

# Прикладная математика

## Лабораторная работа №8

### SVM

М33001  
ВОЕВОДСКИЙ ДМИТРИЙ  
ЕВТУШЕНКО ИВАН  
БЛАЖКОВ АЛЕКСАНДР

2022 год, университет ИТМО

```

1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from tqdm.notebook import tqdm, trange
6
7 np.random.seed(42)

```

Машина опорных векторов отлично показывает себя в случаях бинарной классификации, когда выборка не является линейно разделимой.

Для примера сгенерируем данные с признаками  $(x, y)$  и границей разделения  $x^2 + y^2 = \frac{1}{4}$

```

1 def generate(size=5000, a=-1, b=1, threshold=1/4):
2     coords = np.random.rand(size, 2)
3     coords = coords * abs(b - a) + a
4     marks = np.array(list(map(lambda x: int(x[0] ** 2 + x[1] ** 2 > threshold), coords)))
5     if np.unique(marks).size == 1:
6         return generate(size, a, b, threshold)
7     df = pd.DataFrame(coords, columns=['x', 'y'])
8     df['mark'] = marks
9     return df

```

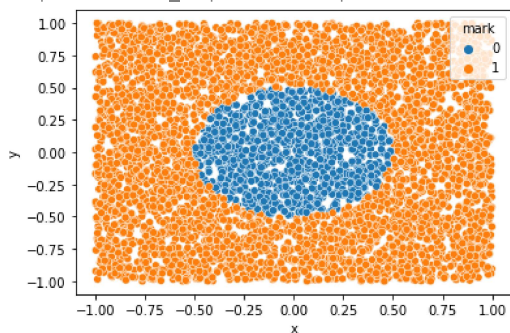
```
1 df = generate()
```

```
1 df.head()
```

	x	y	mark
0	-0.250920	0.901429	1
1	0.463988	0.197317	1
2	-0.687963	-0.688011	1
3	-0.883833	0.732352	1
4	0.202230	0.416145	0

```
1 sns.scatterplot(data=df, x='x', y='y', hue='mark')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f2b0a636af0>



Подготовим вспомогательные методы для получения метрик качества будущих моделей.

```

1 class Metrics:
2     @staticmethod
3     def get_metrics_names():
4         return ['accuracy', 'precision', 'recall', 'f1']
5
6     @staticmethod
7     def get_metrics(y_true, y_pred):
8         return {
9             'accuracy': Metrics.accuracy(y_true, y_pred),
10            'precision': Metrics.precision(y_true, y_pred),
11            'recall': Metrics.recall(y_true, y_pred),

```

```

12         'f1': Metrics.f1(y_true, y_pred)
13     }
14
15     @staticmethod
16     def accuracy(y_true, y_pred):
17         return np.mean(y_true == y_pred)
18
19     @staticmethod
20     def precision(y_true, y_pred):
21         tp = np.sum(y_true * y_pred)
22         fp = np.sum((1 - y_true) * y_pred)
23         return tp / (tp + fp)
24
25     @staticmethod
26     def recall(y_true, y_pred):
27         tp = np.sum(y_true * y_pred)
28         fn = np.sum(y_true * (1 - y_pred))
29         return tp / (tp + fn)
30
31     @staticmethod
32     def f1(y_true, y_pred):
33         p = Metrics.precision(y_true, y_pred)
34         r = Metrics.recall(y_true, y_pred)
35         return 2 * p * r / (p + r)
36
37     @staticmethod
38     def report(y_true, y_pred):
39         print(f'accuracy:  {Metrics.accuracy(y_true, y_pred):.4f}')
40         print(f'precision: {Metrics.precision(y_true, y_pred):.4f}')
41         print(f'recall:    {Metrics.recall(y_true, y_pred):.4f}')
42         print(f'f1:       {Metrics.f1(y_true, y_pred):.4f}')

```

Импортируем необходимые инструменты из библиотеки `sklearn`.

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.svm import NuSVC
3 from sklearn.model_selection import GridSearchCV

```

Разделяем данные на данные для обучения и для тестирования.

```

1 X_train, X_test, y_train, y_test = train_test_split(df[['x', 'y']], df['mark'], test_size=0.2)

```

Определим, какое ядро модели лучше всего решает задачу классификации в нашем случае.

