



Правила остановки последовательного выбора признаков и статистические методы для моделей с бинарным откликом

Научный руководитель: Булинский Александр Вадимович
Мехмат МГУ

Сперва рассмотрим результат из статье [1] Jan Mielniczuk, Pawel Teisseyre.

Будем рассматривать p -тое количество признаков X_1, \dots, X_p и целевая классовая величина(отклик) - Y . Пусть X_S - подмножество X_1, \dots, X_p , $S \subset \{1, \dots, p\}$. Определим $p(x_j) := P(X_j = x_j)$, $x_j \in \mathcal{X}_j$, где \mathcal{X}_j - область значений X_j и $|\mathcal{X}_j|$ - мощность множества(конечно). Область значений классовой величины - \mathcal{Y} (конечно). Совместная вероятность - $p(x_i, x_j) = P(X_i = x_i, X_j = x_j)$. $\hat{p}(x_j)$ - обозначение оценки $p(x_j)$.

Задача - зная значения признаков X_1, \dots, X_p , предположить значение Y . Например, задача можем сводиться к вопросу о $P(Y = y | X_1 = x_1, \dots, X_p = x_p)$

По известному набору наблюдений необходимо найти зависимость отклика от значений признаков.

Однако, наш отклик Y может не зависеть от существенной части признаков, либо же иметь пренебрежительно слабую зависимость. Рассмотрение всех признаков может быть существенно затруднено, например, в силу высоких вычисленных требований в виде экспоненциального роста сложности.

Цель - выбрать наиболее подходящие признаки.

$$\arg \max_{S: |S|=k} I(X_S, Y)$$

Однако, перебор по всем признакам является крайне затратным, так как сложность растёт экспоненциально от числа признаков

Следовательно, можно перейти к последовательному выбору признаков:

$$\arg \max_{j \in S^c} [I(X_{S \cup \{j\}}, Y) - I(X_S, Y)] = \arg \max_{j \in S^c} I(X_j, Y \mid X_S)$$

Разложим условную совместную информацию через информации многостороннего взаимодействия:

$$\begin{aligned}
 & I(X_{S \cup \{j\}}, Y) - I(X_S, Y) \\
 &= I(X_j, Y \mid X_S) = \sum_{k=0}^{|S|} \sum_{\{i_1, \dots, i_k\} \subseteq S} II(X_{i_1}, \dots, X_{i_k}, X_j, Y) \\
 &= I(X_j, Y) + \sum_{i \in S} II(X_i, X_j, Y) + \sum_{i_1, i_2 \in S: i_1 < i_2} II(X_{i_1}, X_{i_2}, X_j, Y) \\
 &+ \dots + II(X_{i_1}, \dots, X_{i_{|S|}}, X_j, Y)
 \end{aligned} \tag{1}$$

Для упрощения вычислений можно взять второй порядок приближения условной совместной информации:

$$\begin{aligned} I(X_j, S) &= I(X_j, Y) + \sum_{i \in S} II(X_i, X_j, Y) \\ &= I(X_j, Y) + \sum_{i \in S} [I(Y, X_j | X_i) - I(Y, X_j)] \\ &= I(X_j, Y) (1 - |S|) + \sum_{i \in S} I(Y, X_j | X_i) \end{aligned} \tag{2}$$

На практике приближение выше второго использовать затруднительно. Для r -того порядка приближения необходимо оценивать $(r + 1)$ -мерные вероятности. Пусть n - число наблюдений, для любого i - $|\mathcal{X}_i| = b$, следовательно, всего b^{r+1} - возможных комбинаций. n/b^{r+1} - в среднем наблюдений на комбинацию.

Например, если $n = 1000$, $b = 5$, $r = 2$, то $n/b^{r+1} = 8$. Если $r = 3$, то $n/b^{r+1} = 1.6$. То есть не удастся приближенно оценить вероятности.

Авторы строят последовательный выбор признаков. Необходима точка остановки, на которой выбор предполагаемых значимых признаков будет окончен.

Предположим выбрали S множество индексов признаков. Тогда выбор признаков останавливается при $I(Y, X_j | X_S) = 0$ для любого $j \in S^c$. Пользуемся приближением второго порядка.

