# Лабораторная работа № 3.

## Подготовка обучающей и тестовой выборки, кросс-валидация и подбор гиперпараметров на примере метода ближайших соседей.

**Цель лабораторной работы: изучение способов подготовки выборки и подбора гиперпараметров на примере метода ближайших соседей.**

### Задание:

1. Выбрать набор данных (датасет) для решения задачи классификации или регрессии.
2. С использованием метода train_test_split разделить выборку на обучающую и тестовую.
3. Обучить модель ближайших соседей для произвольно заданного гиперпараметра К. Оценить качество модели с помощью подходящих для задачи метрик.
4. Произвести подбор гиперпараметра К с использованием GridSearchCV и/или RandomizedSearchCV и кросс-валидации, оценить качество оптимальной модели. Желательно использование нескольких стратегий кросс-валидации.
5. Сравнить метрики качества исходной и оптимальной моделей.

### Выполнение:

В качестве набора данных будем использовать датасет про вино: https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_wine.html#sklearn.datasets.load_wine

### Импорт библиотек

In [86]:

```python
import numpy as np
import pandas as pd
from sklearn.datasets import load_wine
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix, plot_confusion_matrix
from sklearn.metrics import recall_score, precision_score
from sklearn.model_selection import cross_validate
from sklearn.model_selection import KFold, LeaveOneOut, LeavePOut
from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
from typing import Dict
import seaborn as sns
import matplotlib.pyplot as plt
```

### Загрузка данных

In [2]:

```python
wine = load_wine()
```

In [3]:

```python
for x in wine:
    print(x)
```

```
data
target
frame
target_names
DESCR
feature_names
```

In [4]:

```python
# Наименование значений целевого признака
wine['target_names']
```

Out[4]:

```
array(['class_0', 'class_1', 'class_2'], dtype='<U7')
```

In [5]:

```python
# Признаки датасета
wine['feature_names']
```

```
['alcohol',
 'malic_acid',
 'ash',
 'alcalinity_of_ash',
 'magnesium',
 'total_phenols',
 'flavanoids',
 'nonflavanoid_phenols',
 'proanthocyanins',
 'color_intensity',
 'hue',
 'od280/od315_of_diluted_wines',
 'proline']
```

```python
wine['data'].shape
```

```
(178, 13)
```

```python
# Преобразование в pandas DataFrame
w_data = pd.DataFrame(data=np.c_[wine['data'], wine['target']], columns = wine['feature_names']+['target'
```

```python
w_data
```

|  | alcohol | malic_acid | ash | alcalinity_of_ash | magnesium | total_phenols | flavanoids | nonflavanoid_phenols | proanthocyanins | color_intensity |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 14.23 | 1.71 | 2.43 | 15.6 | 127.0 | 2.80 | 3.06 | 0.28 | 2.29 | 5.64 |
| 1 | 13.20 | 1.78 | 2.14 | 11.2 | 100.0 | 2.65 | 2.76 | 0.26 | 1.28 | 4.38 |
| 2 | 13.16 | 2.36 | 2.67 | 18.6 | 101.0 | 2.80 | 3.24 | 0.30 | 2.81 | 5.68 |
| 3 | 14.37 | 1.95 | 2.50 | 16.8 | 113.0 | 3.85 | 3.49 | 0.24 | 2.18 | 7.80 |
| 4 | 13.24 | 2.59 | 2.87 | 21.0 | 118.0 | 2.80 | 2.69 | 0.39 | 1.82 | 4.32 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 173 | 13.71 | 5.65 | 2.45 | 20.5 | 95.0 | 1.68 | 0.61 | 0.52 | 1.06 | 7.70 |
| 174 | 13.40 | 3.91 | 2.48 | 23.0 | 102.0 | 1.80 | 0.75 | 0.43 | 1.41 | 7.30 |
| 175 | 13.27 | 4.28 | 2.26 | 20.0 | 120.0 | 1.59 | 0.69 | 0.43 | 1.35 | 10.20 |
| 176 | 13.17 | 2.59 | 2.37 | 20.0 | 120.0 | 1.65 | 0.68 | 0.53 | 1.46 | 9.30 |
| 177 | 14.13 | 4.10 | 2.74 | 24.5 | 96.0 | 2.05 | 0.76 | 0.56 | 1.35 | 9.20 |

178 rows × 14 columns

```python
# Статистические характеристики датасета
w_data.describe()
```

