

```
1 import numpy as np
2 import pandas as pd
3
4 from sklearn.model_selection import train_test_split,\
5                                     LeaveOneOut,\
6                                     LeavePOut,\
7                                     KFold,\
8                                     ShuffleSplit,\
9                                     GridSearchCV,\
10                                    StratifiedKFold,\
11                                    StratifiedShuffleSplit,\
12                                    LeaveOneGroupOut,\
13                                    LeavePGroupsOut,\
14                                    GroupKFold,\
15                                    GroupShuffleSplit,\
16                                    cross_val_score
17
18 from sklearn import datasets
19 from sklearn.linear_model import LinearRegression
20
21 import matplotlib.pyplot as plt
22 import seaborn as sns
23 from matplotlib.collections import LineCollection
24
25 sns.set(style='whitegrid', font_scale=1.3, palette='Set2')
```

▼ 2. Валидация

В простом случае предполагается, что у нас есть данные для обучения (train set) и данные для тестирования (test set). Ответы для тестовых данных считаются неизвестными, а потому их нельзя использовать для обучения модели.

Главная цель валидации — оценить какое качество модель способна показать на тестовых данных. Отсюда вытекает и главное правило валидации:

Валидационные данные должны быть как можно сильнее похожи на тестовые.

В этом ноутбуке мы обсудим различные стратегии валидации.

▼ 2.1 Валидация на отложенной выборке (holdout validation)

При оценке качества модели нельзя использовать данные, которые использовались для ее обучения, т.к. при таком подходе мы не сможем оценить адекватность модели на новых данных и контролировать переобучение. Для решения данной проблемы существуют подходы, использующие понятие **отложенной выборки**: X_{val} , Y_{val} . Отложенной выборкой называют размеченные данные, которые мы не используем при обучении модели.

В рассматриваемом нами случае ответы на тестовых данных неизвестны, а потому отложенную выборку мы можем сформировать лишь из обучающих данных. Такую выборку обычно называют валидационным множеством (validation set или development set).

В scikit-learn разбиение на обучающую и тестовую выборки можно легко получить с помощью функции `train_test_split`.

Рассмотрим на примере задачи классификации ирисов.

```
1 # загружаем датасет
2 data_full = datasets.load_iris()
3 print("Shape: {}".format(data_full.data.shape))

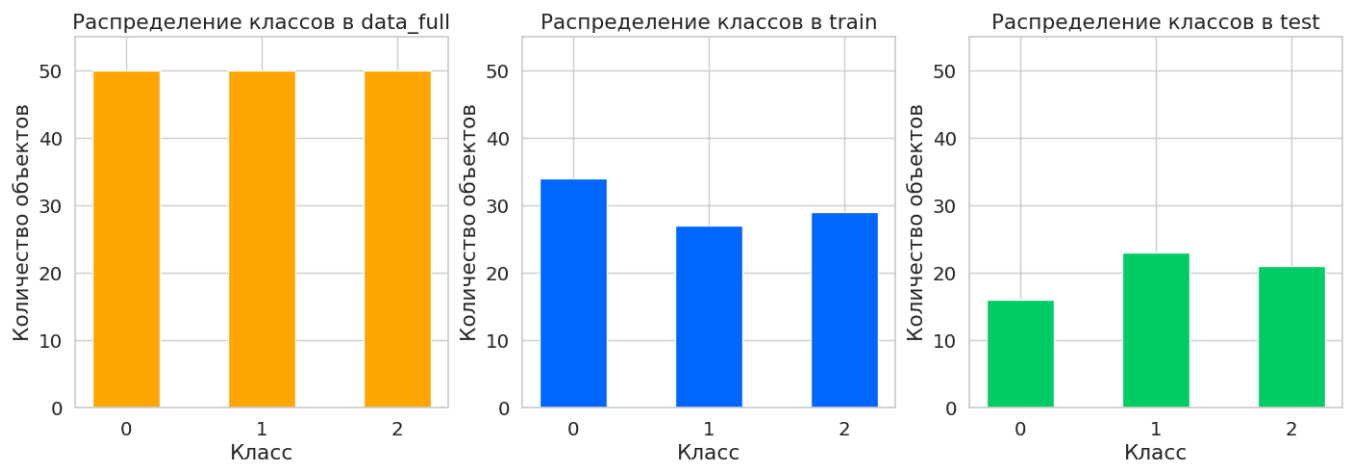
    Shape: (150, 4)

1 X_train, X_test, y_train, y_test = train_test_split(
2     # *arrays: принимает индексруемые объекты с совпадающей shape[0].
3     # Например: list, np.array, pd.DataFrame.
4     data_full.data, data_full.target,
5     test_size=0.4, # доля данных, которые берем в тестовую выборку
6     random_state=0, # фиксируем случайность
7     shuffle=True, # перемешивает данные в случайном порядке
8     stratify=None # если не None, то сохраняет доли классов при разбиении
9 )

1 print("Shape of train data: {} {}".format(X_train.shape, y_train.shape))
2 print("Shape of test data: {} {}".format(X_test.shape, y_test.shape))

    Shape of train data: (90, 4) (90,)
    Shape of test data: (60, 4) (60,)
```

