

# Шпаргалка по Scikit-learn / Sklearn Cheatsheet (XeLaTeX)

Краткий справочник по основным операциям

April 2, 2025

## Contents

1	Основной API: Глаголы ML	1
2	Pipeline: Конвейер Обработки	1
3	ColumnTransformer: Работа с Разными Типами Признаков	1
4	Подбор Гиперпараметров	2
5	Часто Используемые Импорты	3

### Scikit-learn: Ваш Швейцарский Нож в ML

**Scikit-learn (sklearn)** — это фундаментальная библиотека Python для классического машинного обучения. Она предоставляет единообразный интерфейс (API) для большинства алгоритмов и инструментов предобработки, оценки и выбора моделей. Знание её основ — **must-have**. **Ключевая сила — единообразный интерфейс**, позволяющий легко пробовать разные алгоритмы без необходимости изучать новый синтаксис для каждой модели. **Аналогия:** Представь себе ящик с инструментами. Sklearn — это такой ящик, где все инструменты (модели, препроцессоры) имеют похожие ручки (методы 'fit', 'predict', 'transform'), что сильно упрощает работу.

## 1 Основной API: Глаголы ML

### Ключевые Методы Оценщиков (Estimators)

Все объекты sklearn, которые учатся на данных (модели, препроцессоры), называются **оценщиками (estimators)**. У них есть стандартные методы:

- fit(X, y):** "Обучись". Главный метод для тренировки модели. Принимает обучающие данные (X) и целевую переменную (y, если модель с учителем). 'y' \*\*не требуется для большинства препроцессоров\*\* (например, Scaler, Encoder, Imputer): fit(X). **Аналогия:** Это как показать собаке команды (X) и правильные реакции (y), чтобы она научилась. Или как настроить инструмент по эталону (X).
- predict(X):** "Предскажи". Используется *после* fit. Генерирует предсказания для новых данных X. **Аналогия:** Попросить обученную собаку выполнить команду на новых данных.
- predict\_proba(X):** "Оцени вероятности". Используется *после* fit для **классификаторов**. Возвращает вероятности принадлежности к каждому классу (массив [n\_samples, n\_classes]). Очень полезно для оценки уверенности модели и для построения ROC-кривых. **Аналогия:** Спросить у собаки, насколько она уверена, что это именно та команда (например, 80).
- transform(X):** "Преобразуй". Используется *после* fit для **препроцессоров** (например, 'StandardScaler', 'PCA', 'OneHotEncoder'). Применяет выученное преобразование к данным X.

**Аналогия:** Использовать настроенный инструмент (например, линейку, откалиброванную по fit) для измерения новых объектов.

- fit\_transform(X, y=None):** "Обучись и преобразуй". Оптимизированная комбинация fit(X) и transform(X). Часто используется для обучающей выборки, чтобы избежать повторных вычислений. **Важно:** Применять только к **обучающей** выборке, чтобы избежать \*\*утечки информации из тестовой выборки\*\* в процесс обучения препроцессора! Для тестовой выборки используется только transform. **Аналогия:** Найти среднее и стандартное отклонение (fit) и сразу же стандартизировать данные (transform) за один проход.
- score(X, y):** "Оцени качество". Используется *после* fit. Возвращает метрику по умолчанию (ассигуру для классификаторов, R<sup>2</sup> для регрессоров) на данных X, y. Удобно для быстрой проверки.

