

Проект на курсе ML

Состав команды:

- Белоглазов Михаил Юрьевич
- Воронин Георгий Владимирович
- Карпов Денис Денисович
- Перминов Матвей Максимович

Этап 1. Датасет и EDA

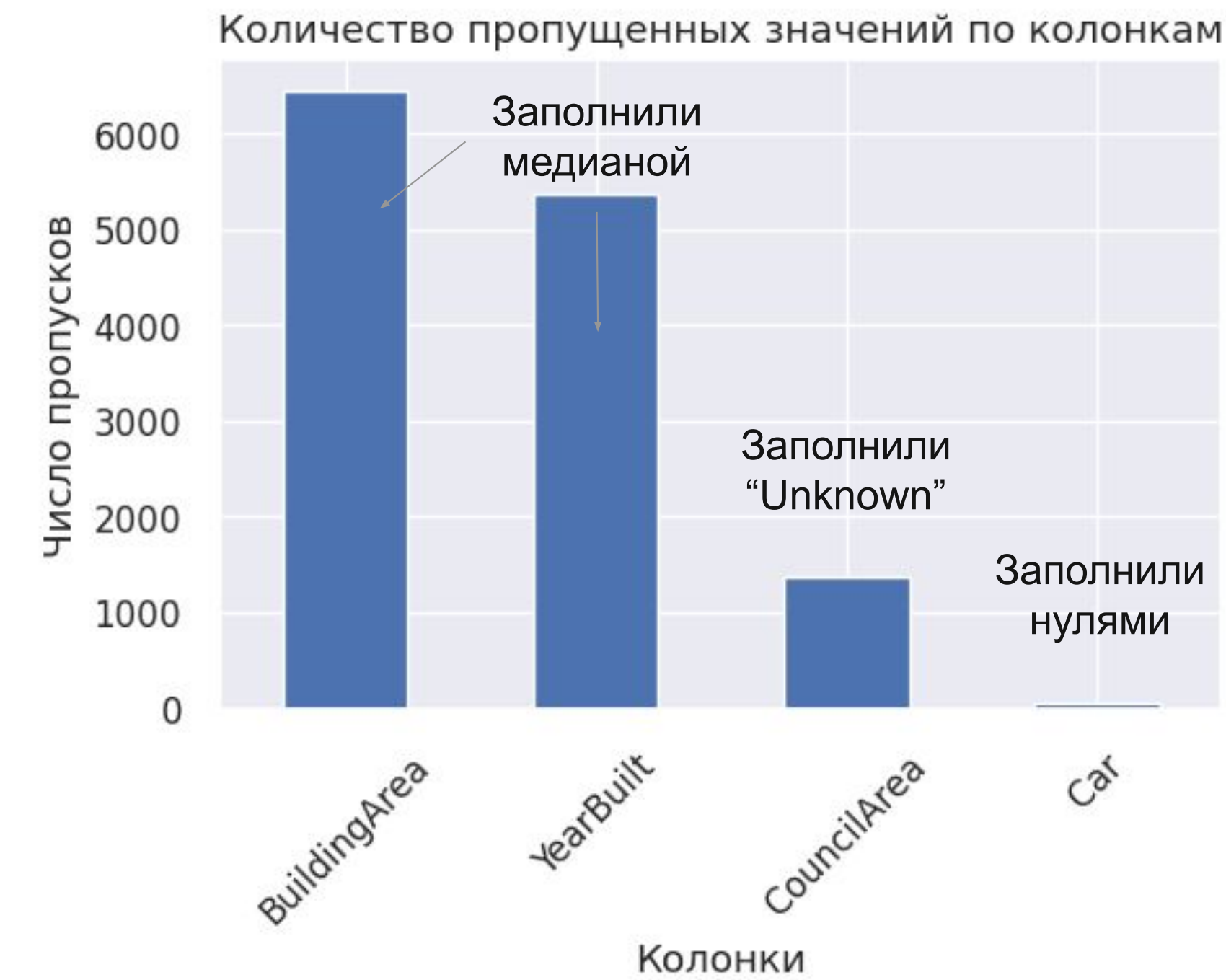
Датасет

Melbourne Housing - содержит данные о продажах разной недвижимости (квартиры, дома, таунхаусы) на аукционе в Мельбурне (Австралия)

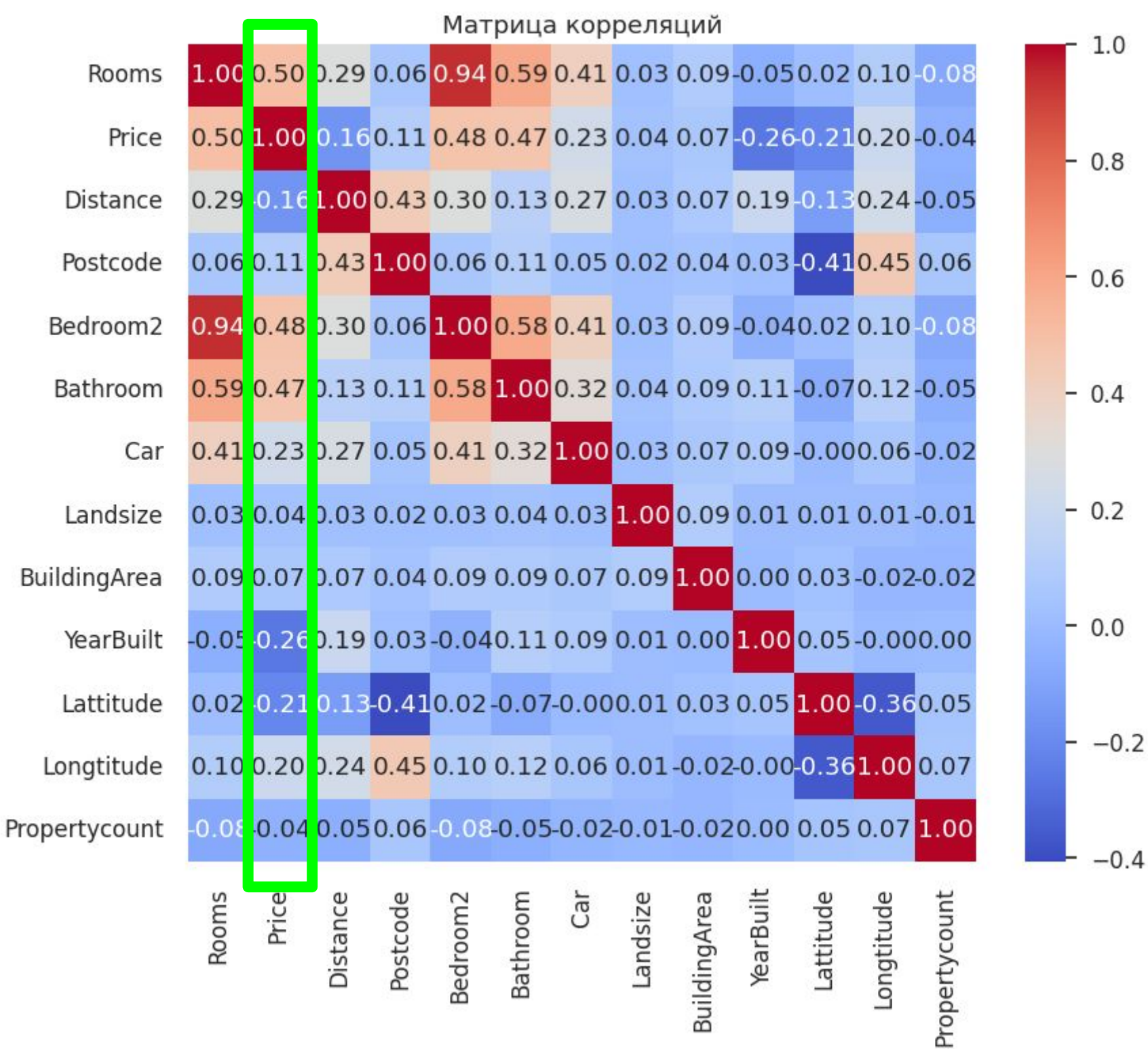
- 13 000 строк
- Целевая переменная: Price (цена, за которую продан дом)
- Есть временные и пространственные признаки