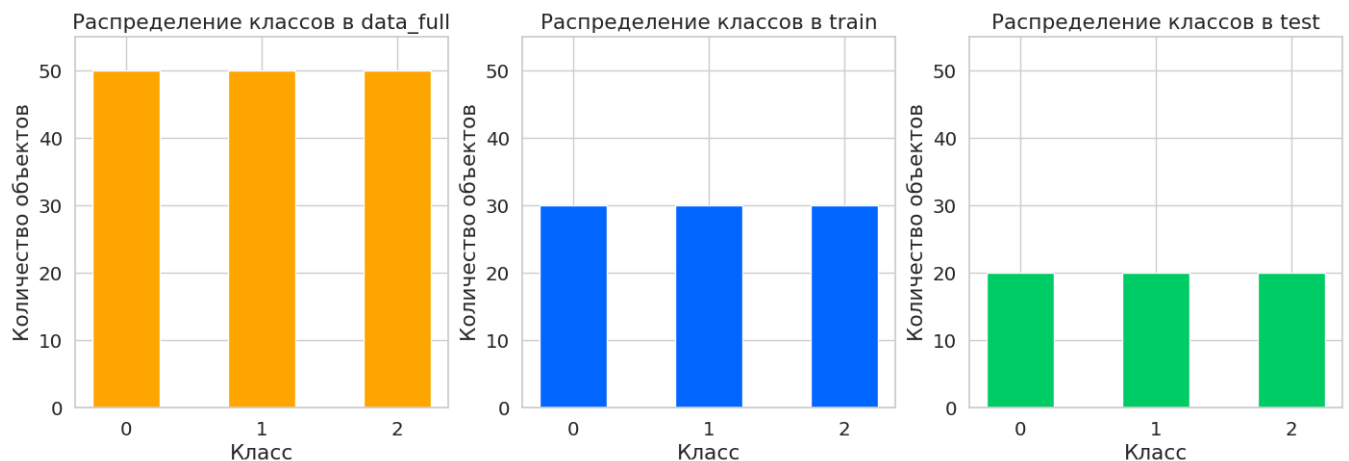
```
1 plt.figure(figsize=(17,5))
2 split_cases = [data_full.target, y_train, y_test]
3 colors = ['orange', '#0066FF', '#00CC66']
4 labels = ['Распределение классов в data_full',
5           'Распределение классов в train',
6           'Распределение классов в test']
7 for i in range(3):
8     plt.subplot(1, 3, i + 1)
9     values, counts = np.unique(split_cases[i], return_counts=True)
10    plt.bar(values, counts, width=0.5, color=colors[i])
11    plt.ylim(0, 55)
12    plt.xticks([0, 1, 2])
13    plt.xlabel('Класс')
14    plt.ylabel('Количество объектов')
15    plt.title(labels[i])
16 plt.show()
```



Видим, что распределение классов в обучающей и тестовой выборках отличаются. Теперь попробуем сделать разбиение с `stratify = data_full.target`. Стратификация — стратегия кросс-валидации, при которой в обучающей и тестовой выборке сохраняется одинаковое распределение целевой переменной, такое же, как во всем датасете. Зачем нужна стратификация узнаем чуть позже.

```
1 X_train, X_test, y_train, y_test = train_test_split(
2     # *arrays: принимает индексруемые объекты с совпадающей shape[0].
3     # Например: list, np.array, pd.DataFrame.
4     data_full.data, data_full.target,
5     test_size=0.4, # доля данных, которые берем в тестовую выборку
6     random_state=0, # фиксируем случайность
7     shuffle=True, # перемешивает данные в случайном порядке
8     # сохраняем доли классов при разбиении как в таргете
9     stratify=data_full.target
10 )
```

```
1 plt.figure(figsize = (17,5))
2 split_cases = [data_full.target, y_train, y_test]
3 colors = ['orange', '#0066FF', '#00CC66']
4 labels = ['Распределение классов в data_full',
5           'Распределение классов в train',
6           'Распределение классов в test']
7 for i in range(3):
8     plt.subplot(1, 3, i + 1)
9     values, counts = np.unique(split_cases[i], return_counts=True)
10    plt.bar(values, counts, width=0.5, color=colors[i])
11    plt.ylim(0, 55)
12    plt.xticks([0, 1, 2])
13    plt.xlabel('Класс')
14    plt.ylabel('Количество объектов')
15    plt.title(labels[i])
16 plt.show()
```



Видим, что после разбиения с `stratify = data_full.target` распределения классов в train и test не отличаются.

Подведем итог для метода отложенной выборки.

Достоинства:

- Быстрый для оценки качества модели. При использовании данной техники разбиения данных для оценки качества модели происходит одна процедура обучения на обучающей выборке, после чего качество модели оценивается на тестовых данных.

Недостатки:

- Результат сильно зависит от способа разбиения. Объекты в train и test могут получиться из разных распределений, если `stratify = False`.
- При обучении модели на обучающей выборке валидационная выборка не используется, то есть мы не задействуем все доступные данные для обучения.
- При оптимизации значения метрики на валидационном множестве модель немного переобучается под него. Таким образом, значение метрики качества на новых данных не будет соответствовать значению метрики на тестовом множестве.

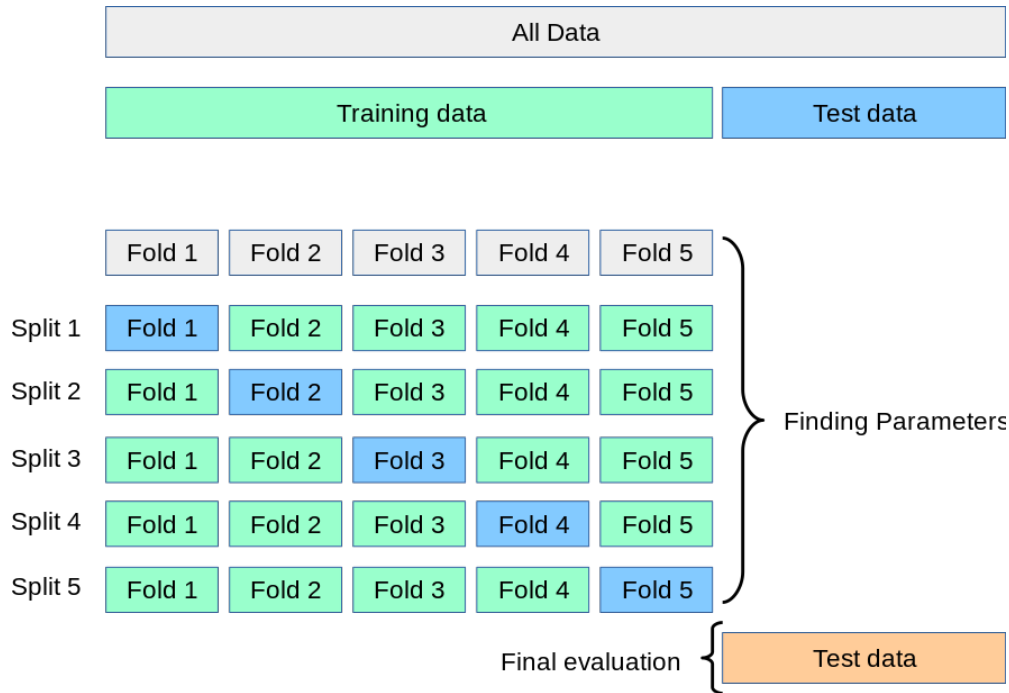
▼ 2.2 Кросс-валидация (cross-validation)

