

# 1. Введение в бустинг и XGBoost

**XGBoost (Extreme Gradient Boosting)** – это реализация градиентного бустинга решающих деревьев, известная высокой скоростью и точностью <sup>1</sup>. В бустинге множество **«слабых» моделей** (обычно неглубоких деревьев) последовательно добавляются к ансамблю, пытаясь исправить ошибки предыдущих моделей <sup>2</sup>. Каждый следующий дерево учится на **остатках (residuals)** – разнице между реальными значениями и предсказаниями текущего ансамбля, что эквивалентно градиентному спуску по функции потерь <sup>3</sup>. В итоге бустинг комбинирует множество простых моделей в одну сильную модель с улучшенной предсказательной способностью.

**Архитектура модели XGBoost** – это ансамбль решающих деревьев (CART), где каждое дерево добавляет свою предсказание к общей сумме <sup>4</sup> <sup>5</sup>. В отличие от случайного леса (bagging) с параллельным обучением независимых деревьев, XGBoost строит деревья *последовательно* (boosting): каждое новое дерево пытается компенсировать ошибки уже построенного ансамбля <sup>2</sup>. Модель оптимизирует специальную цель, включающую функцию потерь (например, MSE для регрессии) и штраф за сложность моделей (регуляризация) <sup>6</sup>. За счёт **L1/L2-регуляризации и контроля структуры деревьев** XGBoost сдерживает переусложнение моделей и улучшает обобщающую способность <sup>6</sup> <sup>7</sup>.

## Преимущества XGBoost:

- Высокая точность на табличных данных за счёт мощного ансамбля деревьев и продвинутой оптимизации.
- Эффективность и масштабируемость: XGBoost реализован с учётом аппаратных оптимизаций (многоядерность, работа с кэшем, приближенные алгоритмы поиска разбиений) <sup>8</sup> <sup>9</sup>, поддерживает параллельное исполнение и работу на GPU.
- Встроенная обработка пропусков: алгоритм автоматически определяет, как учитывать отсутствующие значения, выделяя для них особое направление в узлах дерева <sup>10</sup>. Это позволяет обучать модель даже при наличии пропусков без явной их импутации.
- Регуляризация и усреднение по ансамблю деревьев делают модель устойчивой к переобучению на больших объёмах данных <sup>7</sup>. Низкая скорость обучения (*small learning rate*) и уровень-ориентированный рост деревьев с ранней обрезкой (pruning) помогают избежать переобучения.

**Недостатки:** Модель XGBoost может переобучаться на небольших наборах данных или шумных признаках, если не настроить регуляризацию <sup>11</sup>. Также она менее интерпретируема, чем простые модели, и может требовать значительных вычислительных ресурсов при подборе параметров и на очень больших данных.

# 2. Установка и конфигурация XGBoost (включая GPU)

Установить XGBoost можно через `pip` или `conda`. Самый простой способ – `pip`:

```
pip install xgboost
```

Эта команда установит последнюю версию, которая включает поддержку GPU (если ваша система соответствует требованиям) <sup>12</sup>. Для Anaconda/Miniconda можно использовать:

```
conda install -c conda-forge py-xgboost
```

Conda-установка автоматически подберёт вариант с поддержкой GPU, если видеокарта NVIDIA доступна <sup>13</sup>. На Windows может потребоваться установка Microsoft Visual C++ Redistributable <sup>14</sup>.

**Настройка GPU:** После установки XGBoost может использовать GPU для ускорения обучения. Чтобы задействовать видеокарту, при инициализации модели укажите параметр устройства `device="cuda"` (требуется NVIDIA GPU с поддержкой CUDA) <sup>15</sup>. Например, для scikit-learn API:

```
import xgboost as xgb
model = xgb.XGBRegressor(tree_method="hist", device="cuda") # тренировать с GPU
```

Здесь `tree_method="hist"` включает гистограммный алгоритм построения деревьев, оптимизированный для GPU. Альтернативно, для версий XGBoost < 2.0, можно указать `tree_method="gpu_hist"`, что эквивалентно явному выбору GPU-алгоритма. Проверьте, что XGBoost распознаёт GPU: при запуске обучения в выводе не должно быть ошибок, а загрузка CPU снизится. Учтите, что на MacOS и некоторых других платформах прямая поддержка GPU может отсутствовать <sup>16</sup>.

**Пример установки и проверки GPU:**

```
pip install xgboost # установка пакета
```

```
import xgboost as xgb
from xgboost import XGBRegressor
print(xgb.__version__) # вывод версии XGBoost
model = XGBRegressor(device="cuda")
print(model.get_params()['device']) # должен показать "cuda", если GPU поддерживается
```

Если GPU недоступен, XGBoost по умолчанию будет работать на CPU. В таком случае убедитесь, что `device` установлен в `"cpu"` (или просто не указывать его).

### 3. Подготовка данных для регрессии

Перед обучением модели необходимо подготовить датасет с числовыми признаками и целевой переменной (target). Основные шаги: загрузка данных, обработка пропусков, разбиение на обучающую и тестовую выборки и при необходимости масштабирование признаков.

**Загрузка и обзор данных:** Используйте библиотеки вроде pandas для чтения данных:

```
import pandas as pd
df = pd.read_csv("data.csv")           # загрузка данных
print(df.info())                       # информация о столбцах и типах данных
print(df.isnull().sum())               # количество пропусков в каждом столбце
```

Убедитесь, что все признаки – числовые (int/float). Если есть категориальные признаки, их нужно преобразовать (например, через one-hot encoding) или использовать специальный режим XGBoost для категорий. В нашем случае предполагается, что признаки уже числовые.

**Обработка пропусков:** XGBoost умеет сам обращаться с пропущенными значениями – во время разбиения узла алгоритм может отправлять **missing** значения в одно из дочерних поддеревьев, выбирая направление, которое даёт наилучший прирост качества <sup>10</sup>. Это значит, что *необязательно* заполнять `NaN` перед обучением – можно передать данные с пропусками в `DMatrix` или `XGBRegressor`, и модель сама найдёт оптимальные места для них. Тем не менее, базовые практики обработки пропусков могут улучшить устойчивость модели. Например, можно заполнить пропуски средним или медианой:

```
# Заполнение пропусков средним по колонке
df.fillna(df.mean(), inplace=True)
```

Либо добавить явный признак-флаг отсутствия значения. Но часто достаточно доверить эту задачу самому XGBoost – он **справляется с пропусками автоматически** <sup>10</sup>.

**Масштабирование признаков:** Для деревьев решений обычно *не требуется* масштабировать или нормализовать числовые признаки, потому что алгоритм построения дерева использует только упорядочение значений (пороговые разбиения) и не зависит от абсолютных масштабов <sup>17</sup>. Например, если один признак измеряется в километрах, а другой в метрах, дерево само выберет оптимальные точки разбиения независимо от масштаба. **Нормализация не изменит порядок наблюдений и потому не влияет на разбиения** <sup>17</sup>. Однако исключением может быть ситуация с сильно разреженными признаками или необходимостью сделать веса регуляризации сопоставимыми – в большинстве случаев можно обойтись без стандартизации для XGBoost. Ниже подтверждение этого:

*“...деревья решений не требуют нормализации входных данных; и поскольку XGBoost – ансамбль деревьев, ему тоже не нужна нормализация признаков.”* <sup>18</sup>.

Таким образом, масштабирование – опционально. Вы можете применить `StandardScaler` или `MinMaxScaler` для совместимости с другими моделями в пайплайне, но непосредственно на качество XGBoost это существенно не повлияет.

**Разделение на обучающую и тестовую выборки:** Всегда отделяйте часть данных для оценки модели. Например:

```
from sklearn.model_selection import train_test_split
X = df.drop("target", axis=1)
y = df["target"]
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

Здесь 20% данных выделяется под тест. Также можно отложить часть обучающих данных как валидационную выборку для контроля качества во время обучения (например, при ранней остановке, о чём далее).

**Важность признаков:** После обучения XGBoost позволяет оценить, какие признаки наиболее значительно влияли на предсказание. Модель присваивает каждому признаку **оценку важности** (например, по частоте использования в разбиениях или по суммарному приросту качества) <sup>19</sup>. На этапе подготовки данных полезно проанализировать корреляции и отбросить нерелевантные признаки. Однако даже если вы не уверены в значимости признаков, XGBoost может сам их взвесить – после обучения с помощью метода `feature_importances_` или функции `plot_importance` можно вывести важность и принять решение об отборе признаков. Сильно коррелированные признаки лучше заранее обработать (например, оставить один из пары), чтобы снизить избыточность.

**Пример подготовки данных:**

```
import pandas as pd
from sklearn.model_selection import train_test_split

# 1. Загрузка данных
df = pd.read_csv("diamonds.csv") # датасет цен на бриллианты, например
df.dropna(inplace=True)          # удаляем строки с пропусками (для примера)

# 2. Признаки и целевая переменная
X = df.drop("price", axis=1)
y = df["price"]

# 3. Разбиение на train/test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

(Здесь мы удалили пропуски для простоты, но можно было оставить – XGBoost сам с ними работает.)

