

## ТЕХНОЛОГИЧЕСКАЯ КАРТА ЗАНЯТИЯ

**Тема занятия:** Пайплайн выбора модели

**Аннотация к занятию:** на данном занятии обучающиеся познакомятся с библиотекой Sklearn. Обсудят, как можно перебирать различные гиперпараметры алгоритма.

**Цель занятия:** сформировать у обучающихся представление о библиотеке Sklearn. Рассмотреть конвейер обработки данных с помощью различных моделей.

**Задачи занятия:**

- познакомить обучающихся с модулями grid search и pipeline библиотеки Sklearn;
- рассмотреть, как можно перебирать различные гиперпараметры алгоритма;
- применить полученные знания на практике.

### Ход занятия

Этап занятия	Время	Деятельность педагога	Комментарии, рекомендации для педагогов
<b>Организационный этап</b>	2 мин.	Добрый день! Добро пожаловать на занятие!	Приветствие. Создание в классе атмосферы психологического комфорта
<b>Постановка цели и задач занятия. Мотивация учебной деятельности обучающихся</b>	10 мин.	<p><b>Вопрос для обсуждения</b> Как обучить алгоритм?</p> <p><b>Возможные ответы обучающихся</b> На этом занятии мы обсудим, как можно перебирать различные гиперпараметры алгоритма. Иными словами, как автоматизировать поиск изначальных настроек алгоритма, наиболее подходящих для конкретных данных. На примере реальных данных мы построим конвейер, включающий исследование и перебор уже пройденных алгоритмов и их параметров. Выберем наиболее оптимальную по метрикам качества модель.</p>	Способствовать обсуждению мотивационных вопросов

### Изучение нового материала

50 мин.

В качестве данных мы используем датасет успеваемости учащихся по курсу математики в средней школе. Он состоит из 33 характеристик по типу пола, возраста, образования семьи и прочих. Подобные модели могут быть полезны для раннего выявления отстающих учеников.

Загрузим данные по успеваемости. Для удобства они были скачаны с Kaggle и загружены на удалённый облачный сервер.

```
1 df = pd.read_csv('https://dl.uploadgram.me/62e85bf854ec4h?raw')
2 df.head()
```

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob	...	fa
0	GP	F	18	U	GT3	A	4	4	at_home	teacher	...	
1	GP	F	17	U	GT3	T	1	1	at_home	other	...	
2	GP	F	15	U	LE3	T	1	1	at_home	other	...	
3	GP	F	15	U	GT3	T	4	2	health	services	...	
4	GP	F	16	U	GT3	T	3	3	other	other	...	

5 rows x 33 columns

Отообразим типы данных. Нас интересует целевая переменная: в данном случае можно выбрать одну из трёх.

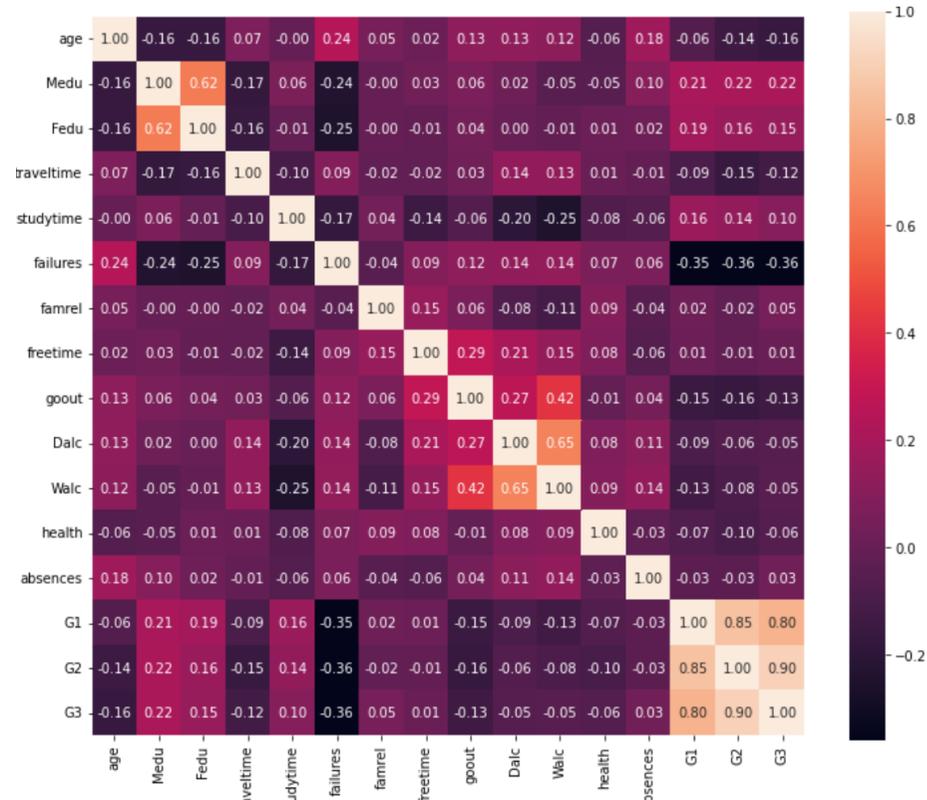
Для справки:  
 Файл для работы:  
[https://drive.google.com/file/d/1Y1bNZwulQeYE9hCWkSoiHy\\_R5qi4eelV/view?usp=sharing](https://drive.google.com/file/d/1Y1bNZwulQeYE9hCWkSoiHy_R5qi4eelV/view?usp=sharing)

```
health      int64
absences    int64
G1          int64
G2          int64
G3          int64
dtype: object
```

```
[ ]  1  df['G1'] = df['G1'].astype(int)
     2  df['G2'] = df['G2'].astype(int)
     3  df['G3'] = df['G3'].astype(int)
```