▼ A. k-Fold Cross Validation

Описанные выше недостатки оказались критичны, поэтому для борьбы с ними придумали кросс-валидацию. **Кросс-валидация** — это метод оценки качества модели, при котором обучающая выборка делится на k частей, или **фолдов**. После чего для каждого из k фолдов проделывается следующая процедура:

- модель обучается на остальных $k - 1$ фолдах, которые вместе формируют обучающую выборку для данной итерации
- обученная модель оценивается на оставшемся k -ом фолде

Таким образом мы получаем k оценок качества. Итоговая метрика считается как среднее полученных оценок. Ниже представлена визуализация рассматриваемой стратегии кросс-валидации для пяти фолдов.



Рассмотрим пример. Будем использовать данные о ценах квартир в Калифорнии.

```
1 housing = datasets.fetch_california_housing()
2
3 X = pd.DataFrame(data=housing['data'], columns=housing['feature_names'])
4 y = housing['target']
```

Описание датасета.

```
1 print(housing['DESCR'])
```

```
.. _california_housing_dataset:
```

California Housing dataset

****Data Set Characteristics:****

:Number of Instances: 20640

:Number of Attributes: 8 numeric, predictive attributes and the target

:Attribute Information:

- MedInc median income in block group
- HouseAge median house age in block group
- AveRooms average number of rooms per household
- AveBedrms average number of bedrooms per household
- Population block group population
- AveOccup average number of household members
- Latitude block group latitude
- Longitude block group longitude

:Missing Attribute Values: None

This dataset was obtained from the StatLib repository.

https://www.dcc.fc.up.pt/~ltorgo/Regression/cal_housing.html

The target variable is the median house value for California districts, expressed in hundreds of thousands of dollars (\$100,000).

This dataset was derived from the 1990 U.S. census, using one row per census block group. A block group is the smallest geographical unit for which the Census Bureau publishes sample data (a block group typically has a population of 600 to 3,000 people).

A household is a group of people residing within a home. Since the average number of rooms and bedrooms in this dataset are provided per household, the columns may take surprisingly large values for block groups with few households and many empty houses, such as vacation resorts.

It can be downloaded/loaded using the

:func:`sklearn.datasets.fetch_california_housing` function.

.. topic:: References

- Pace, R. Kelley and Ronald Barry, Sparse Spatial Autoregressions, Statistics and Probability Letters, 33 (1997) 291-297


```
1 X.head()
```

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Lon
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	

```
1 X.shape
```

```
(20640, 8)
```

При помощи функции [cross_val_score](#) можем получить значение выбранной метрики на всех фолдах.

В качестве примера применим ее к линейной модели.

```
1 model = LinearRegression()
2 scores = cross_val_score(
3     estimator=model, # модель, качество которой хотим оценить
4     X=X, # данные для обучения (не содержат целевую переменную)
5     y=y, # значения целевой переменной
6     cv=5, # количество фолдов
7     scoring='neg_mean_squared_error', # метрика качества
8     n_jobs=-1 # количество ядер для вычислений, -1 - использование всех ядер
9 )
10 scores
```

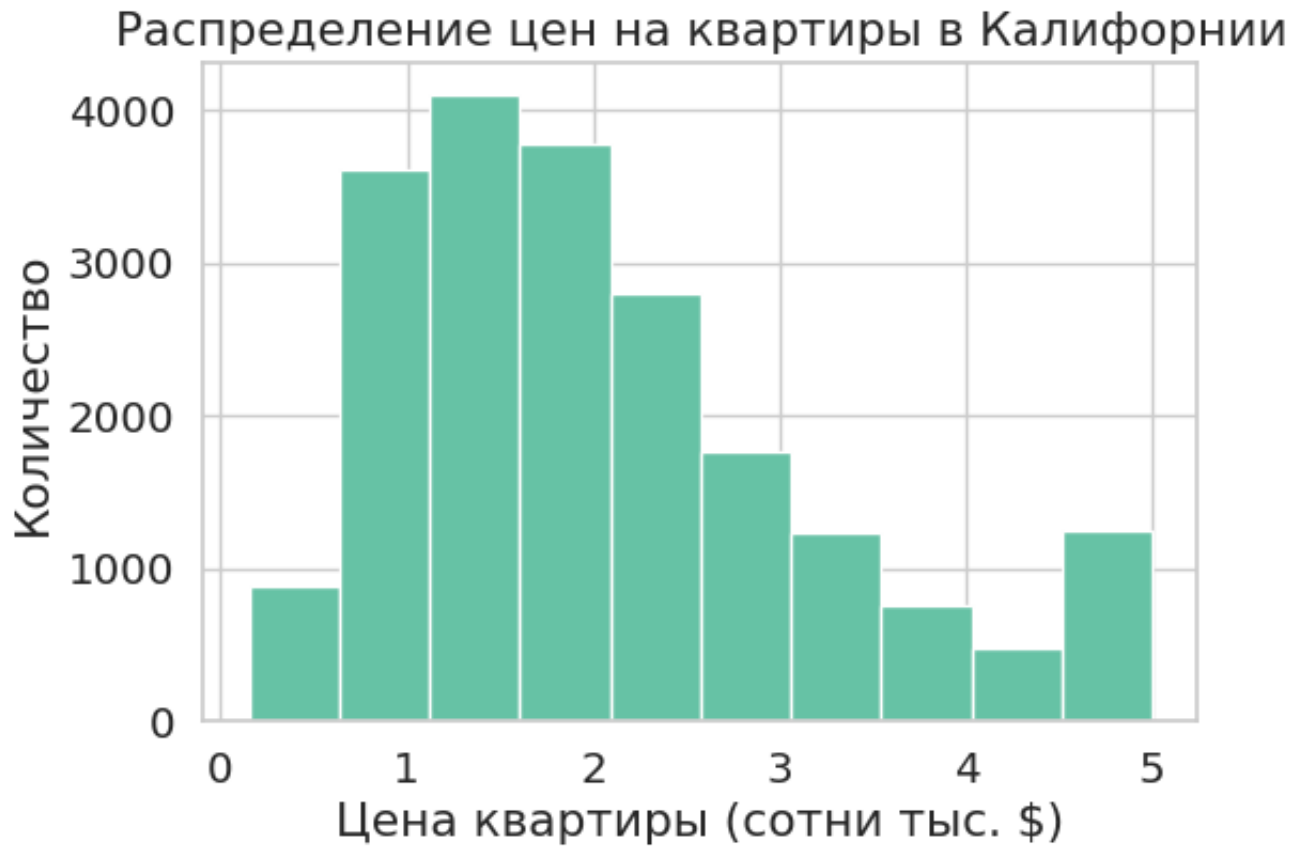
```
array([-0.48485857, -0.62249739, -0.64621047, -0.5431996 , -0.49468484])
```

Стоит отметить, что в качестве scoring мы используем `neg_mean_squared_error`.