- linear
- poly
- rbf
- sigmoid

```

1 fig, axs = plt.subplots(ncols=2, nrows=2, figsize=(15, 15))
2 for i, kernel in enumerate(['linear', 'poly', 'rbf', 'sigmoid']):
3     model = NuSVC(kernel=kernel, nu=0.4, degree=2)
4     model.fit(X_train, y_train)
5     y_pred = model.predict(X_test)
6     print(kernel)
7
8     res_df = pd.DataFrame(X_test, columns=['x', 'y'])
9     res_df['y_pred'] = y_pred
10
11     sns.scatterplot(data=res_df, x='x', y='y', hue='y_pred', ax=axs[i // 2, i % 2])
12     axs[i // 2, i % 2].set_title(kernel)
13     x_min, x_max = X_train['x'].min() - 0.1, X_train['x'].max() + 0.1
14     y_min, y_max = X_train['y'].min() - 0.1, X_train['y'].max() + 0.1
15     xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.01), np.arange(y_min, y_max, 0.01))
16     Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
17     Z = Z.reshape(xx.shape)
18     axs[i // 2, i % 2].contourf(xx, yy, Z, alpha=0.4, cmap='coolwarm')
19
20
21 Metrics.report(y_test, y_pred)
22 print()

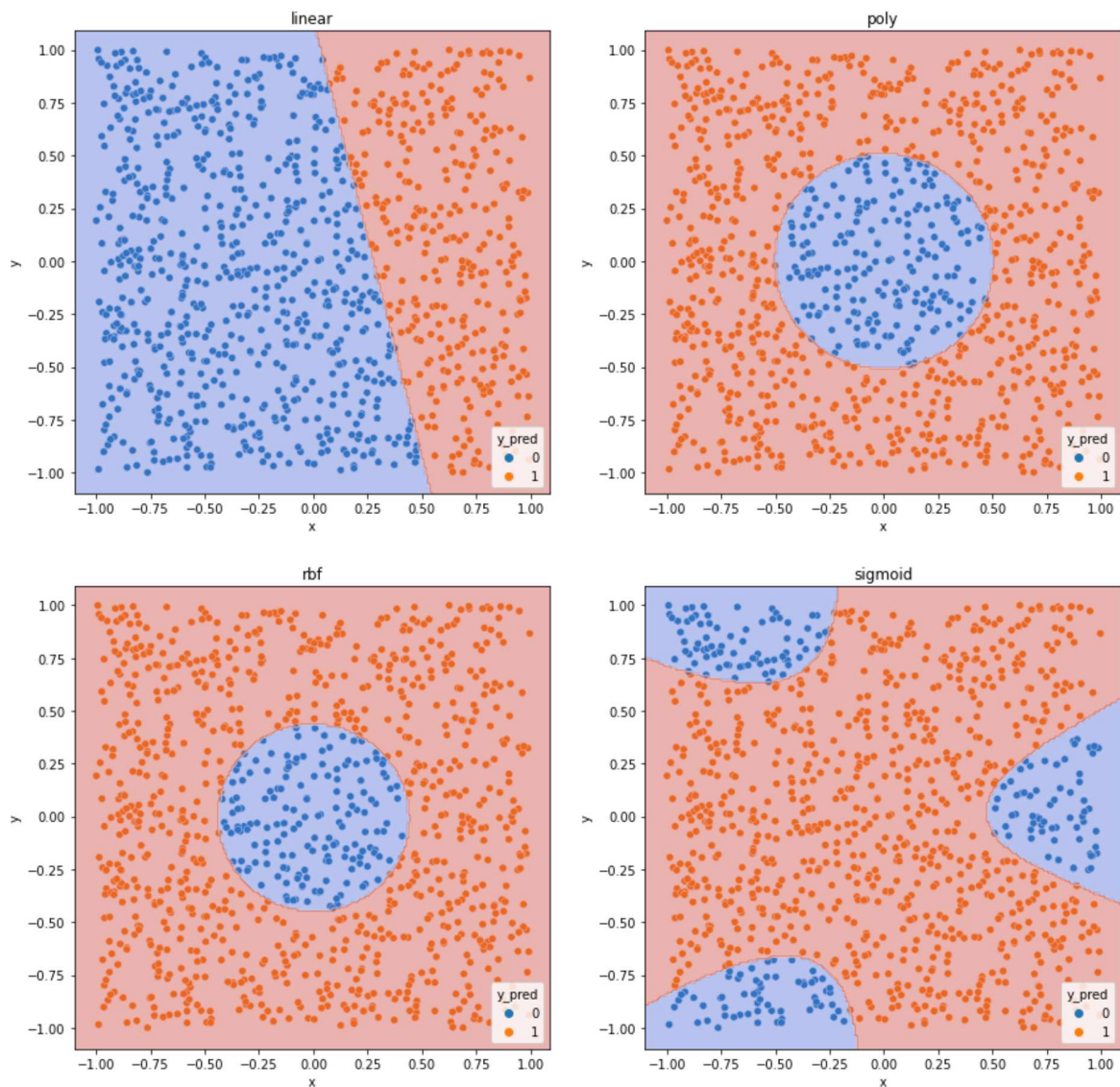
```

```
linear
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:450: UserWarning: X does not have valid feature names, but NuSVC was fitted with
  warnings.warn(
accuracy: 0.4590
precision: 0.9043
recall: 0.3649
f1: 0.5200
```

```
poly
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:450: UserWarning: X does not have valid feature names, but NuSVC was fitted with
  warnings.warn(
accuracy: 0.9930
precision: 1.0000
recall: 0.9913
f1: 0.9956
```

```
rbf
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:450: UserWarning: X does not have valid feature names, but NuSVC was fitted with
  warnings.warn(
accuracy: 0.9680
precision: 0.9617
recall: 1.0000
f1: 0.9805
```

```
sigmoid
/usr/local/lib/python3.8/dist-packages/sklearn/base.py:450: UserWarning: X does not have valid feature names, but NuSVC was fitted with
  warnings.warn(
accuracy: 0.5980
precision: 0.7522
recall: 0.7447
f1: 0.7484
```