## 4. Базовое применение XGBoost для регрессии

Чтобы начать работу с XGBoost, воспользуемся его Python API. XGBoost предоставляет два основных интерфейса: **свой DMatrix API** и совместимый с **scikit-learn API** (`XGBRegressor`). Для большинства задач удобно пользоваться **scikit-learn-совместимым** классом `xgboost.XGBRegressor`, который ведёт себя аналогично другим моделям из `sklearn` (имеет методы `fit`, `predict`, `score` и т.д.) <sup>20</sup>.

**Обучение модели:** Минимальный пример – обучить модель на обучающих данных и сделать прогнозы:

```
import xgboost as xgb
from xgboost import XGBRegressor
from sklearn.metrics import mean_squared_error

# Инициализируем регрессор XGBoost с базовыми гиперпараметрами
model = XGBRegressor(objective="reg:squarederror",
                      n_estimators=100,          # число деревьев
                      learning_rate=0.1,        # скорость обучения
                      random_state=42)

# Обучение модели на тренировочных данных
model.fit(X_train, y_train)

# Предсказание на тестовой выборке
y_pred = model.predict(X_test)

# Оценка качества (RMSE)
rmse = mean_squared_error(y_test, y_pred, squared=False)
print(f"RMSE на тесте: {rmse:.2f}")
```

В этом коде мы явно указали `objective="reg:squarederror"` – цель обучения для задачи регрессии (квадрат ошибки), хотя для `XGBRegressor` можно было не указывать (по умолчанию для регрессии используется именно `reg:squarederror`). Мы задали 100 деревьев (`n_estimators=100`) и скорость обучения 0.1. После вызова `fit` модель строит ансамбль деревьев, оптимизируя MSE на тренировочной выборке. Метод `predict` возвращает предсказанные значения для тестовых наблюдений, которые затем можно сравнить с реальными значениями `y_test`.

**Выводы модели:** После обучения вы можете получить доступ к некоторым внутренним параметрам:

- `model.feature_importances_` – важности признаков (если используется бустер дерева `gbtree`).
- `model.best_score`, `model.best_iteration` – если использовалась ранняя остановка, здесь будут наилучший результат и номер итерации.
- `model.eval_result()` – история метрик по итерациям (если указаны `eval_set`).

**Использование DMatrix (альтернативный подход):** В библиотеке XGBoost есть собственный тип данных `DMatrix`, оптимизированный для хранения признаков и меток. Можно обучить модель через вызов `xgb.train`, передав параметры и `DMatrix`:

```
# Альтернативный способ через DMatrix и xgb.train
dtrain = xgb.DMatrix(X_train, label=y_train)
dtest = xgb.DMatrix(X_test, label=y_test)
params = {"objective": "reg:squarederror", "learning_rate": 0.1, "max_depth":
6}
```

```
bst = xgb.train(params, dtrain, num_boost_round=100, evals=[(dtest, "test")])
y_pred = bst.predict(dtest)
```

Этот способ эквивалентен использованию `XGBRegressor`, но даёт чуть больше контроля (например, возможность задавать callbacks, использовать XGBoost CV и пр.). Для начала работы проще придерживаться `XGBRegressor`, так как он лучше интегрируется в sklearn-пайплайн.

**Пример полного цикла:**

```
# 1. Определяем и обучаем модель
model = XGBRegressor(n_estimators=200, max_depth=5, learning_rate=0.05,
                     objective="reg:squarederror", random_state=42)
model.fit(X_train, y_train, early_stopping_rounds=10,
         eval_set=[(X_test, y_test)], verbose=False)

# 2. Предсказываем и оцениваем
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
r2 = model.score(X_test, y_test)
print(f"Test MSE: {mse:.2f}")
print(f"Test R^2: {r2:.3f}")
```

В этом примере мы задали `early_stopping_rounds=10` с указанием `eval_set` – это включило **раннюю остановку**: модель автоматически прекратит добавлять деревья, если качество (по умолчанию следится метрика на `eval_set` – `rmse`) не улучшается в течение 10 итераций. Это полезно, чтобы не переобучиться чрезмерно большим числом деревьев.

**Совет:** После базового обучения убедитесь, что модель адекватно работает: посмотрите метрику (например, RMSE, MAE) на тестовой выборке, убедитесь, что она разумна, и сравните с тривиальным прогнозом (например, всегда предсказывать среднее). Если базовая модель уже переобучилась (очень маленькая ошибка на трейне и большая на тесте), подумайте о регуляризации или уменьшении сложности (см. Раздел 9).

## 5. Полный обзор гиперпараметров модели (включая оптимизации для GPU)

XGBoost предоставляет множество гиперпараметров, позволяющих точно настроить поведение модели. Их можно разделить на несколько групп: **общие параметры**, **параметры бустера (деревьев)** и **параметры задачи (целевой функции)** <sup>21</sup>. Ниже мы рассмотрим ключевые из них, особенно влияющие на регрессионные задачи, а также параметры, касающиеся производительности на CPU/GPU.

**Общие параметры:**

- **booster** (string, по умолч. 'gbtree'): тип базового алгоритма. `gbtree` – классический ансамбль деревьев; `dart` – вариант с отсевом деревьев (Dropout); `gblinear` – линейная модель (реже используется). В регрессии почти всегда используется `gbtree`.
- **n\_estimators** (int, по умолч. 100): число деревьев (итераций бустинга) в модели. Эквивалент

параметру `num_boost_round`. Больше число деревьев повышает потенциальную сложность модели (риск переобучения), если не компенсируется снижением `learning_rate` или регуляризацией. Обычно подбирается с ранней остановкой или валидацией.

- **learning\_rate** (*float*, алиас: *eta*, по умолч. 0.3): коэффициент обучения, который снижает вклад каждого отдельного дерева. Малые значения (0.01–0.1) делают обучение более “медленным”, требуя больше деревьев, но могут повышать обобщающую способность (каждое дерево – небольшой корректирующий шаг) <sup>22</sup>. Слишком большой `learning_rate` может привести к переобучению, поэтому часто рекомендуют начинать с ~0.1 или меньше <sup>23</sup>.

- **device** (*string*, по умолч. `'cpu'`): устройство вычислений. Может быть `'cpu'` или `'cuda'` (а также `'cuda:<idx>'` для конкретного GPU) <sup>24</sup> <sup>25</sup>. На обучение модели это напрямую не влияет, но определяет, будут ли использоваться GPU-оптимизированные алгоритмы. Для использования GPU также нужно задать соответствующий метод дерева (см. ниже).

### Параметры бустера (деревьев) – основное влияние на модель:

- **max\_depth** (*int*, по умолч. 6): максимальная глубина деревьев <sup>26</sup>. Ограничивает количество разбиений от корня до листа. Более глубокие деревья могут моделировать более сложные зависимости, но склонны к переобучению. Значение 0 означает *без ограничения*, что обычно не используется (деревья будут расти, пока не вычерпается критерий разбиения). Рекомендуемый диапазон – 3 до 10 для начала <sup>27</sup>. Например, **max\_depth=3-5** для простых моделей, **max\_depth=10** для очень сложных зависимостей (но риск переобучения велик).

- **min\_child\_weight** (*int*, по умолч. 1): минимальный вес наблюдений в листе <sup>28</sup>. В случае регрессии это минимальное суммарное количество наблюдений в узле (т.к. вес равен 1 для каждого), требуемое для его разделения. Более высокие значения делают модель более консервативной: узел не разделится, если в нём мало данных. Например, `min_child_weight=5` заставит каждое листовое узел содержать как минимум 5 наблюдений. Увеличение этого параметра помогает бороться с переобучением на мелких выборках, но слишком большое значение может привести к недообучению (модель станет слишком простой). Обычно подбирается вместе с `max_depth`: сначала грубо (например, сетка 3-5-7 по `max_depth` и 1-5-10 по `min_child_weight`) <sup>29</sup>.

- **gamma** (*float*, алиас *min\_split\_loss*, по умолч. 0): порог уменьшения ошибки для разделения узла <sup>30</sup>. Если при расщеплении узла прирост метрики (например, снижение MSE) меньше `gamma`, то разбиение не производится. Таким образом, `gamma` добавляет *штраф* за каждый новый лист – чем выше `gamma`, тем более значимое улучшение нужно для создания разбиения. При `gamma=0` разбиения принимаются при любом положительном улучшении. Увеличение `gamma` делает модель проще и может предотвратить переобучение. Рекомендуется проверять небольшие значения 0, 0.1, 0.2, 0.5 и т.д. на валидации <sup>31</sup>.

- **subsample** (*float*, по умолч. 1.0): доля выборки, используемая при построении каждого дерева <sup>32</sup>. Значение 0.8 означает, что перед построением каждого нового дерева бустинга, XGBoost случайно выберет 80% обучающих наблюдений (стохастический градиентный бустинг). Это помогает уменьшить переобучение, внося разнообразие в ансамбль. Типичный диапазон – 0.5-1.0; по умолчанию 1.0 (используются все данные). При сильном переобучении можно попробовать `subsample=0.8` или 0.6.