Префикс `neg` показывает, что мы оптимизируем $(-1) \cdot \text{MSE}$. Дело в том, что оптимизации в `sklearn` подразумевают **максимизацию** метрики качества.

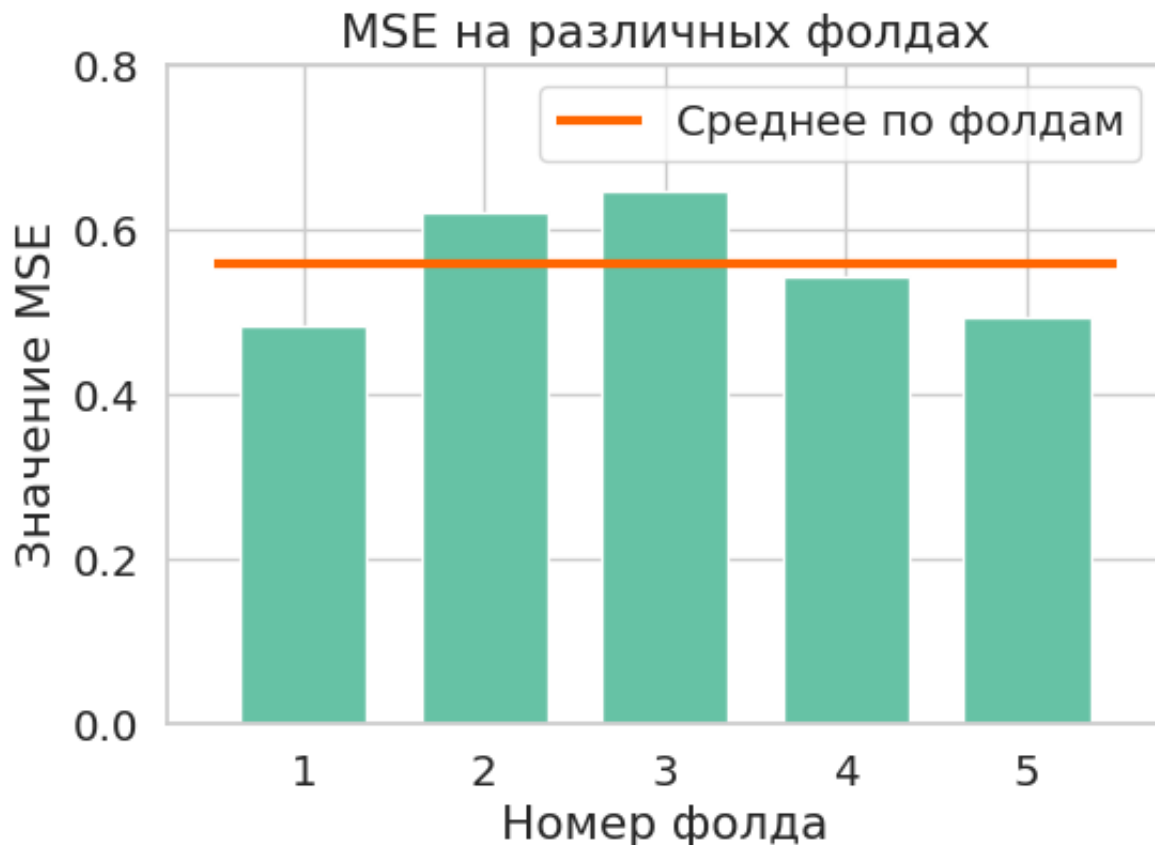
Посмотрим на распределение цен, чтобы понимать в каком масштабе находятся значения MSE.

```
1 plt.figure(figsize=(6, 4))
2 plt.hist(y)
3 plt.xlabel('Цена квартиры (сотни тыс. $)')
4 plt.ylabel('Количество')
5 plt.title('Распределение цен на квартиры в Калифорнии')
6 plt.show()
```



Визуализируем MSE на всех фолдах.

```
1 plt.figure(figsize=(6, 4))
2 plt.bar(range(1, 6), (-1)*scores, width=0.7)
3 plt.hlines(np.mean((-1)*scores), 0.5, 5.5, color='#FF6600', lw=3, label='Сред
4 plt.xlabel('Номер фолда')
5 plt.ylabel('Значение MSE')
6 plt.title('MSE на различных фолдах')
7 plt.ylim((0, 0.8))
8 plt.legend()
9 plt.show()
```



Полезно знать:

- [cross_validate](#) — позволяет задать сразу несколько метрик для подсчета качества модели. Возвращаем значения данных метрик для каждой итерации кросс-валидации в виде словаря.
- [cross_val_predict](#) — возвращает предсказания, полученные для каждого объекта выборки при кросс-валидации.

Выше мы рассмотрели функцию `cross_val_score`, которая имеет аргумент `cv`. По умолчанию данный аргумент использует стратегию кросс-валидации `KFold`, но ему можно передавать и другие стратегии кросс-валидации. Рассмотрим аналогичный способ использования `KFold` кросс-валидации, который на практике является более гибким.

Задаем стратегию кросс-валидации `KFold`.

```
1 kf = KFold(
2     n_splits=2, # количество фолдов
3     shuffle=False # перемешиваем ли данные перед разбиением
4 )
5 kf

KFold(n_splits=2, random_state=None, shuffle=False)
```

В `sklearn` объекты классов, которые соответствуют стратегиям кросс-валидации, обычно имеют два метода:

- `get_n_splits` — возвращает количество итераций, которое необходимо для заданной стратегии кросс-валидации;
- `split` — возвращает генератор индексов для разбиения данных на `train` и `test`.

