinearDiscriminantAnalysis* с lsqr solver и shrinkage

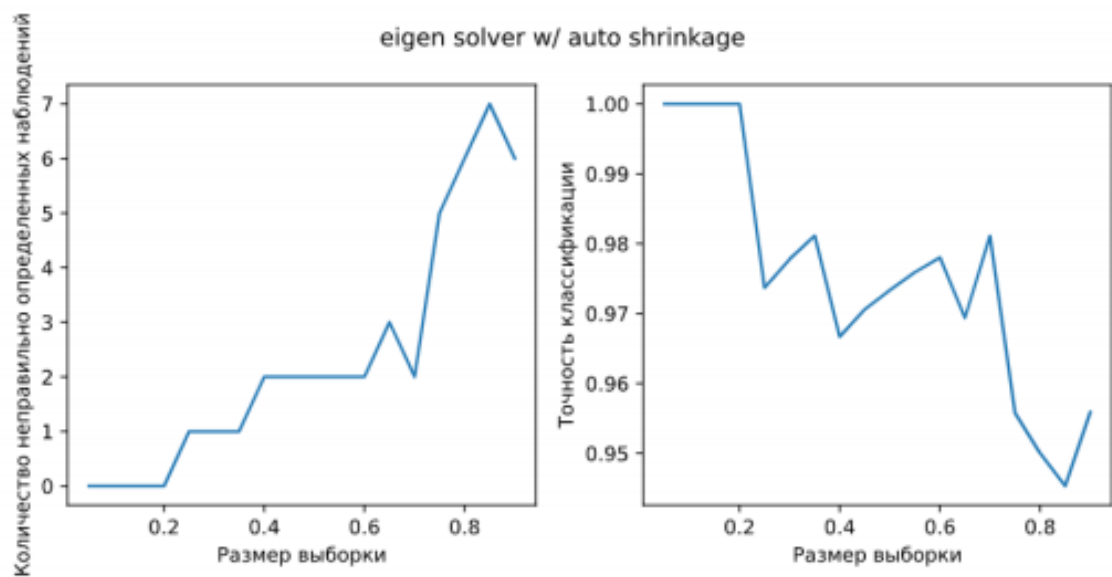


Рисунок 6 – Классификация *LinearDiscriminantAnalysis* с eigen solver и shrinkage

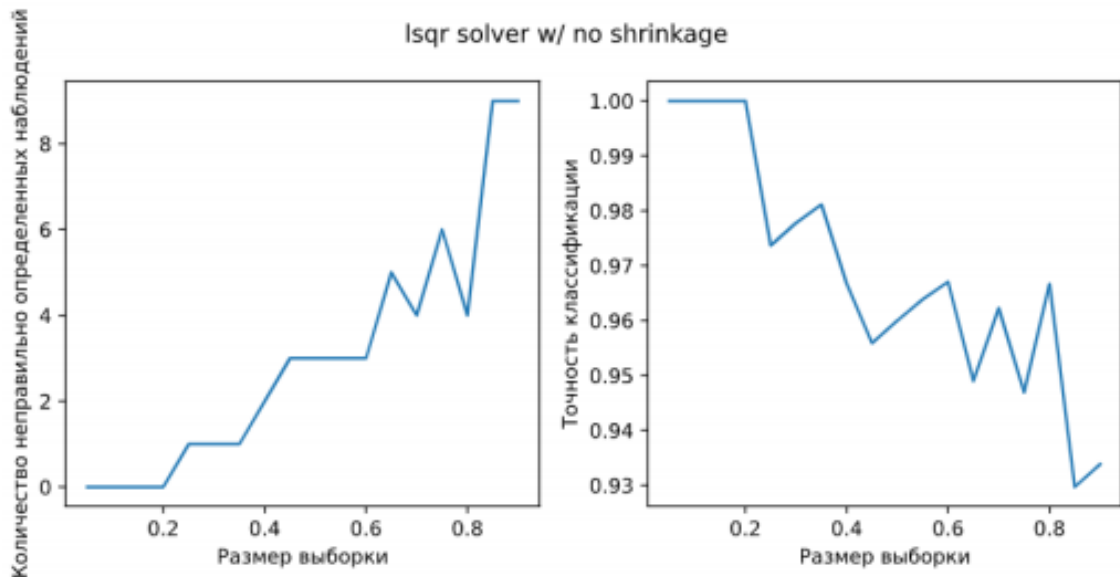


Рисунок 7 – Классификация *LinearDiscriminantAnalysis* с lsqr solver без shrinkage

6. Заданы собственные значения априорных вероятностей классов, результаты представлены в табл. 3 и на рис. 8.