Рассмотрим коррелограмму — карту линейных корреляций признаков по Пирсону. Результаты всех трёх тестов (g1, g2, g3) схожи между собой, будем предсказывать g3. Явных линейных связей с целевой переменной от остальных признаков не наблюдается, что говорит нам о том, что простые линейные модели, вероятно, для этого не подходят.

Для справки: взаимосвязь двух переменных проявляется в совместной вариации: при изменении одного показателя меняется другой. Такая взаимосвязь называется корреляцией.



Подготовка данных

Выделим целевую переменную и разделим выборку на train и test в соотношении 80 на 20:

```
[ ] 1 from sklearn.model_selection import train_test_split
    2
    3 X_train, X_test, y_train, y_test = train_test_split(df.drop(columns = ['G1', 'G2', 'G3']), df['G3'], test_size=0.2, random_state=42)

[ ] 1 categorical_features = df.drop(columns=['G1', 'G2', 'G3']).select_dtypes(include= object).columns.values
    2 categorical_features

array(['school', 'sex', 'address', 'famsize', 'Pstatus', 'Mjob', 'Fjob',
       'reason', 'guardian', 'schoolsup', 'famsup', 'paid', 'activities',
       'nursery', 'higher', 'internet', 'romantic'], dtype=object)
```

Выделим столбцы по типу данных: категориальные и числовые.

```
▶ 1 categorical_features = df.drop(columns=['G1', 'G2', 'G3']).select_dtypes(include= object).columns.values
  2 categorical_features

▶ array(['school', 'sex', 'address', 'famsize', 'Pstatus', 'Mjob', 'Fjob',
       'reason', 'guardian', 'schoolsup', 'famsup', 'paid', 'activities',
       'nursery', 'higher', 'internet', 'romantic'], dtype=object)

▶ 1 numeric_features = df.drop(columns=['G1', 'G2', 'G3']).select_dtypes(exclude= object).columns.values
  2 numeric_features

array(['age', 'Medu', 'Fedu', 'traveltime', 'studytime', 'failures',
       'famrel', 'freetime', 'goout', 'Dalc', 'Walc', 'health',
       'absences'], dtype=object)
```

С последним обновлением библиотека Sklearn позволяет автоматизировать обработку колонок в зависимости от типа данных. Рассмотрим этот функционал:

```

1 from sklearn.pipeline import make_pipeline # конвейер подготовки данных
2 from sklearn.pipeline import Pipeline # конвейер алгоритмов
3 from sklearn.impute import SimpleImputer # пропуски заполняются медианными значениями
4 from sklearn.preprocessing import OneHotEncoder, StandardScaler # onehot encoding и стандартизации
5 from sklearn.compose import ColumnTransformer
6 from sklearn.ensemble import RandomForestRegressor
7 # конвейер подготовки данных числового типа
8 numeric_transformer = make_pipeline(SimpleImputer(strategy="median"),
9                                     StandardScaler())
10
11 # конвейер подготовки данных категориального типа
12 preprocessor = ColumnTransformer(
13     [
14         ('num', numeric_transformer, numeric_features),
15         (
16             'cat',
17             OneHotEncoder(handle_unknown='ignore', sparse=False),
18             categorical_features
19         )
20     ],
21     verbose_feature_names_out = False,
22 )
23 # подготовка и обучение алгоритма
24 alg = Pipeline(steps=[
25     ('enc', preprocessor), # подготовка
26     ('alg', RandomForestRegressor()) # алгоритм
27 ])
28
29 # обучение алгоритма
30 alg.fit(X_train, y_train)

```

В первую очередь определим функцию для обработки колонок (column transformer). Для числовых (num) и категориальных (cat) колонок мы можем определить набор (названий) колонок и функции для их обработки. В нашем случае числовые колонки будут обрабатывать numeric transform, который включает функцию заполнения медианным значением пропусков и стандартизацию. Категориальные колонки же будут заполнены с помощью OneHotEncoder. Аргумент verbose\_features\_names\_out функции ColumnTransformer будет добавлять ко всем именам признаков префикс имени преобразователя, сгенерировав этот признак при значении True. Нам это не нужно, установим его значение на False.