Замечание.

У этого класса нет никакого механизма стратификации.

```
1 kf.get_n_splits()

2

1 kf.split(
2     X=X # данные для разбиения
3 )

<generator object _BaseKFold.split at 0x7bd78ed95150>
```

Приведем пример, демонстрирующий работу метода `split`.

```

1 data = np.array([[81, 27], [26, 45], [83, 64], [25, 98]])
2 data

array([[81, 27],
       [26, 45],
       [83, 64],
       [25, 98]])

1 for train_index, test_index in kf.split(data):
2     print("TRAIN:", train_index, "TEST:", test_index)

TRAIN: [2 3] TEST: [0 1]
TRAIN: [0 1] TEST: [2 3]

```

Случай, когда аргумент `cv` функции `cross_val_score` принимает на вход стратегию кросс-валидации. На выходе функции получаем значения метрик для каждой итерации кросс-валидации.

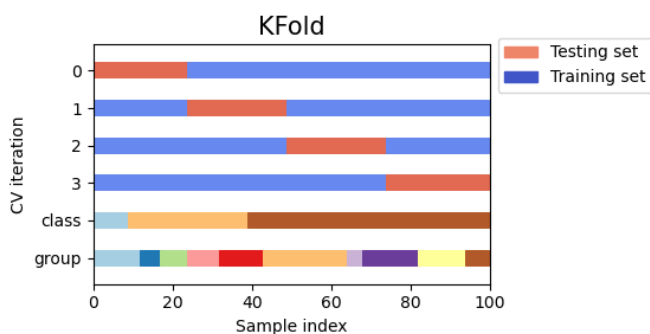
```

1 scores = cross_val_score(estimator=model, X=X, y=y, cv=kf,
2                           scoring='neg_mean_squared_error')
3 scores

array([-0.59144828, -0.54821269])

```

Визуальная интерпретация KFold кросс-валидации:



На данном графике представлена 4-Fold CV (сокращение от Cross-Validation). Каждую горизонтальную полосу стоит понимать как одну и ту же выборку из 100 элементов. По горизонтальной оси показан номер элемента выборки. По вертикальной оси сверху отложены номера фолдов. Снизу отложены разбиения выборки по классу (целевая переменная) и какому-то категориальному признаку (группа). О группах будет сказано чуть позже.

Достоинства:

- Оценивается качество модели, полученное при обучении на всех данных.
- При подборе гиперпараметров можем контролировать переобучение, т.к. выбирается модель, показавшая лучшее качество на отложенных (тестовых) фолдах. Переобучение — это ситуация, когда модель показывает хорошее качество на обучающей выборке, но плохое качество на отложенной выборке.

Недостатки:

- Значительная вычислительная сложность. Вместо одной процедуры обучения приходится обучать модель k раз.
- Никак не учитывает распределение значений целевой переменной.
- Не учитывает разбиение объектов на группы (что это такое разберемся чуть ниже).

▼ B. Leave One Out (LOO)

Данная стратегия кросс-валидации по сути является N Fold CV, где N — количество элементов в обучающей выборке. На каждой итерации мы обучаем модель на $N - 1$ элементах и оцениваем качество на оставшемся элементе.

```
1 X = [1, 2, 3, 4]
2 loo = LeaveOneOut()
3
4 # итерируемся по разбиениям множества индексов
5 for train, test in loo.split(X):
6     print("%s %s" % (train, test))
```

```
[1 2 3] [0]
[0 2 3] [1]
[0 1 3] [2]
[0 1 2] [3]
```

Достоинства:

- На каждой итерации при обучении модели используются все данные, за исключением одного элемента.
- Исследование отдельных объектов. Если на каком-то объекте допускается большая ошибка, может это выброс.
- В некоторых случаях выведены теоретические формулы результата LOO.

Недостатки:

- Огромная вычислительная сложность, не рекомендуется использовать на больших данных.
- Модель, полученная на конкретной итерации, не сильно отличается от моделей, которые получены на других итерациях. Таким образом ошибка сильно зависит от отложенного элемента, вследствие чего среди ошибок на отложенных элементах можно наблюдать высокий разброс.

▼ C. LeavePOut

Данная стратегия кросс-валидации заключается в следующем. Пусть n — размер выборки. При $p = 1$ данная стратегия эквивалентна LOO. Для оценки модели будет обучено C_n^p моделей, где на каждой итерации для обучающей выборки будет взято $n - p$ элементов, а для тестовой p элементов. Пример использования:

```
1 X = [0.76, 0.43, 0.47, 0.82, 0.22] # какая-то выборка размера 5
2 lpo = LeavePOut(p=2) # p - количество элементов в отложенном фолде
3 for train, test in lpo.split(X):
4     print("%s %s" % (train, test))
```

```
[2 3 4] [0 1]
[1 3 4] [0 2]
[1 2 4] [0 3]
[1 2 3] [0 4]
[0 3 4] [1 2]
[0 2 4] [1 3]
[0 2 3] [1 4]
[0 1 4] [2 3]
[0 1 3] [2 4]
[0 1 2] [3 4]
```

Достоинства:

- Является *исчерпывающей* стратегией кросс-валидации для заданного размера тестовой выборки, т.е. проверяет все возможные способы разделения исходной выборки на обучающее и тестовое множества заданного размера.

Недостатки:

- Огромная вычислительная сложность, которая быстро растет с увеличением параметра p . Не рекомендуется использовать на больших данных. Например, при $n = 100$ и $p = 30$ необходимо обучить примерно $3 \cdot 10^{25}$ моделей.
- На некоторых итерациях распределение целевой переменной в обучающей и тестовой выборке может быть слишком разным (нет стратификации).

Замечание.

Важно понимать, что `LeavePOut(p)` не является `KFold(n_splits=n_samples // p)`, т.к. `KFold` создает непересекающиеся тестовые множества.