|  | alcohol | malic_acid | ash | alcalinity_of_ash | magnesium | total_phenols | flavanoids | nonflavanoid_phenols | proanthocyanins | co |
|---|---|---|---|---|---|---|---|---|---|---|
| count | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | 178.000000 | |
| mean | 13.000618 | 2.336348 | 2.366517 | 19.494944 | 99.741573 | 2.295112 | 2.029270 | 0.361854 | 1.590899 | |
| std | 0.811827 | 1.117146 | 0.274344 | 3.339564 | 14.282484 | 0.625851 | 0.998859 | 0.124453 | 0.572359 | |
| min | 11.030000 | 0.740000 | 1.360000 | 10.600000 | 70.000000 | 0.980000 | 0.340000 | 0.130000 | 0.410000 | |
| 25% | 12.362500 | 1.602500 | 2.210000 | 17.200000 | 88.000000 | 1.742500 | 1.205000 | 0.270000 | 1.250000 | |
| 50% | 13.050000 | 1.865000 | 2.360000 | 19.500000 | 98.000000 | 2.355000 | 2.135000 | 0.340000 | 1.555000 | |
| 75% | 13.677500 | 3.082500 | 2.557500 | 21.500000 | 107.000000 | 2.800000 | 2.875000 | 0.437500 | 1.950000 | |
| max | 14.830000 | 5.800000 | 3.230000 | 30.000000 | 162.000000 | 3.880000 | 5.080000 | 0.660000 | 3.580000 | |

```python
# Значения признаков
wine.data[:5]
```

```
array([[1.423e+01, 1.710e+00, 2.430e+00, 1.560e+01, 1.270e+02, 2.800e+00,
        3.060e+00, 2.800e-01, 2.290e+00, 5.640e+00, 1.040e+00, 3.920e+00,
        1.065e+03],
       [1.320e+01, 1.780e+00, 2.140e+00, 1.120e+01, 1.000e+02, 2.650e+00,
        2.760e+00, 2.600e-01, 1.280e+00, 4.380e+00, 1.050e+00, 3.400e+00,
        1.050e+03],
       [1.316e+01, 2.360e+00, 2.670e+00, 1.860e+01, 1.010e+02, 2.800e+00,
        3.240e+00, 3.000e-01, 2.810e+00, 5.680e+00, 1.030e+00, 3.170e+00,
        1.185e+03],
       [1.437e+01, 1.950e+00, 2.500e+00, 1.680e+01, 1.130e+02, 3.850e+00,
        3.490e+00, 2.400e-01, 2.180e+00, 7.800e+00, 8.600e-01, 3.450e+00,
        1.480e+03],
       [1.324e+01, 2.590e+00, 2.870e+00, 2.100e+01, 1.180e+02, 2.800e+00,
        2.690e+00, 3.900e-01, 1.820e+00, 4.320e+00, 1.040e+00, 2.930e+00,
        7.350e+02]])
```

```
# Отмасштабируем признаки, т.к. min и max у некоторых признаков находятся не в одном диапазоне
sc = MinMaxScaler()
sc_wine = sc.fit_transform(wine.data)
```

**Разделение выборки на обучающую и тестовую**

```
X_train, X_test, Y_train, Y_test = train_test_split(sc_wine, wine.target, test_size=0.3, random_state=1)
```

```
# Размер обучающей выборки
X_train.shape, Y_train.shape
```

```
((124, 13), (124,))
```

```
# Размер тестовой выборки
X_test.shape, Y_test.shape
```

```
((54, 13), (54,))
```
Функция train_test_split делит выборку так, чтобы сохранились все классы

```
np.unique(Y_train)
```

```
array([0, 1, 2])
```

```
np.unique(Y_test)
```

```
array([0, 1, 2])
```

**Построим базовую модель на основе ближайших соседей с произвольно заданным гиперпараметром K**

```
# 4 ближайших соседа
cl1 = KNeighborsClassifier(n_neighbors=4)
cl1.fit(X_train, Y_train)
target1 = cl1.predict(X_test)
target11 = cl1.predict(X_train)
len(target1), target1
```

```
(54,
 array([2, 1, 0, 1, 0, 2, 1, 0, 2, 1, 0, 0, 1, 0, 1, 1, 2, 0, 1, 0, 0, 1,
        2, 0, 0, 2, 0, 0, 0, 2, 1, 2, 2, 0, 1, 1, 1, 1, 1, 0, 0, 1, 2, 0,
        0, 0, 0, 0, 0, 0, 1, 2, 2, 0]))
```

```
# 2 ближайших соседа
cl2 = KNeighborsClassifier(n_neighbors=2)
cl2.fit(X_train, Y_train)
target2 = cl2.predict(X_test)
len(target2), target2
```

```
(54,
 array([2, 1, 0, 1, 0, 1, 1, 0, 2, 1, 0, 0, 1, 0, 1, 1, 2, 0, 1, 0, 0, 1,
        2, 0, 0, 2, 0, 0, 0, 2, 1, 2, 2, 0, 1, 1, 1, 1, 1, 0, 0, 1, 2, 0,
        0, 0, 1, 0, 0, 0, 1, 2, 2, 0]))
```

**Метрика Accuracy**

```
accuracy_score(Y_test, target1)
```

```
0.9629629629629629
```

```
accuracy_score(Y_test, target2)
```

```
0.9629629629629629
```

Точность в случае двух и четырех ближайших соседей составляет 96%. Однако эта метрика показывает точность по всем классам в целом, поэтому точность в каждом отдельном классе может отличаться.

```python
def accuracy_score_for_classes(
    y_true: np.ndarray,
    y_pred: np.ndarray) -> Dict[int, float]:
    """
    Вычисление метрики accuracy для каждого класса
    Возвращает словарь: ключ - метка класса,
    значение - Accuracy для данного класса
    """
    # Для удобства фильтрации сформируем Pandas DataFrame
    d = {'t': y_true, 'p': y_pred}
    df = pd.DataFrame(data=d)
    # Метки классов
    classes = np.unique(y_true)
    # Результирующий словарь
    res = dict()
    # Перебор меток классов
    for c in classes:
        # отфильтруем данные, которые соответствуют текущей метке класса в истинных значениях
        temp_data_flt = df[df['t']==c]
        # расчет accuracy для заданной метки класса
        temp_acc = accuracy_score(
            temp_data_flt['t'].values,
            temp_data_flt['p'].values)
        # сохранение результата в словарь
        res[c] = temp_acc
    return res

def print_accuracy_score_for_classes(
    y_true: np.ndarray,
    y_pred: np.ndarray):
    """
    Вывод метрики accuracy для каждого класса
    """
    accs = accuracy_score_for_classes(y_true, y_pred)
    if len(accs)>0:
        print('Метка \t Accuracy')
    for i in accs:
        print('{} \t {}'.format(i, accs[i]))
```

```
print_accuracy_score_for_classes(Y_test, target1)
```

```
Метка    Accuracy
0    1.0
1    0.8947368421052632
2    1.0
```