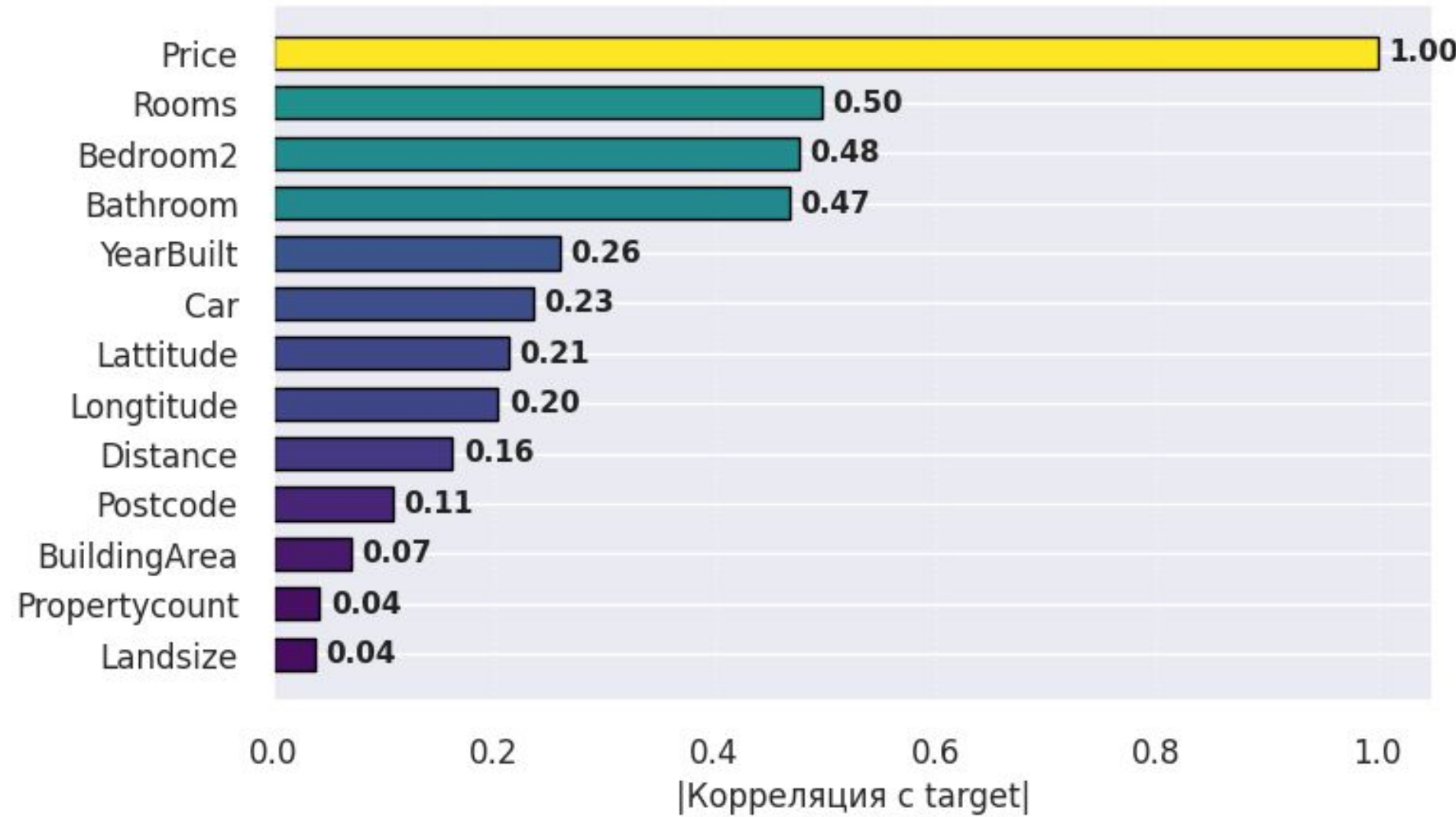
Пропущенные значения



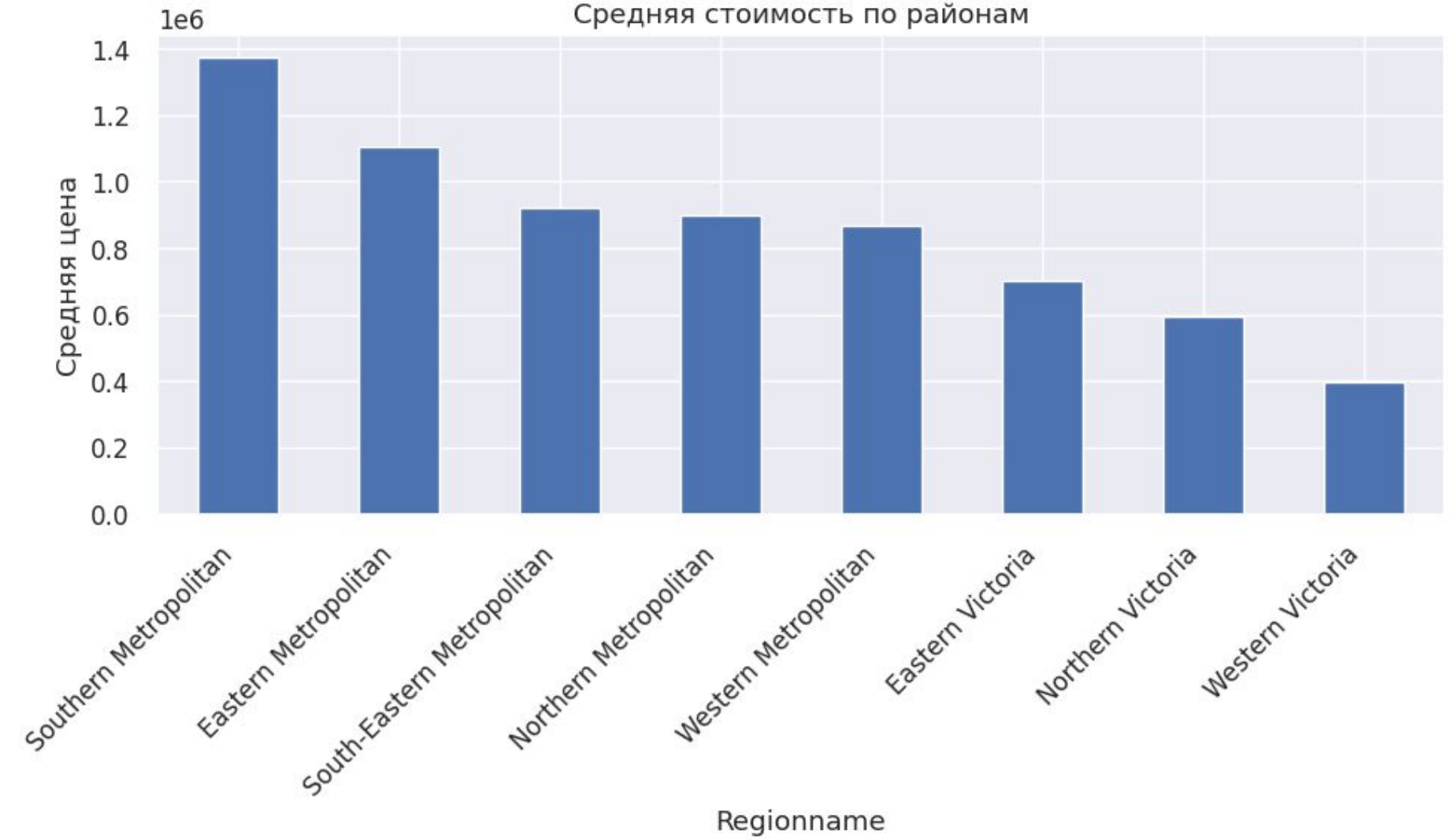
Корреляции



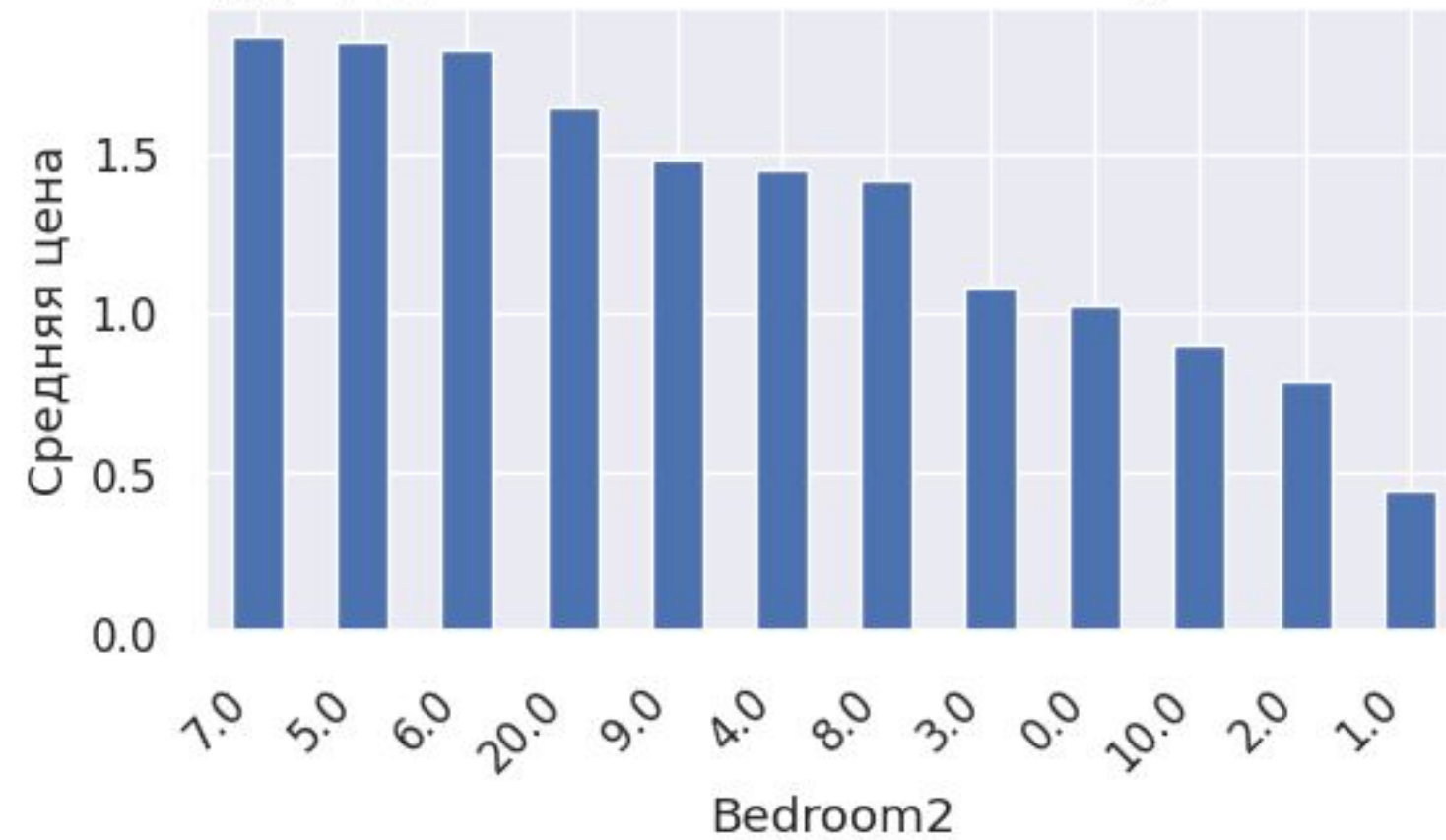
Топ-10 признаков по абсолютной корреляции с target

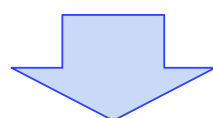


Средняя стоимость по районам



Средняя стоимость по количеству спален





Baseline

| | rmse_train | rmse_val |
|-------------------|---------------|---------------|
| Catboost_baseline | 203991.248249 | 269203.054416 |
| Linreg_baseline | 253387.388905 | 379754.136154 |



Этап 2. Работа с аномалиями и генерация признаков

Анализ таргета



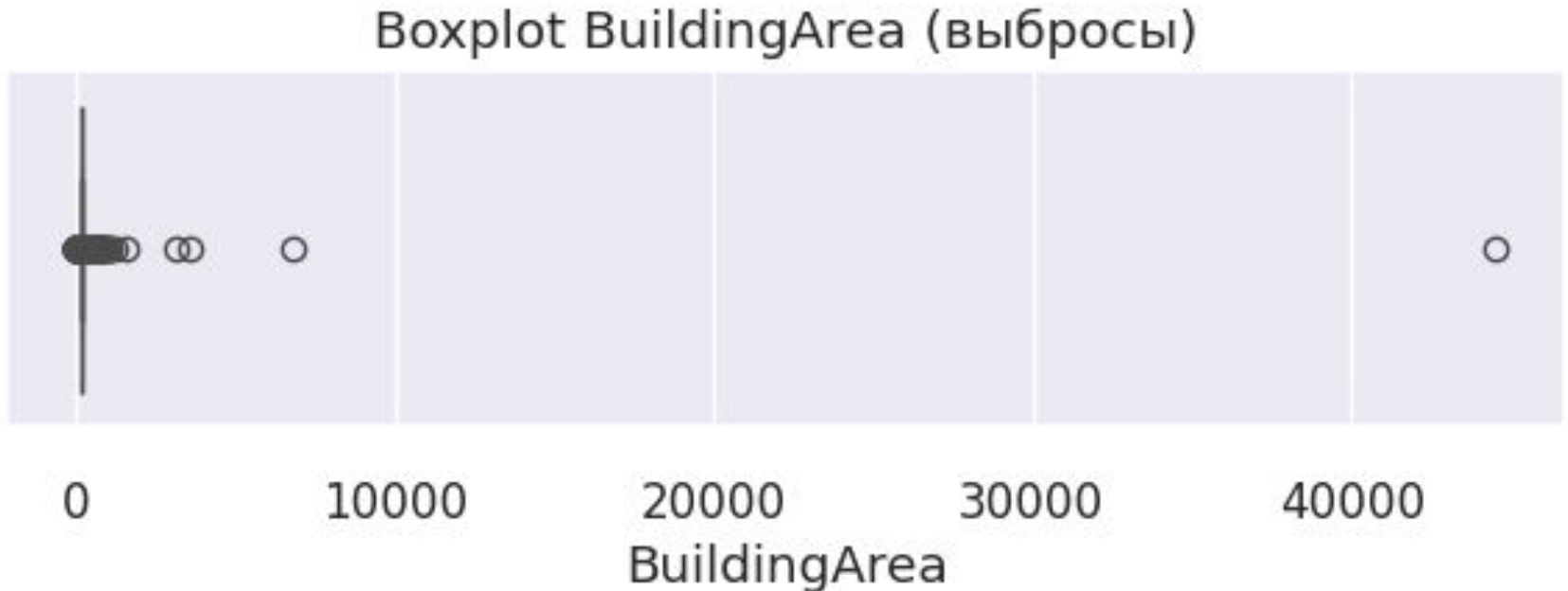
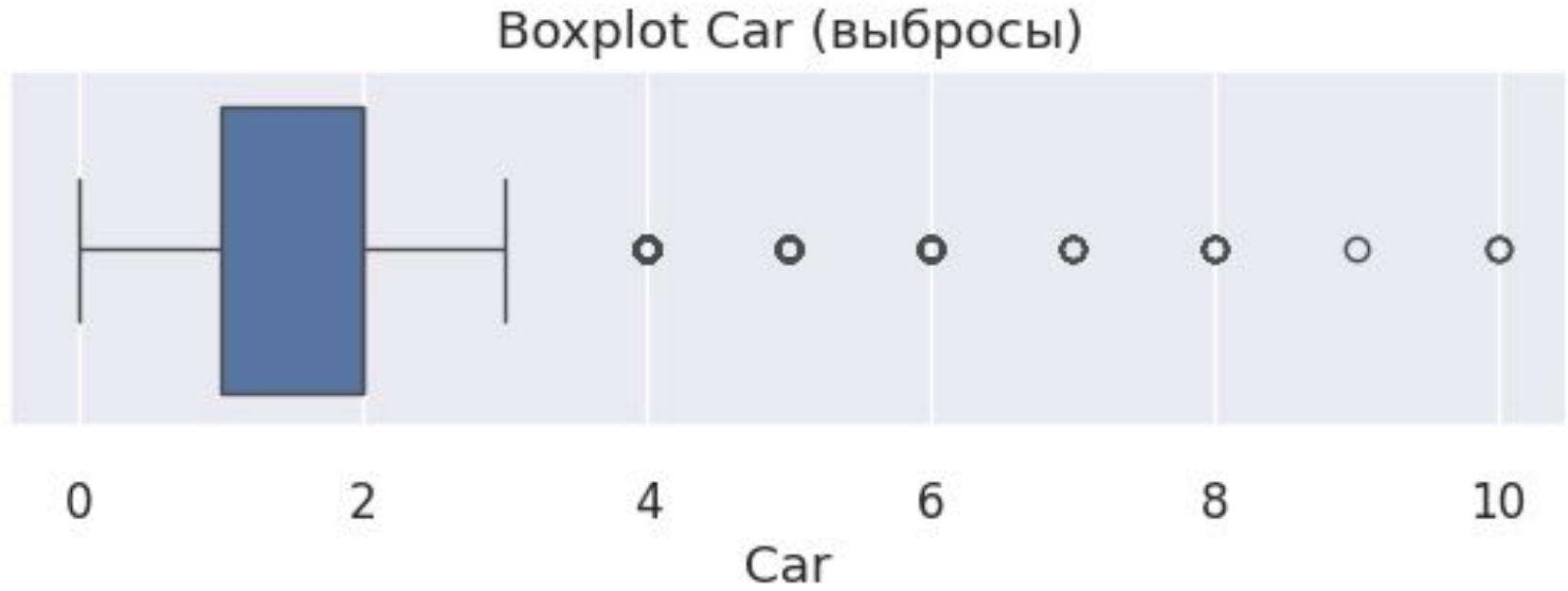
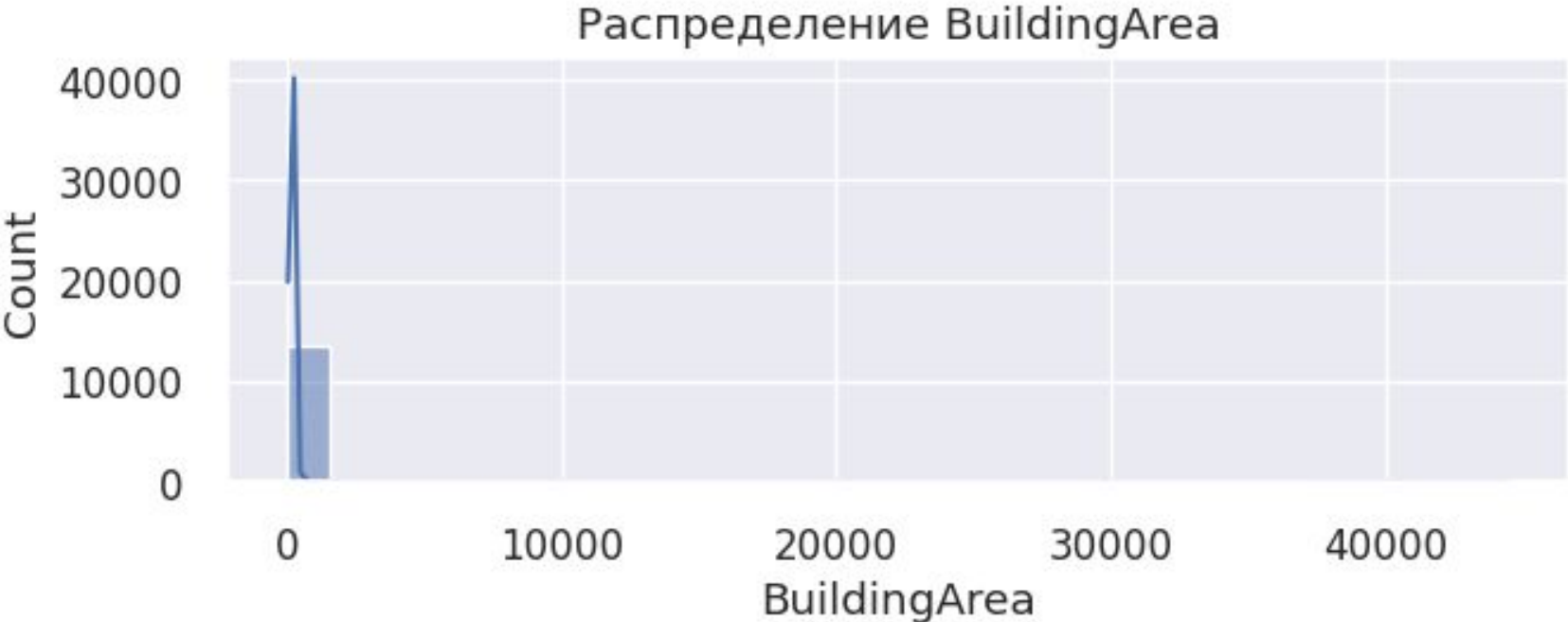
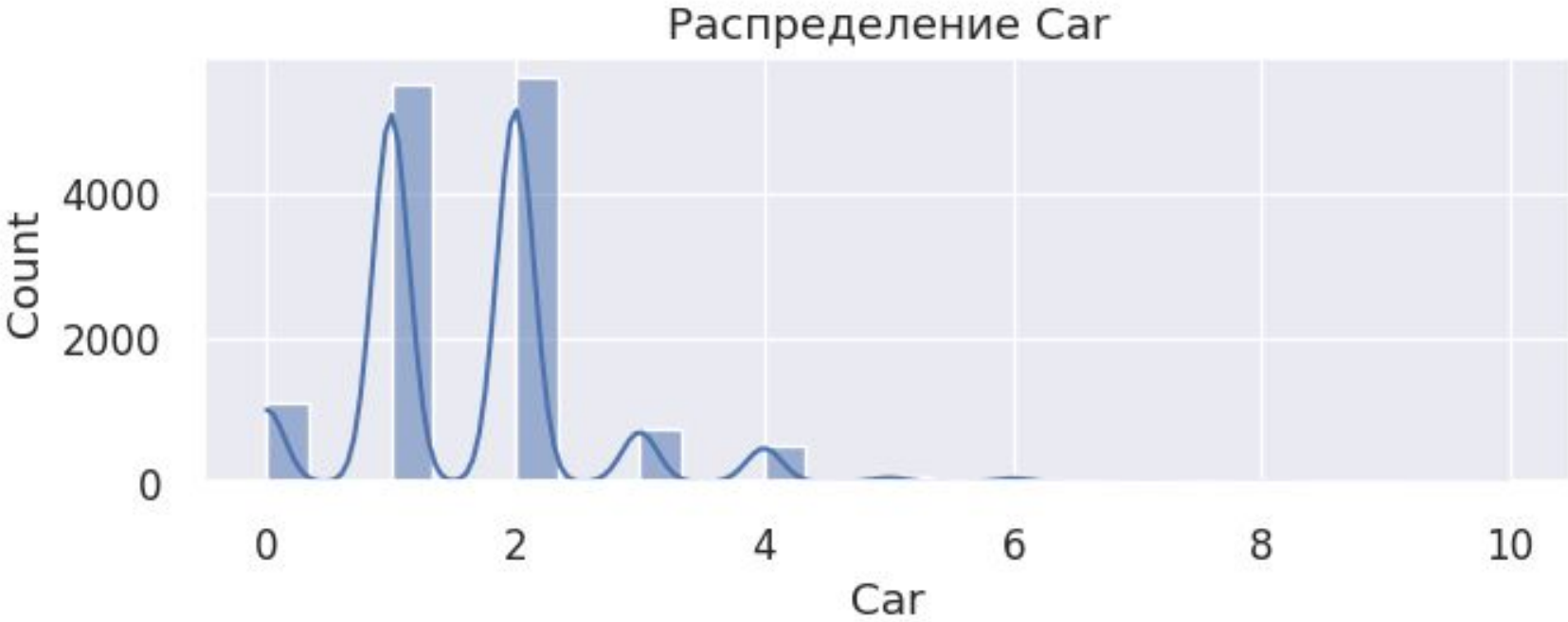
=== Выбросы по Z-score ($>|3|$) ===
Количество выбросов: 232

=== Выбросы по IQR ($1.5 * IQR$) ===
Q1: 650000.0 Q3: 1330000.0 IQR: 680000.0
Нижняя граница: -370000.0 Верхняя граница: 2350000.0
Количество выбросов: 612

