

Анализ ценообразования ноутбуков

Александр Алгазинов

Высшая школа экономики

15 мая 2023

Все, что нужно знать

Постановка задачи

Мы хотим понять, что определяет цену ноутбуков, и насколько точно это можно смоделировать. Для этого мы будем использовать эконометрический анализ и инструменты машинного обучения

Описание данных

Возьмем готовый датасет с Kaggle. Содержит 1300 наблюдений, которые описываются 11 признаками. Особенности: пространственная выборка, содержит выбросы, требует тщательную предобработку

План работы

Выдвижение гипотез; обработка и анализ данных; тестирование гипотез и построение линейных моделей; проведение более глубокого анализа; подведение итогов исследования

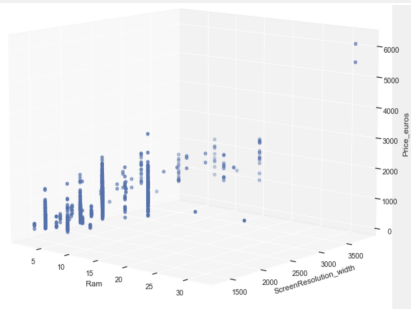
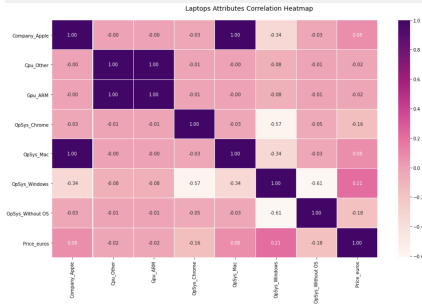
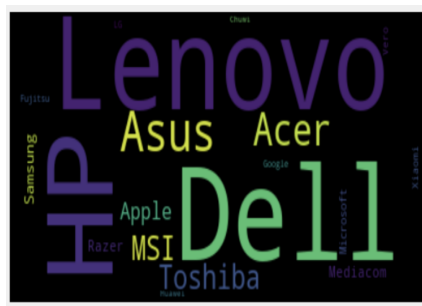
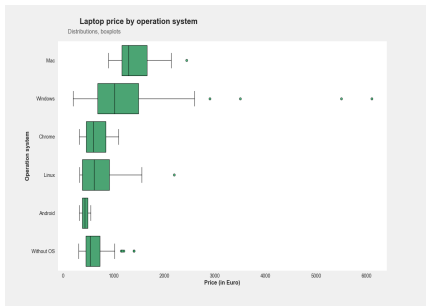
То, во что мы верили изначально

- 1. Линейность и дифференцируемость исходной задачи
- 2. Категориальные и качественные признаки важнее непрерывных
- 3. Некоторые признаки не являются значимыми

Более абстрактные, но все еще интересные идеи

- 1. Комбинации признаков улучшают модель
- 2. Неверность некоторых предпосылок ТГМ
- 3. Inductive bias линейных моделей не очень удачный
- 4. Более робастные линейные модели работают лучше

Анализ



Предпосылки ТГМ

- 1. ✗Отсутствие мультиколлинеарности (избавились удалением признаков, критерий - $VIF > 10$ + матрица корреляций)
- 2. ✗Гомоскедастичность (группы объектов, тесты Уайта и Бройша-Пагана, $p\text{-value} = 0$)
- 3. ✓Остальные предпосылки выполнены

Другие наблюдения

- 1. Модель значима ($p\text{-value} = 0$, $R^2_{adj} = 0.727$)
- 2. Статистически незначимые признаки: Weight ($p\text{-value} = 0.486$), Memory (маленький и отрицательный коэффициент)
- 3. Непрерывные признаки важнее категориальных ($R^2_{adj} = 0.548 > R^2_{adj} = 0.206$)
- 4. Комбинации непрерывных признаков не важны (переменная Ram x Screen не улучшила результаты)

Моделирование (МНК)

Мы оценили четыре модели:

- 1. $\hat{y} = w_0 + x^T w; R_{adj}^2 = 0.728$
- 2. $\log(\hat{y}) = w_0 + x^T w; R_{adj}^2 = 0.790$
- 3. $\hat{y} = w_0 + w_R \log(x_R) + w_S \log(x_S) + \sum_{i \in cat} w_i x_i; R_{adj}^2 = 0.710$
- 4. $\log(\hat{y}) = w_0 + w_R \log(x_R) + w_S \log(x_S) + \sum_{i \in cat} w_i x_i; R_{adj}^2 = 0.796$

Основные выводы:

- 1. Лучше логарифмировать таргет
- 2. Оценивание модели 3. привело к переобучению
- 3. Выбираем модель 4. из-за удобства интерпретации

$$\log(\hat{y}) = 2.004 + 0.345 \log(x_R) + 0.448 \log(x_S) + 0.619 \mathbb{1}[A] + \dots + 0.334 \mathbb{1}[W]$$

Содержательная интерпретация приведена в отчете



Мы разделили выборку на обучающую и тестовую, после чего обучили три модели, измерив качество на тесте

Модель	Функция потерь	MSE	R^2	MAPE
LinReg	MSE	102239	0.732	19.298
QuantReg ($\tau = 0.5$)	MAE	81723	0.786	17.715
Random Forest	MSE	65115	0.829	19.597

Таблица: Сравнение моделей на цифрах

Нас интересует способность модели понимать, что происходит в целом, поэтому Random Forest \succ QuantReg \succ LinReg. Оценка МНК оказалась слишком простой, но все еще хорошей для моделирования задачи

Выводы по итогам работы

- Задача линейна и дифференцируема, несмотря на частичное невыполнение предпосылок ТГМ
- Данные имеют сложную структуру (мультиколлинеарность, разделение объектов на группы)
- Результаты чувствительны к изменениям априорного представления о задаче
- Не все признаки статистически значимы. Более того, некоторые даже мешают оценивать модель
- Непрерывные признаки существенно важнее категориальных

Спасибо за внимание!