### Базовый пример API

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.model_selection import train_test_split
4 import numpy as np
5
6 # Допустим, есть данные X (фичи) и y (цель)
7 # X = np.array(...)
8 # y = np.array(...)
9 # X_train, X_test, y_train, y_test =
   train_test_split(X, y, test_size=0.2)
10
11 # 1. Препроцессор (Scaler)
12 # Масштабирование часто улучшает сходимость и
   производительность моделей, чувствительных к масштабу
   признаков (линейные модели, SVM, нейросети, kNN).
13 scaler = StandardScaler()
14 # Обучаем на ТРЕНИРОВОЧНЫХ данных
15 scaler.fit(X_train)
16 # Преобразуем ТРЕНИРОВОЧНЫЕ и ТЕСТОВЫЕ данные
17 X_train_scaled = scaler.transform(X_train)
18 X_test_scaled = scaler.transform(X_test)
19 # Альтернатива для трейна: X_train_scaled =
   scaler.fit_transform(X_train)
20
21 # 2. Модель (Классификатор)
22 model = LogisticRegression()
23 # Обучаем модель на масштабированных ТРЕНИРОВОЧНЫХ
   данных
24 model.fit(X_train_scaled, y_train)
25
26 # 3. Предсказания
27 # Предсказания классов для теста
28 y_pred = model.predict(X_test_scaled)
29 # Предсказания вероятностей для теста
30 y_pred_proba = model.predict_proba(X_test_scaled)
31
32 print(f"Предсказанные классы: {y_pred[:5]}")
33 print(f"Вероятности классов: \n{y_pred_proba[:5]}")
34 # Оценка качества (ассигура по умолчанию для
   классификатора)
35 # print(f"Точность на тесте:
   {model.score(X_test_scaled, y_test)}")
```

## 2 Pipeline: Конвейер Обработки

### Зачем нужен Pipeline?

Часто рабочий процесс ML включает несколько шагов предобработки (масштабирование, кодирование категорий, извлечение признаков) перед

обучением модели. **Pipeline** позволяет объединить эти шаги в единый объект (оценщик). **Преимущества:**

- Удобство:** Все шаги выполняются одной командой fit и predict/transform.
- Предотвращение утечки данных (Data Leakage):** Гарантирует, что fit препроцессоров вызывается "только" на обучающих фолдах при кросс-валидации, а transform — на обучающих и валидационных/тестовых. Это **критически важно** для получения несмещенной оценки качества модели.
- Легкость настройки:** Гиперпараметры шагов пайплайна можно подбирать с помощью 'GridSearchCV'/'RandomizedSearchCV'.

**Аналогия:** Пайплайн — это как сборочный конвейер на заводе. Сырье (данные) проходит через несколько станций обработки (препроцессоры) и на выходе получается готовый продукт (предсказание модели).

### Пример Pipeline

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.impute import SimpleImputer # Для
   заполнения пропусков
3 from sklearn.preprocessing import StandardScaler,
   OneHotEncoder
4 from sklearn.linear_model import LogisticRegression
5 # Допустим, есть X_train, y_train
6
7 # Создаем пайплайн:
8 # Шаг 1: Заполнение пропусков средним
9 # Шаг 2: Масштабирование числовых признаков
10 # Шаг 3: Обучение модели
11 pipe = Pipeline([
12     ('imputer', SimpleImputer(strategy='mean')), # Имя
   шага, объект
13     ('scaler', StandardScaler()),
14     ('classifier', LogisticRegression())
15 ])
16
17 # Обучаем весь пайплайн как единое целое
18 pipe.fit(X_train, y_train)
19
20 # Делаем предсказания (данные пройдут через imputer ->
   scaler -> predict)
21 # y_pred_pipe = pipe.predict(X_test)
22 # print(f"Точность пайплайна: {pipe.score(X_test,
   y_test)}")
```

## 3 ColumnTransformer: Работа с Разными Типами Признаков

## Обработка Гетерогенных Данных

Реальные данные часто содержат столбцы разных типов: числовые, категориальные. К ним нужно применять разную предобработку (например, масштабирование к числовым, One-Hot Encoding к категориальным). **ColumnTransformer** позволяет применять разные трансформаторы к разным подмножествам столбцов. Сам 'ColumnTransformer' является таким же оценщиком (*estimator*) sklearn и может использоваться как самостоятельно, так и (что чаще всего) в качестве шага внутри 'Pipeline'. **Аналогия:** Представь, что в больницу пришел пациент (строка данных). У него есть разные проблемы (признаки): перелом (числовой признак), аллергия (категориальный). 'ColumnTransformer' — это регистратура, которая направляет пациента к нужным специалистам: травматологу (масштабирование) и аллергологу (One-Hot Encoding).