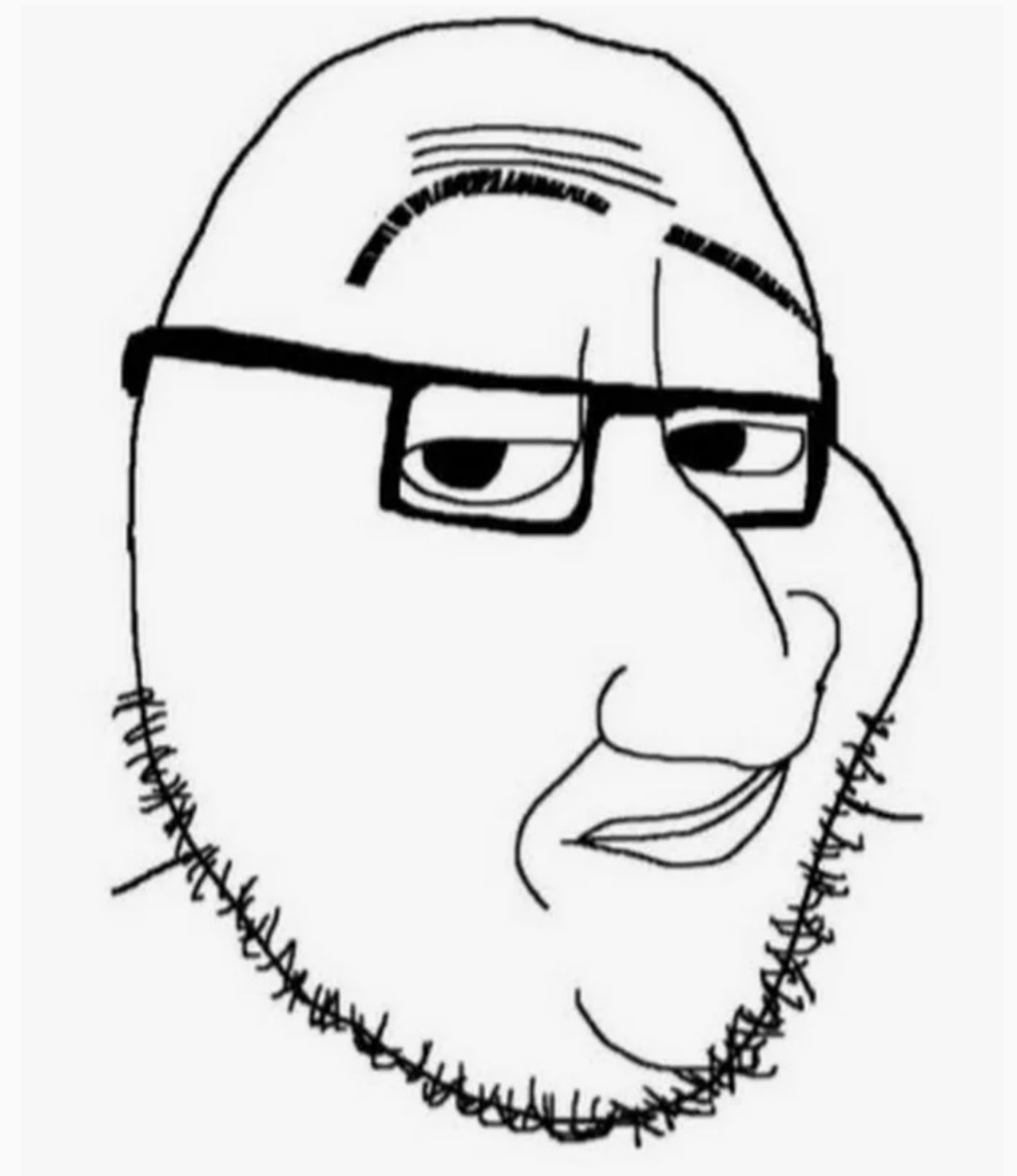
Прости, но ты живешь в США



Анализ Car и BuildingArea

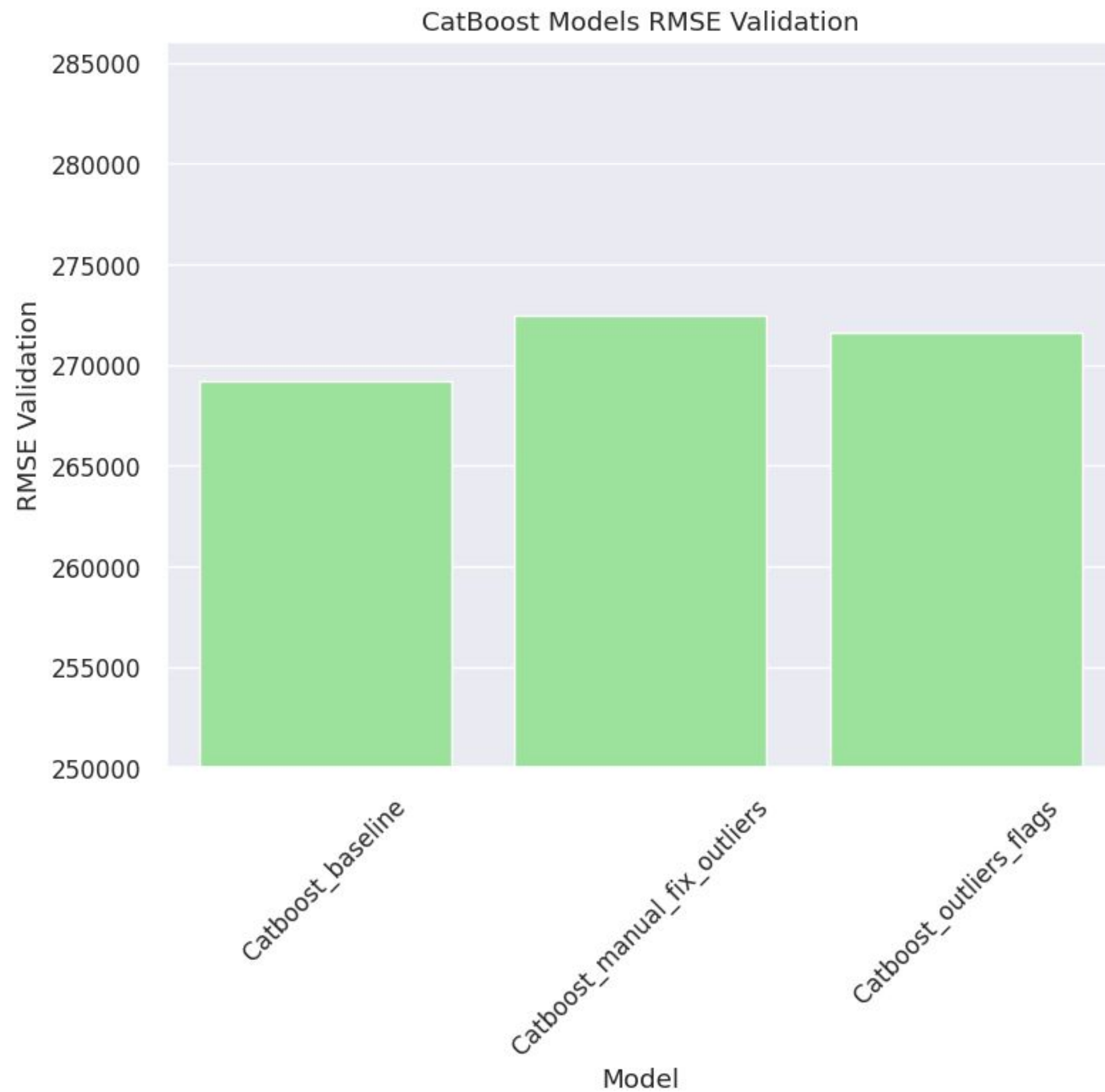


Пример выброса



Поставлю площадь
здания 44 515 м²
(около 6-7 футбольных полей)

Добавим признаки связанные с аномалиями

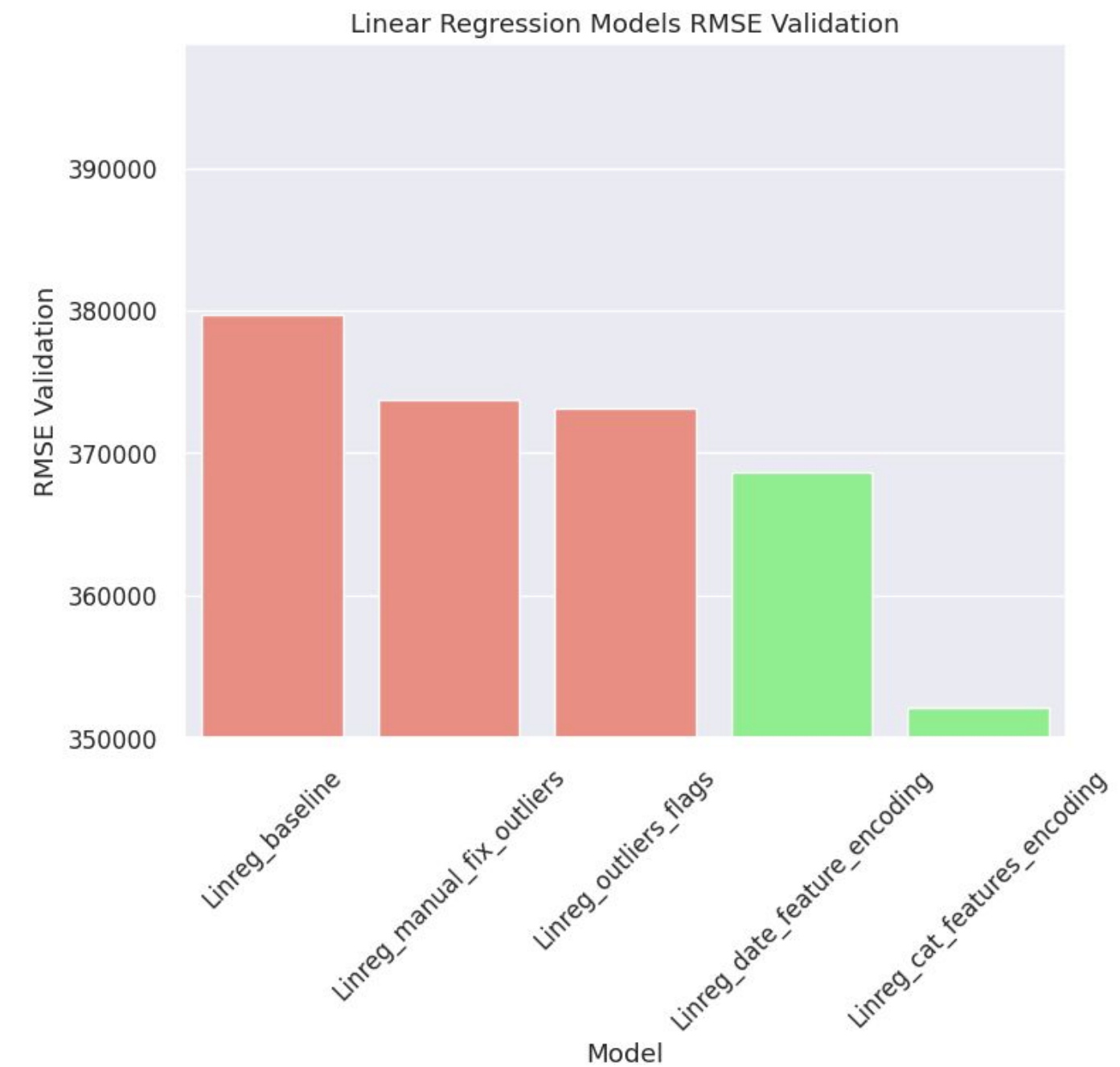
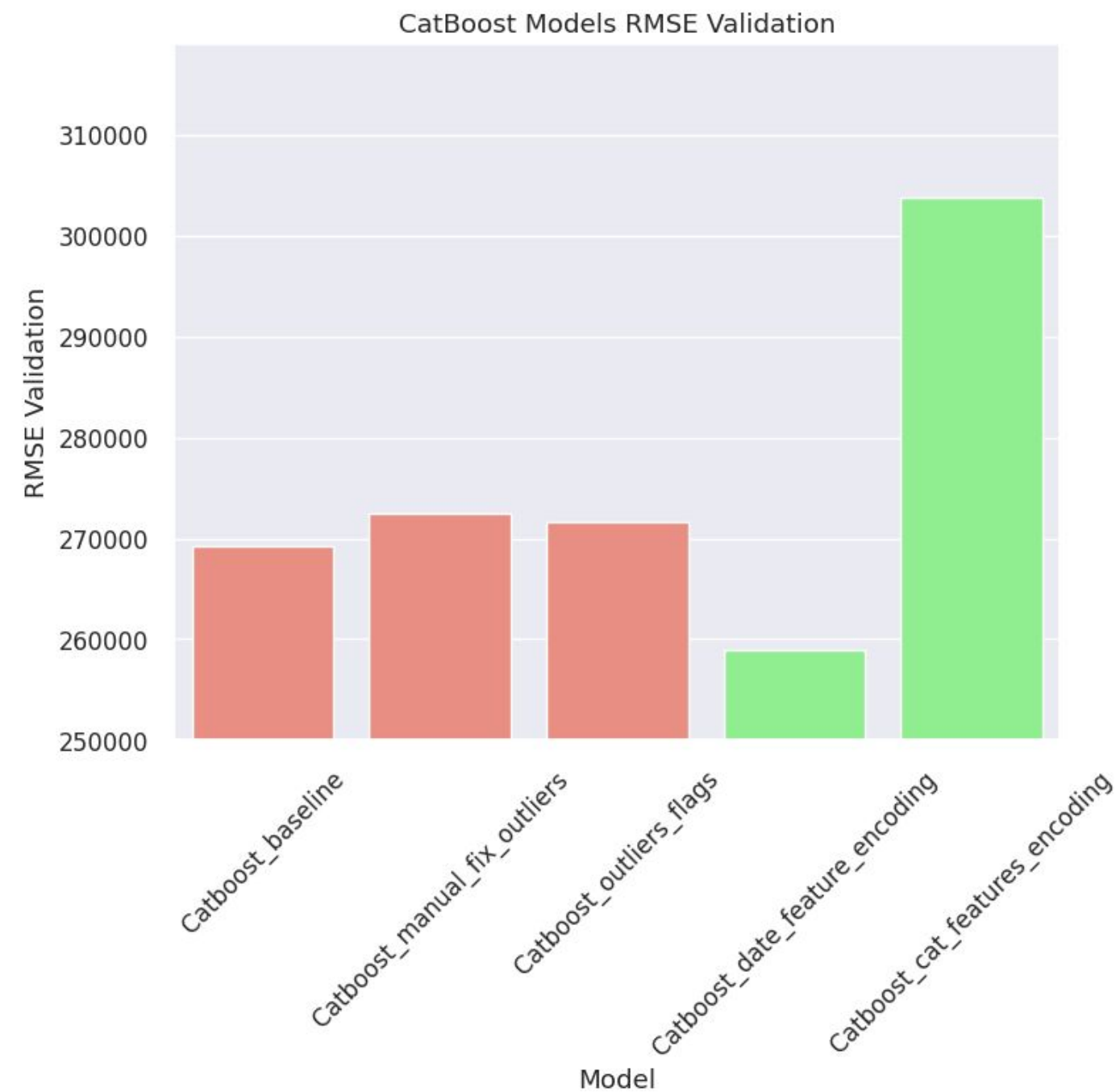