Accuracy для класса 0 и 2 составляет 100%, но для класса 1 проседает до 89%.

```
print_accuracy_score_for_classes(Y_test, target2)
```

```
Метка    Accuracy
0    1.0
1    0.9473684210526315
2    0.9166666666666666
```

Accuracy для класса 0 составляет 100%, но для классов 1 и 2 - 95% и 92% соответственно.

```python
# Конвертация целевого признака в бинарный
def convert_target_to_binary(array:np.ndarray, target:int) -> np.ndarray:
    # Если целевой признак совпадает с указанным, то 1 иначе 0
    res = [1 if x==target else 0 for x in array]
    return res
```

```python
# Допустим целевой признак == 2, будем считать этот случай = 1 в бинарном признаке
bin_Y_train = convert_target_to_binary(Y_train, 2)
list(zip(Y_train, bin_Y_train))[:10]
```

```
[(1, 0),
 (1, 0),
 (0, 0),
 (1, 0),
 (2, 1),
 (1, 0),
 (1, 0),
 (0, 0),
 (2, 1),
 (1, 0)]
```

```
bin_Y_test = convert_target_to_binary(Y_test, 2)
list(zip(Y_test, bin_Y_test))[:10]
```

```
[(2, 1),
 (1, 0),
 (0, 0),
 (1, 0),
 (0, 0),
 (2, 1),
 (1, 0),
 (0, 0),
 (2, 1),
 (1, 0)]
```

```
bin_target1=convert_target_to_binary(target1, 2)
bin_target2=convert_target_to_binary(target2, 2)
```

## Матрица ошибок

```
# для 2 ближайших соседей с целевым признаком == 2
confusion_matrix(bin_Y_test, bin_target2, labels = [0, 1])
```

```
array([[42,  0],
       [ 1, 11]])
```

```
# для 4 ближайших соседей
confusion_matrix(Y_test, target1, labels = [0, 1, 2])
```

```
array([[23,  0,  0],
       [ 2, 17,  0],
       [ 0,  0, 12]])
```

```
# для 4 ближайших соседей
plot_confusion_matrix(cl1, X_test, Y_test,
                      display_labels=wine.target_names, cmap=plt.cm.YlGnBu, normalize='true')
```

```
<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7fc3f80510a0>
```

```
# для 2 ближайших соседей
plot_confusion_matrix(cl2, X_test, Y_test,
                      display_labels=wine.target_names, cmap=plt.cm.YlGnBu, normalize='true')
```

```
<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7fc3f8051d90>
```



## Метрика recall (полнота), precision

1) Доля верно предсказанных классификатором положительных объектов, из всех действительно положительных объектов. (TP/(TP+FN))

```
recall_score(bin_Y_test, bin_target1), recall_score(bin_Y_test, bin_target2)
```

```
(1.0, 0.9166666666666666)
```

2) Доля верно предсказанных классификатором положительных объектов, из всех объектов, которые классификатор верно или неверно определил как положительные. (TP/(TP+FP))

```
precision_score(bin_Y_test, bin_target1), precision_score(bin_Y_test, bin_target2)
```

```
(1.0, 1.0)
```

## Кросс-валидация

```
# Словарь метрик качества
scoring = {'precision': 'precision_weighted',
           'recall': 'recall_weighted',
           'f1': 'f1_weighted'}
```

```
# Кросс-валидация для 5 фолдов, на которые мы разделили обучающую и тестовую выборки
scores = cross_validate(KNeighborsClassifier(n_neighbors = 2),
                        sc_wine, wine.target, scoring = scoring, cv = 5,
                        return_train_score = True)
scores
```

```
{'fit_time': array([0.00216603, 0.00057101, 0.00048399, 0.00048113, 0.00061607]),
 'score_time': array([0.00609899, 0.00359392, 0.003232  , 0.00280595, 0.00400901]),
 'test_precision': array([0.80886752, 0.94910645, 1.        , 1.        , 0.92593407]),
 'train_precision': array([0.97404032, 0.96817443, 0.97404032, 0.98027972, 0.98025504]),
 'test_recall': array([0.80555556, 0.94444444, 1.        , 1.        , 0.91428571]),
 'train_recall': array([0.97183099, 0.96478873, 0.97183099, 0.97902098, 0.97902098]),
 'test_f1': array([0.80041152, 0.94352462, 1.        , 1.        , 0.91265664]),
 'train_f1': array([0.9718937 , 0.96487031, 0.9718937 , 0.9790619 , 0.97905014])}
```

Значения метрик precision, recall и f1 на обучающей выборке стабильнее и ближе к 100%, чем на тестовой выборке.

```
# Кросс-валидация методом LeaveOneOut: в тестовой выборке 1 элемент, а все остальные в обучающей
loo = LeaveOneOut()
```

```
scores1 = cross_validate(KNeighborsClassifier(n_neighbors = 2),
                        sc_wine, wine.target, scoring = scoring, cv = loo,
                        return_train_score = True)
scores1
```

```
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
```

UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted
samples. Use `zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
/Users/kalashnikova/opt/anaconda3/lib/python3.8/site-packages/sklearn/metrics/_classification.py:1221:
UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use `
zero_division` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))