- **colsample\_bytree** (*float*, по умолч. 1.0): доля признаков, случайно выбираемых для каждого дерева <sup>33</sup>. Аналогично Random Forest, где берётся случайное подмножество признаков при построении дерева. Значение 0.8 означает, что каждое дерево видит только 80% случайно выбранных признаков. Это ещё один способ сделать ансамбль более разнообразным и устойчивым к переобучению. Семейство параметров `colsample_by*` также включает `colsample_bylevel` (подвыбор признаков на каждом уровне глубины) и `colsample_bynode` (на каждом разбиении) <sup>34</sup> <sup>35</sup>, но чаще всего достаточно настроить `colsample_bytree`.

- **lambda** (*float*, алиас *reg\_lambda*, по умолч. 1): коэффициент L2-регуляризации на весах листьев <sup>36</sup>. Повышение `lambda` усиливает штраф за крупные значения весов (параметров) в листьях

деревьев, делая модель более гладкой. L2-регуляризация особенно полезна при наличии многих коррелирующих признаков или при переобучении. Обычно `lambda` оставляют `=1`, но в случае сильного переобучения можно увеличить (например, 5 или 10).

- **alpha** (*float*, алиас `reg_alpha`, по умолч. 0): коэффициент L1-регуляризации на веса <sup>37</sup>. Добавляет штраф за количество активных весов (способствует разреженности). L1 может занулить некоторые веса, устраняя незначимые признаки. Полезно в задачах с высокоразмерными данными. Значение `alpha=0.0` (нет L1) по умолчанию, можно попробовать 0.1, 1, 5 при настройке.
- **tree\_method** (*string*, по умолч. `'auto'`): метод построения деревьев <sup>38</sup>. Основные варианты:
  - `'exact'`: перебор всех разбиений (медленно на больших данных).
  - `'hist'`: гистограммный алгоритм (быстрее, использует бинирование признаков; обычно предпочтителен).
  - `'approx'`: устаревший приблизительный метод на основе скетчей (исторически для распределённых вычислений).
  - `'gpu_hist'`: версия гистограммного алгоритма на GPU (для XGBoost < 2.0, сейчас заменён на `device="cuda" + hist`).

Обычно `tree_method="hist"` даёт наилучший баланс скорости и качества на большом датасете <sup>39</sup> <sup>40</sup>. `auto` сам выберет `'hist'` в большинстве случаев. Если используете GPU, установите `device='cuda'` (как выше) – тогда `'hist'` автоматически будет выполняться на GPU <sup>41</sup>.

#### Параметры задачи:

- **objective** (*string*): функция потерь, которую оптимизирует бустинг. Для регрессии по умолчанию `'reg:squarederror'` (минимизация MSE). Также доступны `'reg:linear'` (то же самое, устарело), `'reg:logistic'` (логистическая регрессия по вещественной цели), специальные для выживания, ранжирования и т.д. В большинстве случаев оставляем `reg:squarederror` <sup>42</sup> <sup>43</sup>.
- **eval\_metric** (*string* или список): метрика(и) для оценки качества на валидационных наборах. По умолчанию для `reg:squarederror` – RMSE. Можно задать `'rmse'`, `'mae'`, `'rmsle'` и др. Это не влияет на обучение напрямую, но влияет на то, что отображается и на критерий ранней остановки. В задаче регрессии удобно смотреть RMSE или MAE.
- **early\_stopping\_rounds** (*int*): не гиперпараметр модели, а опция обучения – число итераций, за которое качество на валидации должно улучшиться. Если нет улучшения в течение заданного числа деревьев, обучение прекращается досрочно. Полезно для экономии ресурсов и автоматического выбора оптимального числа деревьев. Например, `early_stopping_rounds=10` остановит обучение, если метрика на `eval_set` не улучшается 10 шагов подряд <sup>44</sup>.

**GPU-специфичные настройки:** Раньше для GPU использовался отдельный параметр `gpu_id` и значение `tree_method 'gpu_hist'`. В современных версиях ( $\geq 2.0$ ) достаточно `device="cuda"`. Другие параметры, влияющие на GPU:

- **max\_bin** (*int*, по умолч. 256): количество бинов для гистограммного алгоритма. На GPU увеличение `max_bin` может улучшить качество (более точные разбиения), но замедляет и увеличивает потребление памяти.
- **sampling\_method** (*string*, по умолч. `'uniform'`): метод сэмплирования при `histogram tree_method`. `'gradient_based'` – продвинутая опция, доступная на GPU, делает сэмплирование пропорционально градиентам (GOSS) <sup>45</sup>, позволяя ставить `subsample` даже 0.1 без потери качества. Обычно не требуется менять.
- **predictor**: XGBoost сам выберет оптимальный предиктор (CPU или GPU). Если нужно, можно указать `'gpu_predictor'` для использования GPU при предсказаниях (в больших задачах).



### Прочие параметры:

- **nthread** (*int*): число потоков CPU. По умолчанию XGBoost задействует все доступные ядра. Можно ограничить, если параллельно идёт другая работа.
- **verbosity** (*int*): уровень логирования (0 – тихо, 1 – предупреждения, 2 – информация). Установите 1 или 2, чтобы видеть предупреждения, например, о неиспользуемых параметрах.

В процессе настройки стоит уделять основное внимание параметрам, сильно влияющим на **компромисс bias-variance**: *n\_estimators*, *learning\_rate*, *max\_depth*, *min\_child\_weight*, *subsample*, *colsample\_bytree*, *gamma*, *reg\_lambda* и *reg\_alpha*. Остальные можно оставить по умолчанию до более тонкой оптимизации.

## 6. Поэтапный подход к настройке модели

Настройка гиперпараметров XGBoost – это поиск баланса между недообучением и переобучением. Рекомендуется подходить к тюнингу **поэтапно**, а не пытаться перебрать сразу все комбинации. Вот пошаговая стратегия, которой можно придерживаться:

**Шаг 1: Определитесь с базовыми параметрами и числом деревьев.** Начните с разумных значений глубины деревьев и других параметров, а затем подберите `n_estimators`. Обычно фиксируют сравнительно высокий *learning\_rate* (например, 0.1) и находят сколько примерно деревьев нужно. Можно использовать встроенную кросс-валидацию `xgb.cv` или `early_stopping_rounds` на валидационном наборе, чтобы оценить оптимальное количество итераций. Например, с `learning_rate=0.1` может оказаться, что оптимально ~100–300 деревьев. Это даст ориентир для дальнейшего тюнинга <sup>46</sup>. Если оптимальное число деревьев слишком велико (>1000), можно увеличить *learning\_rate*; если слишком мало (<10), наоборот снизить его.

**Шаг 2: Подбор параметров деревьев (*max\_depth*, *min\_child\_weight*).** Эти параметры сильно влияют на сложность модели и обычно настраиваются одними из первых <sup>29</sup>. Проведите грубый поиск по сетке: например, *max\_depth* от 3 до 9, *min\_child\_weight* от 1 до 5, перебирая несколько значений. Цель – найти область, где модель начинает переобучаться или недообучаться. Обычно более глубокие деревья требуют также большего *min\_child\_weight*, чтобы не переобучиваться на мелких узлах. Начните с шага 2 (например, проверяя *max\_depth* = 4,6,8 и *min\_child\_weight* = 1,5,9) – выберите комбинацию с наилучшим скользящим контролем качества (напр. на основе RMSE/MAE на валидации). Затем можно сузить поиск вокруг лучшего значения (например, если оптимум был *max\_depth*=6, попробовать 5,6,7 подробнее) <sup>47</sup> <sup>48</sup>.

**Шаг 3: Подбор *gamma* (минимальный прирост для разбиения).** После фиксации глубины и *min\_child\_weight* добавьте регуляризацию разбиений. Проверьте несколько значений *gamma* – 0 (нет порога), 0.1, 0.2, 0.5, 1. Более высокий *gamma* обычно нужен, если даже с подобранной глубиной модель всё ещё переобучивается. Увеличивая *gamma*, вы требуете более существенного снижения ошибки при каждом расщеплении, тем самым сокращая неинформативные разбиения <sup>31</sup>. Выберите наибольшее *gamma*, которое не ухудшает качество на валидации.

**Шаг 4: Подбор *subsample* и *colsample\_bytree*.** Эти параметры отвечают за случайность отбора данных и признаков для каждого дерева. При уже настроенных глубине и *gamma*, попробуйте снизить *subsample* до 0.8 или 0.5 и посмотреть, уменьшается ли переобучение (метрика на валидации улучшается относительно тренировочной). То же с *colsample\_bytree* – попробуйте 0.8, 0.5. Обычно значение около 0.8 для обоих – хороший старт <sup>49</sup> <sup>50</sup>. Если модель недообучается

(т.е. явно не хватает точности на трейне), можно оставить эти параметры =1.0. Если переобучивается – уменьшение до 0.5 поможет. Подберите грубо эти параметры.

