

Семинар: Терминология. Введ sklearn. KNN. Задачи

Задание 1: kNN. Визуализация решающих г kNN.

В этом задании мы изобразим решающую поверхность для
чтобы наглядно увидеть, как классификатор принимает ре
Для простоты будем работать со встроенным в sklearn
содержащим информацию о характеристиках трёх видов в
можно найти [здесь \(https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_wine.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_wine.html)
и [здесь \(https://rdrr.io/cran/rattle.data/man/wine.html\)](https://rdrr.io/cran/rattle.data/man/wine.html).

Загрузим набор данных и сохраним информацию о призна
зависимой переменной – в переменную `y`.

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import matplotlib.pyplot as plt
```

```
In [2]: 1 from sklearn.datasets import load_wine
        2
        3 data = load_wine()
        4 X = pd.DataFrame(data['data'], columns = data[
        5 y = data['target']
        6 X.head(8)
```

```
Out[2]:
```

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	to
0	14.23	1.71	2.43	15.6	127.0	
1	13.20	1.78	2.14	11.2	100.0	
2	13.16	2.36	2.67	18.6	101.0	
3	14.37	1.95	2.50	16.8	113.0	
4	13.24	2.59	2.87	21.0	118.0	
5	14.20	1.76	2.45	15.2	112.0	
6	14.39	1.87	2.45	14.6	96.0	
7	14.06	2.15	2.61	17.6	121.0	

```
In [3]: 1 X.shape
```

```
Out[3]: (178, 13)
```

Задача 1.1. Есть ли в наборе данных пропущенные значения? Есть ли в наборе данных категориальные переменные? Если при помощи OneHot-кодирования.

```
B [4]: 1 X.isna().mean() # Тут мы измеряем соотношение
      2             # что пропущенных значений у нас
```

```
Out[4]: alcohol      0.0
        malic_acid    0.0
        ash           0.0
        alcalinity_of_ash 0.0
        magnesium     0.0
        total_phenols  0.0
        flavanoids     0.0
        nonflavanoid_phenols 0.0
        proanthocyanins 0.0
        color_intensity 0.0
        hue           0.0
        od280/od315_of_diluted_wines 0.0
        proline        0.0
        dtype: float64
```

```
B [5]: 1 X.dtypes == "object" # Тут мы проверяем, есть
      2                       # (тип object), значит нет
      3                       # значит у нас нет категориальных
```

```
Out[5]: alcohol      False
        malic_acid    False
        ash           False
        alcalinity_of_ash False
        magnesium     False
        total_phenols  False
        flavanoids     False
        nonflavanoid_phenols False
        proanthocyanins False
        color_intensity False
        hue           False
        od280/od315_of_diluted_wines False
        proline        False
        dtype: bool
```

```
B [6]: 1 # Нам нечего кодировать (((
```

Задача 1.2. Используя функцию `train_test_split()`, разделив данные на тренировочную и тестовую, и долю тестовой выборки задав `test_size`, разбиение осуществляется случайным образом, не забудьте `np.random.seed()` для воспроизводимости результатов.

```
B [7]: 1 from sklearn.model_selection import train_test
```

```
B [8]: 1 X_train, X_test, y_train, y_test = train_test_
      2 # разбиваем фреймы X и y на test и train с дол
```

```
B [9]: 1 X_train.shape # проверяем размер X_train. Д
```

```
Out[9]: (124, 13)
```

```
B [10]: 1 X_test.shape # Проверяем длину X_test. Она
```

```
Out[10]: (54, 13)
```

Задача 1.3. На тренировочной выборке обучите шесть классификаторов, отличающихся только числом соседей. Для первого классификатора поставьте равным 1, для второго - 3, для третьего – 5, для четвертого – 15 и для шестого – 25 (обратите внимание на параметр `k` в `KNeighborsClassifier`). Для обучения используйте толпы `alcohol` и `magnesium` – и евклидово расстояние. Не забудьте масштабировать данные, например, при помощи модуля `StandardScaler`.

Выведите долю правильных ответов на тренировочной и тестовой выборке для каждого классификатора.