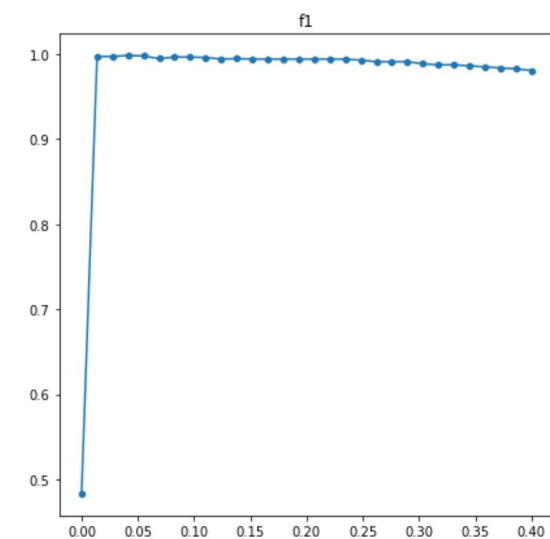
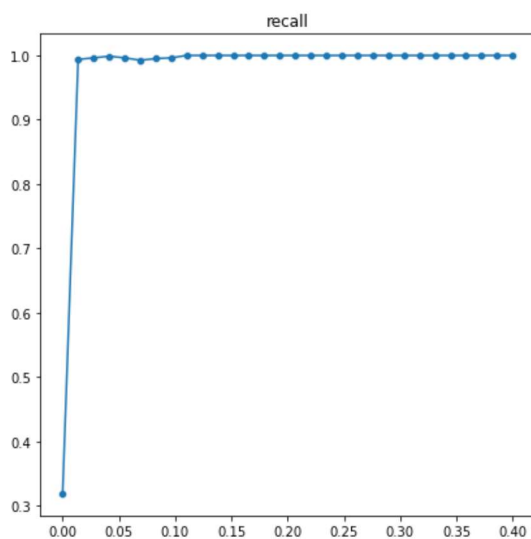
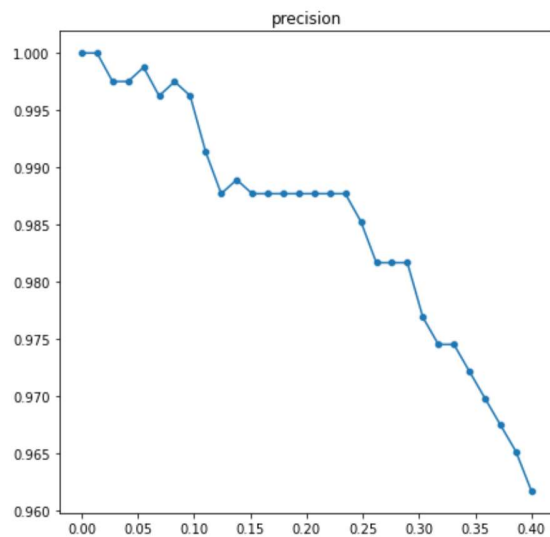
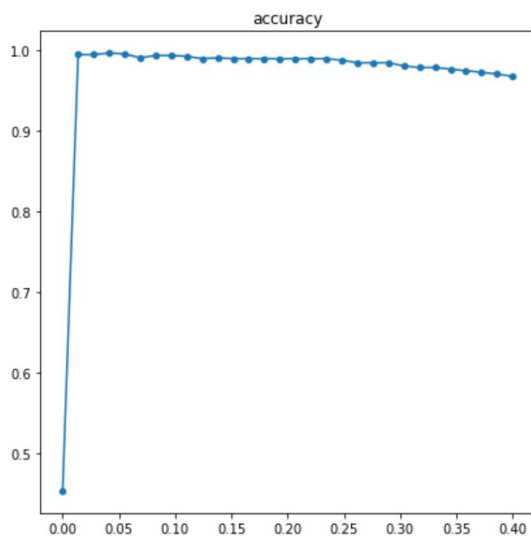