▼ D. ShuffleSplit

Данная стратегия состоит в следующем:

- фиксируем количество итераций `n_splits`, т.е. количество разбиений, которое мы хотим получить.
- фиксируем размер тестовой выборки `test_size`, который будет одинаковым на каждой итерации.
- перемешиваем выборку и делим ее на две части: `train` и `test`. Проделываем это `n_splits` раз.

Пример работы:


```

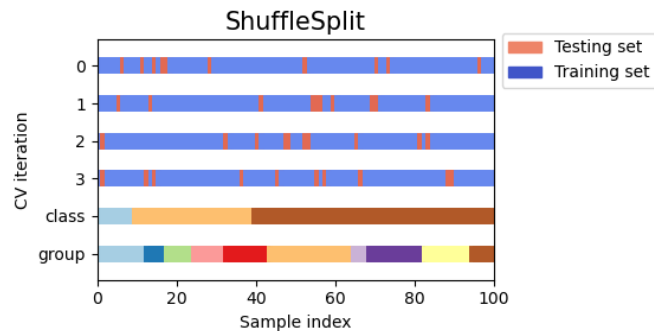
1 X = np.arange(10) # какая-то выборка размера 10
2 ss = ShuffleSplit(
3     # количество итераций перемешивания с разбиением на train и test
4     n_splits=5,
5     # доля объектов, которые хотим класть в test на каждой итерации
6     test_size=0.25,
7     random_state=0
8 )
9
10 for train_index, test_index in ss.split(X):
11     print("%s %s" % (train_index, test_index))

```

```

[9 1 6 7 3 0 5] [2 8 4]
[2 9 8 0 6 7 4] [3 5 1]
[4 5 1 0 6 9 7] [2 3 8]
[2 7 5 8 0 3 4] [6 1 9]
[4 1 0 6 8 9 3] [5 2 7]

```



Визуальная интерпретация:

Достоинства:

- Является хорошей альтернативой KFold, т.к. дает более четкий контроль над количеством итераций и разбиением на train и test.
- Результат разбиения случаен, поэтому не зависит от порядка объектов в данных.

Недостатки:

- На некоторых итерациях распределение целевой переменной в обучающей и тестовой выборке может быть слишком разным, так как разбиение случайно.
- Не учитывает разбиение объектов на группы (что это такое разберемся чуть ниже).

Замечание.

Ключевым отличием KFold от ShuffleSplit является тот факт, что в KFold каждый объект выборки в одной из итераций попадает в тестовый фолд, а в остальных итерациях используется для обучения. В ShuffleSplit разбиение каждой итерации не зависит от предыдущих итераций, объект выборки может как попасть в тестовый фолд, так и не попасть.

▼ 2.3 Стратифицированная Кросс-валидация

▼ A. Stratified KFold

При рассмотрении различных стратегий кросс-валидации выше мы неоднократно отмечали, что многие стратегии не учитывают распределение целевой переменной. В некоторых задачах это может быть критично.

Например, если вам нужно будет предсказать вероятность заболевания у пациента, то в выборке наверняка будет сильный дисбаланс классов.

Таким образом при кросс-валидации объекты выборки из положительного класса могут просто не попасть в фолды для обучения. Для решения данной проблемы используется Stratified KFold. При таком подходе каждый фолд имеет примерно такое же распределение целевого класса, как и во всем датасете. Это нужно, т.к. мы хотим получить значение метрики, которая отражает *реальное* качество модели, а значение метрики сильно зависит от баланса классов.

Например, модель, предсказывающая для всего класс 0, будет иметь хорошее качество на множестве 0, 0, 1, 0, но плохое на 1, 1, 1, 0. Сохраняя баланс классов, нам удастся получить значение метрики, которое более приближено к реальному качеству модели.

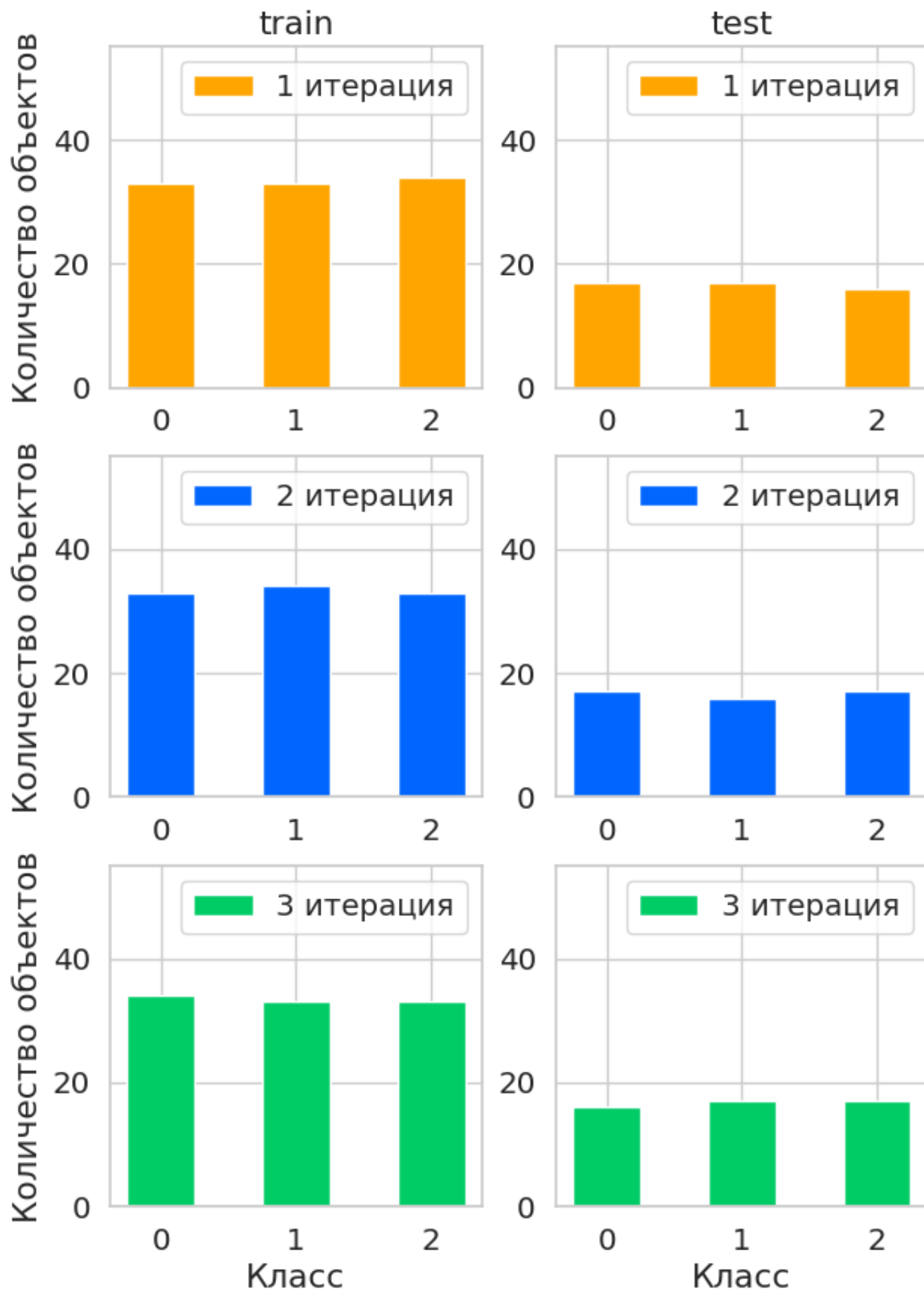
Замечание.

