

Дополнительные материалы

[Материалы](#) > [Анализ и обработка данных](#)

Продолжим работу в том же блокноте

Pipeline и ColumnTransformer

ColumnTransformer

Ранее мы рассмотрели применение пайплайна для последовательного преобразования данных и обучения модели. При этом, обратите внимание, у нас были только количественные данные (собственно, один количественный признак LSTAT).

Что делать, если у нас есть как количественные, так и категориальные признаки, и им соответственно нужны разные преобразования (более того, разным количественным и категориальным признакам также могут понадобиться разные преобразования)?

Здесь выручает **ColumnTransformer**. Он позволяет «прописать» отдельным признакам (т.е. столбцам, columns) свои преобразования, а затем объединить результат и передать в модель. Рассмотрим пример.

```
1 # создадим датасет с данными о клиентах банка
2 scoring = {
3     'Name' : ['Иван', 'Николай', 'Алексей', 'Александра', 'Евгений', 'Елена'],
4     'Age' : [35, 43, 21, 34, 24, 27],
5     'Experience' : [7, 13, 2, np.nan, 4, 12],
6     'Salary' : [95, 135, 73, 100, 78, 110],
7     'Credit_score' : ['Good', 'Good', 'Bad', 'Medium', 'Medium', 'Good'],
8     'Outcome' : [1, 1, 0, 1, 0, 1]
9 }
10
11 scoring = pd.DataFrame(scoring)
12 scoring
```

	Name	Age	Experience	Salary	Credit_score	Outcome
0	Иван	35	7.0	95	Good	1
1	Николай	43	13.0	135	Good	1
2	Алексей	21	2.0	73	Bad	0
3	Александра	34	NaN	100	Medium	1
4	Евгений	24	4.0	78	Medium	0
5	Елена	27	12.0	110	Good	1

```
1 # разобьем данные на признаки и целевую переменную
2 X = scoring.iloc[:, 1 :-1]
3 y = scoring.Outcome
4
5 # поместим название количественных и категориальных признаков в списки
6 num_col = ['Age', 'Experience', 'Salary']
7 cat_col = ['Credit_score']
8
9 # ColumnTransformer позволяет применять разные преобразователи к разным столбцам
10 from sklearn.pipeline import make_pipeline
11 from sklearn.compose import ColumnTransformer
12
13 # создадим объекты преобразователей для количественных
14 from sklearn.impute import SimpleImputer
15 imputer = SimpleImputer(strategy = 'mean')
16
17 from sklearn.preprocessing import StandardScaler
18 scaler = StandardScaler()
19
20 # и категориального признака
21 from sklearn.preprocessing import OrdinalEncoder
22 encoder = OrdinalEncoder(categories = [['Bad', 'Medium', 'Good']])
23
24 # поместим их в отдельные пайплайны
25 num_transformer = make_pipeline(imputer, scaler)
26 cat_transformer = make_pipeline(encoder)
27
28 # поместим пайплайны в ColumnTransformer
29 preprocessor = ColumnTransformer(
30     transformers=[('num', num_transformer, num_col),
31                  ('cat', cat_transformer, cat_col)])
32
33 # создадим объект модели, которая будет использовать все признаки
34 from sklearn.linear_model import LogisticRegression
35
36 model = LogisticRegression()
37
38 # создадим еще один пайплайн, который будет включать объект ColumnTransformer и
39 # объект модели
40 pipe = make_pipeline(preprocessor, model)
41
42 pipe.fit(X, y)
43
44 # сделаем прогноз
45 pipe.predict(X)
```

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

Библиотека joblib

Сохранение пайплайна

Библиотека `joblib` позволяет сохранить пайплайн в файл примерно так, как мы поступали с моделями и библиотекой `pickle`.

```
1 import joblib
2
3 # сохраним пайплайн в файл с расширением .joblib
4 joblib.dump(pipe, 'pipe.joblib')
5
6 # импортируем из файла
7 new_pipe = joblib.load('pipe.joblib')
```

```
8 |
9 | # обучим модель и сделаем прогноз
10 | new_pipe.fit(X, y)
11 | pipe.predict(X)
```

