

# Машинное обучение, DS-поток

## Домашнее задание 7

### Правила:

- Дедлайн **10 апреля 02:00** (по состоянию на момент выдачи). После дедлайна работы не принимаются кроме случаев наличия уважительной причины.
- Выполненную работу нужно отправить на почту `mipt.stats@yandex.ru`, указав тему письма "[ml] Фамилия Имя - задание 7". Квадратные скобки обязательны. Если письмо дошло, придет ответ от автоответчика.
- Прислать нужно ноутбук и его pdf-версию (без архивов). Названия файлов должны быть такими: `7.N.ipynb` и `7.N.pdf`, где `N` - ваш номер из таблицы с оценками.
- Теоретические задачи необходимо оформить в теке/markdown или же прислать фотку в правильной ориентации рукописного решения, **где все четко видно**.
- Решения, размещенные на каких-либо интернет-ресурсах не принимаются. Кроме того, публикация решения в открытом доступе может быть приравнена к предоставлению возможности списать.
- Для выполнения задания используйте этот ноутбук в качестве основы, ничего не удаляя из него.
- Никакой код из данного задания при проверке запускаться не будет.

### Баллы за задание:

- Задача 1 - 3 балла
- Задача 2 - 4 балла
- Задача 3 - 3 балла
- Задача 4 - 3 балла
- Задача 5 - 5 баллов
- Задача 5 - 6 баллов

In [ ]:

```
1 import os
2
3 import numpy as np
4 import scipy.stats as sps
5 import pandas as pd
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8
9 from sklearn.pipeline import Pipeline
10 from sklearn.svm import SVC, SVR
11 from sklearn.metrics import accuracy_score, mean_absolute_error
12 from sklearn.model_selection import train_test_split, RandomizedSearchCV
13 from sklearn.preprocessing import StandardScaler
14 from sklearn.datasets import make_moons, make_blobs, make_circles
15
16 sns.set(font_scale=1.2, palette='Set2')
```

### Задача 1.

Рассмотрим следующие функции

- $K(x, z) = \langle x, z \rangle^k$ , где  $x, z \in \mathcal{X} = \mathbb{R}^2, k \in \mathbb{N}$
- $K(x, z) = \exp(-\gamma \|x - z\|^2)$ , где  $x, z \in \mathcal{X} = \mathbb{R}^2, \gamma > 0$

Для этих ядер найдите спрямляющее пространство  $\mathcal{H}$  и отображение  $\psi : \mathcal{X} \rightarrow \mathcal{H}$ , для которого верно  $K(x, z) = \langle \psi(x), \psi(z) \rangle$ . Какой вид имеет классификатор SVM, основанный на этих ядрах?

## Задача 2.

Какие из следующих функций являются ядрами?

- $K(x, z) = \cos(x - z)$  где  $x, z \in \mathbb{R}$
- $K(x, z) = \cos^2(x - z)$  где  $x, z \in \mathbb{R}$
- $K(x, z) = \sin^2(x - z)$  где  $x, z \in \mathbb{R}$
- $K(x, z) = \min\{x, z\}$ , где  $x, z \in \mathbb{R}_+$

## Задача 3.

Рассмотрите следующую двумерную выборку

In [ ]:

```
1 X, Y = make_circles(n_samples=100)
```

Изобразите эти выборки на плоскости, используя Y в качестве цвета точек.

In [ ]:

```
1
```

Рассмотрим ядро  $K(x, z) = \langle x, z \rangle^2$ . Переведите точки выборки в соответствующее данному ядру спрямляющее пространство и визуализируйте все проекции на пары координат. Рисуйте эти графики в строчку.

In [ ]:

```
1
```

Добавьте к признакам некоторый сдвиг и повторите те же действия. Что можно наблюдать?

In [ ]:

```
1
```

Повторите те же действия с RBF-ядром  $K(x, z) = \exp(-\gamma \|x - z\|^2)$ , рассмотрев проекции на пары координат спрямляющего пространства, которые соответствуют многочленам степени не более 3 исходного пространства. Для простоты рисуйте графики в виде сетки  $k \times k$ .

In [ ]:

```
1
```

Для RBF-ядра повторите те же действия для следующей выборки.