Обработка категориальных признаков



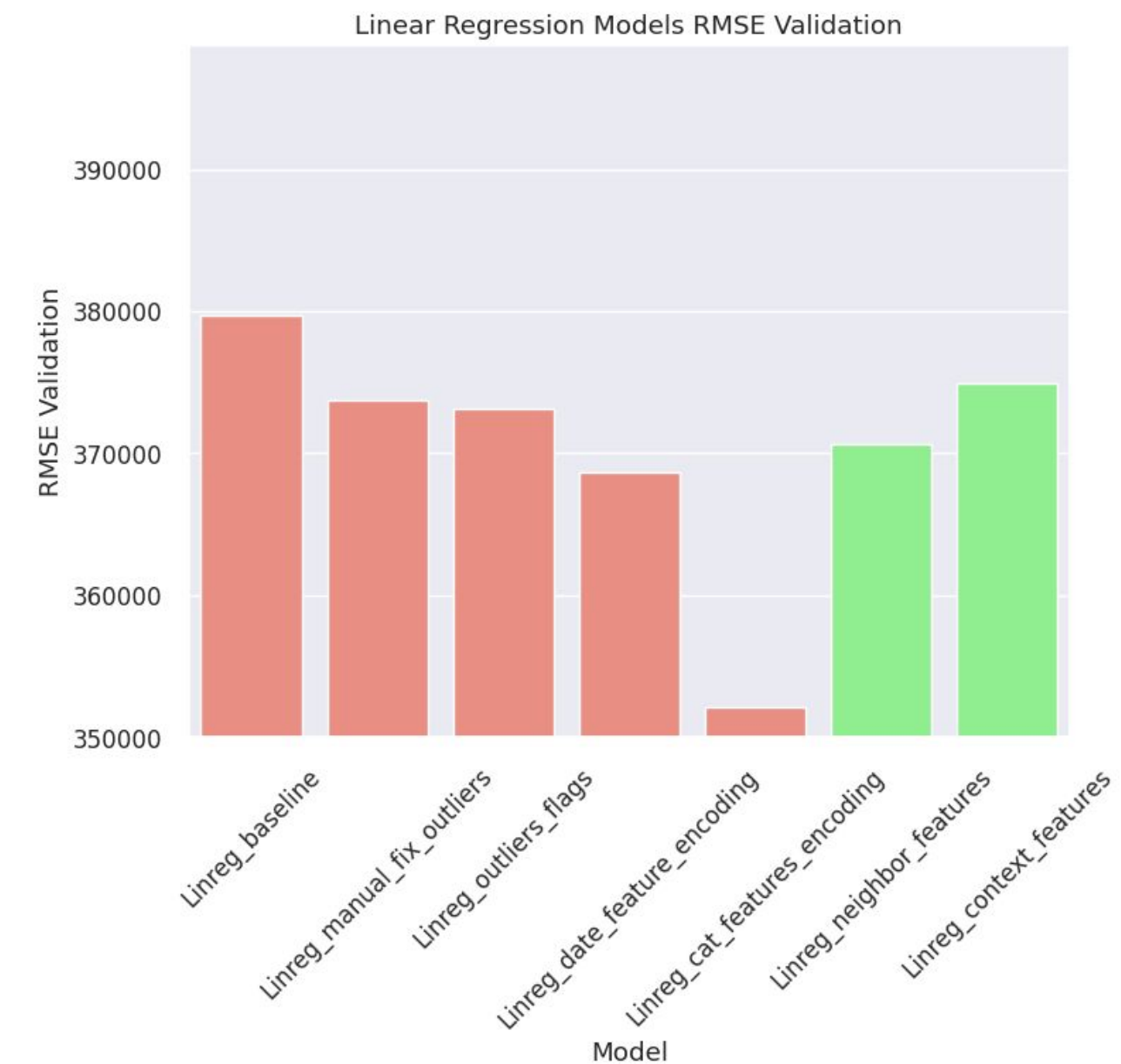
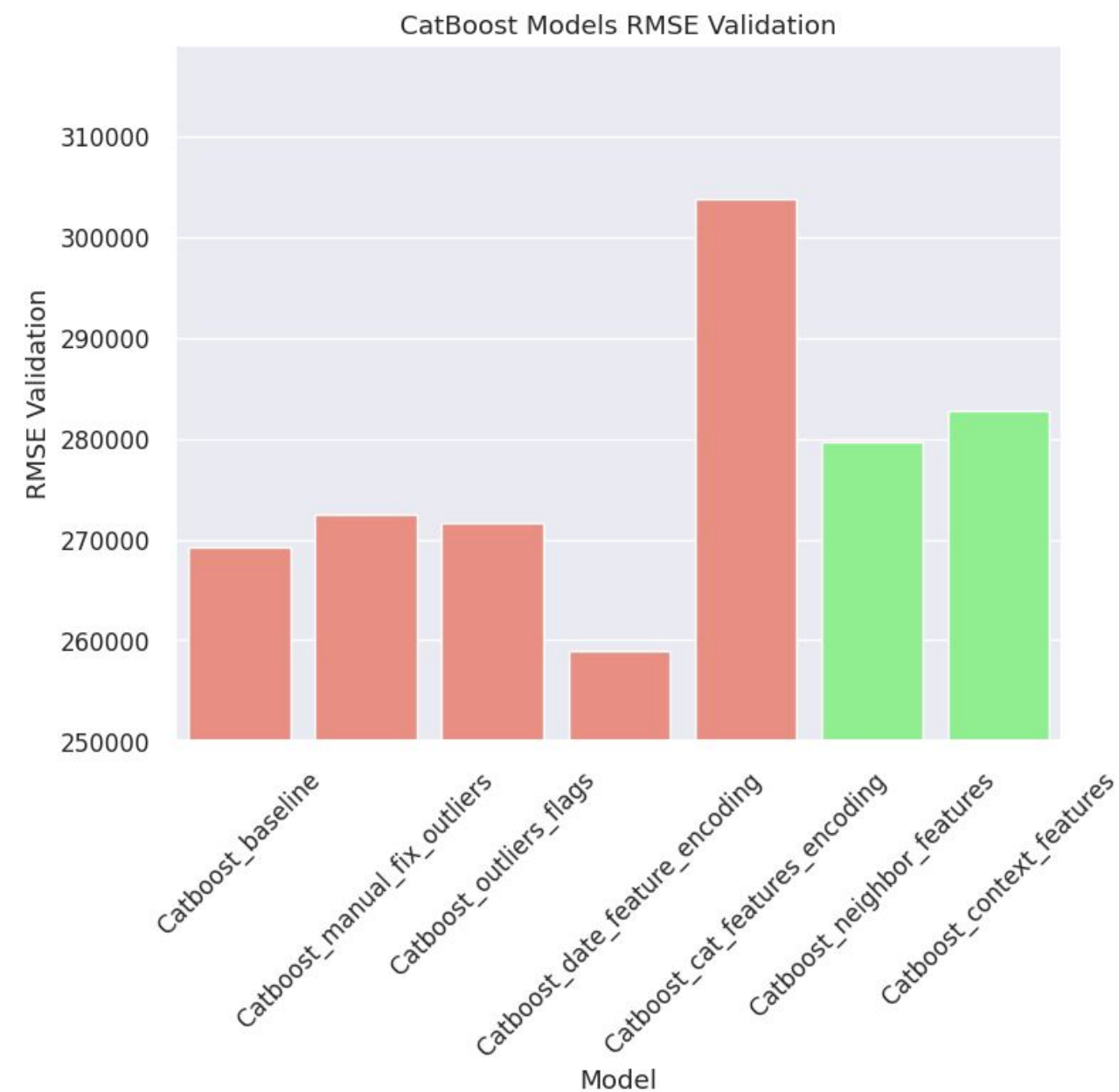
Промежуточные метрики

Обработали
временной признак
Date, разбив его на
отдельные части.



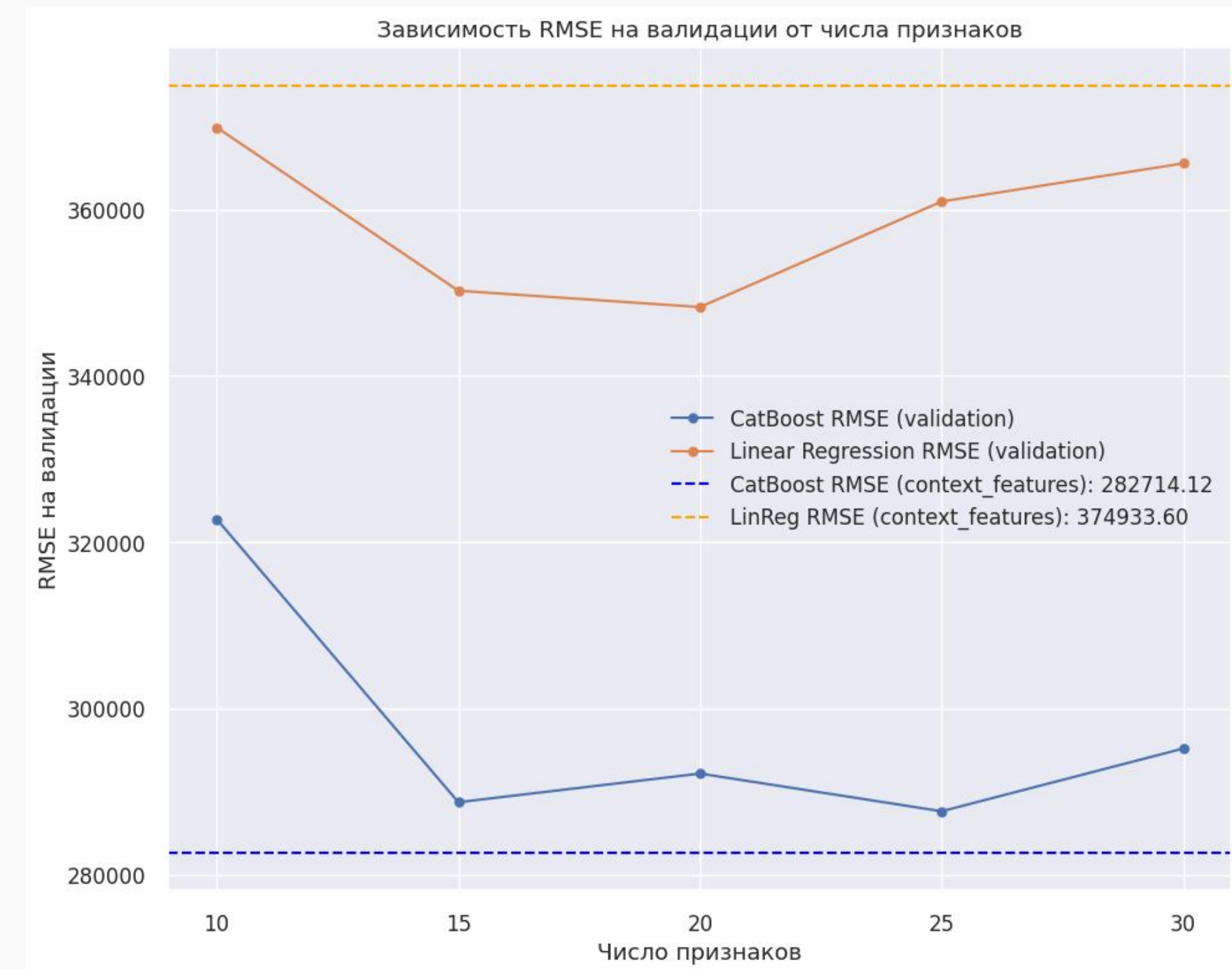
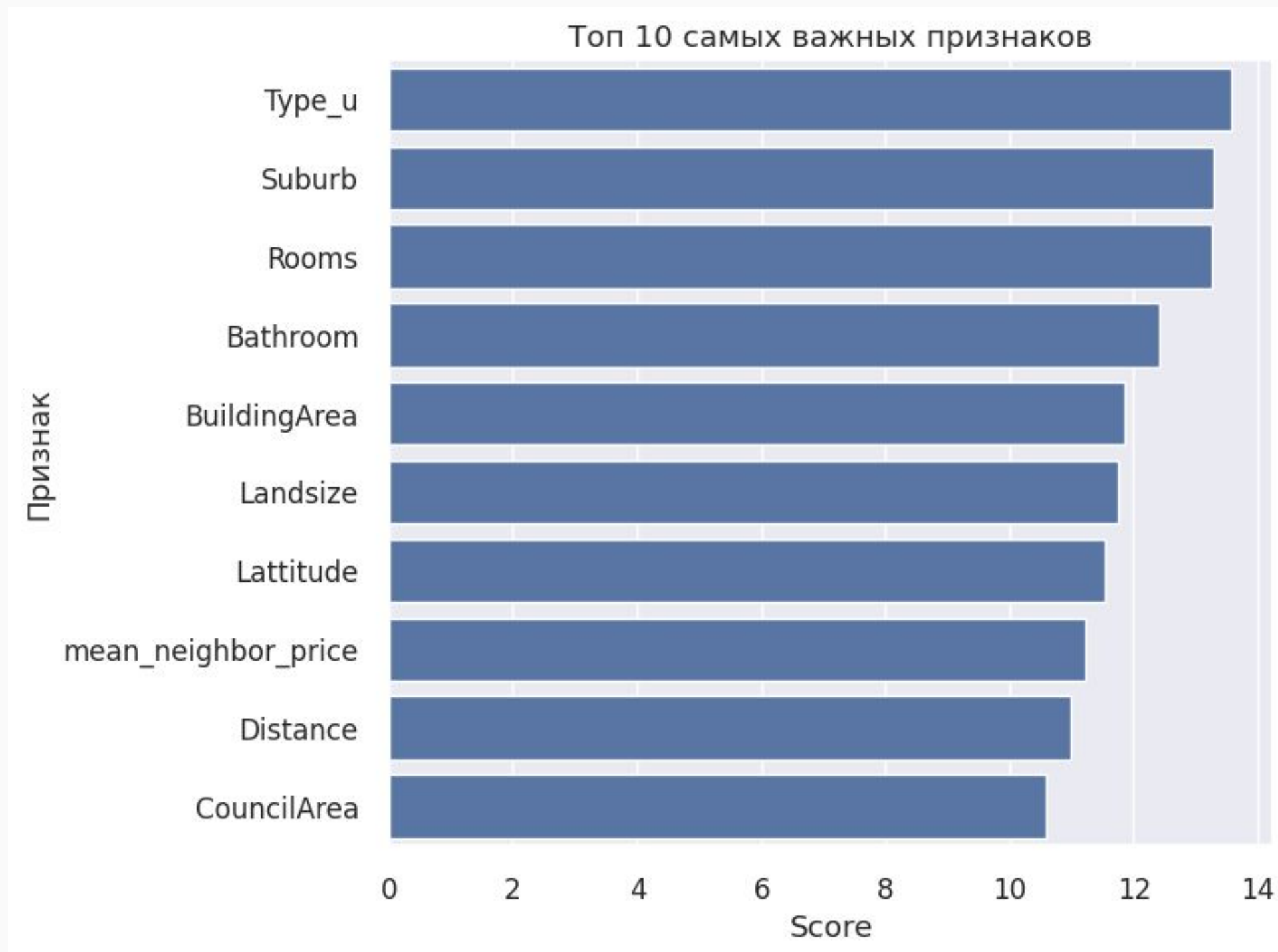
Генерация новых признаков

- Признаки, основанные на ближайших географических соседях (средняя цена, расстояние до ближайшего соседа, среднее расстояние до соседей)
- Контекстные признаки:
- Отношение ванных комнат к общему числу комнат
- Отношение жилой площади к площади участка



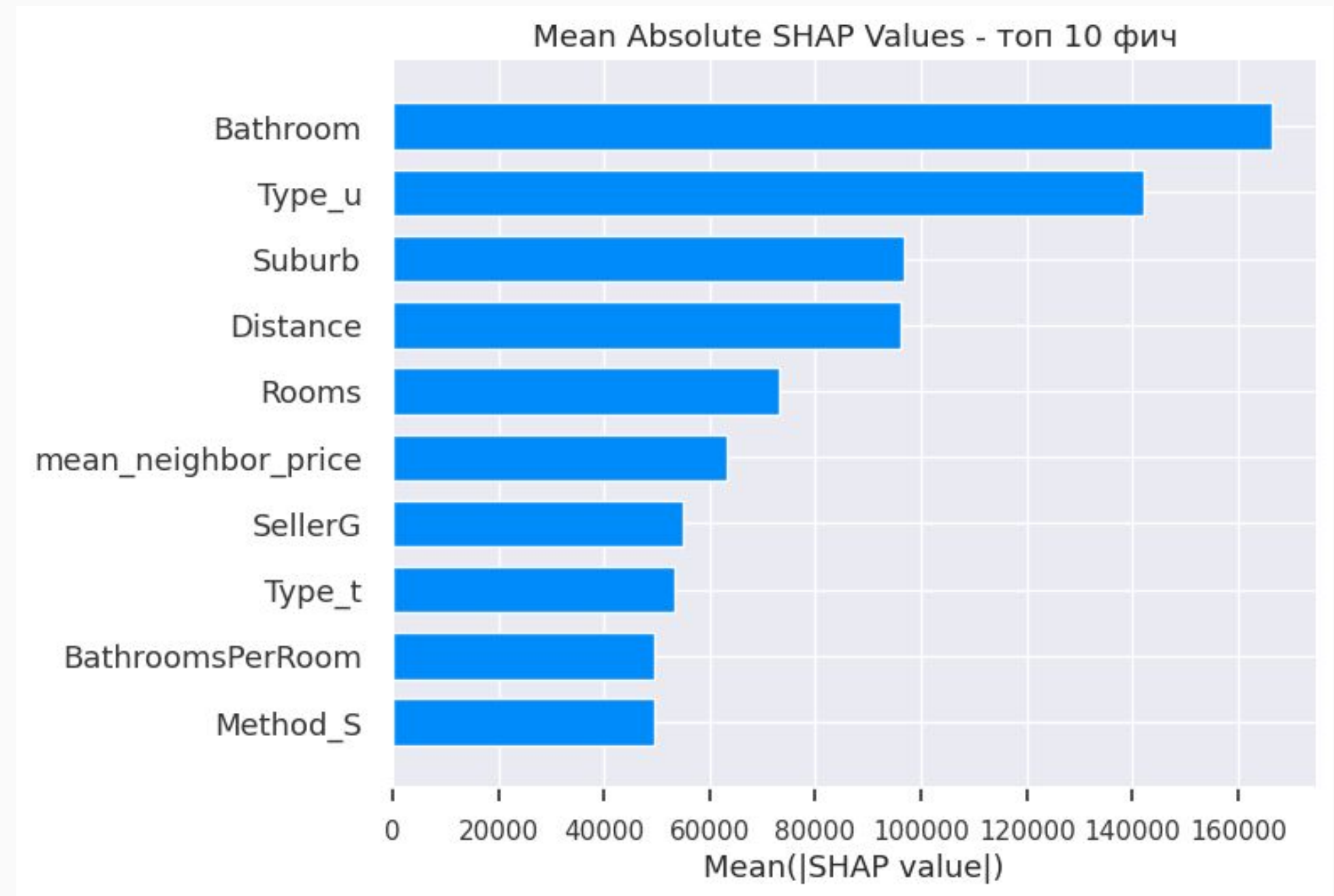
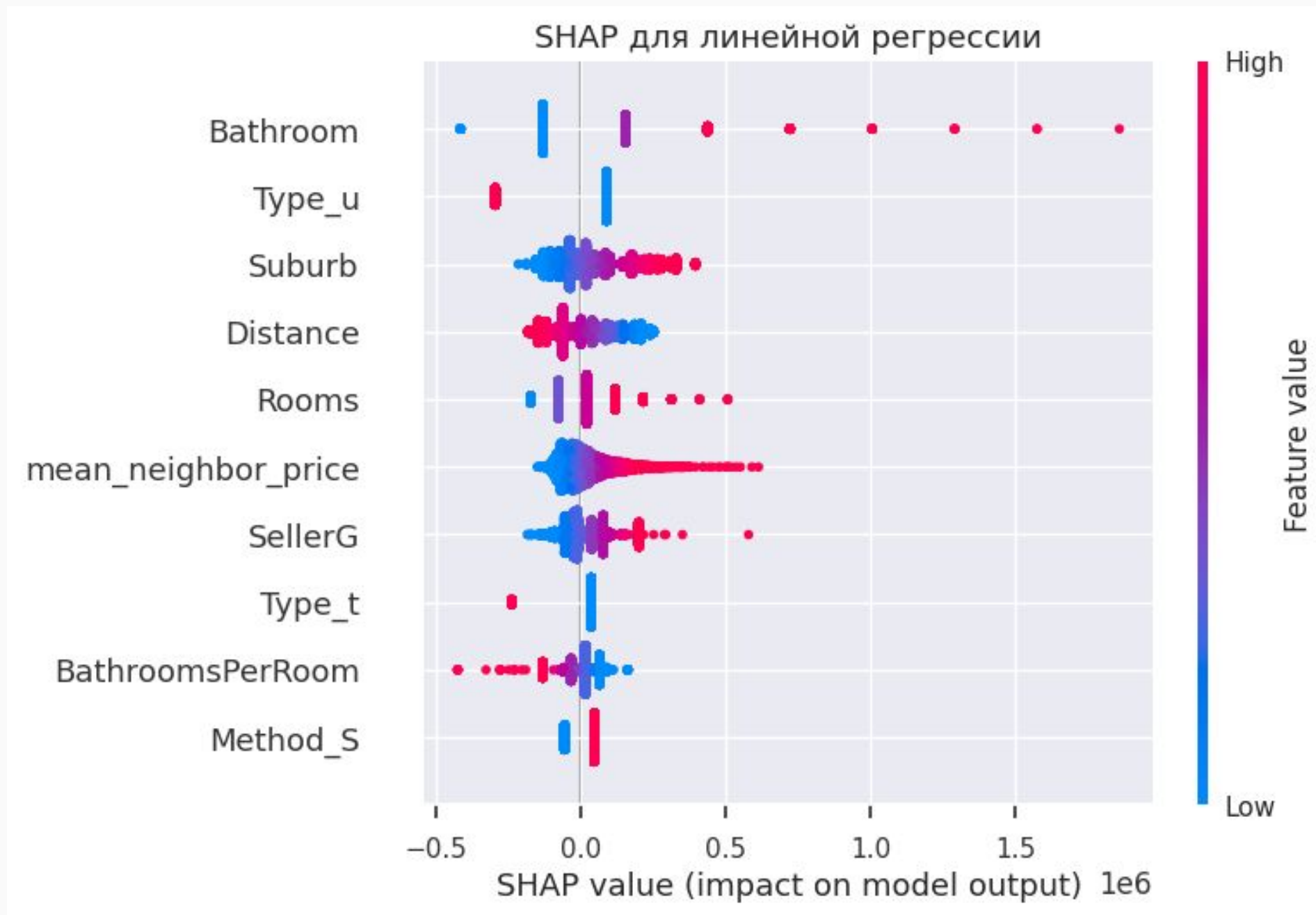
Отбор признаков