Пусть S_k - множество индексов, выбранных на k -том шаге, где $S_0 = \emptyset$. На k -том шаге определим $S_{k+1} = S_k \cup \{j_k\}$ такое, что $j_k = \arg \max_{j \in S_k^c} J(X_j, S_k)$. Таким образом момент остановки в приближении:

$$t := \arg \min_{1 \leq k \leq p} (J(X_j, S_k) = 0), \forall j \in S_k^c.$$

Так как мы не знаем распределений, нужны оценки.

$$\begin{aligned}\hat{J}(X_j, S_k) &= \hat{I}(X_j, Y) + \sum_{i \in S_k} \hat{I}(X_i, X_j, Y) \\ &= \hat{I}(X_j, Y) + \sum_{i \in S_k} \left[\hat{I}(Y, X_j | X_i) - \hat{I}(Y, X_j) \right] \\ &= \hat{I}(X_j, Y) (1 - |S_k|) + \sum_{i \in S_k} \hat{I}(Y, X_j | X_i)\end{aligned}\tag{3}$$

Заметим, что при $J(X_j, S_k) = 0$, может выполняться $\hat{J}(X_j, S_k) > 0$.

Цель - построить приближение распределения $2n\hat{J}(X_j, S_k)$ при гипотезе, что $X_j \perp Y \mid X_{S_k}$; и ввести правило остановки \hat{t} , основанное на кватилиях распределения.

Теорема 1

Пусть X , Y и Z - случайные величины, принимающие значения в $|\mathcal{X}|$, $|\mathcal{Y}|$ и $|\mathcal{Z}|$, соответственно. Предположим, что Y и Z независимы при условии X . Тогда $2n\hat{I}(Y, Z | X) \approx \sum_{i=1}^{|\mathcal{X}|} W_i$, где W_i имеет χ^2 асимптотическое распределение со $(|\mathcal{Y}| - 1)(|\mathcal{Z}| - 1)$ степенями свободы и \approx - означает, что обе части отличаются только на величину, стремящуюся к нулю по вероятности.

□ Для упрощения записи $\hat{p}_{ijk} = \hat{p}(X = x_i, Y = y_j, Z = z_k)$, $\hat{p}_{ij} = \hat{p}(X = x_i, Y = y_j)$ и т.д.. Запишем, используя определение условной совместной информации:

$$2n\hat{I}(Y, Z | X) = 2n \sum_{i,j,k} \hat{p}_{ijk} \log \frac{\hat{p}_{ijk}\hat{p}_i}{\hat{p}_{ij}\hat{p}_{ik}}$$

$$= 2n \sum_i \frac{1}{\hat{p}_i} \sum_{j,k} \hat{p}_{ijk}\hat{p}_i \log \left(1 + \frac{\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}}{\hat{p}_{ij}\hat{p}_{ik}} \right) \quad (4)$$

Используя $\log(1 + x) = x - x^2/2 + O(x^3)$ для малых x , получим:

$$\log \left(1 + \frac{\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}}{\hat{p}_{ij}\hat{p}_{ik}} \right) = \frac{\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}}{\hat{p}_{ij}\hat{p}_{ik}} - \frac{1}{2} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{(\hat{p}_{ij}\hat{p}_{ik})^2}$$

$$+ O \left(\frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^3}{(\hat{p}_{ij}\hat{p}_{ik})^3} \right) \quad (5)$$

Вставим выражение выше в (4), видим, что член содержащий последний член в (5) ограничен некоторой $C > 0$ как

$$\begin{aligned}
 & C \times 2n \sum_i \frac{1}{\hat{p}_i} \sum_{j,k} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^3}{(\hat{p}_{ij}\hat{p}_{ik})^3} \\
 & \leq C \times 2n \sum_i \frac{1}{\hat{p}_i} \times \frac{\max_{j,k} |\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}|}{\min_{j,k} (\hat{p}_{ij}\hat{p}_{ik})^2} \sum_{j,k} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{\hat{p}_{ij}\hat{p}_{ik}} \\
 & = C \times 2 \sum_i \frac{1}{\hat{p}_i^2} \times \frac{\max_{j,k} |\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}|}{\min_{j,k} (\hat{p}_{ij}\hat{p}_{ik})^2} \\
 & \quad \times n \hat{p}_i \sum_{j,k} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{\hat{p}_{ij}\hat{p}_{ik}}
 \end{aligned} \tag{6}$$

$$\begin{aligned}
&= C \times 2 \sum_i \frac{1}{\hat{p}_i^2} \times \frac{\max_{j,k} |\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}|}{\min_{j,k} (\hat{p}_{ij}\hat{p}_{ik})^2} \\
&\quad \times n\hat{p}_i \sum_{j,k} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{\hat{p}_{ij}\hat{p}_{ik}}
\end{aligned} \tag{6}$$