## Пример ColumnTransformer в Pipeline

```
1 from sklearn.compose import ColumnTransformer
2 from sklearn.pipeline import Pipeline
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import StandardScaler,
5   OneHotEncoder
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.model_selection import train_test_split
8 import pandas as pd
9
10 # Допустим, есть DataFrame df с числовыми и
11   категориальными фичами
12 # df = pd.DataFrame(...)
13 # X = df.drop('target', axis=1)
14 # y = df['target']
15 # X_train, X_test, y_train, y_test =
16   train_test_split(X, y, test_size=0.2)
17
18 # Определяем числовые и категориальные столбцы
19 numeric_features =
20   X_train.select_dtypes(include=['int64',
21   'float64']).columns
22 categorical_features =
23   X_train.select_dtypes(include=['object',
24   'category']).columns
25
26 # Создаем пайплайны для каждого типа признаков
27 numeric_transformer = Pipeline(steps=[
28   ('imputer', SimpleImputer(strategy='median')),
29   ('scaler', StandardScaler())])
30
31 categorical_transformer = Pipeline(steps=[
32   ('imputer', SimpleImputer(strategy='most_frequent')),
33   ('onehot', OneHotEncoder(handle_unknown='ignore'))]) #
34   handle_unknown='ignore' важен!
35
36 # Создаем ColumnTransformer
37 preprocessor = ColumnTransformer(
38   transformers=[
39     ('num', numeric_transformer, numeric_features),
40     # Имя, трансформер, колонки
41     ('cat', categorical_transformer,
42     categorical_features)])
43
44 # Создаем основной пайплайн, включающий препроцессор и
45   модель
46 model_pipe = Pipeline(steps=[('preprocessor',
47   preprocessor),
48   ('classifier', LogisticRegression())])
49
50 # Обучаем и используем как обычно
51 model_pipe.fit(X_train, y_train)
52 # y_pred = model_pipe.predict(X_test)
53 # print(f"Точность модели с ColumnTransformer:
54   {model_pipe.score(X_test, y_test)}")
```

## GridSearchCV vs RandomizedSearchCV

Большинство моделей имеют **гиперпараметры** — настройки, которые не выучиваются из данных, а задаются до обучения (например, глубина дерева, коэффициент регуляризации). Их нужно подбирать.

- **GridSearchCV**: "Поиск по сетке". Перебирает **все возможные комбинации** заданных гиперпараметров. Тщательно, но медленно, особенно если параметров много. **Аналогия:** Пекарь пробует испечь хлеб при всех комбинациях температуры (180, 200, 220) и времени (30, 40, 50 минут) — всего 3\*3=9 попыток.
- **RandomizedSearchCV**: "Случайный поиск". Выбирает **случайные комбинации** гиперпараметров из заданных диапазонов или распределений. Работает быстрее GridSearchCV, часто находит достаточно хорошие параметры за меньшее число итераций (n\_iter). Особенно эффективен, когда параметров много или не все параметры одинаково важны (позволяет быстрее найти "достаточно хорошую" область). **Аналогия:** Пекарь пробует случайные 5 комбинаций температуры и времени из возможных диапазонов. Может не найти идеал, но быстро найдет хороший вариант.

Оба инструмента используют **кросс-валидацию (CV)** для оценки качества каждой комбинации параметров на обучающей выборке, чтобы избежать переобучения на конкретном разбиении.