Мы использовали комплексную систему отбора признаков, которая учитывает разные способы отбора с весами

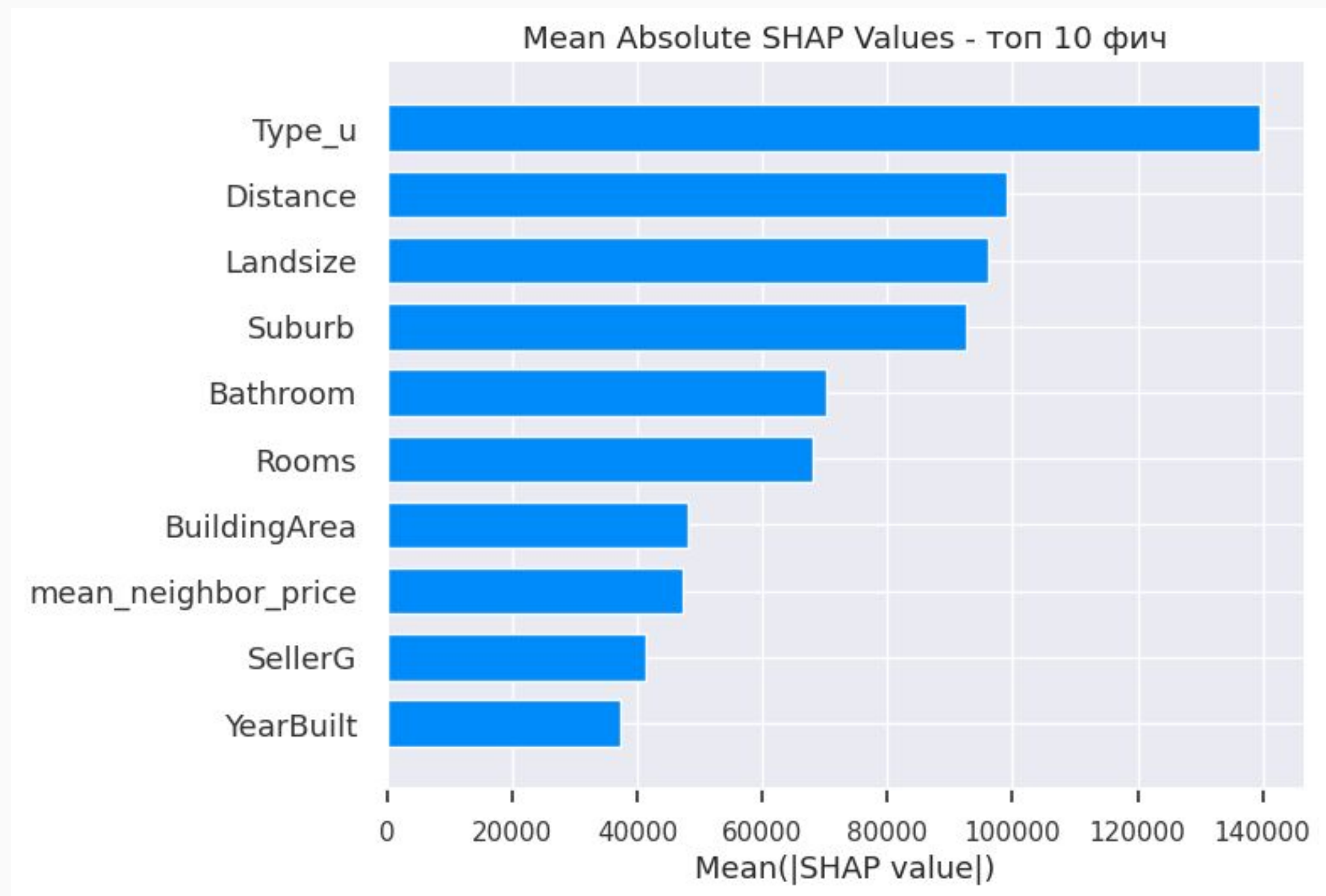
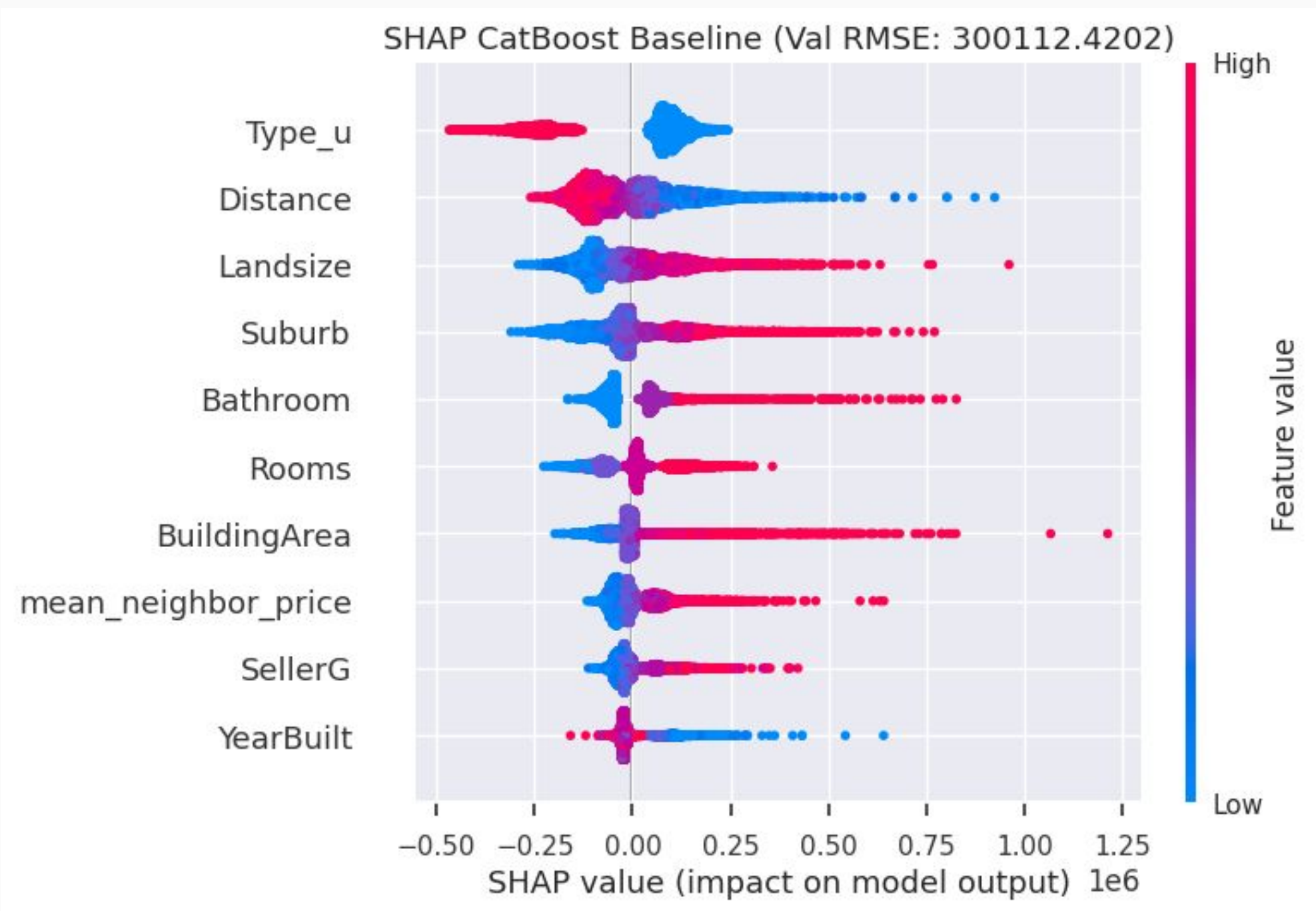