**Шаг 5: Регуляризация по весам (lambda, alpha).** Если после предыдущих шагов разница между ошибкой на трейне и валидации всё ещё большая, можно усилить регуляризацию весов. Попробуйте увеличить `reg_lambda` с 1 до, скажем, 5 или 10. Обратите внимание на качество – сильная L2-регуляризация может немного повысить ошибку на трейне, но заметно снизить на тесте, если была пере настройка. Также можно добавить L1-регуляризацию: `reg_alpha` = 0.1, 1 или 5 – она может обнулить некоторые веса, что полезно при большом числе признаков (feature selection). Настраивайте lambda и alpha по одному: сначала lambda, потом alpha (или наоборот), так как оба сразу труднее искать. Обычно начинают с L2 (lambda), поскольку она более предсказуемо работает.

**Шаг 6: Точная настройка learning\_rate и n\_estimators.** После установки всех вышеперечисленных параметров вернитесь к скорости обучения. Теперь можно попробовать снизить learning\_rate до, например, 0.05 или 0.01 и пропорционально увеличить n\_estimators (например, в 2-3 раза), чтобы добиться лучшего качества. Малый шаг обучения позволит модели учиться более тонким шагом и потенциально достичь более низкой ошибки. Это нужно делать *после* настройки остальных параметров, иначе подбор станет очень долгим. Например, если изначально при 0.1 learning\_rate было оптимально ~200 деревьев, вы можете попробовать learning\_rate=0.05 и до 400 деревьев, используя early\_stopping\_rounds, чтобы модель сама остановилась на нужном количестве (может быть ~350 деревьев, к примеру). Главное – убедиться, что качество на валидации действительно улучшилось. Если прироста нет, возможно, нет смысла уменьшать шаг.

**Шаг 7: Финальная проверка на тесте.** После всей оптимизации, **замерьте финальные метрики на отложенном тестовом наборе**, который не использовался при подборе параметров. Это необходимо для честной оценки обобщающей способности. Если тестовая ошибка значительно хуже валидационной – возможно, некоторые параметры всё же подобраны на переобученную валидацию; тогда стоит пересмотреть раннюю остановку или использовать более строгую кросс-валидацию в процессе тюнинга.

Следуя этому поэтапному процессу (скорее, чередуя ручной анализ с автоматическим перебором нескольких параметров на каждом шаге), вы сузите пространство гиперпараметров. Такой подход эффективнее, чем GridSearch по всем возможным сочетаниям сразу, особенно на больших данных <sup>51</sup> <sup>52</sup>. На каждом шаге вы используете знания о том, как тот или иной параметр влияет на модель: например, если видите, что модель переобучена, сначала пробуете уменьшить глубину или увеличить min\_child\_weight, а не перебираете вслепую десятки комбинаций.

## 7. Grid Search, Random Search, Bayesian Optimization с XGBoost

Помимо пошагового ручного подхода, существуют автоматизированные методы подбора гиперпараметров. Три популярных стратегии: **Grid Search** (полный перебор по сетке значений), **Random Search** (случайный перебор) и **Bayesian Optimization** (баесовская оптимизация). XGBoost поддерживает их опосредованно через инструменты sklearn и сторонние библиотеки.

**Grid Search (полный перебор):** Выбираются несколько значений для каждого параметра и проверяются все комбинации. В sklearn это делается с помощью `GridSearchCV`. Например, подберём `max_depth`, `learning_rate` и `n_estimators`:

```
from sklearn.model_selection import GridSearchCV

param_grid = {
    'max_depth': [3, 5, 7],
    'learning_rate': [0.01, 0.1, 0.2],
    'n_estimators': [100, 200]
}
model = XGBRegressor(objective='reg:squarederror', random_state=42)
grid_search = GridSearchCV(estimator=model, param_grid=param_grid,
                           cv=3, scoring='neg_mean_squared_error', n_jobs=-1)
grid_search.fit(X_train, y_train)
print("Лучшие параметры:", grid_search.best_params_)
print("Лучший RMSE CV:", (-grid_search.best_score_)*0.5)
```

Здесь `cv=3` выполняет 3-кратную кросс-валидацию для каждой комбинации <sup>53</sup> <sup>54</sup>, а `n_jobs=-1` использует все ядра CPU для параллельного перебора. Параметр `scoring='neg_mean_squared_error'` означает, что мы максимизируем отрицательный MSE (т.е. минимизируем MSE). По окончании `grid_search.best_params_` покажет оптимальную комбинацию по CV, а `grid_search.best_score_` – средний MSE (с отрицательным знаком) на ней <sup>54</sup>. Grid Search гарантированно найдёт лучший вариант на заданной сетке, но **комбинаторный взрыв** количества комбинаций ограничивает его: добавление каждого нового параметра или увеличения количества значений ведёт к резкому росту проверок. Поэтому выбирайте сетку грубо и осмысленно.

**Random Search (случайный поиск):** Вместо полного перебора задаётся распределение или список возможных значений для каждого параметра, и случайным образом выбирается некоторое количество комбинаций. Этот метод зачастую находит почти такие же хорошие параметры, как Grid Search, но за меньшее число итераций. Используется `RandomizedSearchCV` из sklearn. Например:

```
from sklearn.model_selection import RandomizedSearchCV
import numpy as np

param_dist = {
    'max_depth': np.arange(3, 10),          # случайное целое от 3 до 9
    'learning_rate': [0.001, 0.01, 0.1, 0.2, 0.3],
    'subsample': [0.5, 0.7, 0.9, 1.0],
    'colsample_bytree': [0.5, 0.7, 1.0],
    'n_estimators': np.linspace(50, 300, 6, dtype=int) # 50, 100, ..., 300
}
random_search = RandomizedSearchCV(estimator=model,
                                    param_distributions=param_dist,
                                    n_iter=20, cv=3,
                                    scoring='neg_mean_absolute_error', random_state=42)
```

```
random_search.fit(X_train, y_train)
print("Лучшие параметры:", random_search.best_params_)
```

Здесь мы указали 20 случайных комбинаций (`n_iter=20`) из заданных сеток. Random Search полезен, когда некоторые параметры не очень чувствительны к точному значению, либо когда у вас ограничено время: вы можете запустить, например, 50 итераций случайного поиска по довольно широкому диапазонам параметров. Часто **случайный поиск** за то же время покрывает более широкое пространство и имеет лучший шанс наткнуться на удачную область, чем grid search с мелкой сеткой на узком диапазоне. Теоретический анализ (Бергстра и Бенджио, 2012) показывает, что Random Search более эффективен, когда значимость разных гиперпараметров неравномерна – он не тратит лишних проверок на маловажные параметры.

**Bayesian Optimization (баесовская оптимизация):** Более продвинутый подход, который строит *surrogate* модель зависимости качества от гиперпараметров и выбирает новые комбинации с учётом предыдущих результатов. Говоря простыми словами, он пытается "угадывать", в каких областях пространства параметров может быть лучше, на основе уже проверенных точек, и проверяет там. Это значительно ускоряет поиск близко к оптимуму по сравнению со случайным блужданием. Для XGBoost популярны библиотеки **Optuna**, **Hyperopt**, **Scikit-Optimize (BayesSearchCV)**. Optuna – один из мощных инструментов для баесовской оптимизации гиперпараметров, интегрируется с XGBoost очень удобно.

Пример использования Optuna для XGBoost:

```
import optuna
from xgboost import XGBRegressor
from sklearn.model_selection import cross_val_score

def objective(trial):
    params = {
        'max_depth': trial.suggest_int('max_depth', 3, 10),
        'min_child_weight': trial.suggest_int('min_child_weight', 1, 10),
        'learning_rate': trial.suggest_loguniform('learning_rate', 1e-3,
3e-1),
        'subsample': trial.suggest_float('subsample', 0.5, 1.0),
        'colsample_bytree': trial.suggest_float('colsample_bytree', 0.5,
1.0),
        'reg_alpha': trial.suggest_float('reg_alpha', 0.0, 5.0),
        'reg_lambda': trial.suggest_float('reg_lambda', 0.0, 5.0),
        'n_estimators': trial.suggest_int('n_estimators', 50, 300)
    }
    model = XGBRegressor(objective='reg:squarederror', **params)
    # Оцениваем качество на кросс-валидации (3-фолд) по RMSE
    scores = cross_val_score(model, X_train, y_train, cv=3,
        scoring='neg_mean_squared_error')
    rmse = (-scores.mean()) ** 0.5
    return rmse

study = optuna.create_study(direction="minimize")
study.optimize(objective, n_trials=50, timeout=600) # 50 испытаний или 10
```

минут

```
print("Лучшие параметры:", study.best_params)
```

Здесь мы определили пространство параметров через методы `trial.suggest_*`. Optuna будет запускать objective-функцию, которая обучает модель с заданными параметрами и возвращает метрику (RMSE). После 50 итераций можно получить `study.best_params_`. Внутри Optuna применяется байесовская стратегия выбора следующих комбинаций (например, алгоритм Tree-structured Parzen Estimator). Таким образом Optuna находит комбинацию, **максимизирующую качество**, значительно быстрее полного перебора <sup>55</sup> <sup>56</sup>.

