

Локальные методы
интерпретации
(одно наблюдение)

LIME

Local interpretable Model-agnostic
explanations

- 1) Выбрать наблюдение x
- 2) Сэмплировать наблюдения из окрестности x
- 3) Взвешивать наблюдения согласно близости к x .
- 4) На полученной датасете строить интерпретируемую модель
- 5) Интерпретация локальной модели

Плюсы

- 1) Независимость от модели
- 2) Универсальность
(таблицы, тексты, картинки)
- 3) Простота

Минусы:

- 1) Понимать как считать и в какой области
- 2) Считать из нормального
- 3) Нестабильность
- 4) Легко манипулировать

Shapley Values

$$\hat{f}(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$\phi_j(f) = \underbrace{\beta_j x_j}_{\leftarrow} - E(\beta_j x_j) = \beta_j x_j - \beta_j E(x_j)$$

$$\sum_{j=1}^p \phi_j(f) = \sum_{j=1}^p (\beta_j x_j - E(\beta_j x_j)) =$$

$$= (\beta_0 + \sum \beta_j x_j) - (\beta_0 + \sum E(\beta_j x_j)) =$$

$$= \hat{f}(x) - E(\hat{f}(x))$$

S - мн-во интересующих нас признаков

$$\phi_j(v) = \sum_{S \subseteq \{1, \dots, p\} \setminus \{j\}} \frac{|S|! (p - |S| - 1)!}{p!} (v(S \cup \{j\}) - v(S))$$

$$Val_x(S) = \int \hat{f}(x_1, \dots, x_p) dP_{x \notin S} - E_x[\hat{f}(x)]$$

4 формулы $S = \{x_1, x_3\}$)

$$\text{Var}_X(S) = \text{Var}_X(\{1, 3\}) =$$

$$= \int_{\mathbb{R}} \int_{\mathbb{R}} \hat{f}(\underbrace{x_1}_{\text{circled}}, x_2, \underbrace{x_3}_{\text{circled}}, x_4) dP_{x_2, x_4} - E_X(\hat{f}(X))$$

$$\text{Var}(S \cup \{j\}) = \text{Var}(S) \Rightarrow \phi_j = 0$$

$$\phi_j = \phi_j^1 + \phi_j^2$$

Вычисление.

M - кол-во итераций

x - наблюдение, X - выборка, f - модель

for min range(M):

1) Выбрать z из (X) случайно

2) Выбирается случайная перестановка точек

$$x_0 = (x_{(1)}, \dots, x_{(j)}, \dots, x_{(p)})$$

$$z_0 = (z_{(1)}, \dots, z_{(j)}, \dots, z_{(p)})$$

3) Создаём два новых наблюдения

$$1) x_{+j} = (x_{(1)}, \dots, x_{(j-1)}, \underbrace{x_{(j)}}_{\text{передаем}}, z_{(j+1)}, \dots, z_{(p)})$$

$$2) x_{-j} = (x_{(1)}, \dots, x_{(j-1)}, z_{(j)}, z_{(j+1)}, \dots, z_{(p)})$$

$$\phi_j^m = \hat{f}(x_{+j}) - \hat{f}(x_{-j})$$

$$\phi_j(x) = \frac{1}{M} \sum_{m=1}^M \phi_j^m$$

← передаем
восстановившем.

Плюсы:

1) Самый основ. метод

Минусы:

- 1) Долго → Всегда использует все фичи
- 2) Всегда одно значение
- 3) Нужно формул к данным
- 4) Переисчисляемые данные

SHAP

Shapley Additive Explanations

$$z' \in \{0, 1\}^M$$

M - max coalition size

$$g(z') = \phi_0 + \sum_{j=1}^M \phi_j z'_j$$

x - все эк-мбы $z' = 1$

Kernel SHAP

- 1) Суммируем координаты $z'_k \in \{0, 1\}^M$,
 $k \in \{1, \dots, K\}$
- 2) Переводим координаты в строки
 $h_x(z'_k)$
- 3) Считаем веса на вх. через
SHAP-kernel
- 4) Вытним мырез.
- 5) Кодаю-ты это наши ϕ