## Пример GridSearchCV (с Pipeline)

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.ensemble import RandomForestClassifier
3 # Используем model_pipe из примера с ColumnTransformer,
4   но заменим LogReg на RF
5 # model_pipe = Pipeline(steps=[('preprocessor',
6   preprocessor),
7   #                               ('classifier',
8   #                               RandomForestClassifier(random_state=42))])
9
10 # Задаем сетку параметров для поиска
11 # Обратите внимание на синтаксис для доступа к
12   параметрам шага пайплайна:
13 # 'имя_шага_имя_параметра'
14 param_grid = {
15   'preprocessor_num_imputer_strategy': ['mean',
16   'median'], # Параметр вложенного пайплайна
17   'classifier_n_estimators': [100, 200], # Параметр
18   модели
19   'classifier_max_depth': [None, 10, 20],
20   'classifier_min_samples_split': [2, 5]
21 }
22
23 # Создаем объект GridSearchCV
24 # cv=5 означает 5-кратную кросс-валидацию
25 # n_jobs=-1 использует все доступные ядра процессора
26 # Важно выбрать метрику (scoring), соответствующую
27   задаче! Не всегда accuracy - лучший выбор.
28 grid_search = GridSearchCV(model_pipe, param_grid,
29   cv=5, scoring='accuracy', n_jobs=-1, verbose=1)
30
31 # Запускаем поиск (может занять время!)
```

## 4 Подбор Гиперпараметров

```

24 grid_search.fit(X_train, y_train)
25
26 # Лучшие параметры и лучший скор
27 print(f"Лучшие параметры: {grid_search.best_params_}")
28 print(f"Лучший скор на CV:
29       {grid_search.best_score_:.4f}")
30
31 # Лучшая модель уже обучена на всех трейн данных с
32   лучшими параметрами
33 best_model = grid_search.best_estimator_
34 # y_pred_best = best_model.predict(X_test)
35 # print(f"Точность лучшей модели на тесте:
36   {best_model.score(X_test, y_test)}")
37
38 # Для RandomizedSearchCV синтаксис похож, но вместо
39   param_grid
40 # используется param_distributions (словари с
41   распределениями scipy.stats)
42 # и добавляется параметр n_iter.

```

## 5 Часто Используемые Импорты

### Джентльменский Набор для Старта

Этот список поможет быстро начать работу над типичной задачей.

### Основные импорты моделей и метрик

```

1 # --- Модели ---
2 # Линейные модели
3 from sklearn.linear_model import LinearRegression #
4   Регрессия
5 from sklearn.linear_model import LogisticRegression #
6   Классификация
7 from sklearn.linear_model import Ridge, Lasso #
8   Регрессия с регуляризацией
9
10 # Метод опорных векторов (SVM)
11 from sklearn.svm import SVC # Классификация
12 from sklearn.svm import SVR # Регрессия
13
14 # Деревья решений
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.tree import DecisionTreeRegressor
17
18 # Ансамбли: Бэггинг (случайный лес)
19 from sklearn.ensemble import RandomForestClassifier
20 from sklearn.ensemble import GradientBoostingRegressor
21
22 # Популярные библиотеки бустинга (устанавливаются
23   отдельно, часто производительнее)
24 # import xgboost as xgb # pip install xgboost
25 # import lightgbm as lgb # pip install lightgbm
26 # import catboost as cb # pip install catboost
27
28 # --- Предобработка ---
29 from sklearn.preprocessing import StandardScaler,
30   MinMaxScaler, RobustScaler
31 from sklearn.preprocessing import OneHotEncoder,
32   OrdinalEncoder
33 from sklearn.impute import SimpleImputer
34 from sklearn.compose import ColumnTransformer

```

```

35 from sklearn.pipeline import Pipeline, make_pipeline #
36   make_pipeline - упрощенное создание Pipeline без явного
37   именования шагов
38
39 # --- Выборка и Оценка Моделей ---
40 from sklearn.model_selection import train_test_split
41 from sklearn.model_selection import cross_val_score,
42   KFold, StratifiedKFold
43 from sklearn.model_selection import GridSearchCV,
44   RandomizedSearchCV
45
46 # --- Метрики ---
47 # Классификация
48 from sklearn.metrics import accuracy_score,
49   precision_score, recall_score, f1_score
50 from sklearn.metrics import roc_auc_score, roc_curve
51 from sklearn.metrics import confusion_matrix,
52   classification_report
53
54 # Регрессия
55 from sklearn.metrics import mean_squared_error,
56   mean_absolute_error, r2_score
57
58 # --- Другое ---
59 import numpy as np
60 import pandas as pd
61 import matplotlib.pyplot as plt
62 import seaborn as sns

```