Стоит отметить, что Optuna можно интегрировать непосредственно с XGBoost через специальный оптимизатор `optuna.integration.XGBoostPruningCallback` для прерывания неудачных проб раньше времени, но это более продвинутые детали. Также библиотека Hyperopt предлагает функцию `fmin` для байесовской оптимизации, а Skopt – класс BayesSearchCV со схожим интерфейсом с GridSearchCV.

**Практические рекомендации:** - Используйте **GridSearchCV** для отладки на небольшом подмножестве данных или если параметров немного. Например, вы точно знаете, что нужно выбрать между глубиной 4, 5, 6 – grid search тут уместен.

- **RandomizedSearchCV** хорош, когда параметров много и вы хотите за разумное время получить неплохой результат, не обязательно идеальный.

- **Байесовская оптимизация** эффективна, когда каждая итерация обучения модели дорогая (долгая) и хочется максимально информативно выбирать следующие комбинации. Она зачастую находит лучше, но требует настройки (например, нужно указывать распределения для параметров).

При любом автоматическом поиске не забывайте: выделяйте независимую тестовую выборку *после* подбора гиперпараметров, чтобы финально проверить модель. Иначе есть риск «утечки» информации через валидацию в процесс оптимизации параметров и, как следствие, переоценки качества.

## 8. Метрики регрессии и оценка качества модели

Для задачи регрессии используются стандартные метрики, оценивающие разницу между предсказанными и фактическими значениями. Наиболее распространённые: **MSE (Mean Squared Error)**, **RMSE (корень из MSE)**, **MAE (Mean Absolute Error)** и **R<sup>2</sup> (коэффициент детерминации)**.

- **MSE** – средний квадрат ошибки:  $\frac{1}{n} \sum_i (y_i - \hat{y}_i)^2$ . Удобен в оптимизации (дифференцируем, используется как loss), но трудно интерпретировать.
- **RMSE** – квадратный корень из MSE:  $\sqrt{\frac{1}{n} \sum (y - \hat{y})^2}$ . Измеряется в тех же единицах, что и цель (например, тысячи долларов), поэтому более понятен. Является стандартной метрикой в XGBoost (по умолчанию оптимизируется RMSE).
- **MAE** – средняя абсолютная ошибка:  $\frac{1}{n} \sum |y - \hat{y}|$ . Менее чувствительна к выбросам (не квадратичная). Иногда предпочтительна, если важна равномерность ошибок, но менее удобна для градиентных методов (абсолют имеет излом в 0).

- **R<sup>2</sup>** – доля дисперсии, объяснённая моделью:  $1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$ . Значение от 0 до 1 (иногда может быть отрицательным, если модель хуже константного среднего). Показывает, насколько модель лучше тривиального предсказания среднего. R<sup>2</sup> удобен для интуитивной интерпретации качества (например, 0.9 значит 90% вариации объяснено моделью).

**Как оценивать модель:** Обычно проводятся: 1. **Валидация при обучении** – если вы используете `eval_set` в `XGBRegressor`, то можете задать `eval_metric` и наблюдать метрику на отложенной выборке по ходу обучения. Например, `eval_metric="rmse"` покажет RMSE на каждой итерации. Можно передать список метрик (например, `["rmse", "mae"]`). Эти метрики служат для контроля и ранней остановки, но не являются итоговой оценкой.

2. **Кросс-валидация** – использование `cross_val_score` из `sklearn` или `xgb.cv` для получения среднего качества на нескольких фолдах. Это даёт более надёжную оценку, особенно на небольших датасетах. Например:

```
from sklearn.model_selection import KFold, cross_val_score
cv = KFold(n_splits=5, shuffle=True, random_state=42)
scores = cross_val_score(model, X, y, cv=cv,
                          scoring='neg_mean_absolute_error')
print("Средний MAE на 5-Фолд CV:", -scores.mean())
```

Здесь мы получили средний MAE по 5 блокам (нужно взять отрицательное значение, так как `cross_val_score` возвращает отрицательную ошибку, потому что большие значения метрики считаются лучше в `sklearn`). Кросс-валидация помогает убедиться, что модель не переобучена на конкретном разбиении.

3. **Тестовая выборка** – окончательная проверка на данных, не участвовавших ни в обучении, ни в подборе параметров. Здесь вычисляем все интересующие метрики и убеждаемся, что они сопоставимы с ожиданиями.

**Пример расчёта метрик:**

```
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
rmse = mse ** 0.5
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)
print(f"Тест MSE: {mse:.2f}")
print(f"Тест RMSE: {rmse:.2f}")
print(f"Тест MAE: {mae:.2f}")
print(f"Тест R^2: {r2:.3f}")
```

Например, вывод может быть: "Тест RMSE: 5.42, Тест MAE: 3.12, R<sup>2</sup>: 0.87" – это говорит, что средняя ошибка ~5.42 (в тех же единицах, что целевая переменная), среднее отклонение 3.12, и модель объясняет ~87% дисперсии данных.

При сравнении моделей (например, разных настроек XGBoost или XGBoost vs другие алгоритмы) выбирайте метрику, релевантную бизнес-цели. RMSE популярен, но не робастен к выбросам – если задача критична к большим ошибкам, RMSE важен; если важна средняя точность – MAE.  $R^2$  просто дает общее понимание качества.

**ВАЖНО:** Используйте метрики и графики ошибок, чтобы **диагностировать модель**. Например, посмотрите распределение ошибок (Residuals). Если оно сильно скошено или имеет странные паттерны, модель систематически ошибается на некоторых диапазонах. Об этом – в следующем разделе.

## 9. Отладка модели: диагностика переобучения/недообучения, использование кросс-валидации

После того, как модель обучена и базовые метрики рассчитаны, важно понять, *насколько хорошо модель обобщает* и нет ли проблем переобучения или недообучения.

**Признаки переобучения (overfitting):** модель очень хорошо показывает себя на обучающей выборке, но на валидации или тесте ошибка значительно выше. Например, RMSE на трейне 1.0, а на тесте 5.0 – явный сигнал переобучения. XGBoost – мощный алгоритм, и без должной регуляризации он легко может переобучиться, особенно если число деревьев большое или деревья слишком глубокие <sup>11</sup>. Для диагностики: - Сравните метрики *train vs test*: `model.evals_result()` (если вы обучали с `eval_set`) даст историю значений на тренировке и валидации по итерациям. Если к концу обучения ошибка на трейне значительно ниже, чем на валидации, – модель переобучивается.

- Постройте **кривые обучения**: график ошибки в зависимости от количества деревьев. Это можно сделать, сохранив `eval_result` или запуская `xgboost.cv`. Если на графике ошибка на обучении продолжает снижаться, а на валидации начала расти – классический признак переобучения.

- Посмотрите на **важности признаков**: если модель при переобучении начинает придавать огромный вес какому-то одному неочевидному признаку или комбинации (особенно если признак может содержать шум/выбросы), это тоже индикатор.

### Что делать с переобучением:

1. Уменьшить сложность модели: снизить *max\_depth*, увеличить *min\_child\_weight*, увеличить *gamma* (т.е. потребовать более значимого разбиения) – все эти меры приведут к более простым деревьям.
2. Уменьшить число деревьев (*n\_estimators*) или использовать *early\_stopping\_rounds*, чтобы не добавлять лишние деревья после того, как валидационная ошибка перестала улучшаться. Early stopping – один из самых эффективных способов предотвратить сильное переобучение: модель сама остановится на итерации, где ошибка минимальна <sup>44</sup>.
3. Добавить или усилить **регуляризацию**: `reg_lambda`, `reg_alpha` – как описано ранее, L2/L1 штрафуют слишком большие веса. Повышение их значений сделает модель более *консервативной*, что может снизить переобучение. Также *learning\_rate* поменьше заставит модель учиться медленнее и, возможно, аккуратнее (но нужно больше деревьев, не забываем).
4. Использовать **subsampling**: `subsample`/`colsample_bytree` < 1: добавление *стохастичности* в бустинг (каждое дерево видит неполные данные или признаки) помогает смягчить переобучение. Например, `subsample=0.8`, `colsample_bytree=0.8` часто рекомендуется по умолчанию.

5. Пересмотреть данные: иногда переобучение связано с выбросами или шумом в данных. Модель может подгоняться под шум (особенно MSE штрафует большие отклонения квадратично). Стоит проверить распределение ошибок: если несколько точек имеют огромные ошибки, можно рассмотреть изменение функции потерь (например, Huber loss или quantile loss, которые поддерживаются XGBoost, как `reg:pseudohubererror` и `reg:quantileerror`), либо убрать/обработать выбросы в данных.

**Признаки недообучения (underfitting):** модель плохо обучилась и на трейне, и на тесте (ошибка высокая везде, мало лучше тривиального прогноза). Пример:  $R^2$  на трейне 0.3, на тесте 0.25. Причины: модель слишком простая (например, `max_depth` слишком маленький, или слишком сильная регуляризация). Решения: увеличить глубину, разрешить больше деревьев, уменьшить регуляризацию, добавить признаки или иным способом усложнить модель, пока она не начнет лучше подгоняться под обучающие данные. Но убедитесь, что не зашли слишком далеко и не возникло переобучения.