Out[68]:

```
{'fit_time': array([0.00428391, 0.00054288, 0.00053024, 0.00052381, 0.00046587,
       0.00046206, 0.00046301, 0.00051832, 0.00052094, 0.00049686,
       0.00048685, 0.00075078, 0.00058413, 0.00166392, 0.00051618,
       0.00044298, 0.00045013, 0.00040388, 0.00042892, 0.00043511,
       0.00043201, 0.00060582, 0.00048304, 0.00047469, 0.00041389,
       0.00046682, 0.00039101, 0.00042486, 0.00039315, 0.00039411,
       0.00038099, 0.0003798 , 0.00037718, 0.00037694, 0.00038719,
       0.00039387, 0.00037789, 0.00037503, 0.00037694, 0.00038385,
       0.00038624, 0.00038409, 0.00038505, 0.00038004, 0.00039697,
       0.00037694, 0.00038195, 0.00037599, 0.00037789, 0.00038624,
       0.00039506, 0.00039887, 0.00037813, 0.00037718, 0.00037599,
       0.0003829 , 0.00037479, 0.00037885, 0.00040698, 0.00041509,
       0.00040317, 0.00041294, 0.00041199, 0.00047207, 0.00044608,
       0.00042701, 0.00043201, 0.00044727, 0.00046301, 0.00049305,
       0.00046587, 0.00046277, 0.00071812, 0.00048804, 0.00046706,
       0.00046062, 0.00070286, 0.0004878 , 0.00046492, 0.00046706,
       0.00073409, 0.00048709, 0.00047088, 0.00049782, 0.00108004,
       0.00062895, 0.00048494, 0.00046515, 0.0004437 , 0.00041199,
       0.00059199, 0.00044203, 0.00043702, 0.00039387, 0.00043488,
       0.00039101, 0.00043178, 0.00037408, 0.00037718, 0.00037217,
       0.0003767 , 0.00037503, 0.00037527, 0.00037408, 0.00037193,
       0.00037098, 0.00038195, 0.00052714, 0.00058508, 0.00081801,
       0.00047684, 0.00045896, 0.00044918, 0.00037885, 0.00042892,
       0.00037694, 0.00047421, 0.00041723, 0.00037789, 0.00049782,
       0.00051999, 0.0004921 , 0.00039601, 0.00038099, 0.000458  ,
       0.00043702, 0.00041699, 0.00037909, 0.00038028, 0.00037885,
       0.00037384, 0.00054407, 0.00042009, 0.00037622, 0.00037599,
       0.00037503, 0.00038624, 0.00037408, 0.00037193, 0.00037098,
       0.00044799, 0.00037503, 0.00037289, 0.00059009, 0.00046515,
       0.00195694, 0.00046325, 0.00075102, 0.00046802, 0.00054097,
       0.00048399, 0.00055313, 0.00045013, 0.00049496, 0.000489  ,
       0.00051308, 0.0004499 , 0.00047708, 0.00045085, 0.00046611,
       0.0004909 , 0.0004518 , 0.00047779, 0.00045109, 0.00046682,
       0.0004549 , 0.00044799, 0.00045729, 0.00041485, 0.00045586,
       0.00052404, 0.00052214, 0.00044918, 0.00046492, 0.00060487,
       0.00048923, 0.00045896, 0.00069809]),
 'score_time': array([0.0070622 , 0.00250316, 0.00240588, 0.00216222, 0.00206423,
       0.00209498, 0.00207996, 0.00300789, 0.00216389, 0.00213695,
```

```
        0.00229692, 0.00461626, 0.00214386, 0.00362396, 0.0028758 ,
        0.00205612, 0.00219321, 0.0023582 , 0.00201678, 0.00250793,
        0.00196505, 0.0027771 , 0.00206375, 0.00205016, 0.002038  ,
        0.00204301, 0.00234985, 0.00200605, 0.00197077, 0.00194192,
        0.00187302, 0.00185823, 0.00187874, 0.00205493, 0.001899  ,
        0.00192499, 0.00186491, 0.00186777, 0.00194311, 0.00188017,
        0.00187373, 0.00190377, 0.00195789, 0.00189185, 0.00186992,
        0.00195503, 0.00186706, 0.00185609, 0.00186801, 0.00187969,
        0.00202227, 0.00187898, 0.00197458, 0.00185394, 0.00187111,
        0.00187182, 0.00187206, 0.00193   , 0.00192404, 0.00194716,
        0.00190091, 0.00207996, 0.01243591, 0.00211906, 0.00209618,
        0.00237083, 0.00202274, 0.0020659 , 0.00222492, 0.0028379 ,
        0.00205421, 0.00341201, 0.00256586, 0.00274992, 0.00205994,
        0.00277519, 0.00248408, 0.00217414, 0.00204492, 0.00267196,
        0.00240898, 0.00208783, 0.00202703, 0.00340295, 0.004035  ,
        0.00239015, 0.00261307, 0.00230694, 0.00216293, 0.00236726,
        0.00195909, 0.00209713, 0.00196314, 0.00194287, 0.00194907,
        0.00190401, 0.00191212, 0.00184178, 0.00183392, 0.00183201,
        0.00182033, 0.00182486, 0.00183082, 0.00184512, 0.0018723 ,
        0.00183988, 0.00276709, 0.00209904, 0.00238895, 0.00362992,
        0.00210214, 0.00217605, 0.00209808, 0.00187111, 0.00192189,
        0.00199604, 0.00195789, 0.00194502, 0.00310206, 0.00207829,
        0.002141  , 0.00281119, 0.00190687, 0.00294495, 0.00197196,
        0.00193   , 0.00189209, 0.00183773, 0.00183606, 0.00185609,
        0.0019238 , 0.00228405, 0.00191307, 0.00186992, 0.00184202,
        0.00182796, 0.00197291, 0.00182676, 0.00184083, 0.0018239 ,
        0.00184011, 0.00184202, 0.00195003, 0.00223708, 0.00206685,
        0.00279999, 0.00201273, 0.00285006, 0.0026269 , 0.00223374,
        0.00234079, 0.00235415, 0.00240517, 0.00242591, 0.00206327,
        0.00209093, 0.00247788, 0.00209308, 0.00254703, 0.00204992,
        0.00217199, 0.00254202, 0.002141  , 0.00219703, 0.00202227,
        0.00241566, 0.00218606, 0.00201559, 0.00236511, 0.00198817,
        0.00227308, 0.00218606, 0.00216174, 0.00207114, 0.0029881 ,
        0.00253391, 0.00211191, 0.00212884]),
 'test_precision': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 0., 1., 1.,
        1., 1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0.,
        1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1.]),
 'train_precision': array([0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97885912, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97885912, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97399336,
        0.97399336, 0.97399336, 0.97399336, 0.97399336, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97883598, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97883598, 0.97395833, 0.97883598, 0.97395833,
        0.97395833, 0.96923077, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97883598, 0.97395833, 0.97883598, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833, 0.97395833,
        0.97395833, 0.97395833, 0.97395833, 0.97395833]),
 'test_recall': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 0., 1., 1.,
        1., 1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0.,
        1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
        1., 1., 1., 1., 1., 1., 1., 1.]),
 'train_recall': array([0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
```

```
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97740113, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97740113, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97740113, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97740113, 0.97175141, 0.97740113, 0.97175141,
          0.97175141, 0.96610169, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97740113, 0.97175141, 0.97740113, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141, 0.97175141, 0.97175141,
          0.97175141, 0.97175141, 0.97175141]),
 'test_f1': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
          1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
          1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
          1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1., 1., 0., 1., 1.,
          1., 1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 0., 1.,
          1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
          1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0.,
          1., 1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
          1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
          1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1., 1.,
          1., 1., 1., 1., 1., 1., 1., 1.]),
 'train_f1': array([0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97745026, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97745026, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97181958,
          0.97181958, 0.97181958, 0.97181958, 0.97181958, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97743927, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97743927, 0.97180245, 0.97743927, 0.97180245,
          0.97180245, 0.9661629 , 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97180245, 0.97743927, 0.97180245, 0.97743927, 0.97180245,
          0.97180245, 0.97180245, 0.97180245, 0.97180245, 0.97180245,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009, 0.97181009, 0.97181009,
          0.97181009, 0.97181009, 0.97181009])}
```