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

Кэширование функции

В качестве небольшого дополнения рассмотрим возможность `joblib` кэшировать, например, созданную нами функцию. Кэширование существенно ускоряет время ее исполнения.

```
1 | import time
2 |
3 | # напишем функцию, которая принимает список чисел
4 | # и выдает их квадрат
5 | def square_range(start_num, end_num):
6 |
7 |     res = []
8 |     # пройдемся по заданному перечню
9 |     for i in range(start_num, end_num):
10 |         res.append(i ** 2)
11 |         # искусственно замедлим исполнение
12 |         time.sleep(0.5)
13 |
14 |     return res
15 |
16 | start = time.time()
17 | res = square_range(1, 21)
18 | end = time.time()
19 |
20 | # посмотрим на время исполнения и финальный результат
21 | print(end - start)
22 | print(res)
```

```
1 | 10.014686584472656
```

```
2 | [1, 4, 9, 16, 25, 36, 49, 64, 81, 100, 121, 144, 169, 196, 225, 256, 289, 324, 361, 400]
```

Поместим функцию в кэш и вызовем из кэша в первый раз.

```
1 | # определим, куда мы хотим сохранить кэш
2 | location = '/content/'
3 |
4 | # используем класс Memory
5 | memory = joblib.Memory(location, verbose = 0)
6 |
7 | def square_range_cached(start_num, end_num):
8 |
9 |     res = []
10 |    # пройдемся по заданному перечню
11 |    for i in range(start_num, end_num):
12 |        res.append(i ** 2)
13 |        # искусственно замедлим исполнение
14 |        time.sleep(0.5)
15 |
16 |    return res
17 |
18 | # поместим в кэш
19 | square_range_cached = memory.cache(square_range_cached)
20 |
21 | start = time.time()
22 | res = square_range_cached(1, 21)
23 | end = time.time()
```

```
24 |
25 | print(end - start)
26 | print(res)
```

```
1 | 10.015617370605469
2 | [1, 4, 9, 16, 25, 36, 49, 64, 81, 100, 121, 144, 169, 196, 225, 256, 289, 324, 361, 400]
```

Вызовем из кэша еще раз.

```
1 | start = time.time()
2 | res = square_range_cached(1, 21)
3 | end = time.time()
4 |
5 | print(end - start)
6 | print(res)
```

```
1 | 0.0011799335479736328
2 | [1, 4, 9, 16, 25, 36, 49, 64, 81, 100, 121, 144, 169, 196, 225, 256, 289, 324, 361, 400]
```

Параллелизация

Параллелизация (parallelization) или параллельное исполнение кода на нескольких процессорах (CPU) может существенно увеличить скорость выполнения операций.

Посмотрим, сколько процессоров доступно в Google Colab.

```
1 | n_cpu = joblib.cpu_count()
2 | n_cpu
```

```
1 | 2
```

Напишем медленную функцию.

```
1 | def slow_square(x):
2 |     time.sleep(1)
3 |     return x ** 2
```

Применим эту функцию к числам от 0 до 9.

```
1 | %time [slow_square(i) for i in range(10)]
```

```
1 | CPU times: user 54.5 ms, sys: 7.15 ms, total: 61.7 ms
2 | Wall time: 10 s
3 | [0, 1, 4, 9, 16, 25, 36, 49, 64, 81]
```

Выполним параллелизацию.

```
1 | from joblib import Parallel, delayed
2 |
3 | # функция delayed() разделяет исполнение кода на несколько задач (функций)
4 | delayed_funcs = [delayed(slow_square)(i) for i in range(10)]
5 |
6 | # класс Parallel отвечает за параллелизацию
7 | # если указать n_jobs = -1, будут использованы все доступные CPU
8 | parallel_pool = Parallel(n_jobs = n_cpu)
9 |
10 | %time parallel_pool(delayed_funcs)
```

```
1 | CPU times: user 56.4 ms, sys: 47.8 ms, total: 104 ms
2 | Wall time: 6.01 s
3 | [0, 1, 4, 9, 16, 25, 36, 49, 64, 81]
```

```
1 | # для наглядности выведем задачи, созданные функцией delayed()
2 | delayed_funcs
```

```

1 [(<function __main__.slow_square(x)>, (0,), {}),
2  (<function __main__.slow_square(x)>, (1,), {}),
3  (<function __main__.slow_square(x)>, (2,), {}),
4  (<function __main__.slow_square(x)>, (3,), {}),
5  (<function __main__.slow_square(x)>, (4,), {}),
6  (<function __main__.slow_square(x)>, (5,), {}),
7  (<function __main__.slow_square(x)>, (6,), {}),
8  (<function __main__.slow_square(x)>, (7,), {}),
9  (<function __main__.slow_square(x)>, (8,), {}),
10 (<function __main__.slow_square(x)>, (9,), {})]

```

Встраивание функций и классов в sklearn

Библиотека sklearn предоставляет инструменты для создания собственных функций и классов, которые затем можно встроить в классы sklearn в соответствии с их парадигмой (т.е. с методами `.fit()`, `.transform()` и т.д.).

