

ПРАВИТЕЛЬСТВО РОССИЙСКОЙ ФЕДЕРАЦИИ

**Федеральное государственное автономное
образовательное учреждение высшего образования**

**Национальный исследовательский университет
«Высшая школа экономики»**

Факультет экономических наук
Образовательная программа «Экономика»

КУРСОВАЯ РАБОТА

«Методы интерпретации моделей машинного обучения»

Студентка группы БЭК171
Махнева Елизавета Александровна

Научный руководитель:
Соколов Евгений Андреевич

Москва 2020

Содержание

| | | |
|----------|---|-----------|
| 1 | Введение | 3 |
| 2 | Интерпретация - подумать над названием | 4 |
| 2.1 | Зачем она нужна | 4 |
| 2.2 | Что можно интерпретировать | 4 |
| 2.3 | Методы | 4 |
| 3 | Методы и их принципы | 5 |
| 3.1 | PDP (partial dependence plot) | 5 |
| 3.1.1 | Принцип работы | 5 |
| 3.1.2 | Реализация | 6 |
| 3.2 | LIME + реализация | 6 |
| 3.3 | SHAP + реализация | 6 |
| 4 | Примеры | 7 |
| 4.1 | Первый прикольный пример | 7 |
| 4.2 | Второй прикольный пример | 7 |
| 5 | Данные и модели | 8 |
| 5.1 | Данные | 8 |
| 5.2 | Модели | 8 |
| 5.3 | Попытка интерпретации | 8 |
| 6 | Анализ результатов | 9 |
| 7 | Заключение | 10 |

1 Введение

2 Интерпретация - подумать над названием

2.1 Зачем она нужна

2.2 Что можно интерпретировать

2.3 Методы

3 Методы и их принципы

3.1 PDP (partial dependence plot)

PDP (Partial Dependence Plot, график частичной зависимости) – график, который показывает зависимость прогноза модели от значения отдельного признака. С его помощью мы можем понять, как некоторый признак влияет на предсказание. Данный график можно изобразить для двух либо трех признаков из имеющихся.

Идея: визуализация – это отличный способ интерпретации. Если мы хотим понять, как признаки влияют на результат, можно посмотреть, как меняется прогноз от изменения одного признака при прочих равных. В идеальной ситуации мы бы построили график зависимости результата от всех признаков и меняли бы только один признак. Однако мы сталкиваемся с проблемой: если признаков больше двух, построить график не получится. Поэтому чтобы сохранить возможность визуализации, можно анализировать зависимость результата от одного признака без учета влияния остальных, построив график зависимости от одного признака. Аналогично можно изучать влияние одновременно двух признаков, построив трехмерный график.

<Пример 2мерного и 3мерного графиков>

3.1.1 Принцип работы

Обозначения:

$X = (x_1, \dots, x_d)$ – матрица признаков

x_1, x_2 – векторы исследуемых признаков

$X_b = (x_3, \dots, x_d)$ – векторы остальных признаков

$a(x_1, \dots, x_d)$ – предсказания модели как функция от признаков

Нам нужно получить функцию зависимости предсказания от одного-двух признаков при зафиксированных остальных: $g(x_1, x_2) = a(x_1, x_2 | x_3, \dots, x_d)$. Но если x_1 и/или x_2 зависимы с признаками из X_b , то возникает проблема. При изменении анализируемого признака меняется и зависимый с ним, который мы не рассматриваем – мы не сможем рассмотреть чистый предельный эффект одного признака, на него всегда будет наложен эффект другого предиктора. Поэтому одной из предпосылок метода является независимость исследуемых признаков от остальных.

Но даже с предпосылкой о независимости признаков функция $g(x_1, x_2)$ не будет показывать точный результат, так как предельные эффекты предикторов разные для разных объектов выборки. Поскольку нашей задачей является посмотреть влияние выбранных признаков в целом, мы рассмотрим, как влияют анализируемые признаки на среднее предсказание. То есть найдем матожидание предсказания модели при фиксированных исследуемых признаках (как констант с точки зрения матожидания):

$$\bar{g}(x_1, x_2) = \mathbb{E}(a(x_1, x_2, X_b) | X_b)$$

Таким образом, мы получим функцию, которая показывает предельные эффекты признаков для среднего предсказания. Но чтобы найти матожидание, мы должны знать истинные распределения признаков. Поскольку нам недоступна данная информация, можно воспользоваться методом Монте-Карло, чтобы примерно оценить искомую функцию:

$$\hat{g}(x_1, x_2) = \frac{1}{n} \sum_i^n a(x_1, x_2, X_b^{(i)}),$$

где $X_b^{(i)}$ – i строка матрицы X_b

Результат: функция показывает, как исследуемые признаки в среднем влияют на результат работы модели. Мы можем построить ее график, чтобы более наглядно посмотреть на влияние предикторов на предсказание.

3.1.2 Реализация

3.2 LIME

3.2.1 Принцип работы

3.2.2 Реализация

3.3 SHAP + реализация

4 Примеры

4.1 Первый прикольный пример

4.2 Второй прикольный пример

5 Данные и модели

5.1 Данные

5.2 Модели

5.3 Попытка интерпретации

6 Анализ результатов

7 Заключение

Список литературы

- [1] [Interpretable Machine Learning](#) | Christoph Molnar | Christoph Molnar | 2020 | all pages