Этап 3. Интерпретация и диагностика моделей

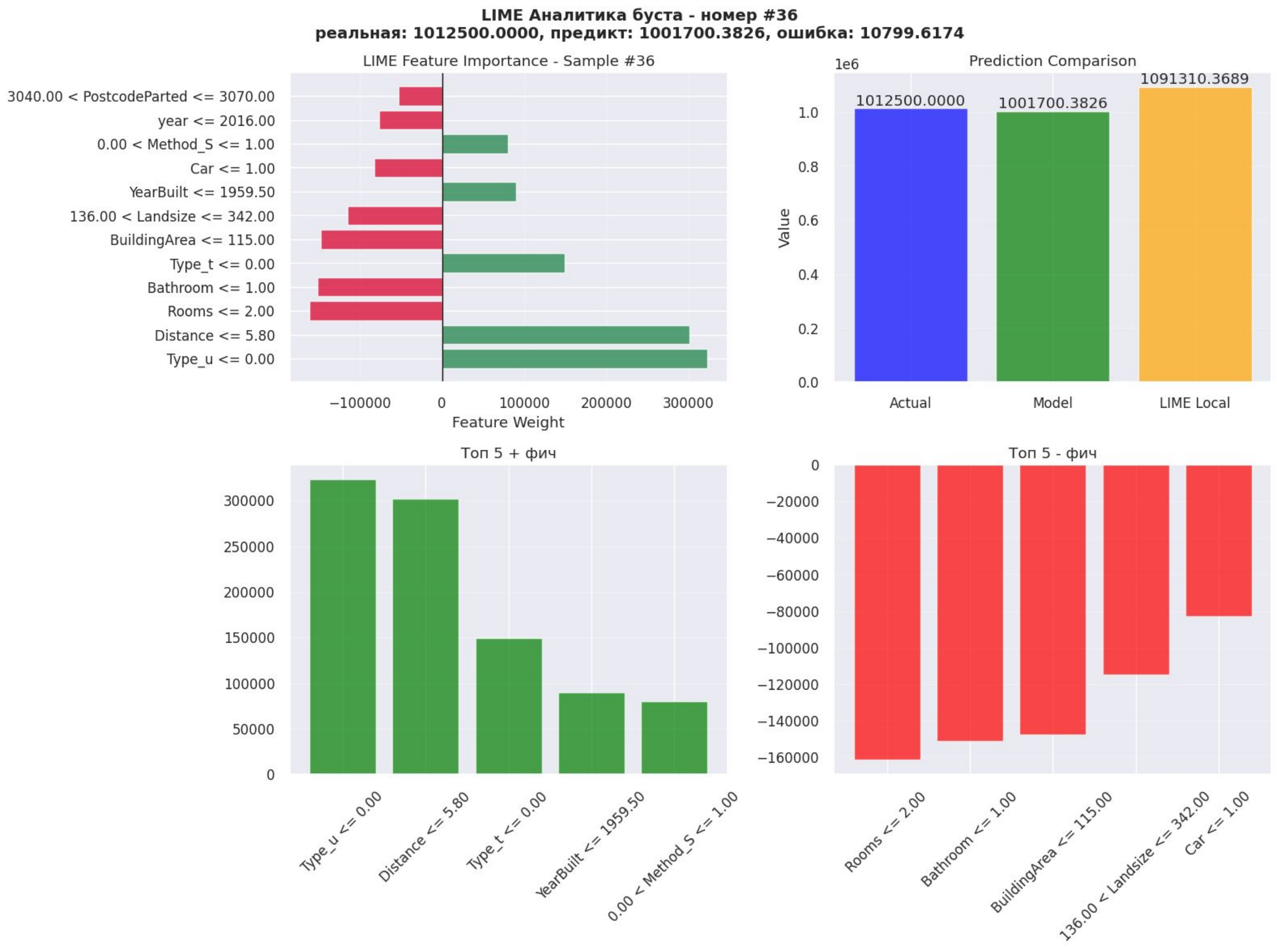
SHAP для линрега



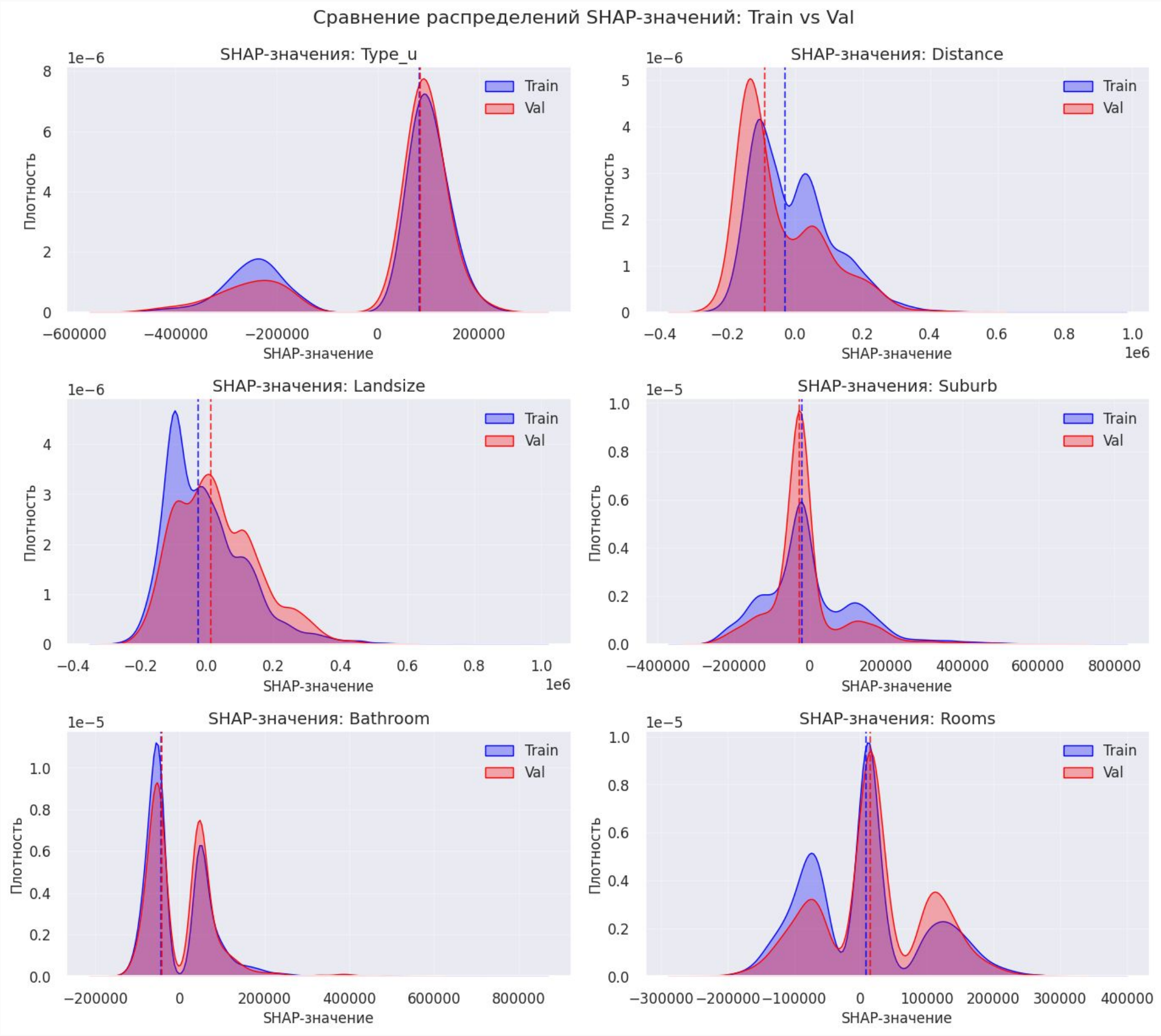
SHAP для кэтбуста



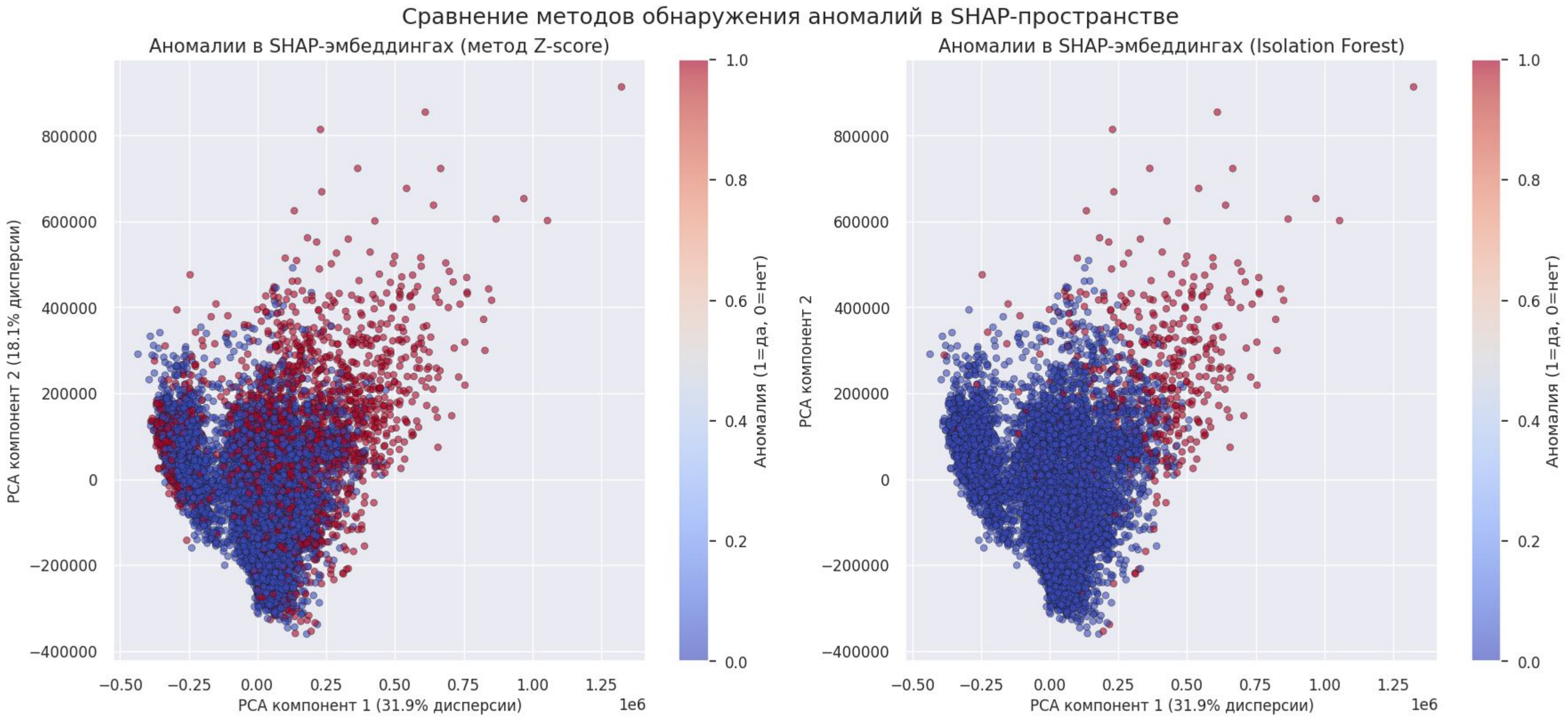
LIME



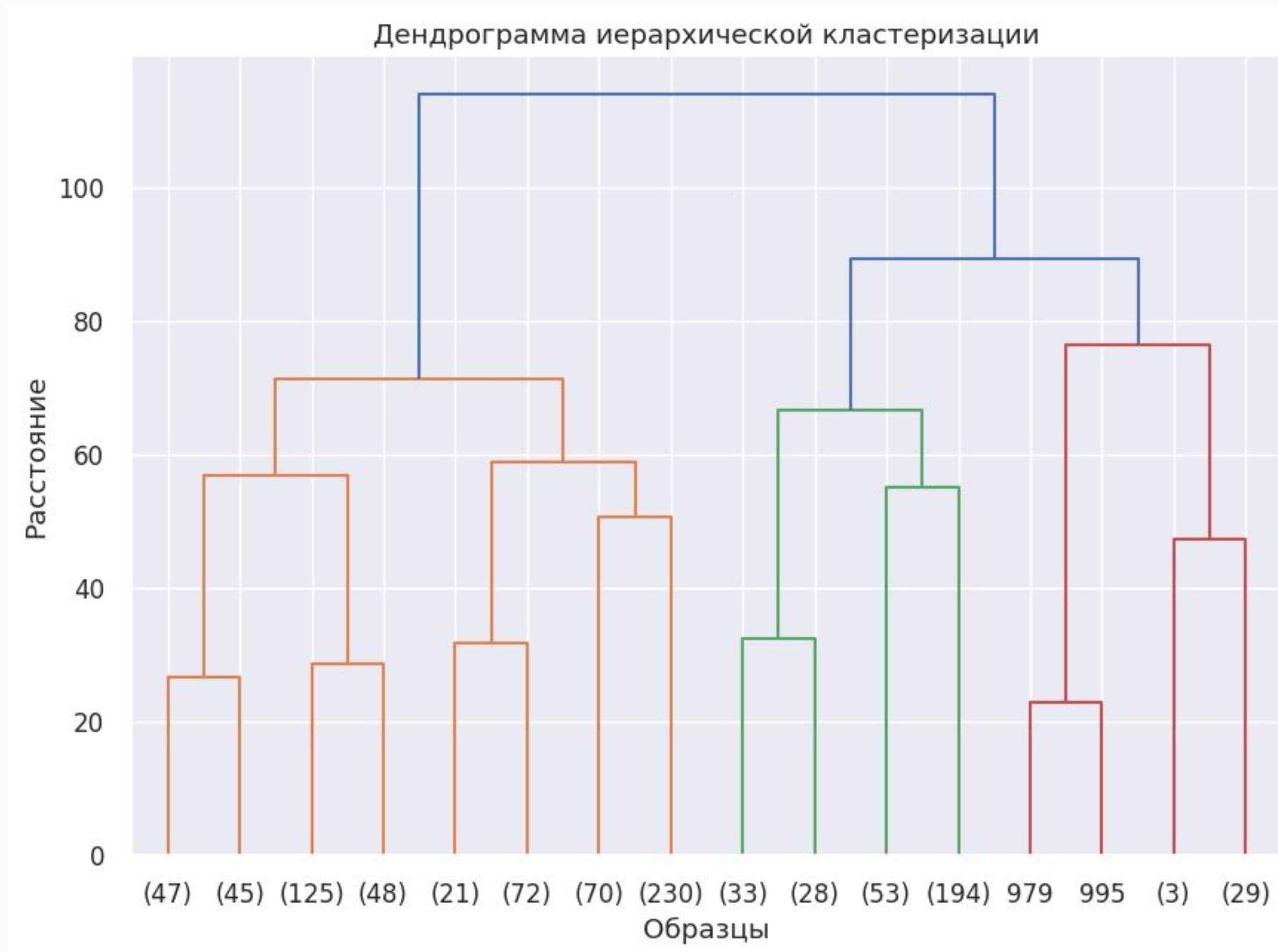
Сдвиги для SHAP значений



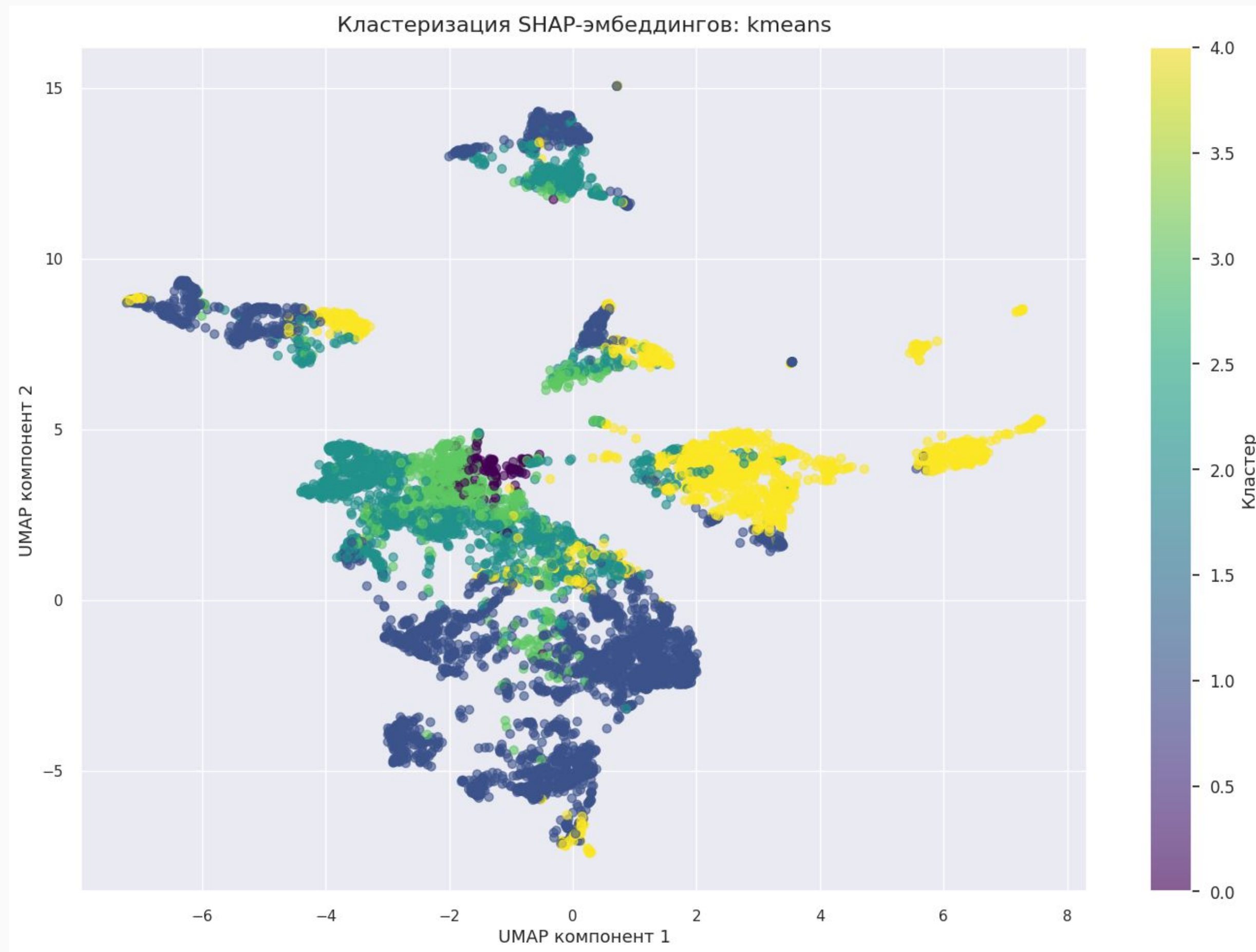
SHAP-аномалии



Дендрограмма и кластеры

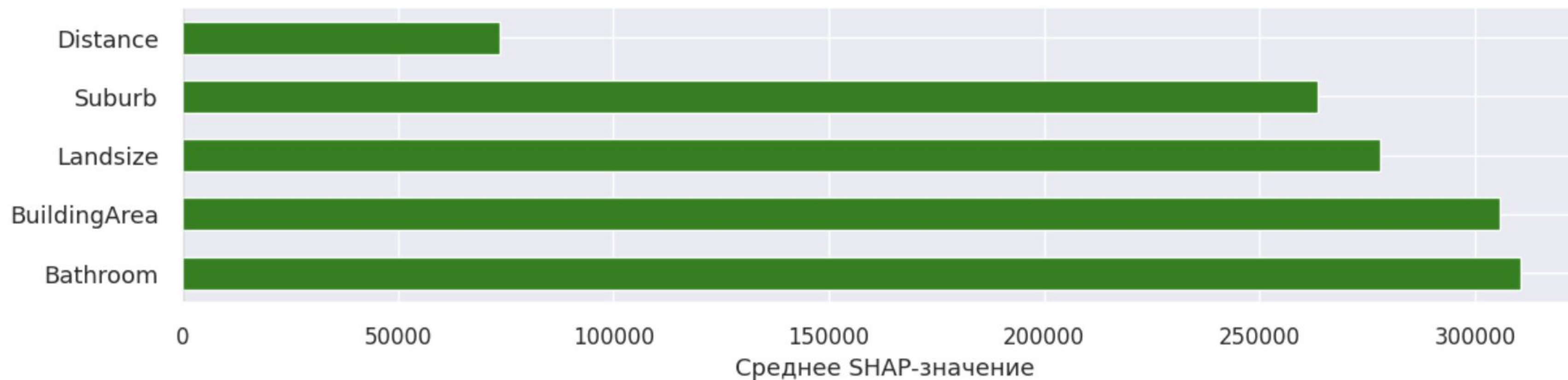


Кластеры Kmeans, Umap

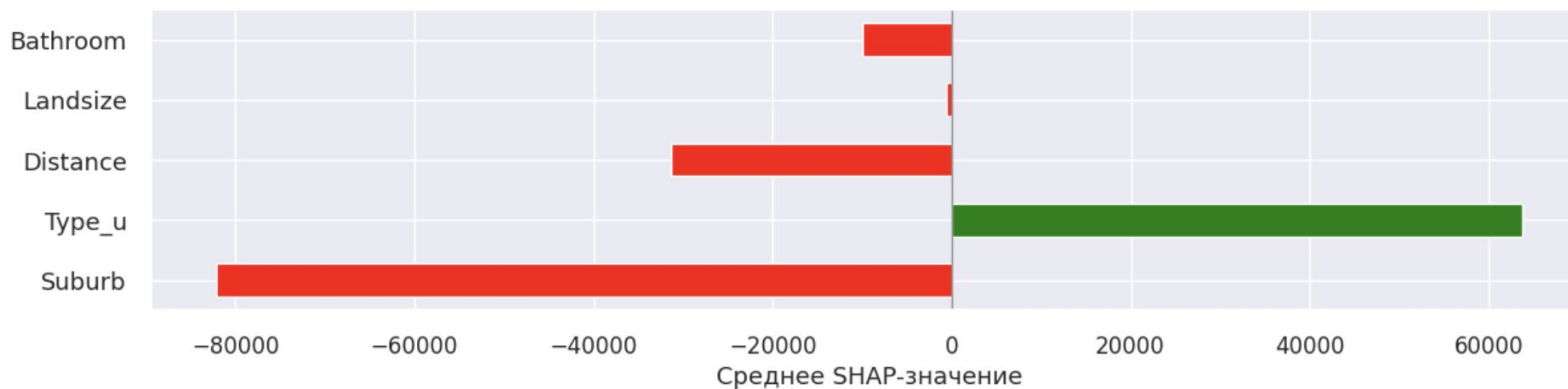


Кластеры

Кластер 0 (размер: 128)