**Кросс-валидация для диагностики:** Используйте k-fold CV на этапе оценки, чтобы убедиться, что качество модели стабильно на разных подвыборках данных. Если качество сильно варьируется от фолда к фолду, модель может быть неустойчивой – возможно, признаков мало или они неинформативны, или модель слишком гибкая. В таких случаях, помимо настройки параметров, можно попробовать *бэггинг* – усреднение нескольких моделей, или, напротив, упростить модель.

**Пример: контроль переобучения с кросс-валидацией и ранней остановкой:**

```
# Используем XGBoost CV для наглядности
import xgboost as xgb
dtrain = xgb.DMatrix(X_train, label=y_train)
params = {"objective": "reg:squarederror", "max_depth": 6, "learning_rate": 0.1}
cv_results = xgb.cv(params, dtrain, num_boost_round=1000, nfold=5,
                    metrics="rmse", early_stopping_rounds=10, seed=42)
print("Лучшее число деревьев по CV:", cv_results.shape[0])
print("RMSE на последней итерации:", cv_results["test-rmse-mean"].min())
```

Здесь `xgb.cv` разбивает данные на 5 фолдов и тренирует 1000 деревьев с ранней остановкой. Выход в `cv_results` покажет, на каком раунде остановилось (например, 120 деревьев) и минимальный средний RMSE на тестовых фолдах. Этот подход в целом совпадает с `manual early_stopping` на `eval_set`, но CV более устойчив к случайным разбиениям. Если CV и `train` сильно расходятся, это намёк на переобучение.

**Диагностика по графикам:** Постройте график фактических vs предсказанных значений:

```
import matplotlib.pyplot as plt
plt.scatter(y_test, y_pred, alpha=0.5)
plt.xlabel("Actual values")
plt.ylabel("Predicted values")
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()],
         color='red')
plt.show()
```



Если модель идеальна, точки лягут близко к красной линии  $y = \hat{y}$ . Рассеивание вокруг линии показывает погрешности. Тренд, отличающийся от 45°, укажет на смещение.

Также посмотрите **график резидуалов**:

```
residuals = y_test - y_pred
plt.scatter(y_pred, residuals, alpha=0.5)
plt.axhline(0, color='red')
plt.xlabel("Predicted")
plt.ylabel("Residual (actual - pred)")
plt.show()
```

В идеале остатки распределены случайно вокруг нуля (красная линия). Если видна структура (например, для малых предсказаний остатки в основном положительные, а для больших – отрицательные), это признак, что модель систематически недоучивает какую-то нелинейность или эффект (возможно, нужно добавить признак или усложнить модель). Если остатки имеют несколько выбросов далеко от 0 – эти точки модель плохо предсказывает, можно их изучить отдельно (возможно, аномалии).

**Вывод:** отладка модели заключается в том, чтобы с помощью кросс-валидации и анализа ошибок убедиться, что модель имеет хорошую обобщающую способность. Если что-то неладно – возвращаемся к настройке гиперпараметров или предобработке данных. XGBoost предоставляет средства для отслеживания обучения (eval\_set, лог метрик), которыми стоит пользоваться для этой цели. В целом градиентный бустинг довольно *устойчив*, если задать разумные ограничения и иметь достаточно данных: «Boosting is resilient and robust method that prevents and curbs overfitting quite easily» <sup>57</sup> (усилиями вроде небольшого шага обучения, субсемплинга и регуляризации, как мы обсудили).

## 10. Визуализация деревьев, важности признаков и ошибок

Визуализация помогает лучше понять работу модели XGBoost, особенно важна для объяснения принятия решений и отладки. Основные аспекты для визуализации: **структура отдельных деревьев, значимость признаков, ошибки модели (резидуалы)**.

**Визуализация деревьев:** XGBoost позволяет вывести каждое решающее дерево из ансамбля. Для этого есть метод `xgboost.plot_tree(booster, num_trees=...)`, который рисует схему дерева с разбиениями. Если вы использовали `XGBRegressor`, можно получить booster через `model.get_booster()`. Пример – нарисуем первое дерево модели:

```
import matplotlib.pyplot as plt
xgb.plot_tree(model, num_trees=0)
plt.show()
```

Этот график может быть большим, поэтому можно ограничить глубину при выводе или увеличить размер холста (`plt.figure(figsize=(width, height))`). В узлах дерева будут

показаны условия разбиения (напр. `feature_2 < 1.5`), значения весов на листьях и некоторые статистики. Структура дерева даст представление, какие признаки и пороги использует модель на ранних уровнях (первые разбиения – самые информативные глобально). Однако для моделей с сотнями деревьев рассматривать каждое – тяжёлая задача, лучше сфокусироваться на совокупной информации (например, важности признаков).

Если требуется отдельный файл, можно выгрузить дерево в формат Graphviz DOT:

```
model.get_booster().dump_model('tree_dump.txt', with_stats=True)
```

или

```
xgb.to_graphviz(model, num_trees=0, rankdir='LR').render('tree0')
```

При наличии установленного Graphviz этот код сохранит визуализацию первого дерева в файл `tree0.pdf`.

**Важность признаков:** XGBoost считает несколько видов важности признаков: по умолчанию `feature_importances_` в `XGBRegressor` выдаёт значения важности по частоте (frequency/weight – сколько раз признак использован для разбиений). Также доступны важности по суммарному приросту качества (gain) и по покрытию (cover) – их можно получить через `model.get_booster().get_score(importance_type='gain')` и др. На практике **gain-важность** более информативна – она показывает, насколько данный признак в среднем улучшает качество, когда используется в дереве.

Для визуализации удобно использовать `xgboost.plot_importance`:

```
ax = xgb.plot_importance(model, importance_type='gain')
plt.show()
```

Получится бар-график с признаками и их относительной важностью <sup>19</sup>. С его помощью можно быстро увидеть топ-5 или топ-10 наиболее влиятельных признаков. Имейте в виду, что если признаки скоррелированы, важность может распределиться между ними произвольно. Также низкая важность не всегда означает бесполезность признака (возможно, он не нужен в присутствии других). Тем не менее, такой анализ помогает **сократить признаки** или, по крайней мере, объяснить, какие факторы влияют на предсказания. Например, Neptune.ai показал график важности, где ведущими факторами цены бриллианта оказались вес (carat) и глубина огранки

<sup>58</sup>.

Кроме того, существуют продвинутые способы интерпретации вроде **SHAP values** (SHapley Additive exPlanations) – XGBoost имеет встроенную поддержку расчёта вкладов признаков через `model.predict(X, pred_contribs=True)`. SHAP даёт для каждого объекта вклад каждого признака в отклонение от среднего прогноза. Графики SHAP позволяют понять локально, почему модель так решила для конкретного примера. Это полезно, если вам нужно интерпретировать модель на уровне отдельных предсказаний.

**Визуализация ошибок (резидуалов):** В разделе 9 мы уже рассмотрели графики *предсказание vs реальность* и *предсказание vs остаток*. Добавим к этому: - Постройте **гистограмму ошибок**:

```
plt.hist(y_test - y_pred, bins=30)
plt.xlabel("Error (actual - pred)")
plt.ylabel("Frequency")
plt.show()
```

Это покажет распределение ошибок модели. В идеале оно должно быть приблизительно симметричным вокруг нуля (немного смещённым вправо, если у вас MSE-оптимизация, т.к. модель склонна чуть недопредсказывать большие значения, но в целом без сильных выбросов). Если гистограмма ошибок имеет длинный хвост – есть наблюдения, где модель сильно ошиблась; можно идентифицировать их и проанализировать отдельно (возможно, особенности данных, которых нет в обучении, или просто сложные случаи).

- **Парные графики «признак – ошибка»:** Для каждого важного признака можно построить график остатка от значения признака:

```
plt.scatter(X_test['feature_name'], y_test - y_pred)
plt.axhline(0, color='red')
plt.xlabel("Feature value")
plt.ylabel("Residual")
```

Например, если увидим, что для какого-то признака при высоких значениях остатки систематически положительные (реальные выше предсказанных) – модель недооценивает цель на высоких значениях этого признака, возможно, связь не линеаризована или нужен полином/взаимодействие. Такие визуализации помогают **внести улучшения в фичи**, что зачастую эффективнее, чем тонкая настройка модели.

**Пример визуализации (код):**

```
# 1. График важности признаков
xgb.plot_importance(model, max_num_features=10)
plt.title("Top-10 важных признаков")
plt.show()

# 2. Визуализация одного из деревьев (например, 0-го)
plt.figure(figsize=(12, 8))
xgb.plot_tree(model, num_trees=0, rankdir="LR")
plt.title("Дерево 0 модели XGBoost")
plt.show()

# 3. Распределение ошибок
errors = y_test - y_pred
plt.hist(errors, bins=20)
plt.xlabel("Error")
plt.ylabel("Count")
plt.show()
```

(График дерева может быть очень загруженным; параметр `rankdir="LR"` рисует дерево слева-направо для компактности.)