```
B [11]: 1 from sklearn.preprocessing import StandardScaler
      2 from sklearn.neighbors import KNeighborsClassifier
```

```
B [12]: 1 X_train_1_3 = X_train[['alcohol', 'magnesium']]
      2 X_test_1_3 = X_test[['alcohol', 'magnesium']]
```

```
B [29]: 1 scaler = StandardScaler()
      2 scaler.fit(X_train_1_3)
      3 scaled_X_train_1_3 = scaler.transform(X_train_1_3)
      4 scaled_X_test_1_3 = scaler.transform(X_test_1_3)
      5 scaled_X_train_1_3[:15] # вывожу первые 15,
```

```
Out[29]: array([[ -1.10453814, -1.29401378],
                [ -0.60884898,  4.7936089 ],
                [  1.17054803,  0.66003794],
                [ -1.37144769, -1.21885795],
                [ -0.44361925, -0.31698792],
                [  0.20458966, -1.51948129],
                [  0.12832978, -0.9182346 ],
                [  1.11970812,  1.26128463],
                [  1.05615822, -0.61761126],
                [ -1.10453814, -2.12072797],
                [ -0.9520184 , -0.16667625],
                [  2.28902615, -0.54245543],
                [ -0.8884685 , -1.36916962],
                [  1.6026873 ,  0.96066128],
                [ -0.83762859, -0.76792293]])
```

```

B [18]: 1 num_neighbours = [1, 3, 5, 10, 15, 25]
        2
        3 tab = pd.DataFrame(columns=['NN', 'Train', 'Te
        4
        5 classifiers = []
        6
        7 for i in num_neighbours:
        8     clf = KNeighborsClassifier(n_neighbors=i,
        9
        10
        11     clf.fit(scaled_X_train_1_3, y_train)
        12
        13     y_train_predicted = clf.predict(scaled_X_t
        14     y_test_predicted = clf.predict(scaled_X_te
        15
        16     raw = pd.Series({'NN': i, 'Train': np.mean
        17     # здесь мы считаем точность предсказаний u
        18     tab = pd.concat([tab, raw.to_frame().T], i
        19     classifiers.append(clf)
        20 tab

```

Out[18]:

	NN	Train	Test
0	1.0	0.991935	0.685185
1	3.0	0.814516	0.685185
2	5.0	0.798387	0.722222
3	10.0	0.741935	0.777778
4	15.0	0.677419	0.814815
5	25.0	0.709677	0.796296

```

B [24]: 1 ##### Пример (не

```

Задача 1.4. Установите библиотеку `mlxtend` командой ни
можно установить из терминала при помощи `pip` или `con`
(<http://rasbt.github.io/mlxtend/installation/>).

B [15]: 1 !pip install mlxtend