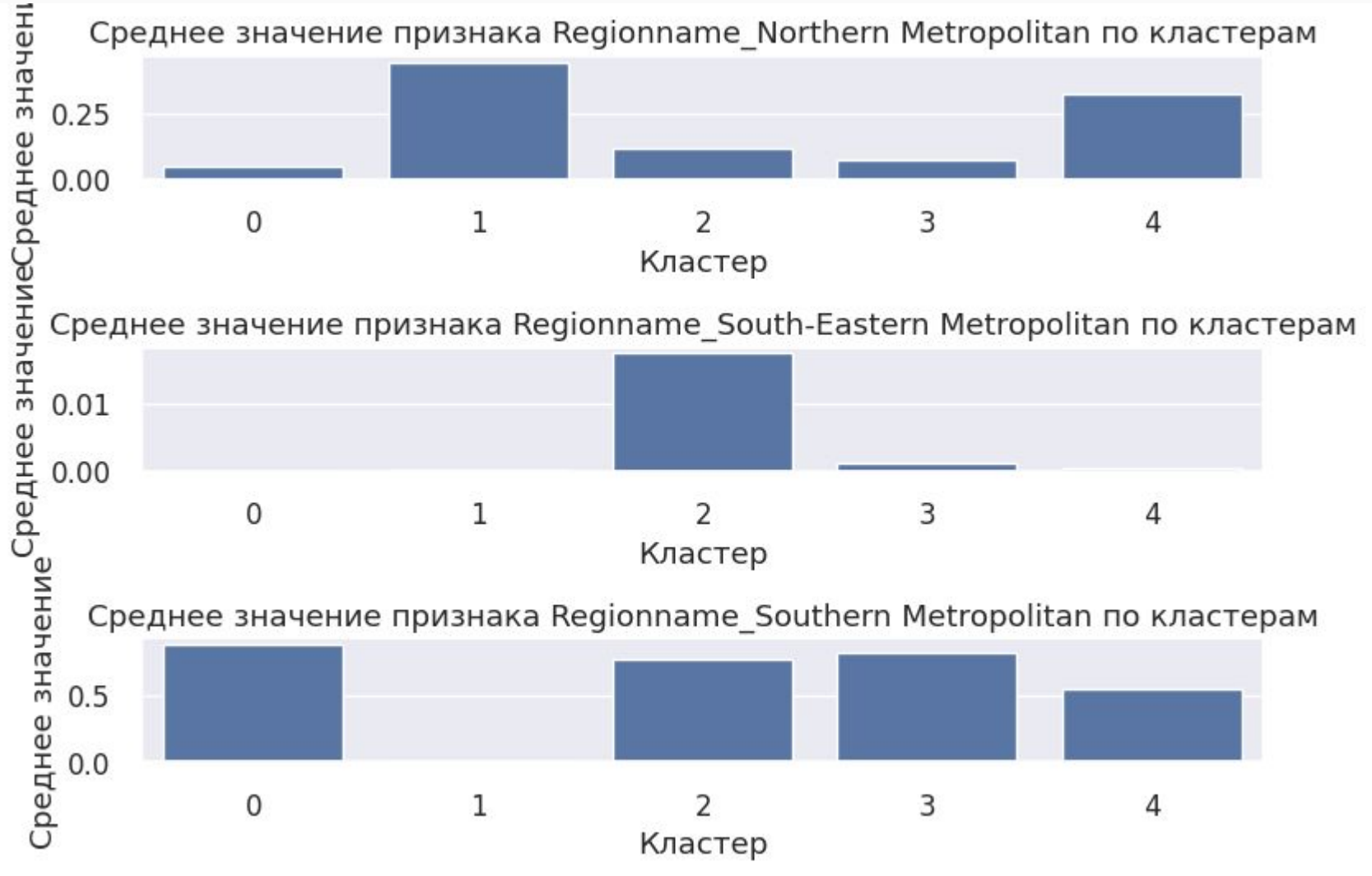
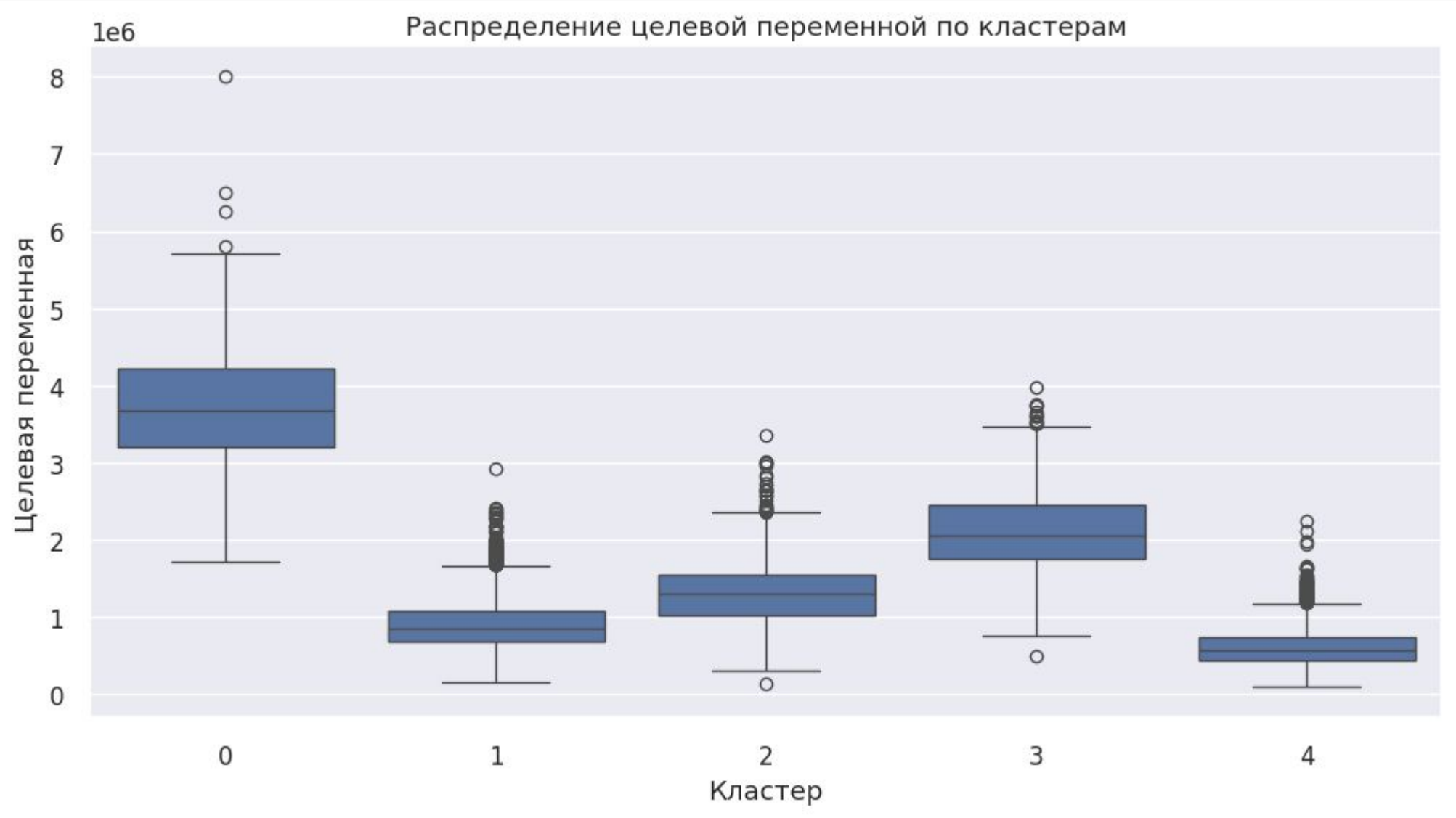
Кластер 1 (размер: 3691)



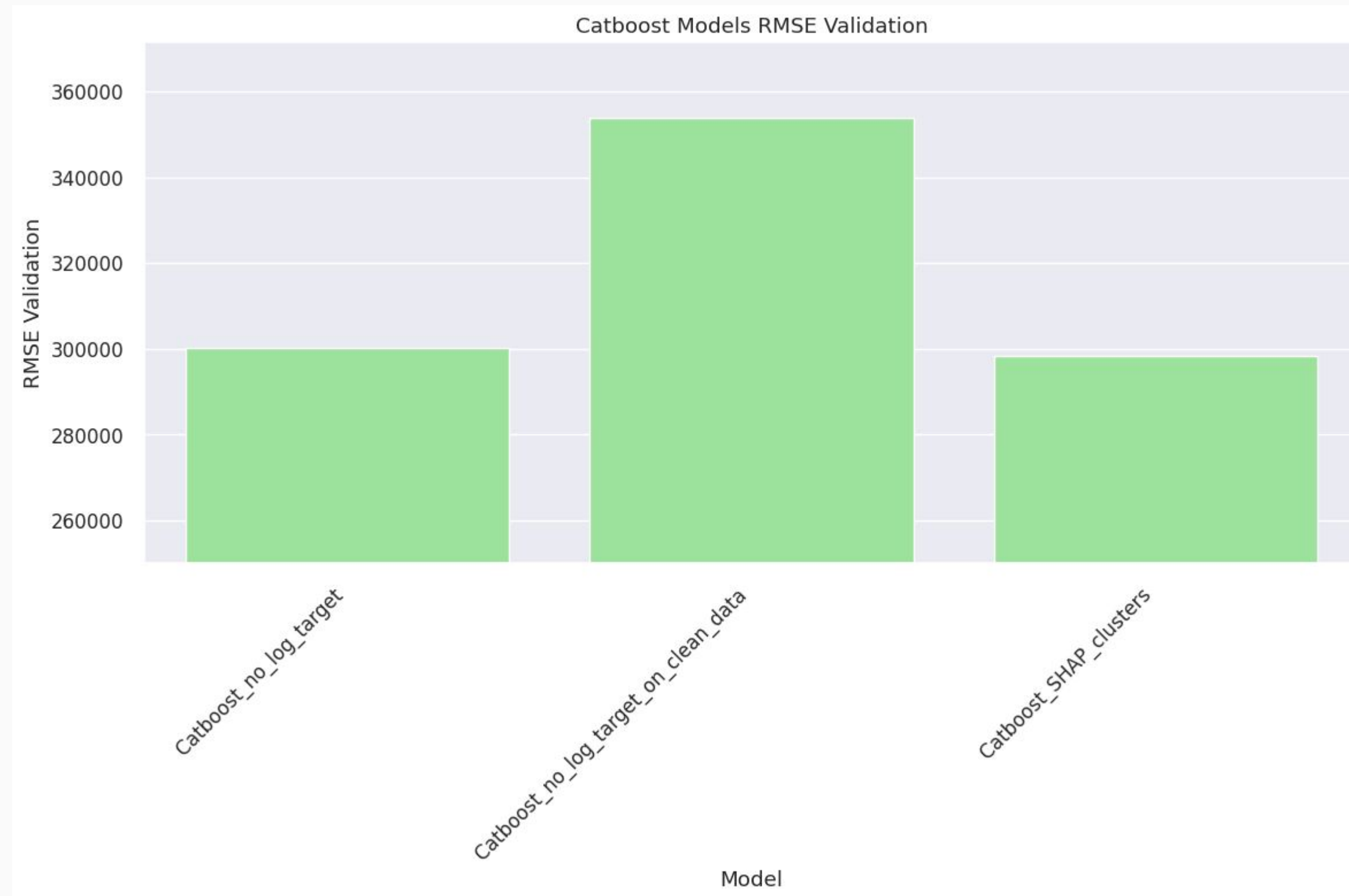
Кластеры



Профили кластеров

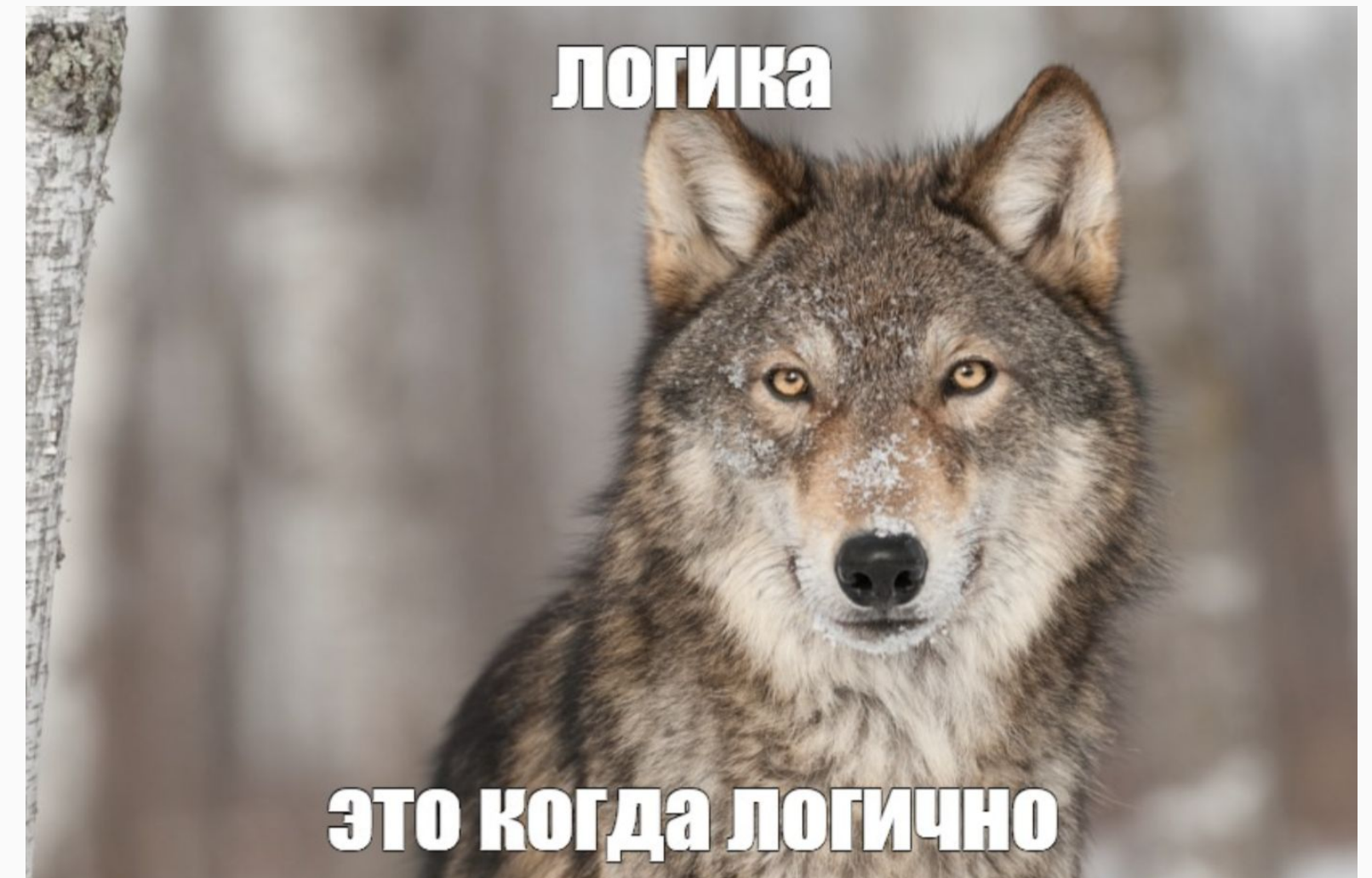
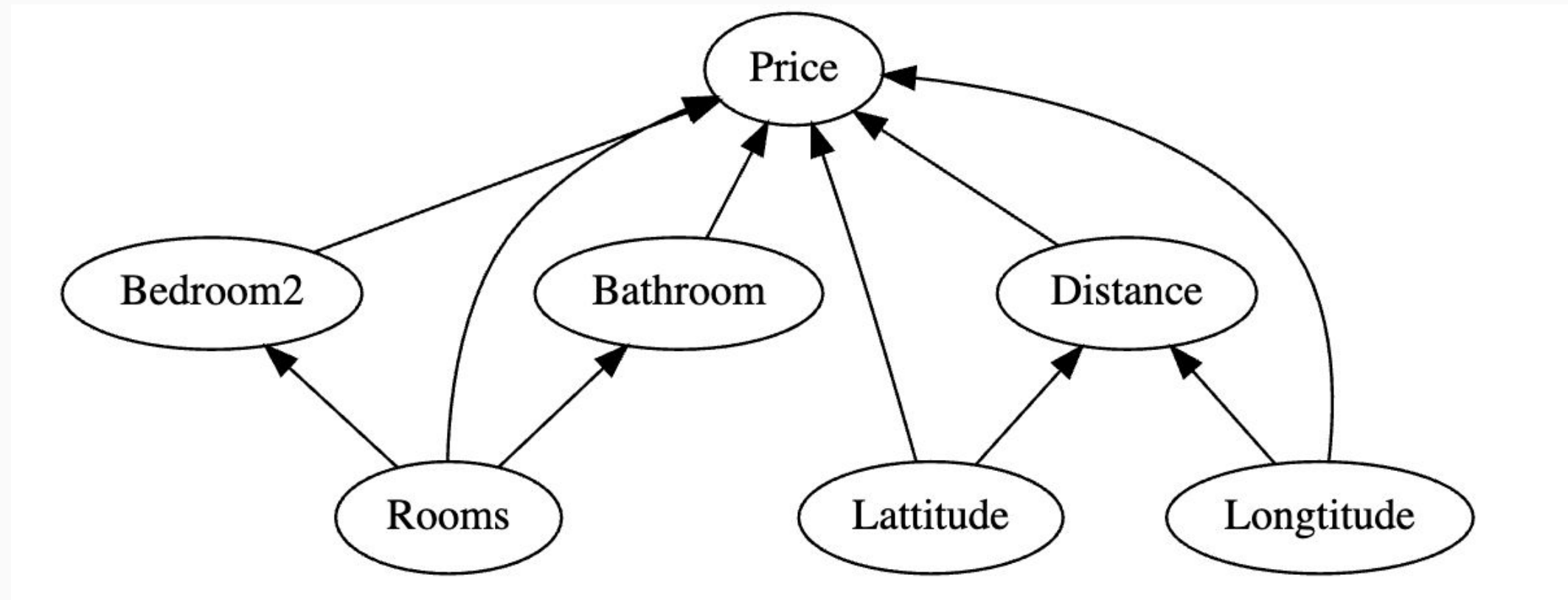