## Подбор гиперпараметра K с импользованием GridSearch и кросс-валидации

```python
# создаем массив со значениями гиперпараметра (кол-во ближайших соседей)
n_range = np.array(range(5,55,5))
tuned_params = [{'n_neighbors': n_range}]
tuned_params
```

```
[{'n_neighbors': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45, 50])}]
```

```
# Подбор гиперпараметра с использованием метода кросс-валидации LeavePOut(2) - 2 параметра в тестовой вы
# остальные в обучающей, и метрики качества accurancy
lpo = LeavePOut(2)
clf_gs = GridSearchCV(KNeighborsClassifier(), tuned_params, cv = lpo, scoring = 'accuracy')
clf_gs.fit(X_train, Y_train)
```

```
GridSearchCV(cv=LeavePOut(p=2), estimator=KNeighborsClassifier(),
             param_grid=[{'n_neighbors': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45, 50])}],
             scoring='accuracy')
```

```
clf_gs.cv_results_
```

```
{'mean_fit_time': array([0.00036827, 0.0003652 , 0.00036578, 0.00036527, 0.00037244,
        0.0004025 , 0.00038617, 0.00038437, 0.00040303, 0.00038359]),
 'std_fit_time': array([4.75352266e-05, 3.67897686e-05, 4.90215461e-05, 4.54180927e-05,
        5.98046029e-05, 8.55278823e-05, 4.55882250e-05, 4.40962847e-05,
        8.89617720e-05, 4.31835096e-05]),
 'mean_score_time': array([0.00061715, 0.00063268, 0.00061605, 0.00062711, 0.00063152,
        0.00070614, 0.00066285, 0.0006728 , 0.00070313, 0.00066211]),
 'std_score_time': array([8.77901857e-05, 1.62070956e-03, 6.98159128e-05, 8.17194101e-04,
        8.95293201e-05, 9.73635963e-04, 8.33577222e-05, 1.15841650e-03,
        6.98868543e-04, 7.86441299e-05]),
 'param_n_neighbors': masked_array(data=[5, 10, 15, 20, 25, 30, 35, 40, 45, 50],
             mask=[False, False, False, False, False, False, False, False,
                   False, False],
        fill_value='?',
             dtype=object),
 'params': [{'n_neighbors': 5},
 {'n_neighbors': 10},
 {'n_neighbors': 15},
 {'n_neighbors': 20},
 {'n_neighbors': 25},
 {'n_neighbors': 30},
 {'n_neighbors': 35},
 {'n_neighbors': 40},
 {'n_neighbors': 45},
 {'n_neighbors': 50}],
 'split0_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split1_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split2_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split3_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split4_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split5_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split6_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split7_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split8_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split9_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split10_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split11_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split12_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split13_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split14_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split15_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split16_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split17_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split18_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split19_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split20_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split21_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split22_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split23_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split24_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split25_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split26_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split27_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split28_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split29_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split30_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split31_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split32_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split33_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split34_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split35_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split36_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split37_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split38_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split39_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split40_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split41_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split42_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split43_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split44_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split45_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split46_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split47_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split48_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.])
```