```
Requirement already satisfied: mlxtend in c:\users
e-packages (0.23.0)
Requirement already satisfied: scipy>=1.2.1 in c:\
b\site-packages (from mlxtend) (1.10.1)
Requirement already satisfied: numpy>=1.16.2 in c:
ib\site-packages (from mlxtend) (1.24.3)
Requirement already satisfied: pandas>=0.24.2 in c
\lib\site-packages (from mlxtend) (1.5.3)
Requirement already satisfied: scikit-learn>=1.0.2
nda3\lib\site-packages (from mlxtend) (1.3.0)
Requirement already satisfied: matplotlib>=3.0.0 i
a3\lib\site-packages (from mlxtend) (3.7.1)
Requirement already satisfied: joblib>=0.13.2 in c
\lib\site-packages (from mlxtend) (1.2.0)
Requirement already satisfied: contourpy>=1.0.1 in
3\lib\site-packages (from matplotlib>=3.0.0->mlxt
Requirement already satisfied: cyclor>=0.10 in c:\
b\site-packages (from matplotlib>=3.0.0->mlxtend)
Requirement already satisfied: fonttools>=4.22.0 i
a3\lib\site-packages (from matplotlib>=3.0.0->mlxt
Requirement already satisfied: kiwisolver>=1.0.1 i
a3\lib\site-packages (from matplotlib>=3.0.0->mlxt
Requirement already satisfied: packaging>=20.0 in
\lib\site-packages (from matplotlib>=3.0.0->mlxten
Requirement already satisfied: pillow>=6.2.0 in c:
ib\site-packages (from matplotlib>=3.0.0->mlxtend)
Requirement already satisfied: pyparsing>=2.3.1 in
3\lib\site-packages (from matplotlib>=3.0.0->mlxt
Requirement already satisfied: python-dateutil>=2.
onda3\lib\site-packages (from matplotlib>=3.0.0->m
Requirement already satisfied: pytz>=2020.1 in c:\
b\site-packages (from pandas>=0.24.2->mlxtend) (20
Requirement already satisfied: threadpoolctl>=2.0.1
onda3\lib\site-packages (from scikit-learn>=1.0.2-
Requirement already satisfied: six>=1.5 in c:\user
te-packages (from python-dateutil>=2.7->matplotlib
0)
```

Если всё прошло успешно, то в выводе команды выше вы увидите "successfully installed", а следующая ячейка выполнится без ошибок.

B [20]: 1 from mlxtend.plotting import plot_decision_regions
2 import matplotlib.pyplot as plt
3 import matplotlib.gridspec as gridspec
4 import itertools

Задача 1.5. Библиотека mlxtend позволяет достаточно просто строить решающие поверхности обученных классификаторов. Изучите (http://rasbt.github.io/mlxtend/user_guide/plotting/plot_decision_regions/) и найдите, как можно построить несколько графиков решающих поверхностей (decision region grids). Постройте такую сетку графиков для классификаторов.

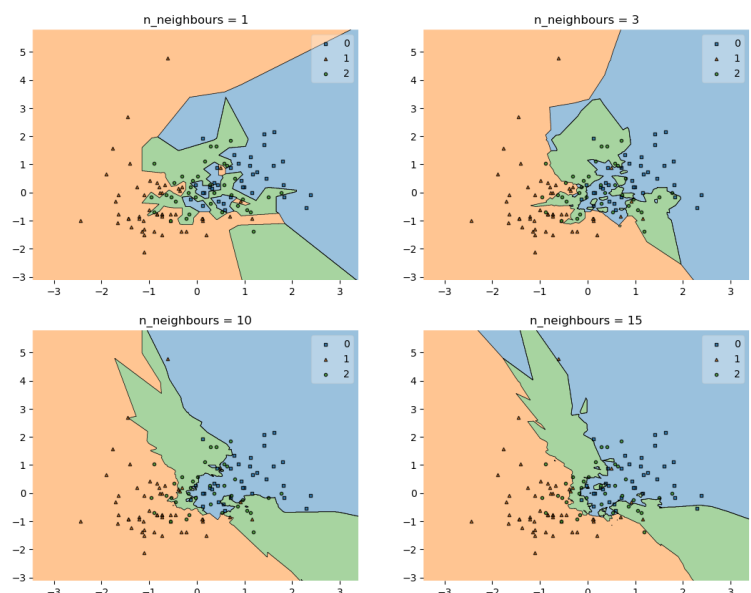
Подсказки:

1. Вы можете использовать готовый код, приведённый в Jupyter Notebook, и адаптировать его для нашего случая.
2. Вам могут понадобиться дополнительные библиотеки, например из документации.
3. Обратите внимание на то, как нужно изменить параметры `itertools.product()` для нашего числа классификаторов.
4. В функции `plot_decision_region()` используйте `y_train` из `X_train`. Возможно, их придётся перевести в формат `array`.
5. Если в задаче 1.3 вы сохранили обученные классификаторы, то не нужно заново обучать их заново.
6. Построение графика может занять некоторое время – подождать!

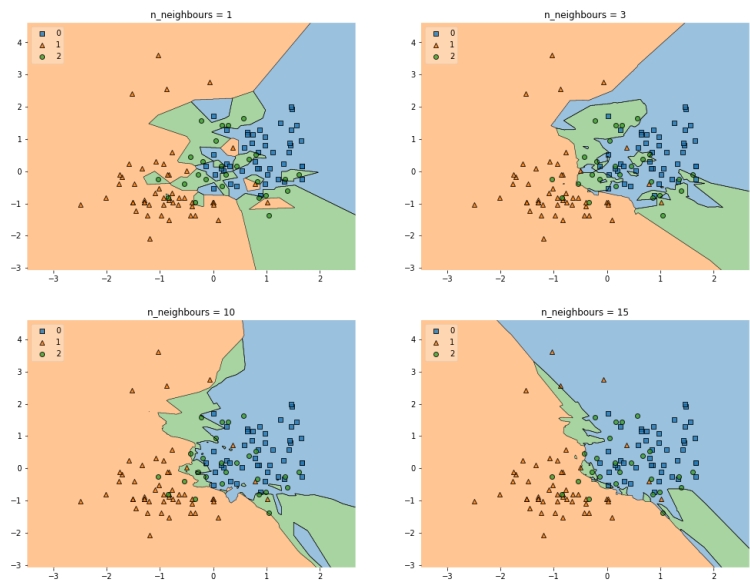
```

