Оценка взаимодействия признаков с помощью ансамблей консервативных деревьев

Громаков И.А.¹ Ланге А.М.²

¹ МФТИ

² Сколковский институт науки и технологий (Сколтех)

Мотивация

• Поиск gene-gene взаимодействий (эпистазис)

- Поиск drug-drug взаимодействий (выявление синергических или антагонистических эффектов между различными препаратами)
- Поиск взаимодействия биомаркеров (комбинации клинических параметров, которые лучше предсказывают риск осложнений)

Обозначения и задача

- X пространство объектов
- Y множество ответов
- $f_j = X \to D_j$, j = 1, 2, ..., p признаки объектов
- $F = \{f_j \mid j \in \{1, 2, ..., p\}\}$ множество признаков
- $\mathscr{F} \subseteq F$ подмножество признаков, $|\mathscr{F}| \ge 2$
- $S(\mathscr{F})$ мера взаимодействия \mathscr{F}

Требуется:

- Для всех выбранных подмножеств \mathscr{F} определить, есть ли у них взаимодействие: $\forall k \leq k_0, \ k \in \{2, 3, ..., |F|\}$, где k_0 заданное ограничение, равное максимальной глубине деревьев
- Найти $T_k = \{ \mathscr{F} | S(\mathscr{F}) \ge S_0^k, k = |\mathscr{F}| \}$
- Ранжировать степень взаимодействия у разлиных групп признаков $\mathscr{F} {\in} T_k$ в каждом найденном T_k

Взаимодействие

Пусть $a(x) = a(x_1, x_2, ..., x_p)$ - функция, аппроксимирующая целевую функцию, а $(x_1, x_2, ..., x_p) - p$ признаков модели

Определение: вещественные признаки \mathscr{F} будем называть взаимодействующими, если

$$E_{x} \left[\frac{\partial a(x)}{\prod_{x \in \mathscr{T}} \partial x} \right]^{2} > 0$$

Определение: признаки \mathscr{F} будем называть *не взаимодействующими*, если a(x) можно представить в виде суммы $N \leq |\mathscr{F}|$ функций, каждая из которых не зависит от хотя бы одной переменной $x \in \mathscr{F}$:

$$a(x) = \sum_{k: x_k \in \mathscr{F}} f_{x_{-k}}(x_{-k}),$$

где x_{-k} означает все признаки кроме x_k

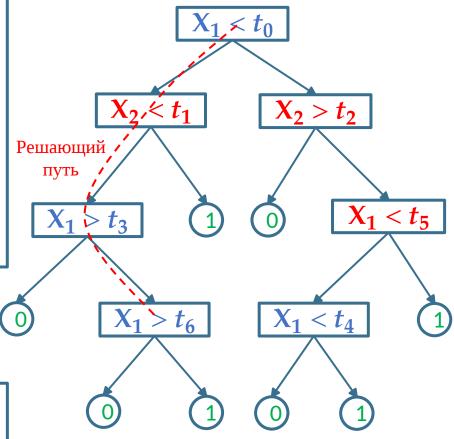
Решающие деревья

Решающее дерево:

Это алгоритм классификации a(x), задающийся деревом с корнем $v_0 \in V$ и множеством внутренних вершин $V = V_{\mathsf{BHYTP}} \cup V_{\mathsf{ЛИСТ}}$, в котором каждой внутренней вершине $v \in V_{\mathsf{BHYTP}}$ сопоставлен дискретный признак (предикат) $\beta_v \colon X \to D_v$, правило перехода в дочерние вершины $S_v \colon D_v \to V$, и каждой листовой вершине $\forall v \in V_{\mathsf{ЛИСТ}}$ ставится в соответствие $y_v \in Y$

$$Gain(\beta, U) = \Phi(U) - \Phi(U \mid \beta),$$
 где U – рассматриваемая выборка объектов

Даже если целевая функция аддитивна, то решающие пути все равно содержат разные признаки. Такие паттерны могут возникать как следствие случайных вариаций или рекурсивной природы алгоритмов построения дерева



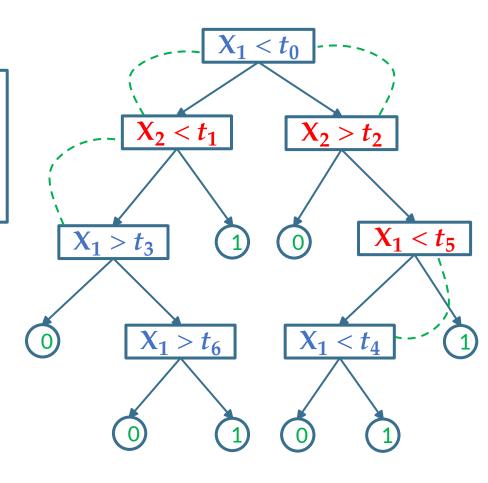
Консервативные решающие деревья

Консервативное решающее дерево:

Это модификация CART, вводящая новый гиперпараметр $rit_alpha^1 \in (-\infty, 1)$. Использование rit_alpha дает преимущество при ветвлении тем признакам, которые уже были выбраны ранее в дереве

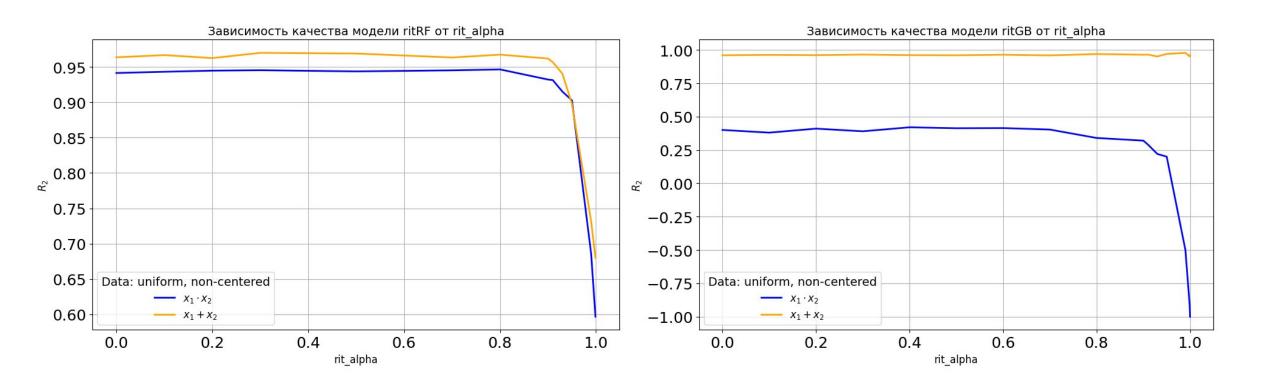
$$B = Gain(\beta_{f_{best}}, U)(1 - rit_alpha)$$

Разные признаки в решающем пути – потенциальное взаимодействие



¹ RIT - Redundancy Insensitive Trees

Random Forest / Gradient Boosting



Подбор гиперпараметра

Пусть

- $\mathcal{M}_{\alpha} = \{M_{\alpha}^{(1)}, ..., M_{\alpha}^{(N)}\}$ множество N моделей градиентного бустинга, обученных с параметром $rit_alpha = \alpha$ на случайных подвыборках обучающей выборки одинакового размера
- score(M) метрика качества модели M
- S_{α} выборка с посчитанными метриками качества моделей

Формальная постановка задачи:

 $\{H_0:$ средние значения $score\left(M_{lpha_0}^{(i)}
ight)$ и $score\left(M_{-\infty}^{(j)}
ight)$ совпадают, $H_1:$ качество при $lpha_0$ статистически значимо ниже, чем при $lpha_0=-\infty$

Для проверки используется односторонний критерий Манна-Уитни (U-test) на уровне значимости $\delta = 0.05$:

$$p = MannWhitneyU(S_{\alpha_0}, S_{-\infty}, alternative = \prime less')$$

Если $p > \delta$, то гипотеза H_0 **не отвергается**, и параметр α_0 может быть принят как допустимый, не ухудшающий качество модели

Статистика взаимодействий

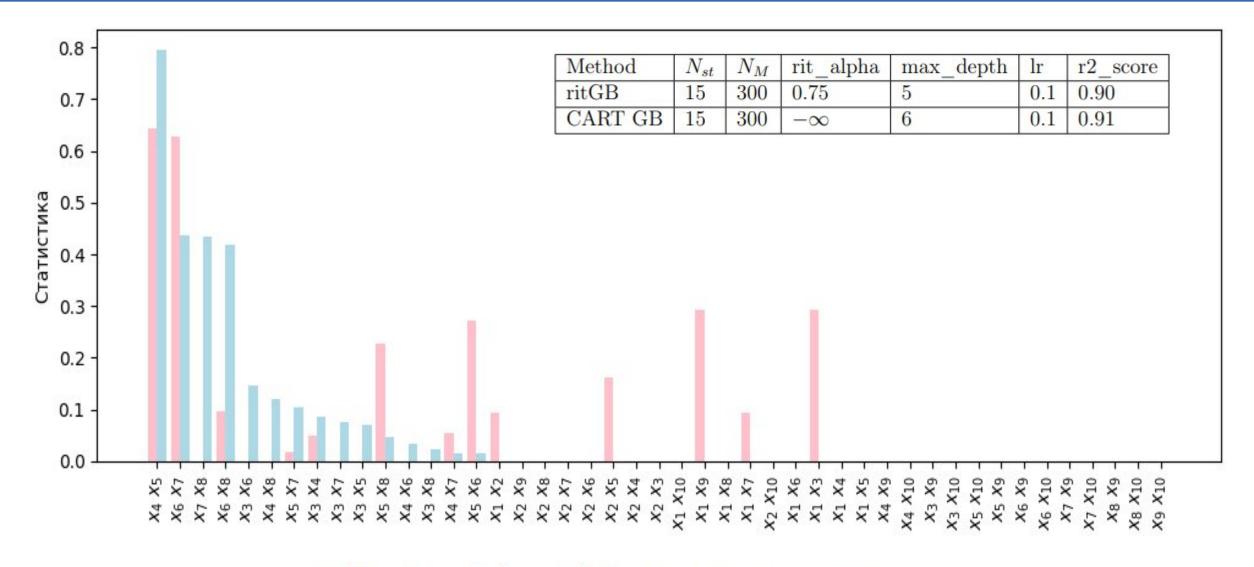
Предлагается подобрать гиперпараметры для N_{st} моделей с найденным rit_alpha и для каждой модели найти нормализованную статистику $S_n(M_\alpha, \mathscr{F})$, а затем усреднить результат:

$$S(\mathcal{F}) = \frac{1}{N_{st}} \sum_{j=1}^{N_{st}} S_n(M_\alpha^j, \mathcal{F}) \qquad S_n(M_\alpha, \mathcal{F}) = \frac{S(M_\alpha, \mathcal{F})}{(\prod_{x \in \mathcal{F}} S^{|\mathcal{F}|}(x))^{1/|\mathcal{F}|}}$$

$$S^{K}(x_{j}) = \frac{1}{N_{M_{\alpha}}} \sum_{i \in M_{\alpha}} \frac{1}{|V_{\text{JIICT}}^{i}|} \sum_{v \in V_{\text{JIICT}}^{i}} \frac{N_{v}^{i}}{N^{i}} \frac{m_{v}^{i,K}(x_{j})}{l_{v}^{i}} \quad S(M_{\alpha}, \mathcal{F}) = \frac{1}{N_{M_{\alpha}}} \sum_{i \in M_{\alpha}} \frac{1}{|V_{\text{JIICT}}^{i}|} \sum_{v \in V_{\text{JIICT}}^{i}} \frac{N_{v}^{i}}{N^{i}} \frac{n_{v}^{i}(\mathcal{F})}{l_{v}^{i}}$$

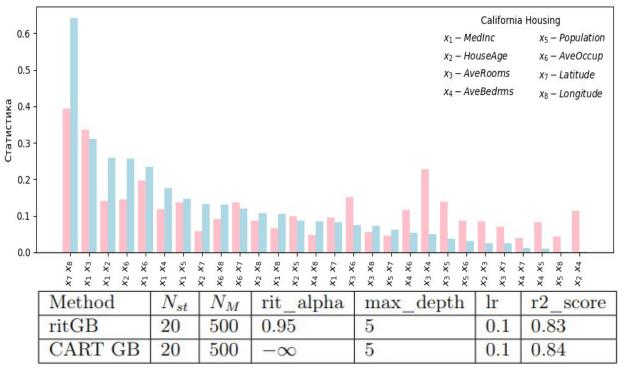
где l_v — длины путей от корня до листа v, $n_v(\mathscr{F})$ — индикатор множества уникальных признаков решающего пути с листовой вершиной v, $m_v^K(x)$ — индикатор наличия признака x в множестве уникальных признаков \mathscr{F} : $|\mathscr{F}| = K$ на решающем пути, N_v — количество объектов из обучающей выборки, которые попали в листовую вершину v, N — общее количество объектов в обучающей выборке, N_{M_α} — количество предикторов модели M_α

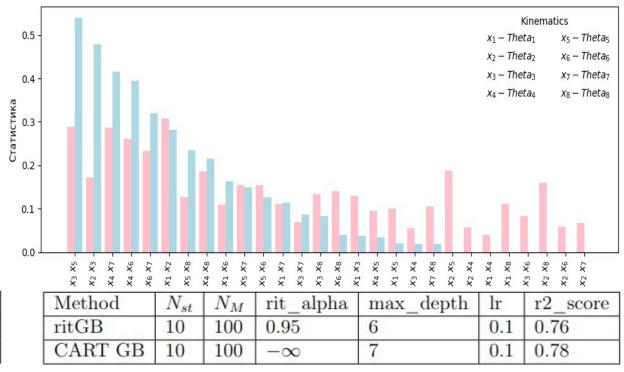
Синтетические данные



$$F(x) = 2x_1 + 3x_2^2 + \sin(x_3) + 5x_4x_5 + 4x_6x_7x_8 + 0.5x_9 + \varepsilon$$

Реальные данные





Выводы

- Метод способен выявлять взаимодействия малых порядков (парные, тройные)
- Метод может быть использован как дополнение к комплексному анализу взаимодействий признаков