Встраивание функций

Напишем собственную функцию, которая будет кодировать категориальную переменную в соответствии с переданным ей словарем.

```

1 # напомним простой encoder
2 # будем передавать в функцию данные, столбец, который нужно кодировать,
3 # и схему кодирования (map)
4 def encoder(df, col, map_dict):
5     df_map = df.copy()
6     df_map[col] = df_map[col].map(map_dict)
7     return df_map

```

```

1 # зададим схему кодирования столбца Credit_score
2 map_dict = {'Bad' : 0,
3             'Medium' : 1,
4             'Good' : 2}

```

Теперь импортируем класс `FunctionTransformer`, которому при создании соответствующего объекта передадим созданную функцию и ее параметры в виде словаря (параметр `kw_args`).

```

1 from sklearn.preprocessing import FunctionTransformer
2
3 encoder = FunctionTransformer(func = encoder,
4                               kw_args = {'col' : 'Credit_score',
5                                           'map_dict' : map_dict})

```

`FunctionTransformer` автоматически создаст стандартные методы класса sklearn, в частности, метод `.fit_transform()`.

```

1 encoder.fit_transform(X)

```

	Age	Experience	Salary	Credit_score
0	35	7.0	95	2
1	43	13.0	135	2
2	21	2.0	73	0
3	34	NaN	100	1
4	24	4.0	78	1

5	27	12.0	110	2
---	----	------	-----	---

Встраивание классов

Помимо функций в парадигму sklearn можно встраивать и классы. Для этого новый класс должен наследовать классы BaseEstimator и TransformerMixin. Создадим такой же encoder, но на этот раз в виде класса.

```
1 # класс BaseEstimator создает методы .get_params() и .set_params()
2 # класс TransformerMixin создает .fit_transform()
3 from sklearn.base import BaseEstimator, TransformerMixin
4
5 class Encode(BaseEstimator, TransformerMixin):
6
7     # при создании объекта класса передаем столбец и схему кодирования
8     def __init__(self, col, map_dict):
9         self.col = col
10        self.map_dict = map_dict
11
12    def fit(self, df, y = None):
13        return self
14
15    def transform(self, df, y = None):
16        df_map = df.copy()
17        # применим преобразование
18        df_map[self.col] = df_map[self.col].map(self.map_dict)
19        return df_map
20
21    def inverse_transform(self, df, y = None):
22        df_inv_map = df.copy()
23        # поменяем ключи и значения местами
24        inv_map = {v: k for k, v in self.map_dict.items()}
25        # применим обратное преобразование
26        df_inv_map[self.col] = df_inv_map[self.col].map(inv_map)
27        return df_inv_map
```

Создадим объект класса Encode и применим метод `.fit_transform()`, который мы не объявляли явным образом.

```
1 encoder = Encode('Credit_score', dict(Bad = 0, Medium = 1, Good = 2))
2
3 X_trans = encoder.fit_transform(X)
4 X_trans
```

	Age	Experience	Salary	Credit_score
0	35	7.0	95	2
1	43	13.0	135	2
2	21	2.0	73	0
3	34	NaN	100	1
4	24	4.0	78	1
5	27	12.0	110	2

Попробуем метод `.inverse_transform()`.

```
1 | encoder.inverse_transform(X_trans)
```

	Age	Experience	Salary	Credit_score
0	35	7.0	95	Good
1	43	13.0	135	Good
2	21	2.0	73	Bad
3	34	NaN	100	Medium
4	24	4.0	78	Medium
5	27	12.0	110	Good

Такие вновь созданные классы sklearn, в частности, можно встраивать в pipeline.

```
1 | imputer = SimpleImputer(strategy = 'mean')
2 | scaler = StandardScaler()
3 |
4 | encoder = Encode('Credit_score', dict(Bad = 0, Medium = 1, Good = 2))
5 |
6 | num_transformer = make_pipeline(imputer, scaler)
7 | cat_transformer = make_pipeline(encoder)
8 |
9 | preprocessor = ColumnTransformer(
10 |     transformers=[('num', num_transformer, num_col),
11 |                   ('cat', cat_transformer, cat_col)])
12 |
13 | model = LogisticRegression()
14 |
15 | pipe = make_pipeline(preprocessor, model)
16 | pipe.fit(X, y)
17 | pipe.predict(X)
```

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

Подведем итог

В дополнительных материалах мы изучили инструмент ColumnTransformer, некоторые возможности библиотеки joblib, а также встраивание собственных функций и классов в парадигму sklearn.