**Визуализация кривых обучения:** Стоит упомянуть, что еще один полезный график – динамика метрики в процессе обучения. Если вы сохраняли `eval_metric` на валидации, можно построить график `iterations` vs `metric`. Это поможет увидеть, на каком примерно количестве деревьев началось переобучение (кривая теста перестала улучшаться или пошла вверх). Пример:

```
results = model.evals_result()
plt.plot(results['validation_0']['rmse'], label='Validation RMSE')
plt.plot(results['training']['rmse'], label='Training RMSE')
plt.legend()
plt.xlabel("Boosting Round")
plt.ylabel("RMSE")
plt.show()
```

(Работает, если при `fit` вы передавали `eval_set=[(X_test, y_test)]` и `eval_metric="rmse"` – тогда `model.evals_result()` заполнится.)

Резюмируя: **визуализация – мощный инструмент** для интерпретации и отладки XGBoost-модели. Она превращает “черный ящик” ансамбля деревьев в понятные артефакты: деревья как правила принятия решений, важности признаков как обобщённый вклад, графики ошибок как отражение недочётов модели. Эти средства помогут вам как эксперту по ML понимать, что происходит внутри модели и как её улучшить.

## 11. Использование XGBoost с pandas и scikit-learn (интерфейс API, совместимость)

Одним из преимуществ XGBoost является полноценная совместимость с экосистемой `scikit-learn` и `pandas`. Класс `xgboost.XGBRegressor` реализует стандартный интерфейс `estimator`’а `scikit-learn` <sup>20</sup>, поэтому вы можете легко интегрировать его в ML-пайплайны наряду с другими преобразованиями.

**Входные данные:** `XGBRegressor` принимает на вход данные в виде numpy массивов или `pandas DataFrame`. Если вы передаёте `DataFrame`, XGBoost сохранит информацию о названиях столбцов (feature names), что полезно для интерпретации и сохранения модели. Например:

```
import pandas as pd
from xgboost import XGBRegressor

# Допустим, у нас данные уже в pandas DataFrame X_train, X_test
model = XGBRegressor()
model.fit(X_train, y_train)          # можно передавать DataFrame напрямую
y_pred = model.predict(X_test)
```

Внутри XGBoost преобразует DataFrame в DMatrix, но при сохранении модели запомнит имена признаков, что потом отразится при печати дерева или важностей.

**Совместимость со sklearn Pipeline и трансформерами:** Вы можете использовать XGBRegressor внутри `sklearn.pipeline.Pipeline` или `FeatureUnion` наряду с преобразованиями признаков. Например, если нужно стандартизовать некоторые признаки или сделать one-hot для категорий до подачи в модель, всё это можно оформить как Pipeline:

```
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler

pipeline = Pipeline([
    ('scaler', StandardScaler()),          # масштабирование признаков
    ('regressor', XGBRegressor(n_estimators=100))
])
pipeline.fit(X_train, y_train)
y_pred = pipeline.predict(X_test)
```

Здесь `StandardScaler` применится к данным, а затем `XGBRegressor` обучится на нормализованных признаках. (Хотя, как мы отмечали, для деревьев нормализация не обязательна, но в составе общего конвейера, либо если сочетать с другими моделями, это может быть нужно.)

**GridSearchCV/RandomizedSearchCV:** Использование `XGBRegressor` с этими классами ничем не отличается от использования, например, `RandomForestRegressor`. Вы можете передавать параметры XGB в `param_grid`, как мы делали в разделе 7. Стоит упомянуть, что в `sklearn 1.0+` реализована поддержка **пользовательских ранних остановок** через встроенные callbacks. В сочетании с `XGBRegressor` это позволяет проводить, например, `RandomizedSearchCV` с ранней остановкой для каждой комбинации (не тратя время на заведомо переобученные большие числа деревьев). Делается это через класс `EarlyStopping` из `sklearn.model_selection`, но на практике можно обойтись встроенным `early_stopping_rounds` в XGBoost, если грамотно обрабатывать.

**Кросс-валидация:** Помимо описанного `xgb.cv`, можно использовать `sklearn.model_selection.cross_val_score` прямо с `XGBRegressor`:

```
from sklearn.model_selection import cross_val_score
scores = cross_val_score(XGBRegressor(n_estimators=100), X, y, cv=5,
    scoring='neg_mean_squared_error')
print(f"Средний RMSE (5-fold CV): {(-scores.mean())*0.5:.2f}")
```

Это удобно и показывает, что `XGBRegressor` полностью соответствует интерфейсу: у него есть методы `fit`, `predict`, и также реализован `score` (который по умолчанию для регрессора возвращает  $R^2$ ). Можно указать `scoring = 'neg_mean_squared_error', 'neg_mean_absolute_error'` или использовать свои метрики.

**Pandas interoperability:** Если вы хотите воспользоваться преимуществом `pandas` (например, удобная фильтрация, обработка NaN), делайте это перед передачей данных в модель.

XGBRegressor вернёт предсказания в виде numpy array (который легко снова преобразовать в Series, если нужно).

**Совместимость версий:** Убедитесь, что версии sklearn и xgboost актуальны, чтобы избежать предупреждений. Например, старые версии XGBoost могли не поддерживать последние изменения sklearn. На момент 2025 года XGBoost 1.6+ хорошо интегрируется со sklearn 1.x.

**Пример с sklearn API:**

```
from sklearn.datasets import fetch_california_housing
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error

# Загрузим датасет калифорнийской недвижимости из sklearn
X, y = fetch_california_housing(return_X_y=True, as_frame=True)
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)

# Обучение XGBRegressor
model = XGBRegressor(n_estimators=200, max_depth=4, learning_rate=0.1)
model.fit(X_train, y_train)

# Оценка
y_pred = model.predict(X_test)
print(f"Test RMSE: {mean_squared_error(y_test, y_pred, squared=False):.2f}")
```

Обратите внимание: мы загрузили данные как pandas DataFrame (`as_frame=True`), разбили их и передали модели. Модель успешно обучилась. **Таким образом, XGBoost легко работает “из коробки” с pandas DataFrame и всеми инструментами sklearn** <sup>59</sup> <sup>60</sup> .

**Дополнительно:** XGBoost совместим с `sklearn.inspection.permutation_importance` (для оценки значимости признаков перемешиванием), а также с такими библиотеками как SHAP (которые могут брать модель типа XGBRegressor и данные, и возвращать объяснения). Также вы можете сохранить XGBRegressor вместе с предобработкой с помощью `sklearn.pipeline.Pipeline` и потом сериализовать пайплайн через joblib.

Подытоживая, интеграция XGBoost в существующий ML-процесс не требует специальных шагов – разработчики позаботились о том, чтобы интерфейс был интуитивно понятен пользователям sklearn. На практике вы получите преимущества XGBoost (скорость и точность) в сочетании с удобством pandas/sklearn (легкость предобработки, оценки, прототипирования).

## 12. Экспорт модели, сохранение и повторное использование

После того, как вы обучили модель XGBoost и удовлетворены её качеством, скорее всего, вы захотите сохранить модель на диск, чтобы затем загрузить для предсказаний (например, в продакшене) или для дальнейшего анализа. XGBoost предлагает несколько способов сериализации модели:

## 1. Сохранение встроенным методом XGBoost (рекомендуется):

Каждый объект `Booster` или `XGBRegressor` имеет метод `save_model(filename)`, который сохранит модель в файл. Начиная с версий 1.0+, XGBoost использует формат JSON/UBJ для сохранения моделей <sup>61</sup>, который стабилен и совместим между версиями (в отличие от устаревшего бинарного). Рекомендуется сохранять с расширением `.json` – тогда XGBoost сохранит только сами деревья и параметры задачи, без лишнего (то есть “модель”, а не снимок сессии) <sup>62</sup> <sup>63</sup>. Пример:

```
# Предположим, model - это обученный XGBRegressor
model.save_model("xgb_model.json")
```

Этот файл содержит структуру всех деревьев, веса листьев, параметры обучения (например, `objective`), имена признаков и т.д. Он **независим от Python** – то есть вы можете загрузить его, например, в R или C++ версии XGBoost, и сделать прогнозы (очень удобно для продакшена, где может быть другая технология). Загрузка модели:

```
model2 = xgb.XGBRegressor()
model2.load_model("xgb_model.json")
# model2 теперь содержит ту же модель, можно предсказывать:
y_pred2 = model2.predict(X_test)
```

Методы `save_model` / `load_model` – предпочтительный способ сохранения, поскольку XGBoost гарантирует поддержку формата модели между версиями <sup>64</sup>. В частности, JSON/UBJ формат будет поддерживаться и в будущих версиях (разработчики обещают обратную совместимость для моделей) <sup>65</sup>.

Стоит различать `save_model` и `dump_model`: второй используется для получения *читаемого описания* модели (например, в текстовом файле с весами и сплитами), но не для повторного использования модели в коде <sup>66</sup>. Для восстановления модели применяйте именно `save_model` / `load_model` <sup>67</sup>.