Стратификация работает только для классификации.

Рассмотрим применение Stratified KFold к задаче классификации ирисов. Построим графики распределения классов на каждой итерации.

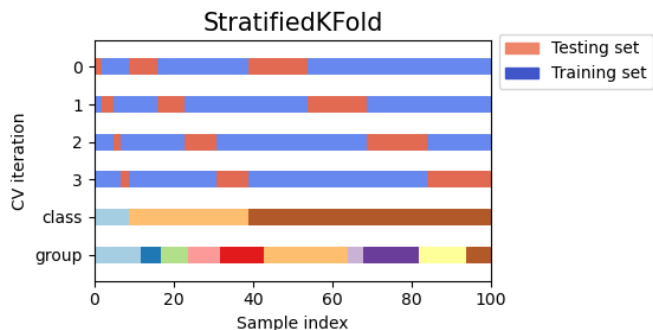
```
1 iris = datasets.load_iris()
2 X = pd.DataFrame(data=iris['data'], columns=iris['feature_names'])
3 y = iris['target']

1 skf = StratifiedKFold(n_splits=3, shuffle=True)
2
3 colors = ['orange', '#0066FF', '#00CC66']
4
5 plt.figure(figsize=(7, 10))
6 for i, (train, test) in enumerate(skf.split(X, y)):
7     plt.subplot(3,2,2*i + 1)
8     if i == 0:
9         plt.title('train')
10    values, counts = np.unique(y[train], return_counts=True)
11    plt.bar(values, counts, width=0.5, color=colors[i],
12           label='{} итерация'.format(i + 1))
13    plt.legend()
14    plt.ylim(0, 55)
15    plt.xticks([0, 1, 2])
16    if 2*i + 1 == 5:
17        plt.xlabel('Класс')
18    plt.ylabel('Количество объектов')
19    plt.subplot(3,2, 2*i + 2)
20    if i == 0:
21        plt.title('test')
22    values, counts = np.unique(y[test], return_counts=True)
23    plt.bar(values, counts, width=0.5, color=colors[i],
24           label='{} итерация'.format(i + 1))
25    plt.legend()
26    plt.ylim(0, 55)
27    plt.xticks([0, 1, 2])
28    if 2*i + 2 == 6:
29        plt.xlabel('Класс')
30
31 plt.show()
```



Ниже представлена визуальная интерпретация Stratified KFold. Стоит обратить внимание, что на каждой итерации доли каждого класса в train и test такие же, как в полном датасете.

Визуальная интерпретация:



Из построенного графика и картинки видим, что при разбиении выборки на train и test на каждой итерации кросс-валидации распределение целевого класса остается примерно одинаковым.

Достоинства:

- Учитывает распределение целевого класса при разбиении на обучающую и тестовую выборку.
- Оценивается качество модели, полученное при обучении на всех данных.
- При подборе гиперпараметров можем контролировать переобучение, т.к. выбирается модель, показавшая лучшее качество на отложенных (тестовых) фолдах.

Недостатки:

- Если не использовать `shuffle`, то результат сильно зависит от порядка объектов в данных.
- Не учитывает разбиение объектов на группы (что это такое разберемся чуть ниже).

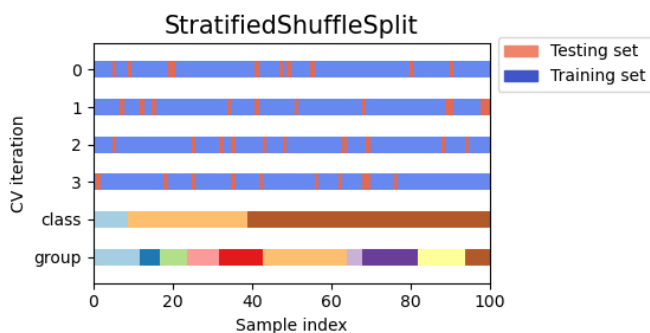
▼ B. Stratified Shuffle Split

Данная стратегия кросс-валидации делает то же самое, что и Shuffle Split, только учитывает распределение целевой переменной. Несмотря на случайность разбиения, в каждой итерации распределение на train и test такое же, как и во всем датасете.

```
1 sss = StratifiedShuffleSplit(n_splits=5, test_size=0.2,
2                             random_state=0)
3 X = np.array([[1, 2], [3, 4], [1, 2], [3, 4], [1, 2], [3, 4]])
4 y = np.array([0, 0, 0, 1, 1, 1])
5
6 for train_index, test_index in sss.split(X, y):
7     print("TRAIN:", train_index, "TEST:", test_index)
8     X_train, X_test = X[train_index], X[test_index]
9     y_train, y_test = y[train_index], y[test_index]
```

```
TRAIN: [2 5 1 3] TEST: [0 4]
TRAIN: [0 4 3 1] TEST: [2 5]
TRAIN: [0 4 3 1] TEST: [2 5]
TRAIN: [5 4 1 2] TEST: [0 3]
TRAIN: [1 5 2 4] TEST: [0 3]
```

Визуальная интерпретация:



Достоинства:

- Учитывает распределение целевого класса.
- Результат разбиения случаен, поэтому не зависит от порядка объектов в данных.