Теперь определим сам пайплайн, включающий себя как этап кодирования, описанный с помощью функции `ColumnTransformer`, так и сам алгоритм. В нашем случае будем использовать решающий лес (`RandomForestRegressor`). Удобство такой инициализации в том, что достаточно просто подставить вместо решающего леса любой другой алгоритм, представленный в `Sklearn`.

#### Обучение алгоритмов

Получим прогнозы на тесте и померяем метрики качества (коэффициент детерминации и медианную абсолютную ошибку). Видно, что пока алгоритм справляется относительно плохо.

```

1 # предсказания
2 preds = alg.predict(X_test)
3
4 # коэффициент детерминации и медианна абсолютная ошибка
5 r2_score(y_test, preds), median_absolute_error(y_test, preds)

```

```
(0.29472844450348523, 2.4000000000000004)
```

```

[ ] 1 # среднее и стандартное отклонение целевой переменной
     2 y_test.mean(), y_test.std()

```

```
(10.772151898734178, 4.55718471092093)
```

Попробуем задать сетку гиперпараметров алгоритма, которые мы можем перебрать для повышения его качества. Рассмотрим разное количество деревьев, минимальное количество примеров в листе и максимальное количество признаков.

Для справки: в задачах машинного обучения для оценки качества моделей и сравнения различных алгоритмов используются метрики.

```
[ ] 1 parameters = {
2     # название этапа подготовки_гиперпараметр: [сетка (список) параметров]
3     'alg_n_estimators': [100, 300, 500, 1000],
4     'alg_min_samples_leaf': [1, 2, 4, 6],
5     'alg_max_features': ['auto', 'sqrt', 'log2']
6 }
```



Передадим алгоритм, сетку гиперпараметров, количество фолдов кросс-валидации и потоков для параллельной обработки в функцию GridSearchCV, которая найдёт нам наиболее оптимальные с точки зрения метрик качества параметры.

```
1 from sklearn.model_selection import GridSearchCV
2 # сетка гиперпараметров
3 gs = GridSearchCV(alg, # выбор алгоритма или пайплайна (подготовка + алгоритм)
4                   parameters, # сетка гиперпараметров
5                   cv=5, # количество фолдов кросс-валидации
6                   verbose=2, # отображать процесс перебора гиперпараметров
7                   n_jobs=2) # параллелизм вычислений
8
9 gs.fit(X_train, y_train)
10 gs.best_params_
```

```
Fitting 5 folds for each of 48 candidates, totalling 240 fits
{'alg__max_features': 'auto',
 'alg__min_samples_leaf': 1,
 'alg__n_estimators': 1000}
```

Используя лучший алгоритм, предскажем метрики на тестовой выборке. Ошибка незначительно изменилась в лучшую сторону. Целесообразно попробовать разные алгоритмы, подставив их вместо RandomForestRegressor) и поменяв соответствующую алгоритму сетку гиперпараметров.

#### Модуль кросс-валидации

```
[ ] 1 preds = gs.best_estimator_.predict(X_test)
     2
     3 r2_score(y_test, preds), median_absolute_error(y_test, preds)
```

```
(0.29946759506767107, 2.3699999999999999)
```

<b>Закрепление изученного материала</b>	15 мин.	<b>Вопросы для обсуждения</b> <ul style="list-style-type: none"> <li>• Для чего нужна библиотека Sklearn?</li> <li>• Как обучить алгоритм?</li> </ul>	Педагог организует беседу по вопросам
<b>Этап подведения итогов занятия (рефлексия)</b>	8 мин.	<b>Вопросы для обсуждения</b> <ul style="list-style-type: none"> <li>• Чему я научился?</li> <li>• С какими трудностями я столкнулся?</li> <li>• Какие вопросы остались? Что осталось непонятным?</li> </ul>	Педагог способствует размышлению обучающихся над вопросами
<b>Информация о домашнем задании, инструктаж по его применению</b>	5 мин.	Дома повторите основные функции библиотеки Sklearn	

### Рекомендуемые ресурсы для дополнительного изучения:

1. Введение в машинное обучение с помощью Scikit-learn. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/post/264241/>
2. Библиотека Scikit-learn в Python. [Электронный ресурс] – Режим доступа: <https://pythonim.ru/libraries/biblioteka-scikit-learn-v-python>
3. От эскиза до релиза: пайплайн. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/company/lightmap/blog/550968/>