 'split48_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split49_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split50_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split51_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split52_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split53_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split54_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split55_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split56_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split57_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split58_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split59_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split60_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split61_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split62_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split63_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split64_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split65_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5, 0.5]),
 'split66_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split67_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split68_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split69_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split70_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split71_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split72_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split73_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split74_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split75_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split76_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split77_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split78_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split79_test_score': array([1. , 0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split80_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split81_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split82_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split83_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split84_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split85_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split86_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split87_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split88_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split89_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split90_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split91_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split92_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split93_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split94_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split95_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split96_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 0.5, 1. , 1. , 0.5]),
 'split97_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split98_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split99_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split100_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split101_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split102_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split103_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split104_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split105_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split106_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split107_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split108_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split109_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split110_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split111_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split112_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split113_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split114_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split115_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split116_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split117_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split118_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split119_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split120_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split121_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split122_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split123_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split124_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split125_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split126_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0. ]),
 'split127_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split128_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split129_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split130_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split131_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split132_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split133_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split134_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split135_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split136_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split137_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split138_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split139_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split140_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split141_test_score': array([0. , 0. , 0.5, 0.5, 0. , 0. , 0. , 0. , 0. , 0. ]),

 'split142_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split143_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split144_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split145_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split146_test_score': array([0. , 0.5, 0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split147_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split148_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split149_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split150_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split151_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split152_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split153_test_score': array([0. , 0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split154_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split155_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split156_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split157_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split158_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split159_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split160_test_score': array([0.5, 0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split161_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split162_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split163_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split164_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split165_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split166_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split167_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split168_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split169_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split170_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split171_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split172_test_score': array([0., 0., 0., 0., 0., 0., 0., 0., 0., 0.]),
 'split173_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split174_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split175_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split176_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split177_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split178_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split179_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split180_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split181_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split182_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split183_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split184_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split185_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split186_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split187_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0. , 0. , 0. ]),
 'split188_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split189_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split190_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split191_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split192_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split193_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split194_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split195_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split196_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split197_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split198_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split199_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split200_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split201_test_score': array([0.5, 0. , 0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split202_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split203_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split204_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split205_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split206_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split207_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split208_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split209_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split210_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split211_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split212_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split213_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split214_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split215_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split216_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split217_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split218_test_score': array([0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split219_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split220_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split221_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split222_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split223_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split224_test_score': array([0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split225_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split226_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split227_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split228_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split229_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split230_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split231_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split232_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split233_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split234_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split235_test_score': array([0. , 0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),

 split235_test_score : array([0. , 0. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split236_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split237_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split238_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split239_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split240_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split241_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split242_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split243_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split244_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split245_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split246_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split247_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split248_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split249_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split250_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split251_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split252_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split253_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split254_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split255_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split256_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split257_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split258_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split259_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split260_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split261_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split262_test_score': array([0.5, 0.5, 1. , 1. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split263_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split264_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split265_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split266_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split267_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split268_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split269_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split270_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split271_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split272_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split273_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split274_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split275_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split276_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split277_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split278_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split279_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split280_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split281_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split282_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split283_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split284_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split285_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split286_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split287_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split288_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split289_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split290_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split291_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split292_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split293_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split294_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split295_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split296_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split297_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split298_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split299_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split300_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split301_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split302_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split303_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split304_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split305_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split306_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split307_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split308_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5, 0.5]),
 'split309_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split310_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split311_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split312_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split313_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split314_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split315_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split316_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split317_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split318_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split319_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split320_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split321_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split322_test_score': array([1. , 0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split323_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split324_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split325_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split326_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split327_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split328_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),

 'split329_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split330_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split331_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split332_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split333_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split334_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split335_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split336_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split337_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split338_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split339_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split340_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split341_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split342_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split343_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split344_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split345_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split346_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split347_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split348_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split349_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split350_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split351_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split352_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split353_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split354_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split355_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split356_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split357_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split358_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split359_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split360_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split361_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split362_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split363_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split364_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split365_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split366_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split367_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split368_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split369_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split370_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split371_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split372_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split373_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split374_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split375_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split376_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split377_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split378_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split379_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split380_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split381_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split382_test_score': array([0.5, 0.5, 1. , 1. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split383_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split384_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split385_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split386_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split387_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split388_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split389_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split390_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split391_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split392_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split393_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split394_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split395_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split396_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split397_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split398_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split399_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split400_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split401_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split402_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split403_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split404_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split405_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split406_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split407_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split408_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split409_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split410_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split411_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split412_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split413_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split414_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split415_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split416_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split417_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split418_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split419_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split420_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split421_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),