В [24]: 1 gs = gridspec.GridSpec(ncols=3, nrows=2) # создаем сетку
2 fig = plt.figure(figsize=(20,10)) # создаем фигуру
3 font = {'size': 10} # задаем шрифт
4 lines = {'markersize': 3} # задаем размер маркера
5
6 plt.rc('font', **font) # применяем шрифт
7 plt.rc('lines', **lines) # применяем размер маркера
8
9 labels = ['n_neighbours = {}'.format(i) for i in range(1, 16)]
10
11 for clf, lab, grd in zip(classifiers, labels,
12 # zip объединяет объекты с соответствующим индексом
13 #
14 ax = plt.subplot(gs[grd[0], grd[1]]) # берем подфигуру
15 # по заданной сетке
16 # из
17
18 fig = plot_decision_regions(X = scaled_X_train,
19 # в X передаем данные (см. задание 1.3), к
20 # метки и в clf наши классификаторы.
21
22 plt.title(lab) # даем заголовок
23
24 plt.show() # выводим

```



В []: 1 ##### Пример (не



Задача 1.6 Прокомментируйте результаты, полученные в 3 число соседей оптимально использовать для обучения кла ваш выбор при помощи описания геометрии данных и полу поверхности.

Type *Markdown* and LaTeX: α^2

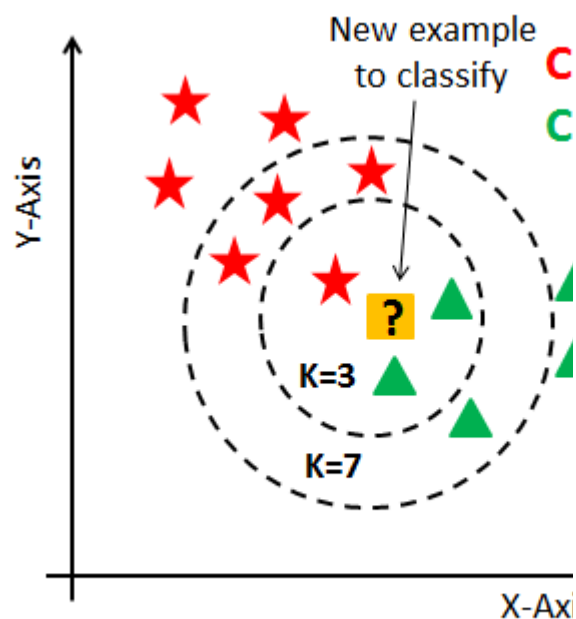
В задании 1.3 мы получили таблицу, отражающую долю пр test выборках, в зависимости от изменения гиперпараметр количество соседей. Видно, что на train выборке точность (что не равно 1), а в дальнейшем она начинает падать, потому ч своей группы объекты (шумовые) попадают уже в другие гр ситуация противоположная. С увеличением числа соседей, определенное количество объектов, потому что модель становится не Падать точность начинает, когда число соседей становится то класс забирает себе чужие области, потому что его элем train выборке. Например, если взять число соседей равное выборке, то все последующие элементы будут относиться к элементов которого больше.

В задании 1.5 видно, что разделяющая поверхность с увеличением числа соседей становится проще, так как на нее меньше действуют шумовые точки. Взгляд, самыми лучшими тут будут поверхности с 10 и 15 соседей.

Для того, чтобы оценить качество модели, нам нужно посмотреть на тестовой выборке, потому что именно они показывают, как модель работает на новых данных, которые она еще не видела. Исходя из результатов в задании 1.3, наиболее подходящее количество соседей для K-NN — это 15. Этот ответ подтверждает и то, что разделяющая поверхность выглядит разумно.

Задание 2. KNN своими руками (дополните)

В данном задании мы попробуем реализовать алгоритм KNN. В данном случае мы попробуем сделать KNN для классификации.