Анализ выбросов и кластеризация не помогают улучшить качество модели

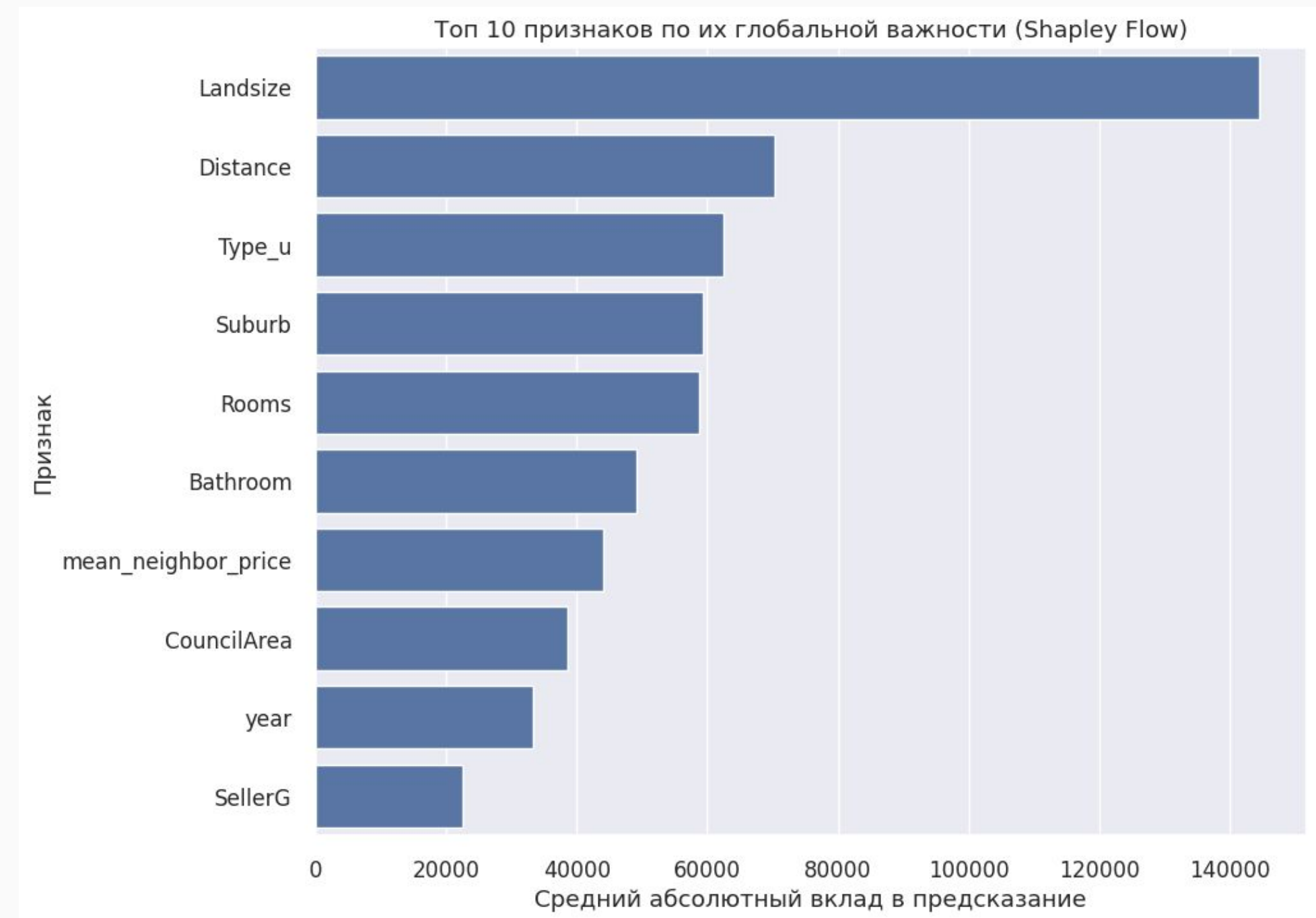
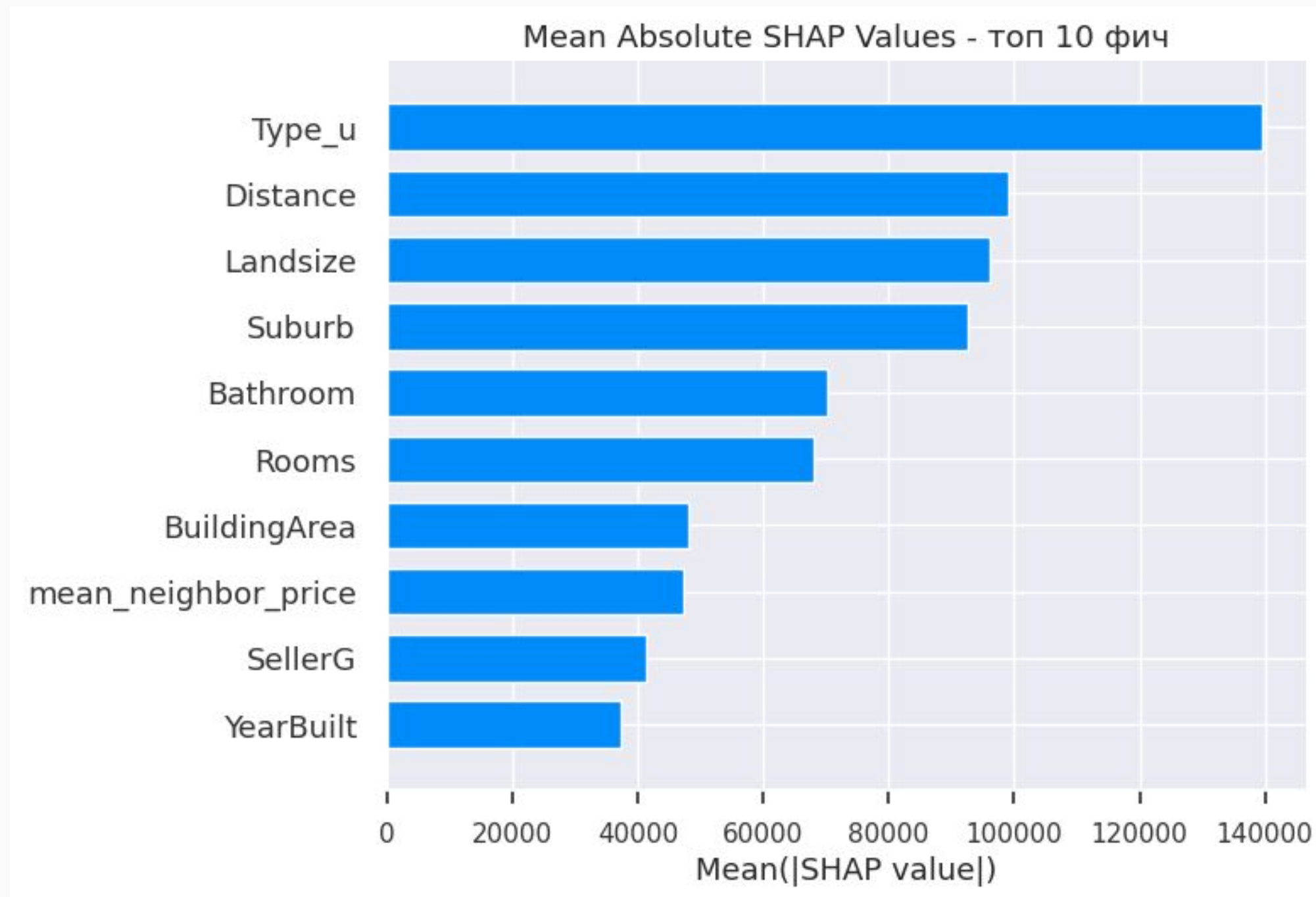


Shapley Flow

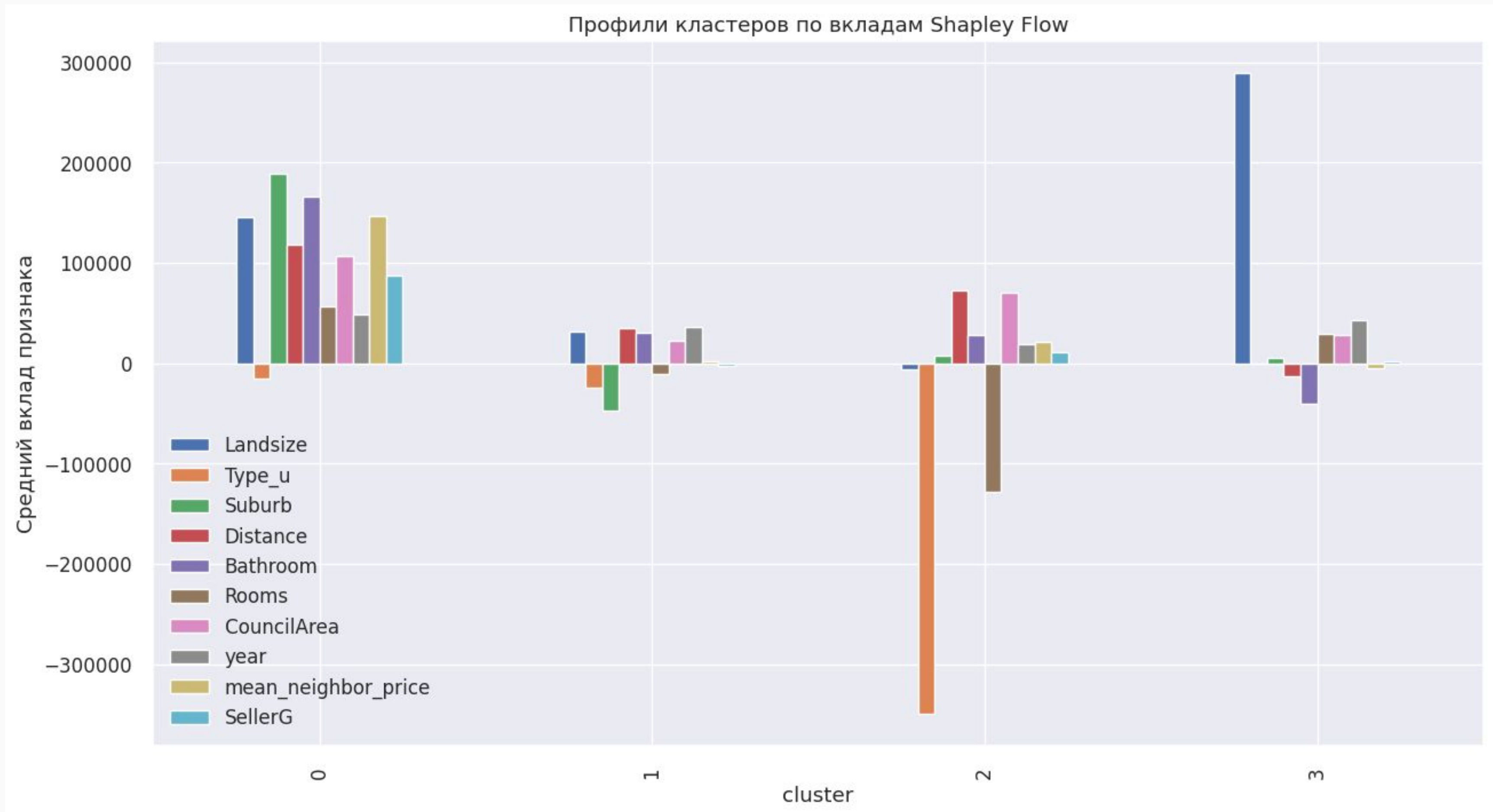
Для метода Shapley Flow нужно задать граф связей между переменными, мы используем самые очевидные связи



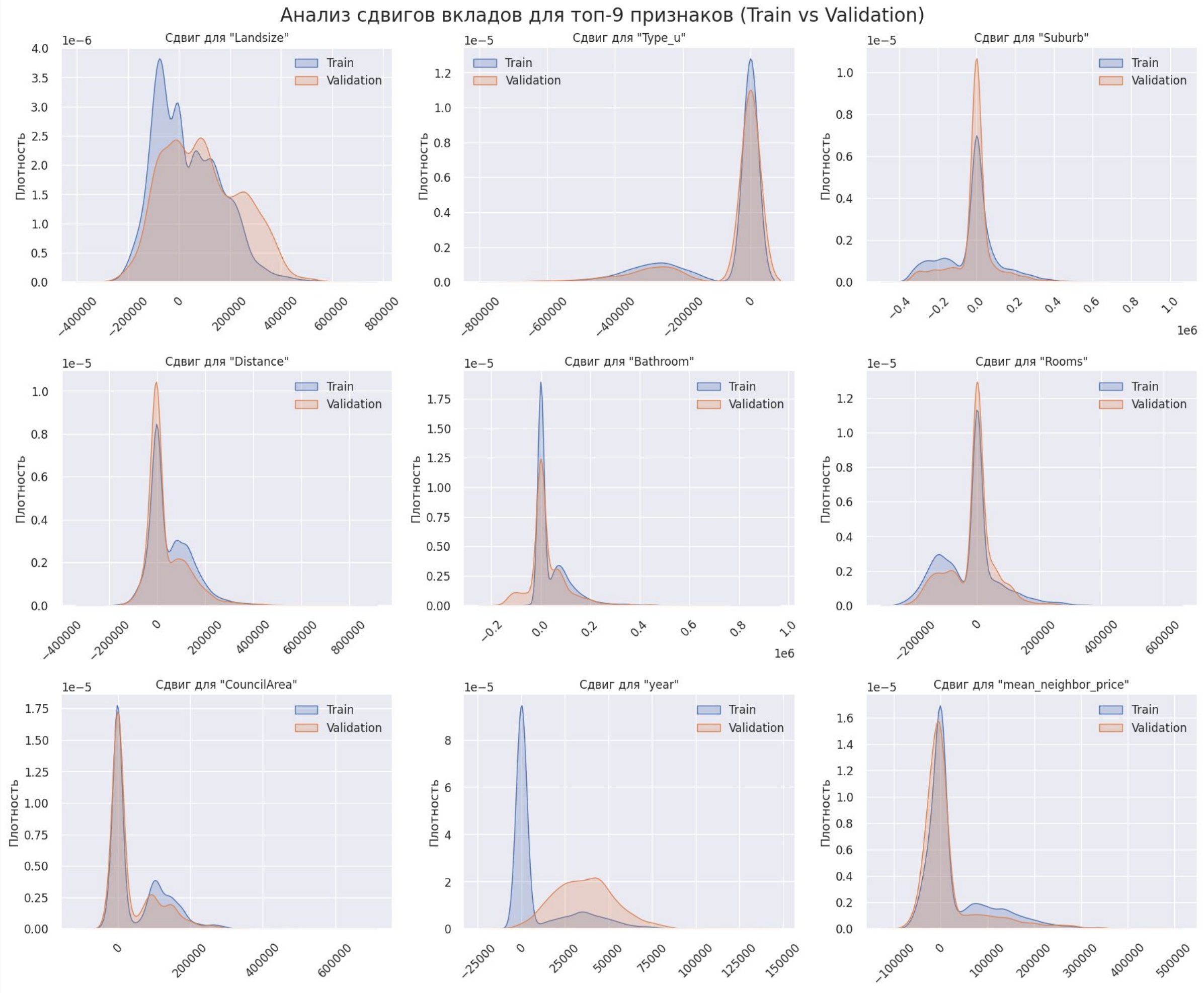
Сравнение самых важных признаков для SHAP и Shapley Flow




Кластеризация с помощью Shapely Flow



Сдвиги для Shapley Flow





**Спасибо за
внимание!**