Недостатки:

- Не учитывает разбиение объектов на группы (что это такое разберемся чуть ниже).

▼ 2.4 Групповая кросс-валидация

На практике могут возникнуть ситуации, когда в таблице появляется понятие **группы**. Это категориальная переменная, которая обладает свойством целостности. Что это означает? Понятнее всего будет показать на примере.

Пример:

Представьте, что у вас есть данные с записями показателя уровня сахара в крови, причем на одного человека приходится по несколько записей. Вы хотите получить модель, которая сможет предсказывать уровень сахара на новых людях.

Таким образом, при обучении модели нужно избегать ситуации, когда и в train, и в test попадают записи, соответствующие одному и тому же человеку. В данном примере **группой** записей является человек (например, его id).

Данные стратегии кросс-валидации работают аналогично уже рассмотренным выше методам, поэтому рассмотрим их вкратце.

▼ A. Group KFold

Разбиение на фолды происходит так, что на каждой итерации в train или test нет представителей одной и той же группы. Другими словами, представители одной группы попадают либо в train фолды, либо в test фолд.

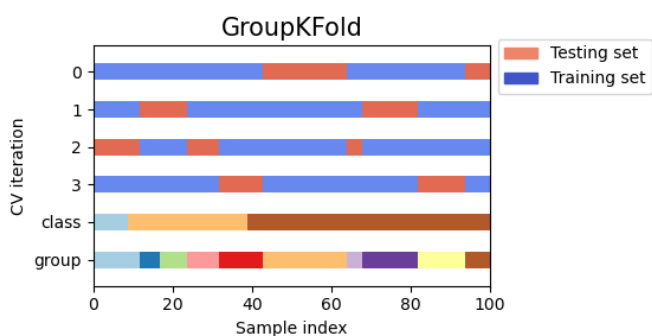
```

1 X = [0.1, 0.2, 2.2, 2.4, 2.3, 4.55, 5.8, 8.8, 9, 10]
2 y = ["a", "b", "b", "b", "c", "c", "c", "d", "d", "d"]
3 groups = [1, 1, 1, 2, 2, 2, 3, 3, 3, 3]
4
5 gkf = GroupKFold(n_splits=3)
6
7 for train, test in gkf.split(X, y, groups=groups):
8     print("%s %s" % (train, test))

[0 1 2 3 4 5] [6 7 8 9]
[0 1 2 6 7 8 9] [3 4 5]
[3 4 5 6 7 8 9] [0 1 2]

```

Визуальная интерпретация:



▼ B. Leave One Group Out

В данной стратегии на каждой итерации в качестве тестовой выборки используются только **все** представители какой-то одной группы. Таким образом, количество итераций равно количеству групп в данных.

```

1 X = [1, 5, 10, 50, 60, 70, 80]
2 y = [0, 1, 1, 2, 2, 2, 2]
3 groups = [1, 1, 2, 2, 3, 3, 3]
4 logo = LeaveOneGroupOut()
5 for train, test in logo.split(X, y, groups=groups):
6     print("%s %s" % (train, test))

[2 3 4 5 6] [0 1]
[0 1 4 5 6] [2 3]
[0 1 2 3] [4 5 6]

```


▼ C. Leave P Groups Out

В данной стратегии рассматриваются все возможные разбиения на обучающее и тестовое множества, где тестовое множество состоит из P полных групп. Таким образом, это аналог стратегии LeavePOut, оперирующий целыми группами.

```
1 X = np.arange(6)
2 y = [1, 1, 1, 2, 2, 2]
3 groups = [1, 1, 2, 2, 3, 3]
4 lpgo = LeavePGroupsOut(n_groups=2)
5 for train, test in lpgo.split(X, y, groups=groups):
6     print("%s %s" % (train, test))
```

```
[4 5] [0 1 2 3]
[2 3] [0 1 4 5]
[0 1] [2 3 4 5]
```

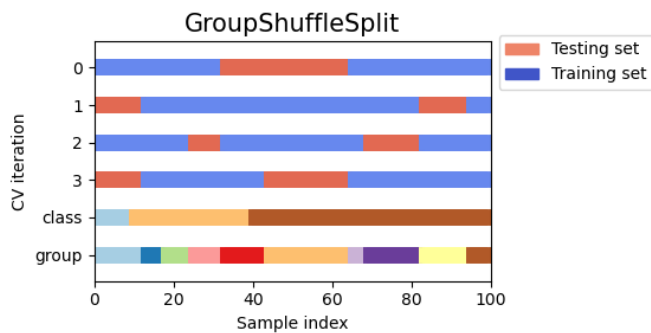
▼ D. Group Shuffle Split

В данной стратегии на каждой итерации данные разбиваются на train и test случайным образом, но так, чтобы представители одной группы не попадали одновременно и в train, и в test.

```
1 X = [0.1, 0.2, 2.2, 2.4, 2.3, 4.55, 5.8, 0.001]
2 y = ["a", "b", "b", "b", "c", "c", "c", "a"]
3 groups = [1, 1, 2, 2, 3, 3, 4, 4]
4 gss = GroupShuffleSplit(n_splits=4, test_size=0.5, random_state=0)
5 for train, test in gss.split(X, y, groups=groups):
6     print("%s %s" % (train, test))
```

```
[0 1 2 3] [4 5 6 7]
[2 3 6 7] [0 1 4 5]
[2 3 4 5] [0 1 6 7]
[4 5 6 7] [0 1 2 3]
```

Визуальная интерпретация:



▼ Итоги

Итак, выше мы рассмотрели множество стратегий кросс-валидации и узнали, в каких случаях какой из них отдавать предпочтение. О том, как использовать кросс-валидацию при настройке гиперпараметров модели, вы узнаете в ноутбуке "Поиск гиперпараметров". В качестве итога можно сказать, что кросс-валидация — это отличный метод для оценки качества вашей модели, потому что при обучении вы можете задействовать все имеющиеся данные. К сожалению, его немаловажным минусом является количество процедур обучения, так как часто требуется большое количество вычислительных мощностей, особенно если вы работаете с большими датасетами.