 'split422_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split423_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split424_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split425_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split426_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split427_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split428_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5, 0.5]),
 'split429_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split430_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split431_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split432_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split433_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split434_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split435_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split436_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split437_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split438_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split439_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split440_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split441_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split442_test_score': array([1. , 0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split443_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split444_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split445_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split446_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split447_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split448_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split449_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split450_test_score': array([1. , 1. , 1. , 1. , 0.5, 1. , 1. , 1. , 1. , 0.5]),
 'split451_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split452_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split453_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split454_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split455_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split456_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split457_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split458_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split459_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split460_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split461_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split462_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split463_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split464_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split465_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split466_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split467_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split468_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split469_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split470_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split471_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split472_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split473_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split474_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split475_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split476_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split477_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split478_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split479_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split480_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split481_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split482_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split483_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split484_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split485_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split486_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split487_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split488_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split489_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split490_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split491_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split492_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split493_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split494_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split495_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split496_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split497_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split498_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split499_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split500_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split501_test_score': array([0.5, 0.5, 1. , 1. , 0.5, 0.5, 0.5, 0.5, 1. , 1. ]),
 'split502_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split503_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split504_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split505_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split506_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split507_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split508_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split509_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split510_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split511_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split512_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split513_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split514_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split515_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),

 'split516_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split517_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split518_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split519_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split520_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split521_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split522_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split523_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split524_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split525_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split526_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split527_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split528_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split529_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split530_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split531_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split532_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split533_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split534_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split535_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split536_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split537_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split538_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split539_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split540_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split541_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split542_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split543_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split544_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split545_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split546_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split547_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5, 0.5]),
 'split548_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split549_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split550_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split551_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split552_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split553_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split554_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split555_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split556_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split557_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split558_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split559_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split560_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split561_test_score': array([1. , 0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split562_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split563_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split564_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split565_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split566_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split567_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split568_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split569_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split570_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split571_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split572_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split573_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split574_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split575_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split576_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split577_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split578_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split579_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split580_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split581_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split582_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split583_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split584_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split585_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split586_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split587_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split588_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split589_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split590_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split591_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split592_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split593_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split594_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split595_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split596_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split597_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split598_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split599_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split600_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split601_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split602_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split603_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split604_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split605_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split606_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split607_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split608_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),

 'split609_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split610_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split611_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split612_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split613_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split614_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split615_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split616_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split617_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split618_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split619_test_score': array([0.5, 0.5, 1. , 1. , 0.5, 0.5, 0.5, 0.5, 0.5, 0. ]),
 'split620_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split621_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split622_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split623_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split624_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split625_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split626_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split627_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split628_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split629_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split630_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split631_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split632_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split633_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split634_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split635_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split636_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split637_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split638_test_score': array([1. , 0.5, 1. , 1. , 1. , 0.5, 0.5, 0.5, 0. , 0.5]),
 'split639_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split640_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split641_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split642_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split643_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split644_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split645_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split646_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split647_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split648_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split649_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split650_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0. ]),
 'split651_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split652_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split653_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split654_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split655_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split656_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split657_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split658_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split659_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split660_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split661_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split662_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split663_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split664_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split665_test_score': array([1. , 1. , 1. , 1. , 1. , 0.5, 0.5, 0.5, 0. , 0. ]),
 'split666_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split667_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split668_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split669_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split670_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split671_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split672_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split673_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split674_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split675_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split676_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split677_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split678_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split679_test_score': array([1. , 0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split680_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split681_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split682_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split683_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split684_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split685_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split686_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split687_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split688_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split689_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split690_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split691_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split692_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split693_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split694_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split695_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split696_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split697_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
 'split698_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split699_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split700_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split701_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
 'split702_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5])

'split702_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split703_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split704_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split705_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split706_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split707_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split708_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split709_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split710_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
'split711_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split712_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split713_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split714_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5]),
'split715_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split716_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split717_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split718_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split719_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split720_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split721_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split722_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5]),
'split723_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split724_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split725_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split726_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split727_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split728_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split729_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split730_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split731_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split732_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split733_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split734_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split735_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split736_test_score': array([0.5, 0.5, 1. , 1. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
'split737_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split738_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split739_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split740_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split741_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
'split742_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split743_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split744_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split745_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split746_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split747_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split748_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
'split749_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split750_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split751_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split752_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split753_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split754_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split755_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
'split756_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split757_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split758_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split759_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split760_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split761_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split762_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split763_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split764_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split765_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split766_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split767_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
'split768_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split769_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split770_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split771_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split772_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split773_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split774_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split775_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split776_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split777_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split778_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split779_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split780_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split781_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split782_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5, 0.5]),
'split783_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split784_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split785_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split786_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split787_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split788_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split789_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split790_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split791_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split792_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split793_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split794_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split795_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),

 'split796_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split797_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split798_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split799_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split800_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split801_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split802_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split803_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split804_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split805_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split806_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split807_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split808_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split809_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split810_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split811_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split812_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split813_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split814_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split815_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split816_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split817_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split818_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split819_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split820_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split821_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split822_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split823_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split824_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split825_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split826_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split827_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split828_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split829_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split830_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split831_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split832_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split833_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split834_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split835_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split836_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split837_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split838_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split839_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split840_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split841_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split842_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split843_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split844_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split845_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split846_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split847_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split848_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split849_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split850_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split851_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split852_test_score': array([0.5, 0.5, 1. , 1. , 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split853_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split854_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split855_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split856_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split857_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split858_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split859_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split860_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split861_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split862_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split863_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split864_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split865_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split866_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split867_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split868_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split869_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split870_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split871_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split872_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split873_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split874_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split875_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split876_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split877_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split878_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split879_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split880_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split881_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split882_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split883_test_score': array([0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5, 0.5]),
 'split884_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split885_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split886_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split887_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split888_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split889_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.])