In [ ]:

```
1 X, Y = make_moons(n_samples=100)
```

In [ ]:

```
1
```

Сделайте выводы.

## Задача 4.

Даны три случайные двумерные выборки для бинарной классификации:

- С линейно разделимыми классами
- С изогнутыми классами в виде лун
- С вложенными друг в друга классами

Код для генерации каждой выборки приводится далее.

In [ ]:

```
1 random_state = 20200329
2 n_samples = 500
3
4 X, Y = np.zeros((3, n_samples, 2)), np.zeros((3, n_samples))
5 X[0], Y[0] = make_blobs(n_samples=n_samples, n_features=2,
6                           centers=2, random_state=random_state)
7 X[1], Y[1] = make_moons(n_samples=n_samples,
8                           noise=0.2, random_state=random_state)
9 X[2], Y[2] = make_circles(n_samples=n_samples,
10                            noise=0.1, random_state=random_state)
```

Визуализация этих выборок

In [ ]:

```
1 plt.figure(figsize=(16, 5))
2
3 for i, name in enumerate(['Линейно разделимые', 'Изогнутые', 'Вложенные']):
4     plt.subplot(1, 3, i+1)
5     plt.scatter(X[i, :, 0], X[i, :, 1], c=Y[i], cmap='spring')
6     plt.title(name)
7 plt.show()
```

Для решения задачи вам выдается функция, аналогичная функции с семинара, но с некоторыми изменениями

In [ ]:

```
1  def draw_graphics(models, X, Y, X_test, Y_test, point_size=35,
2                      ncol=3, margin=False, params=['C']):
3      ...
4      Визуализирует решающие правила для каждой модели.
5      models --- все обученные SVM-модели, которые нужно визуализировать.
6      X --- объекты для визуализации (предполагается обучающая выборка)
7      Y --- ответы для визуализации (предполагается обучающая выборка)
8      X_test --- объекты, на которых необходимо посчитать качество
9      Y_test --- соответствующие им ответы.
10     point_size --- размер точки на графике
11     ncol --- количество колонок у таблицы графиков
12     margin --- если True, то визуализируется решающая функция,
13             иначе решающее правило
14     params --- список параметров SVM, которые нужно напечатать на графике
15     ...
16
17     # определение количества строк таблицы графиков в зависимости
18     # от количества колонок и количества моделей
19     n_rows = (len(models) + ncol-1) // ncol
20
21     plt.figure(figsize=(16, 5*n_rows))
22     for i_model, model in enumerate(models):
23         plt.subplot(n_rows, 3, i_model+1)
24
25         # Визуализация опорных векторов model.support_vectors_
26         plt.scatter(
27             model.support_vectors[:, 0], model.support_vectors[:, 1],
28             edgecolor='black', s=1.5*point_size,
29             alpha=0.5, zorder=10, linewidths=7
30         )
31
32         # Визуализация остальных точек
33         plt.scatter(
34             X[:, 0], X[:, 1], c=Y, zorder=10, s=point_size, alpha=0.7,
35             cmap=plt.cm.spring, edgecolor='black', linewidths=1.5
36         )
37
38         # Определение границ графика
39         x_min = X[:, 0].min() - 1.5
40         x_max = X[:, 0].max() + 1.5
41         y_min = X[:, 1].min() - 1.5
42         y_max = X[:, 1].max() + 1.5
43
44         # Сетка точек в пространстве
45         XX, YY = np.mgrid[x_min:x_max:500j, y_min:y_max:500j]
46         # Значения решающей функции для этой сетки.
47         # Чтобы их получить, нужно передать матрицу размера (N, 2)
48         Z = model.decision_function(np.c_[XX.ravel(), YY.ravel()])
49         # Ответы -- вектор. Переводим их к размеру сетки
50         Z = Z.reshape(XX.shape)
51
52         # Отрисовка решающей функции
53         if margin:
54             plt.pcolormesh(XX, YY, Z, cmap=plt.cm.RdBu)
55         else:
56             plt.pcolormesh(XX, YY, Z > 0, cmap=plt.cm.spring)
57
58         # Отрисовка разделяющей прямой и разделяющей полосы
59         plt.contour(
```

```

60         XX, YY, Z, colors=['k', 'k', 'k'],
61         linestyle=['--', '-', '-'],
62         levels=[-1, 0, 1]
63     )
64
65     plt.xlim(x_min, x_max)
66     plt.ylim(y_min, y_max)
67     plt.xticks(())
68     plt.yticks(())
69
70     # Вычисление качества
71     score = accuracy_score(Y_test, model.predict(X_test))
72     # Значения гиперпараметров
73     params_line = ', '.join([
74         name + '=' + str(model.get_params()[name]) for name in params
75     ])
76
77     plt.title(params_line + ', {}, Acc. = {:.1f}%'.format(
78         model.get_params()['kernel'], score*100
79     ))
80
81     plt.tight_layout()
82     plt.show()