В силу условной независимости и сходимости $\hat{p}_{ij} \rightarrow p_{ij} > 0$, $\hat{p}_{ik} \rightarrow p_{ik} > 0$ имеем

$$|\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}| \rightarrow |p_{ijk}p_i - p_{ij}p_{ik}| = 0 \tag{7}$$

при $n \rightarrow \infty$.

В силу (7) и (10) ниже, последний член в (6) - сумма двух членов таких, что первый член сходится к нулю, и второй член имеет распределение хи-квадрат. Следовательно, из теоремы Slutsky следует, что оценка в (6) сходится по вероятности к 0, когда $n \rightarrow \infty$. Таким образом $2n\hat{I}(Y, Z | X)$ примерно равен

$$\begin{aligned}
 & 2n \sum_i \frac{1}{\hat{p}_i} \sum_{j,k} [\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik} + \hat{p}_{ij}\hat{p}_{ik}] \\
 & \times \left[\frac{\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}}{\hat{p}_{ij}\hat{p}_{ik}} - \frac{1}{2} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{(\hat{p}_{ij}\hat{p}_{ik})^2} \right] \\
 & = 2n \sum_i \frac{1}{\hat{p}_i} \sum_{j,k} \left(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik} + \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{\hat{p}_{ij}\hat{p}_{ik}} \right. \\
 & \quad \left. - \frac{1}{2} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{\hat{p}_{ij}\hat{p}_{ik}} + \frac{1}{2} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^3}{(\hat{p}_{ij}\hat{p}_{ik})^2} \right)
 \end{aligned}$$

$$\begin{aligned}
&= 2n \sum_i \frac{1}{\hat{p}_i} \sum_{j,k} \left(\hat{p}_{ijk} \hat{p}_i - \hat{p}_{ij} \hat{p}_{ik} + \frac{(\hat{p}_{ijk} \hat{p}_i - \hat{p}_{ij} \hat{p}_{ik})^2}{\hat{p}_{ij} \hat{p}_{ik}} \right. \\
&\quad \left. - \frac{1}{2} \frac{(\hat{p}_{ijk} \hat{p}_i - \hat{p}_{ij} \hat{p}_{ik})^2}{\hat{p}_{ij} \hat{p}_{ik}} + \frac{1}{2} \frac{(\hat{p}_{ijk} \hat{p}_i - \hat{p}_{ij} \hat{p}_{ik})^3}{(\hat{p}_{ij} \hat{p}_{ik})^2} \right) \\
&\approx n \sum_i \frac{1}{\hat{p}_i} \sum_{j,k} \left(2(\hat{p}_{ijk} \hat{p}_i - \hat{p}_{ij} \hat{p}_{ik}) + \frac{(\hat{p}_{ijk} \hat{p}_i - \hat{p}_{ij} \hat{p}_{ik})^2}{\hat{p}_{ij} \hat{p}_{ik}} \right) \\
&= n \sum_i \frac{1}{\hat{p}_i} \sum_{j,k} \frac{(\hat{p}_{ijk} \hat{p}_i - \hat{p}_{ij} \hat{p}_{ik})^2}{\hat{p}_{ij} \hat{p}_{ik}}
\end{aligned} \tag{8}$$

где приближение в (8) получено аналогично как в (6), и последнее равенство следует из то, что $\sum_{j,k}(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}) = 0$. Последнее полученное выражение равно

$$\begin{aligned} n \sum_i \hat{p}_i \sum_{j,k} \frac{(\hat{p}_{ijk}/\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}/\hat{p}_i^2)^2}{\hat{p}_{ij}\hat{p}_{ik}/\hat{p}_i^2} \\ = \sum_i n_i \sum_{j,k} \frac{(\hat{p}_{ijk}/\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik}/\hat{p}_i^2)^2}{\hat{p}_{ij}\hat{p}_{ik}/\hat{p}_i^2} =: \sum_{i=1}^{|\mathcal{X}|} W_i \end{aligned} \quad (9)$$

где $n_i = n\hat{p}_i$.

