

Шпаргалка по Scikit-learn / Sklearn Cheatsheet (XeLaTeX)

Краткий справочник по основным операциям

Содержание

1	Основной API: Глаголы ML	1
2	Pipeline: Конвейер Обработки	1
3	ColumnTransformer: Работа с Разными Типами Признаков	2
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". Использу-

ется *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 (ассигасу для классификаторов, R² 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 =
10 # train_test_split(X, y, test_size=0.2)
11
12 # 1. Препроцессор (Scaler)
13 # Масштабирование часто улучшает сходимость и
14 # производимость моделей, чувствительных к масштабу
15 # признаков (линейные модели, SVM, нейросети, kNN).
16 scaler = StandardScaler()
17 # Обучаем на ТРЕНИРОВОЧНЫХ данных
18 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 =
20 # scaler.fit_transform(X_train)
21
22 # 2. Модель (Классификатор)
23 model = LogisticRegression()
24 # Обучаем модель на масштабированных ТРЕНИРОВОЧНЫХ
25 # данных
26 model.fit(X_train_scaled, y_train)
27
28 # 3. Предсказания
29 # Предсказания классов для теста
30 y_pred = model.predict(X_test_scaled)
31 # Предсказания вероятностей для теста
32 y_pred_proba = model.predict_proba(X_test_scaled)
33
34 print(f"Предсказанные классы: {y_pred[5]}")
35 print(f"Вероятности классов: \n{y_pred_proba[5]}")
36 # Оценка качества (accuracy по умолчанию для
37 # классификатора)
38 # print(f"Точность на тесте:
39 # {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,
  OneHotEncoder
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.model_selection import train_test_split
7 import pandas as pd
```

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

4 Подбор Гиперпараметров

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 $3 \times 3 = 9$ 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,
  но заменим LogReg на RF
4 # model_pipe = Pipeline(steps=[('preprocessor',
  preprocessor),
5                               ('classifier',
  RandomForestClassifier(random_state=42))])
6
7 # Задаем сетку параметров для поиска
8 # Обратите внимание на синтаксис для доступа к
  параметрам шага пайплайна:
9 # 'имя_шага__имя_параметра'
10 param_grid = {
11     'preprocessor__num__imputer__strategy': ['mean',
12     'median'], # Параметр вложенного пайплайна
13     'classifier__n_estimators': [100, 200], # Параметр
14     модели
15     'classifier__max_depth': [None, 10, 20],
16     'classifier__min_samples_split': [2, 5]
17 }
18
19 # Создаем объект GridSearchCV
20 # cv=5 означает 5-кратную кросс-валидацию
21 # n_jobs=-1 использует все доступные ядра процессора
22 # Важно выбрать метрику (scoring), соответствующую
23 задаче! Не всегда accuracy - лучший выбор.
24 grid_search = GridSearchCV(model_pipe, param_grid,
25 cv=5, scoring='accuracy', n_jobs=-1, verbose=1)
26
27 # Запускаем поиск (может занять время!)
28 grid_search.fit(X_train, y_train)
29
30 # Лучшие параметры и лучший скор
31 print(f"Лучшие параметры: {grid_search.best_params_}")
32 print(f"Лучший скор на CV:
33 {grid_search.best_score_:.4f}")
34
35 # Лучшая модель уже обучена на всех трейн данных с
36 лучшими параметрами
37 best_model = grid_search.best_estimator_
38 # y_pred_best = best_model.predict(X_test)
39
40 # print(f"Точность лучшей модели на тесте:
41 {best_model.score(X_test, y_test)}")
42
43 # Для RandomizedSearchCV синтаксис похож, но вместо
44 param_grid
45 # используется param_distributions (словари с
46 распределениями scipy.stats)
47 # 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 #
  Регрессия
4 from sklearn.linear_model import LogisticRegression #
  Классификация
5 from sklearn.linear_model import Ridge, Lasso #
  Регрессия с регуляризацией
6
7 # Метод опорных векторов (SVM)
8 from sklearn.svm import SVC # Классификация
9 from sklearn.svm import SVR # Регрессия
10
11 # Деревья решений
12 from sklearn.tree import DecisionTreeClassifier
13 from sklearn.tree import DecisionTreeRegressor
14
15 # Ансамбли: Бэггинг (случайный лес)
16 from sklearn.ensemble import RandomForestClassifier
17 from sklearn.ensemble import RandomForestRegressor
18
19 # Ансамбли: Бустинг
20 from sklearn.ensemble import GradientBoostingClassifier
21 from sklearn.ensemble import GradientBoostingRegressor
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,
  MinMaxScaler, RobustScaler
30 from sklearn.preprocessing import OneHotEncoder,
  OrdinalEncoder
31 from sklearn.impute import SimpleImputer
32 from sklearn.compose import ColumnTransformer
33 from sklearn.pipeline import Pipeline, make_pipeline #
  make_pipeline - упрощенное создание Pipeline без явного
  именованя шагов
34
35 # --- Выборка и Оценка Моделей ---
36 from sklearn.model_selection import train_test_split
37 from sklearn.model_selection import cross_val_score,
  KFold, StratifiedKFold
38 from sklearn.model_selection import GridSearchCV,
  RandomizedSearchCV
39
40 # --- Метрики ---
41 # Классификация
42 from sklearn.metrics import accuracy_score,
  precision_score, recall_score, f1_score
43 from sklearn.metrics import roc_auc_score, roc_curve
44 from sklearn.metrics import confusion_matrix,
  classification_report
45 # Регрессия
46 from sklearn.metrics import mean_squared_error,
  mean_absolute_error, r2_score
47
48 # --- Другое ---
49 import numpy as np
50 import pandas as pd
51 import matplotlib.pyplot as plt
52 import seaborn as sns
```