

# SHAP, LIME

## Интерпретируемость

Мир  $\rightarrow$  Данные  $\rightarrow$  Модель  $\rightarrow$  Интерпр.  $\rightarrow$  Profit!

1) Модель интерпретируема

а) Линейн  $\rightarrow$  Коэф.

$$\hat{\beta}_j, \hat{\beta}_j x_{ij}$$
$$y_i = \sum_j (\beta_j x_{ij}) + \beta_0$$

б) Логит  $\frac{\partial P}{\partial x_j}$ , отношение шансов

в) Деревья, разбиения

Model-agnostic methods

1) Global 2) Local

Global methods

PDP, ALE, Global surrogate

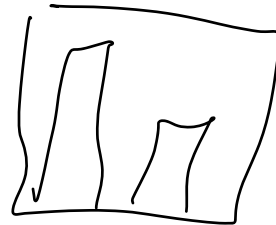
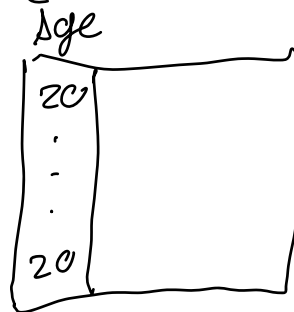
PDP - partial dependence plot

$\hat{f}$  - наша модель

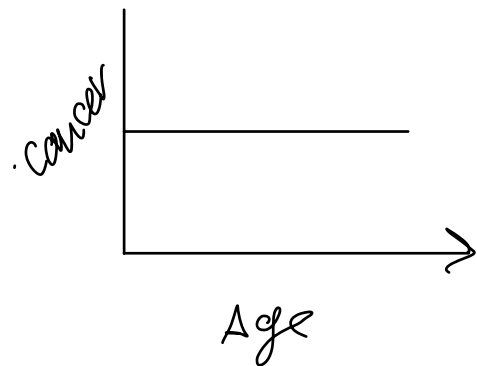
$S$  - мн-во фичей, которые мы рассматрив.

$C$  - мн-во остальных фичей

$$\hat{f}_S(x_S) = E_{x_C} [\hat{f}(x_S, x_C)] = \int \hat{f}(x_S, x_C) dP(x_C)$$



$$\hat{f}_S(x_S) = \frac{1}{n} \sum_{i=1}^n \hat{f}(x_S, x_C^{(i)})$$



$$I(x_S) = \sqrt{\frac{1}{K-1} \sum_{k=1}^K (\hat{f}_S(x_S^{(k)}) - \frac{1}{K} \sum_{k=1}^K \hat{f}_S(x_S^{(k)}))^2}$$

кон-во уника  
значений  $x_S$

$$I(x_S) = \max_k (f_S(x_S^{(k)})) - \min_k (f_S(x_S^{(k)})) / n_{cats}$$

## ALE - Accumulated Local Effects

PDP:

$$\hat{f}_{s, PDP}(x_s) = E_{x_c} [\hat{f}(x_s, x_c)]$$

M-plot

$$\hat{f}_{s, M}(x_s) = E_{x_c | X_s} [f(x_s, x_c) | X_s = x_s]$$

ALE:

$$\hat{f}_{s, ALE}(x_s) = \int_{z_{0,s}}^{x_s} E_{x_c | X_s = z_s} [\hat{f}'(x_s, x_c) | X_s = z_s] dz_s$$

- + Неинвариантность
- + Большое количество
- + Интерпретация
- + Декомпозиция

$$f = f_1 + f_2 \dots$$

- Неустойчивость
- Не декомпозируется (...)

- сложнее интерпретировать

Global surrogate

$f(x)$  - целевая ф-я (black box)

$$\hat{f}(x) = g(x)$$

- 1) Получить прогноз на всей выборке  $\hat{f}(x)$
- 2) Выбрать интерпр. модель  $g(x)$
- 3) Оценить  $\hat{f}(x) = g(x)$
- 4) Оценить качество  $g(x)$
- 5) Интерпретировать  $g(x)$

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^n (\hat{y}_g^{(i)} - \hat{y}_f^{(i)})^2}{\sum_{i=1}^n (\hat{y}_f^{(i)} - \bar{\hat{y}}_f^{(i)})^2}$$

+ Гибкость

+ Интерпретируемость

+  $R^2$  просто интерпретировать

- Неясен cut-off для  $R^2$

- Все минусы модели  $g$

Local

ICE-plots