```

Разбейте выборки на обучающую и тестовую в соотношении 7:3.

Далее для каждой выборки и для каждого случая постройте сетку графиков:

- Линейное ядро, три значения  $C$ , сетка графиков 1x3
- Полиномиальное ядро степени 2, три значения  $C$  и три значения  $\gamma$ , сетка графиков 3x3
- RBF-ядро степени 2, три значения  $C$  и три значения  $\gamma$ , сетка графиков 3x3
- RBF-ядро степени 2, три значения  $C$  и три значения  $\gamma$ , сетка графиков 3x3, визуализировать саму решающую функцию

Значения гиперпараметров подберите так, чтобы визуально получились "хороший", "средний" и "плохой" результаты. Какие значения гиперпараметров приводят к хорошему качеству визуально? По значению метрики? Какая форма разделяющей поверхности возможна для квадратичного ядра? Сделайте выводы.

## Задача 5.

Оптимизационная задача в методе SVM-регрессии имеет вид.

$$\begin{cases} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n (\xi_i^+ + \eta_i^+) \longrightarrow \max_{\theta, \theta_0, \xi, \eta} \\ Y_i - \varepsilon - \eta_i \leq \theta^T X_i + \theta_0 \leq Y_i + \varepsilon + \xi_i, \quad i = 1, \dots, n, \end{cases}$$

где  $\xi_i^+ = \max \{0, \xi_i\}$  и  $\eta_i^+ = \max \{0, \eta_i\}$ .

Сведите задачу к двойственной. Какие объекты являются опорными?

## Задача 6.

Рассмотрим [данные \(https://archive.ics.uci.edu/ml/datasets/Forest+Fires\)](https://archive.ics.uci.edu/ml/datasets/Forest+Fires) о лесных пожарах. Необходимо по различным показателям предсказать площадь лесного пожара. Описание данных приведено по ссылке.

1. Загрузите данные и постройте гистограмму целевого признака. Что можно сказать по гистограмме? Какое преобразование лучше совершить над ним?

In [ ]:

1

Сделайте это преобразование и еще раз посмотрите на гистограмму. Стало ли лучше?

In [ ]:

1

2. Оставьте только 4 признака: 'temp', 'RH', 'wind', 'rain', которые отвечают за температуру воздуха, влажность, скорость ветра и количество осадков соответственно.

In [ ]:

1

Разбейте выборку на обучающую и тестовую.

In [ ]:

1

Обучите линейную SVM-регрессию с  $\epsilon = 1$ , предварительно отмасштабирав признаки.

In [ ]:

1

Какова доля опорных векторов? Визуализируйте гистограмму значений двойственных коэффициентов. Что можно по ней сказать?

In [ ]:

1

Посчитайте ошибки предсказаний для каждого объекта обучающей выборки. Какие значения ошибок имеют опорные векторы? Остальные?

In [ ]:

1

3. Подберите оптимальные гиперпараметры SVM: величины  $C$ ,  $\epsilon$ , ядро, а также параметры ядра. Перед обучением необходимо провести стандартизацию признаков.

Для реализации используйте `Pipeline`. Учтите, что в таком случае для сетки гиперпараметров их имена должны иметь вид `model__param`, где `model` -- имя модели в `Pipeline`, а `param` -- имя ее гиперпараметра.

In [ ]:

1	
---	--

Посчитайте MAE-ошибку предсказания и доверительный интервал для нее. Перед этим не забудьте сделать обратное преобразование.

In [ ]:

1	
---	--

**4.** Выполните п.3 для всех признаков и сравните качество полученных моделей. Является ли разница статистически значимой?

In [ ]:

1	
---	--