Заметим, что W_i есть хи-квадрат статистика для тестирования независимости Y и Z при $X = x_i$, которая при предположении независимости Y и Z при X имеет асимптотическое χ^2 распределение со $(|\mathcal{Y}| - 1)(|\mathcal{Z}| - 1)$ степенями свободы

$$W_i = n_i \sum_{j,k} \frac{(\hat{p}_{ijk}\hat{p}_i - \hat{p}_{ij}\hat{p}_{ik})^2}{\hat{p}_{ij}\hat{p}_{ik}} \approx \chi^2_{(|\mathcal{Y}|-1)(|\mathcal{Z}|-1)} \quad (10)$$

в силу (Теорема 6.9, Shao Mathematical Statistics). Автором было показано, что число степеней свободы предельного χ^2 распределения равно $p - s - 1$, где $p = |\mathcal{Y}||\mathcal{Z}|$ и $s = (|\mathcal{Y}| - 1) + (|\mathcal{Z}| - 1)$, при предположении условной независимости $(|\mathcal{Y}| - 1) + (|\mathcal{Z}| - 1)$ условных вероятностей Y и Z при $X = x_i$ нужно определить, чтобы определить условное распределение (Y, Z) при $X = x_i$. Таким образом общее число степеней свободы равно $|\mathcal{Y}||\mathcal{Z}| - (|\mathcal{Y}| - 1) - (|\mathcal{Z}| - 1) - 1 = (|\mathcal{Y}| - 1)(|\mathcal{Z}| - 1)$. \square

Далее перейдём к теме потенциального применения правил остановки. Далее материал из статьи [2].

X_1^j, \dots, X_n^j - случайные величины признаков j -того наблюдения. Y^j - соответственный бинарный отклик (-1 или 1). Пусть

$$\xi = (\xi^1, \dots, \xi^N)$$

где

$$\xi^j = (X^j, Y^j), \quad j = 1, \dots, N$$

Будем считать ξ^1, \dots, ξ^N - н.о.р. случайные вектора. Введём (X, Y) независимый с ξ и имеющий распределение ξ^1 . Все случайные вектора рассматриваются на (Ω, \mathcal{F}, P) , E - интегрирование по P .

Пусть $X := \{0, 1, 2\}^n$ - пространство возможных значений объясняющих переменных. Функция $f : X \rightarrow \{-1, 1\}$ называется теоритической функцией предсказания. Определим сбалансированную и нормированную ошибку предсказания для f как

$$\text{Err}(f) := \mathbb{E}|Y - f(X)|\psi(Y)$$

где $\psi : \{-1, 1\} \rightarrow \mathbf{R}_+$ - штрафная функция. Следовательно,

$$\begin{aligned} \text{Err}(f) = & 2\psi(-1)\mathbf{P}(f(X) = 1, Y = -1) \\ & + 2\psi(1)\mathbf{P}(f(X) = -1, Y = 1) \end{aligned} \tag{1}$$

Очевидно $\text{Err}(f)$ зависит от распределения (X, Y) , но опустим обозначение. Ссылаясь на идеи предыдущих работ авторы устанавливают

$$\psi(y) = \frac{1}{4P(Y=y)}, \quad y \in \{-1, 1\}$$

где случаи $P(Y = -1) = 0$ и $P(Y = 1) = 0$ исключаются. Тогда

$$\text{Err}(f) = \frac{1}{2}P(f(X) = 1 \mid Y = -1) + \frac{1}{2}P(f(X) = -1 \mid Y = 1) \quad (2)$$

Если $P(Y = -1) = P(Y = 1) = 1/2$, то называем сбалансированным: $\text{Err}(f) = E|Y - f(X)|/2$.

Следовательно, $\text{Err}(f)$ равно ошибке классификации $P(Y \neq f(X))$. В общем,

$$\text{Err}(f) = \frac{1}{2} E |Y^* - f(X^*)|$$

где (X^*, Y^*) имеет распределение

$$P(X^* = x, Y^* = y) = \frac{1}{2} P(X = x | Y = y)$$

$$(x, y) \in X \times \{-1, 1\}$$

Можно заметить, что оптимальная теоритическая функция предсказания минимизирующая сбалансированную ошибку предсказания:

$$f^*(x) = \begin{cases} 1, p(x) > P(Y = 1) \\ -1, \text{ иначе} \end{cases}$$

где

$$p(x) = P(Y = 1 \mid X = x), x \in X \quad (4)$$

Тогда каждый $x \in \mathcal{X}$ классифицируется как вероятный, если $f^*(x) = 1$, иначе при $f^*(x) = -1$.