Таблица 3 – Результаты классификации *LinearDiscriminantAnalysis*

Априорные вероятности классов	Количество неправильно определенных наблюдений	Точность классификации
[0.38666667, 0.26666667, 0.34666667]	3	0.987
[0.15, 0.7, 0.15]	5	0.987

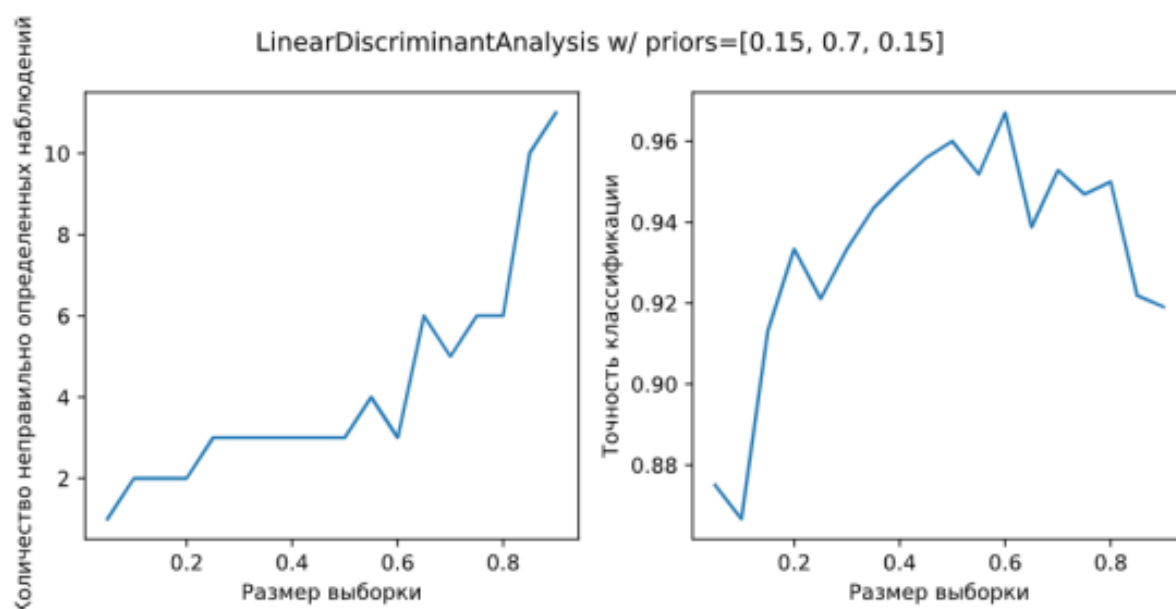


Рисунок 8 – Классификация *LinearDiscriminantAnalysis* с заданными априорными вероятностями

### Метод опорных векторов

1. Проведена классификация наблюдений с помощью метода опорных векторов на тех же данных. Выявлено 4 неправильно классифицированных наблюдения.
2. Точность классификации получена с помощью функции `score()` и составляет 96%.



3. Атрибут *support\_* хранит индексы опорных векторов, *support\_vectors\_* – сами опорные вектора, *n\_support\_* – количество опорных векторов для каждого класса.
4. Построен график зависимости неправильно классифицированных наблюдений и точности классификации от размера тестовой выборки. График представлен на рис. 9.