## 2. Сериализация средствами Python (pickle/joblib):

Поскольку `XGBRegressor` – это Python-объект, его можно сохранить стандартными способами, например:

```
import joblib
joblib.dump(model, "model.joblib")
# ... позже
model_loaded = joblib.load("model.joblib")
```

или с помощью `pickle`:

```
import pickle
pickle.dump(model, open("model.pkl", "wb"))
model_loaded = pickle.load(open("model.pkl", "rb"))
```

Это сохранит весь объект вместе с внутренним бустером. Такой подход тоже работает (и часто используется при простом прототипировании), однако у него есть **недостатки**: файлы могут получиться большими, и главное – бинарная совместимость pickle может нарушиться при обновлении версии XGBoost или Python. Например, модель, сериализованная joblib в XGBoost 1.3, может не загрузиться в XGBoost 1.6, так как внутреннее устройство объекта изменилось. Поэтому **официальная рекомендация**: использовать `save_model` для долговременного хранения, а pickle/joblib – только для короткоживущих моделей или если у вас зафиксирована версия окружения

64 .

**3. Экспорт в формат PMML, ONNX и др.:** Есть проекты, позволяющие конвертировать модель XGBoost в межплатформенные форматы: например, ONNX (Open Neural Network Exchange) или PMML (Predictive Model Markup Language). ONNX поддерживает бустинговые модели, и существуют конверторы (например, `onnxmltools`), которые могут перевести XGBoost-модель в onnx-файл для использования в разных фреймворках или на мобильных устройствах. Этот путь однако извилистый и не всегда 100% поддерживает все фишки, но опция есть.

**Загрузка модели и прогнозы:** После `load_model` вы получаете модель, готовую делать `predict`. Если вы сохраняли XGBRegressor, можно загружать через `XGBRegressor.load_model`; если вы использовали низкоуровневый Booster (например, получили через `model.get_booster()`), то загружать надо в Booster. Пример:

```
bst = model.get_booster()
bst.save_model("xgb_model.json")

# позже
bst2 = xgb.Booster()
bst2.load_model("xgb_model.json")
y_pred = bst2.predict(xgb.DMatrix(X_new))
```

Обратите внимание, `Booster.predict` требует `DMatrix`. Если хотим продолжить пользоваться `scikit-обёрткой`:

```
xgb_reg = xgb.XGBRegressor()
xgb_reg.load_model("xgb_model.json")
y_pred = xgb_reg.predict(X_new)
```

Это более удобно.

**Версионирование моделей:** Если вы проводите эксперименты, можно сохранять несколько версий модели с разными параметрами (например, `model_v1.json`, `model_v2.json`...). Платформы типа MLflow или neptune.ai могут автоматически сохранять артефакты модели по завершении тренировки. Neptune, кстати, отмечает: “Neptune автоматически сохраняет текущую версию модели по окончании тренировки, так что вы всегда можете загрузить предыдущие версии” 68 .

**Применение сохранённой модели на новом данных:** Убедитесь, что на этапе инференса вы применяете те же преобразования к данным, что и на этапе обучения. Например, если вы нормировали признаки или делали encoding – при загрузке модели это не «помнится», поэтому нужно воспроизвести преобразования. Хорошей практикой будет оборачивать модель и



препроцессинг в единый Pipeline и сохранять/загружать весь Pipeline (при условии, что препроцессинг сериализуем). Либо хранить параметры масштабирования отдельно.

**Контроль совместимости:** Как было упомянуто, лучше всего сохранять в `.json`. XGBoost гарантирует, что модель, сохранённая в JSON, будет читаться более новыми версиями библиотеки (в разумных пределах). Однако обратное не всегда верно: новая модель может не загрузиться старой библиотекой, если добавились новые фичи. Поэтому планируйте совместное обновление кода и моделей.

**Размер модели:** Файлы XGBoost моделей могут быть довольно большими, если много деревьев. Можно уменьшить размер, сохранив в UBJ (это бинарный аналог JSON) – достаточно указать расширение `.ubj` в `save_model`:

```
model.save_model("model.ubj")
```

UBJ экономит место и быстрее грузится (бинарник), при тех же гарантиях совместимости.

**Пример сохранения/загрузки:**

```
model = XGBRegressor(n_estimators=100)
model.fit(X_train, y_train)
# Сохраняем модель
model.save_model("final_model.json")

# ... позже или в другом скрипте/окружении ...
model_loaded = XGBRegressor()
model_loaded.load_model("final_model.json")
predictions = model_loaded.predict(X_test)
```

Как показывает практика: *"It's officially recommended to use the `save_model()` and `load_model()` functions to save and load models."*<sup>67</sup> Причём, *"`dump_model()` используется для интерпретации/визуализации, а не для сохранения обученного состояния"*<sup>66</sup>. Если соблюдаете эти рекомендации, проблем с переносом модели не возникнет.

После загрузки стоит сделать небольшую проверку: например, сравнить несколько предсказаний `model.predict` до сохранения и после загрузки на одних и тех же данных – они должны совпадать точно. Также, если у вас в данных были категориальные признаки, убедитесь, что порядок/имена столбцов совпадают, иначе могут быть смещения.

**Повторное использование модели:** Загруженная модель работает так же, как свеженатренированная, но имейте в виду: если вы планируете **дообучивать (continue training)** модель, XGBoost это позволяет через передачу сохранённого бустера в параметр `xgb_model` метода `fit`. Например:

```
model = XGBRegressor(n_estimators=50)
model.fit(X_train, y_train)
model.save_model("model50.json")
```

```
# позже, хотим добавить ещё 50 деревьев
model2 = XGBRegressor(n_estimators=100)
model2.fit(X_train, y_train, xgb_model="model150.json")
```

Здесь model2 начнёт обучение с уже имеющимися 50 деревьями и построит ещё 50 (итого станет 100). Это удобно, если вы недообучили модель и хотите продолжить с того же места, не начиная заново. Но будьте осторожны: если данные те же, это эквивалентно просто больше итераций. Если данные новые (скажем, дозагрузка), то дообучение модели на новых данных – более сложная тема (например, можно сделать доп. boosting на новых данных, но лучше обычно объединить датасеты и переобучить заново с early stopping).

**Вывод:** сохранение и загрузка моделей XGBoost – довольно прямой процесс. Следуя рекомендованным практикам (использование `save_model` / `load_model`), вы можете легко перенести обученную модель между средами и языками, а также хранить её в файловом хранилище для последующего использования без необходимости повторно обучать. Это заключительный шаг в нашем руководстве – теперь у вас есть готовая, настроенная модель XGBoost, которую вы можете применять к новым данным, отслеживать её работу и при необходимости обновлять. Успехов в ваших задачах регрессии с XGBoost! 69 70

---

1 2 3 What is the XGBoost algorithm and how does it work?

<https://www.analyticsvidhya.com/blog/2018/09/an-end-to-end-guide-to-understand-the-math-behind-xgboost/>

4 5 Introduction to Boosted Trees — xgboost 3.0.2 documentation

<https://xgboost.readthedocs.io/en/stable/tutorials/model.html>

6 7 8 9 10 11 XGBoost - GeeksforGeeks

<https://www.geeksforgeeks.org/machine-learning/xgboost/>

12 13 14 16 Installation Guide — xgboost 3.0.2 documentation

<https://xgboost.readthedocs.io/en/stable/install.html>

15 41 XGBoost GPU Support — xgboost 3.1.0-dev documentation

<https://xgboost.readthedocs.io/en/latest/gpu/>

17 18 decision trees - Is it necessary to normalize data for XGBoost? - Data Science Stack Exchange

<https://datascience.stackexchange.com/questions/60950/is-it-necessary-to-normalize-data-for-xgboost>

19 23 51 52 53 54 57 58 59 60 68 XGBoost: Everything You Need to Know

<https://neptune.ai/blog/xgboost-everything-you-need-to-know>

20 42 43 How to Use XGBoost XGBRegressor | XGBoosting

<https://xgboosting.com/how-to-use-xgboost-xgbregressor/>

21 22 24 25 26 28 30 32 33 34 35 36 37 38 39 40 45 XGBoost Parameters — xgboost 3.0.2 documentation

<https://xgboost.readthedocs.io/en/stable/parameter.html>

27 29 31 46 47 48 49 50 XGBoost Parameters Tuning: A Complete Guide with Python Codes

<https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/>

44 Avoid Overfitting By Early Stopping With XGBoost In Python

<https://www.machinelearningmastery.com/avoid-overfitting-by-early-stopping-with-xgboost-in-python/>

55 Optimizing XGBoost: A Guide to Hyperparameter Tuning | by RITHP

<https://medium.com/@rithpansanga/optimizing-xgboost-a-guide-to-hyperparameter-tuning-77b6e48e289d>

56 XGBoost Hyperparameter Tuning With Optuna (Kaggle Grandmaster ...

<https://forecastegy.com/posts/xgboost-hyperparameter-tuning-with-optuna/>

61 62 63 65 Introduction to Model IO — xgboost 3.1.0-dev documentation

[https://xgboost.readthedocs.io/en/latest/tutorials/saving\\_model.html](https://xgboost.readthedocs.io/en/latest/tutorials/saving_model.html)

64 66 67 69 70 How to Save and Load XGBoost Models

<https://stackabuse.com/bytes/how-to-save-and-load-xgboost-models/>