Как видим, правильная классификация может быть осуществлена с помощью ядер `poly` и `rbf`. Далее будем рассматривать ядро `rbf`, так как именно это нужно делать по плану лабораторной работы!

Double-click (or enter) to edit

Определим, как метрики качества зависят от нижней границы доли опорных векторов  $\nu$ .

```
1 nus = np.linspace(1e-10, 0.4, 30)
2 metrics = {}
3 for metric in Metrics.get_metrics_names():
4     metrics[metric] = []
5
6 for nu in nus:
7     model = NuSVC(kernel='rbf', nu=nu)
8     model.fit(X_train, y_train)
9     y_pred = model.predict(X_test)
10    calc_metrics = Metrics.get_metrics(y_test, y_pred)
11    for metric in calc_metrics.keys():
12        metrics[metric].append(calc_metrics[metric])
13
14 fig, axs = plt.subplots(ncols=2, nrows=2, figsize=(15, 15))
15 for i, metric in enumerate(metrics.keys()):
16     sns.lineplot(x=nus, y=metrics[metric], ax=axs[i // 2, i % 2])
17     sns.scatterplot(x=nus, y=metrics[metric], ax=axs[i // 2, i % 2])
18     axs[i // 2, i % 2].set_title(metric)
```



Как видно по графикам, с увеличением доли опорных векторов качество модели снижается. Скорее всего, это связано с т.н. переобучением. Переобучение в этом случае, это не "обучение заново", как могли подумать некоторые читатели, а феномен "хороших" показателей метрик на данных для обучения, и не таких "хороших" показателей на тестовых данных.