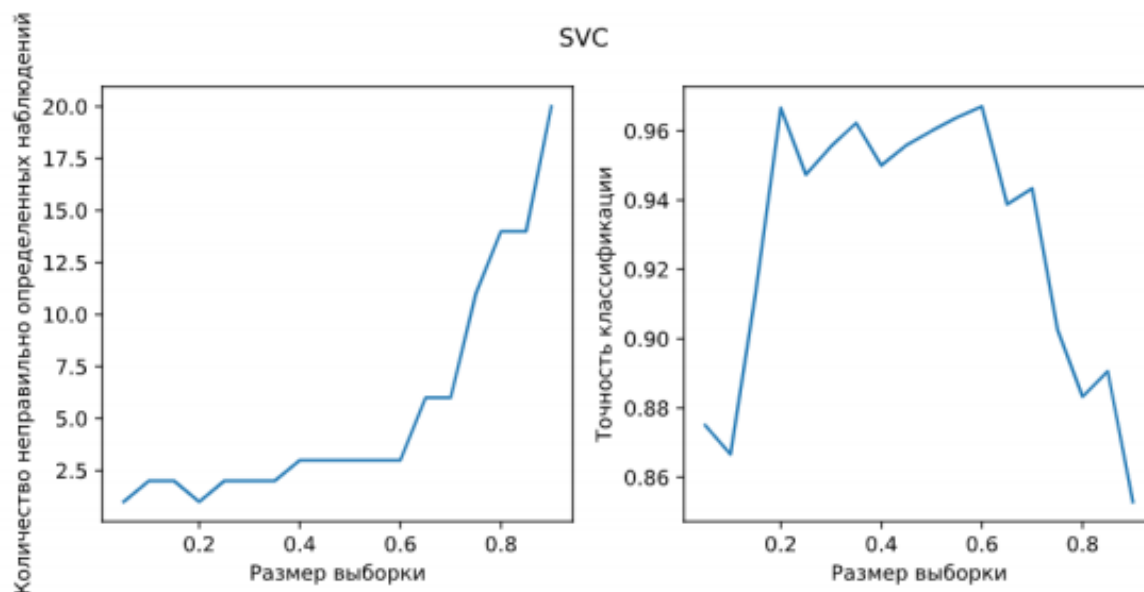


Рисунок 9 – Классификация *SVC*

5. Исследована работа метода опорных векторов при различных значениях параметров *kernel*, *degree*, *max\_iter*. Результаты представлены в табл. 4-6.

Таблица 4 – Результаты классификации *SVC*

Тип ядра	Количество неправильно определенных наблюдений	Точность классификации
linear	2	0.973
poly	6	0.987
rbf (default)	4	0.960
sigmoid	54	0.386

Таблица 5 – Результаты классификации *SVC*

Степень	Количество неправильно определенных наблюдений	Точность классификации
1	5	0.960
2	6	0.960
3	6	0.987
4	5	0.987
5	3	0.987

Таблица 6 – Результаты классификации *SVC*

Количество итераций	Количество неправильно определенных наблюдений	Точность классификации
Без ограничений	4	0.960
1	9	0.960
2	8	0.987
3	5	0.973
4	3	0.973
5	1	0.973
6	3	0.973

6. Классификация методами *NuSVC* и *LinearSVC* представлены на рис. 10-11.