 split889_test_score : array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split890_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split891_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split892_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split893_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split894_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split895_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split896_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split897_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split898_test_score': array([1. , 1. , 1. , 1. , 1. , 1. , 1. , 0.5, 0.5, 0.5]),
 'split899_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split900_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split901_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split902_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split903_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split904_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split905_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split906_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split907_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split908_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split909_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split910_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split911_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split912_test_score': array([1. , 0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split913_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split914_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split915_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split916_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split917_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split918_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split919_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split920_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split921_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split922_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split923_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split924_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split925_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split926_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split927_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split928_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split929_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split930_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split931_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split932_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split933_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split934_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split935_test_score': array([0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split936_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split937_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split938_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split939_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split940_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split941_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split942_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split943_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split944_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split945_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split946_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split947_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split948_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split949_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split950_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split951_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split952_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split953_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split954_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split955_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split956_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split957_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split958_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split959_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split960_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split961_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split962_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split963_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split964_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split965_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split966_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split967_test_score': array([0.5, 0.5, 1. , 1. , 0.5, 0.5, 0.5, 0.5, 1. , 1. ]),
 'split968_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split969_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split970_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split971_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split972_test_score': array([0.5, 1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split973_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split974_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split975_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split976_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split977_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split978_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split979_test_score': array([0.5, 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
 'split980_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split981_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
 'split982_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),

```
'split983_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split984_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split985_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split986_test_score': array([1. , 0.5, 1. , 1. , 1. , 1. , 1. , 1. , 1. , 1. ]),
'split987_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split988_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split989_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split990_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split991_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split992_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
'split993_test_score': array([1., 1., 1., 1., 1., 1., 1., 1., 1., 1.]),
...}
```

In [77]:

```python
# Лучшая модель
clf_gs.best_estimator_
```

Out[77]:

```
KNeighborsClassifier(n_neighbors=20)
```

In [78]:

```python
# Лучшее значение метрики
clf_gs.best_score_
```

Out[78]:

```
0.9832153160241279
```

In [81]:

```python
# Лучшее значение параметров
clf_gs.best_params_
```

Out[81]:

```
{'n_neighbors': 20}
```

In [82]:

```python
# Изменение качества на тестовой выборке в зависимости от K-соседей
plt.plot(n_range, clf_gs.cv_results_['mean_test_score'])
```

Out[82]:

```
[<matplotlib.lines.Line2D at 0x7fc3f8553ee0>]
```



In [85]:

```python
lpo = LeavePOut(2)
clf_rs = RandomizedSearchCV(KNeighborsClassifier(), tuned_params, cv = lpo, scoring = 'accuracy')
clf_rs.fit(X_train, Y_train)
```

Out[85]:

```
RandomizedSearchCV(cv=LeavePOut(p=2), estimator=KNeighborsClassifier(),
                   param_distributions=[{'n_neighbors': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45, 50])}]
,
                   scoring='accuracy')
```

In [87]:

```python
clf_rs.best_score_, clf_rs.best_params_
```

Out[87]:

```
(0.9832153160241279, {'n_neighbors': 20})
```

In [88]:

```python
# Подбор гиперпараметра с использованием метода кросс-валидации KFold - деление выборок на 5 фолдов,
# и метрики качества accurancy
clf_gs1 = GridSearchCV(KNeighborsClassifier(), tuned_params, cv = 5, scoring = 'accuracy')
clf_gs1.fit(X_train, Y_train)
```

Out[88]:

```
GridSearchCV(cv=5, estimator=KNeighborsClassifier(),
             param_grid=[{'n_neighbors': array([ 5, 10, 15, 20, 25, 30, 35, 40, 45, 50])}],
             scoring='accuracy')
```

In [89]:

```python
clf_gs1.best_score_, clf_gs1.best_params_
```

```
(0.976, {'n_neighbors': 15})
```

```
plt.plot(n_range, clf_gs1.cv_results_['mean_test_score'])
```

```
[<matplotlib.lines.Line2D at 0x7fc3fd244430>]
```



## Сравнение метрик качества исходной и оптимальной моделей

```
clf_gs.best_estimator_.fit(X_train, Y_train)
opt_target_train = clf_gs.best_estimator_.predict(X_train)
opt_target_test = clf_gs.best_estimator_.predict(X_test)
```

```
# Точность для оптимальной модели
accuracy_score(Y_train, opt_target_train), accuracy_score(Y_test, opt_target_test)
```

```
(0.9838709677419355, 0.9814814814814815)
```

```
# Точность для исходной модели
accuracy_score(Y_train, target11), accuracy_score(Y_test, target1)
```

```
(0.967741935483871, 0.9629629629629629)
```

Точность у оптимальной модели выше и разница в точности между обучающей и тестовой выборками уменьшилась с 0.0048 до 0.0023

```
# Precision для оптимальной модели
precision_score(Y_train, opt_target_train, average='macro'), precision_score(Y_test, opt_target_test, ave
```

```
(0.981981981981982, 0.9861111111111112)
```

```
# Precision для исходной модели
precision_score(Y_train, target11, average='macro'), precision_score(Y_test, target1, average='macro')
```

```
(0.9653499653499654, 0.9733333333333333)
```

Разница метрики между выборками уменьшилась с 0.0079 до 0.0041

```
# Полнота для оптимальной модели, учитывается вес классов
recall_score(Y_train, opt_target_train, average='weighted'), recall_score(Y_test, opt_target_test, averag
```

```
(0.9838709677419355, 0.9814814814814815)
```

```
# Полнота для исходной модели, учитывается вес классов
recall_score(Y_train, target11, average='weighted'), recall_score(Y_test, target1, average='weighted')
```

```
(0.967741935483871, 0.9629629629629629)
```

Разница полноты между выборками уменьшилась с 0.0047 до 0.0023