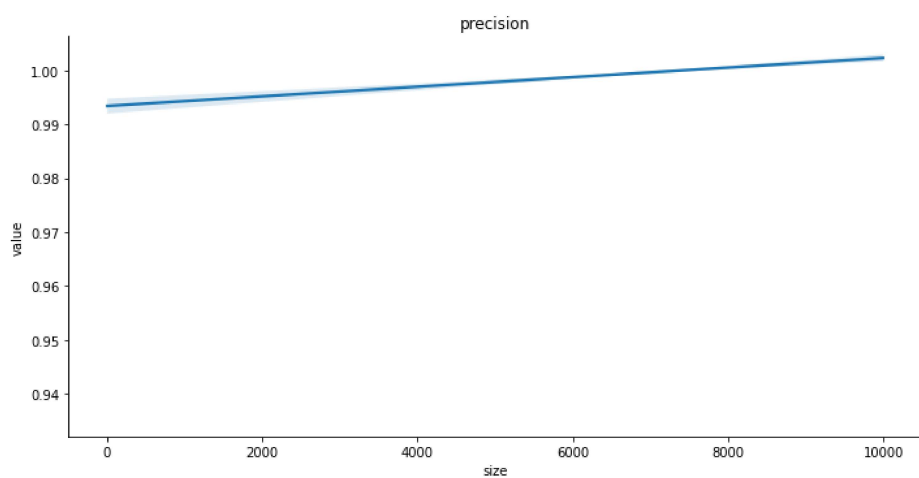
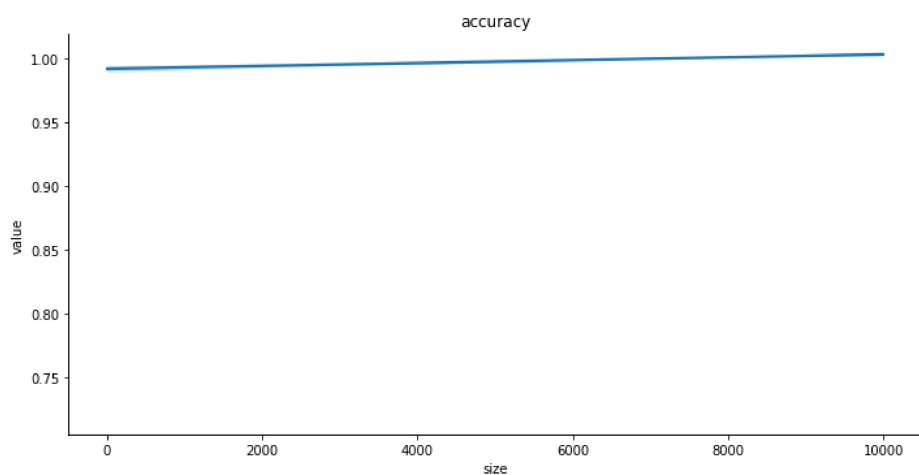
```
1 def smooth_line(y, window=10):
2     return np.convolve(y, np.ones(window) / window, mode='valid')
```

Посмотрим, как размер обучающей выборки влияет на метрики качества.

```
1 metrics = {}
2
3 import warnings
4 warnings.filterwarnings("ignore")
5
6 test = generate(50)
7 X_test, y_test = test[['x', 'y']], test['mark']
8 # nus = []
9
10 for metric in Metrics.get_metrics_names():
11     metrics[metric] = []
12
13 sizes = range(10, 10000, 10)
14
15 for size in tqdm(sizes):
16     model = NuSVC(kernel='rbf')
17     train = generate(size)
18
19     # search = GridSearchCV(model, {'nu': np.linspace(1e-10, 0.4, 30)}, cv=5, scoring='f1', refit=True, n_jobs=-1, verbose=0).fit(train[['x', 'y']])
20     # nus.append(search.best_params_['nu'])
21
22     model = NuSVC(kernel='rbf', nu=0.1)
23
24     X_train, y_train = train[['x', 'y']], train['mark']
25     model.fit(X_train, y_train)
26
27     y_pred = model.predict(X_test)
28
29     calc_metrics = Metrics.get_metrics(y_test, y_pred)
30
31     for metric in calc_metrics.keys():
32         metrics[metric].append(calc_metrics[metric])
```

100% 999/999 [05:37<00:00, 1.08it/s]

```
1 for i, metric in enumerate(metrics.keys()):
2     data = pd.DataFrame({'size': sizes, 'value': metrics[metric]})
3     sns.lmplot(data=data, x='size', y='value', height=5, aspect=2, markers='')
4     plt.title(metric)
```



Невероятно, но факт: чем больше размер тестовых данных, тем лучше показатели метрик. Однако в нашем случае ощутимого роста, по крайней мере, визуально не наблюдается. Это связано с тем, что в сгенерированных данных отсутствует шум.