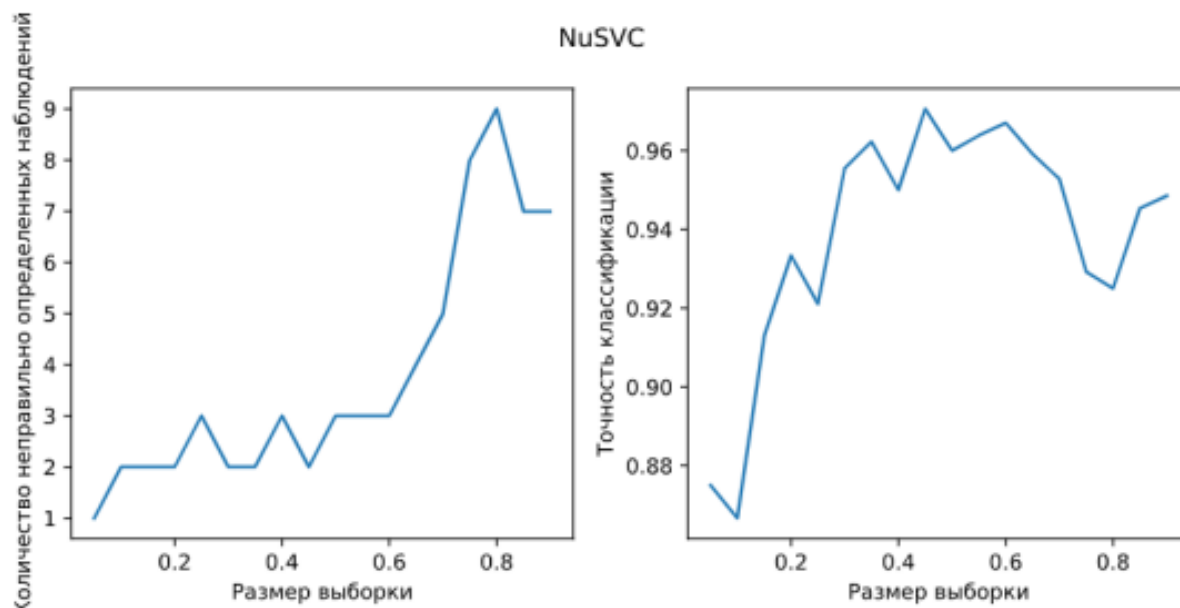


Рисунок 10 – Классификация *NuSVC*

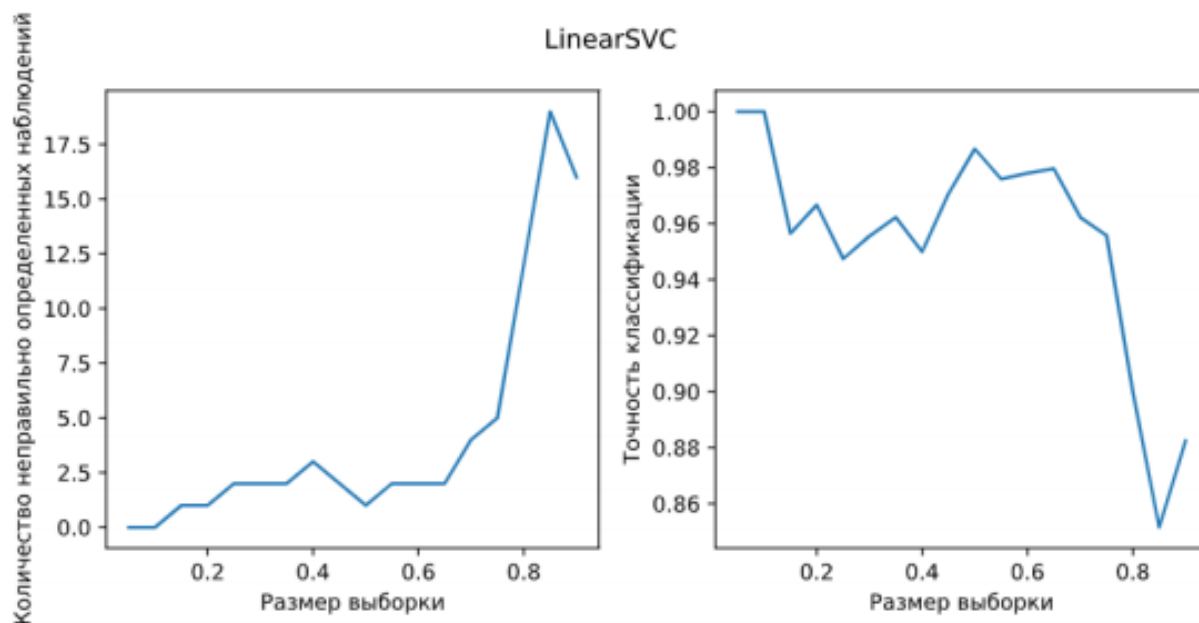


Рисунок 10 – Классификация *LinearSVC*

*NuSVC* подобен *SVC*, но использует параметр для управления количеством опорных векторов.

*LinearSVC* аналогично *SVC* с линейным ядром, но лучше масштабируется для большого числа выборок.

## **Выводы**

В ходе лабораторной работы рассмотрены такие методы классификации модуля *Sklearn*, как *LinearDiscriminantAnalysis*, *SVC*, *NuSVC* и *LinearSVC*.