```
В [40]: 1 import numpy as np
        2 from collections import Counter  #не пригодило
```



```
B [25]: 1 import numpy as np
2 from collections import Counter
3 import statistics as st
4
5 class KNN:
6     def __init__(self, k:int): # конструктор
7         self.k = k
8
9     def fit(self, X, Y): # функция обуч
10         self.X_train = X # Значит нам н
11         self.y_train = Y
12         self.m, self.n = self.X_train.shape #
13
14     def predict(self, X_test): # метод предск
15         self.X_test = X_test # сохраняем те
16         self.m_test, self.n = self.X_test.shap
17
18
19         y_predict = np.zeros(self.m_test)
20
21         for i in range(self.m_test):
22             distances = np.zeros(self.m)
23
24             for j in range(self.m): # 6
25                 distances[j]=np.sqrt((np.sum(n
26
27             neighbors = self.y_train[distances
28
29             y_predict[i]= st.mode(neighbors)
30
31         return y_predict
```

```

В [26]: 1 # Не меняйте файл!
2 def test_knn(KNN):
3     knn = KNN(k=1)
4     X_train = np.array([[1, 1], [2, 2]])
5     y_train = np.array([0, 1])
6     X_test = np.array([[1.5, 1.5]])
7     knn.fit(X_train, y_train)
8     assert knn.predict(X_test) == [0]
9
10    knn = KNN(k=3)
11    X_train = np.array([[1, 1], [2, 2], [3, 3]])
12    y_train = np.array([0, 0, 0, 0, 1, 1, 1, 1])
13    X_test = np.array([[9.5, 9.5]])
14    knn.fit(X_train, y_train)
15    assert knn.predict(X_test) == [1]
16
17    knn = KNN(k=3)
18    X_train = np.array([[1, 1], [2, 2], [3, 3]])
19    y_train = np.array([0, 0, 0, 0, 1, 1, 1, 1])
20    X_test = np.array([[5.5, 5.5]])
21    knn.fit(X_train, y_train)
22    assert knn.predict(X_test) == [1]
23
24    knn = KNN(k=3)
25    X_train = np.array([[1, 1], [2, 2], [3, 3]])
26    y_train = np.array([0, 0, 0, 0, 1, 1, 1, 1])
27    X_test = np.array([[15, 15]])
28    knn.fit(X_train, y_train)
29    assert knn.predict(X_test) == [1]
30
31    knn = KNN(k=3)
32    X_train = np.array([[1, 1], [2, 2], [3, 3]])
33    y_train = np.array([0, 0, 0, 0, 1, 1, 1, 1])
34    X_test = np.array([[5, 5], [2, 2]])
35    knn.fit(X_train, y_train)
36    assert all(knn.predict(X_test) == [1, 0])

```

```

В [27]: 1 # Если тесты эти пройдены, то все верно!
2 test_knn(KNN)

```

```

В [28]: 1 # Все тесты пройдены.

```