$p(x)$ и $P(Y = 1)$ неизвестны. Необходимо найти приближение f^* , используя

$$f_{PA} = f_{PA}(x, \xi(S))$$

со значениями в $\{-1, 1\}$, которая зависит от $x \in \mathcal{X}$ и наблюдений

$$\xi(S) = \{\xi^j, j \in S\}$$

где

$$S \subset \{1, \dots, N\} \tag{5}$$

Авторы вводят оценки $p(x)$ и $P(Y = 1)$ как

$$\hat{p}(x, \xi(S)) = \frac{\sum_{j \in S} I\{Y^j = 1, X^j = x\}}{\sum_{j \in S} I\{X^j = x\}}, x \in \mathcal{X} \quad (6)$$

и

$$\hat{P}_S(Y = 1) = \frac{1}{\#S} \sum_{j \in S} I\{Y^j = 1\} \quad (7)$$

где $\#D$ - мощность множества D .

Устанавливается аналогично для f_{PA} :

$$\begin{aligned} & \text{Err}(f_{PA}(\cdot, \xi(S))) \\ &= \frac{1}{2} \sum_{y \in \{-1, 1\}} \mathbb{P}(f_{PA}(X, \xi(S)) \neq y \mid Y = y) \end{aligned} \quad (9)$$

Распределение (X, Y) неизвестно, строится оценка $\hat{\text{Err}}(f_{PA}(\cdot, \xi(S)))$

Строится оценка для f_{PA} с помощью K -свёрточной кросс-валидации ($K > 1$):

$$\begin{aligned} & \hat{\text{Err}}_K(f_{PA}(\cdot, \xi), \xi) \\ &= \frac{1}{2} \sum_{y \in \{-1, 1\}} \frac{1}{K} \sum_{k=1}^K \frac{\sum_{(k)} I \left\{ f_{PA} \left(X^j, \xi \left(\overline{S}_k \right) \right) \neq y, Y^j = y \right\}}{\sum_{(k)} I \{ Y^j = y \}} \end{aligned}$$

где сумма $\sum_{(k)}$ берётся по всем j принадлежащим

$$\begin{aligned} S_k = \{ & (k-1)[N/K] + 1, \dots \\ & k[N/K] I\{k < K\} + NI\{k = K\} \} \end{aligned} \quad (11)$$

$$C_{k_1, \dots, k_r}(x) = \{u = (u_1, \dots, u_n) \in X : u_{k_i} = x_{k_i}, i = 1, \dots, r\}.$$

Теорема 2

Предположим существует подмножество $U \subset X$ и подмножество $\{k_1, \dots, k_r\} \subset \{1, \dots, n\}$ такие, что выполнено:

1. Для любого $x \in \mathcal{X}$ и каждого конечномерного вектора v с компонентами в $X \times \{-1, 1\}$, функции $f_{PA}(\cdot, v)$ and f постоянны на $C_{k_1, \dots, k_r}(x)$.
2. Для любого $x \in U$ и каждого $W_N \subset \{1, \dots, N\}$ с $\#W_N \rightarrow \infty$, верно $f_{PA}(x, \xi(W_N)) \rightarrow f(x)$ п.н. при $N \rightarrow \infty$.
3. $P(Y = 1 \mid X_{k_1} = x_{k_1}, \dots, X_{k_r} = x_{k_r}) = P(Y = 1)$, если $x \in X \setminus U$.
4. f постоянна на $X \setminus U$.

Тогда $\text{Err}_K(\hat{f}_{PA}(\cdot, \xi), \xi) \rightarrow \text{Err}(f)$ п.н., $N \rightarrow \infty$.

Замечание. Если заменить условие 3 в Теореме 1 более усиленным

3') $P(Y = 1 \mid X = x) = P(Y = 1)$ для любого $x \in X \setminus U$, то можно брать $\{k_1, \dots, k_r\} = \{1, \dots, n\}$, чтобы убрать условие 1.

□. Доказательство основано на следующей лемме:

Лемма 1

Пусть $\left\{ \left(Z_j^{(m)}, Y_j^{(m)} \right), \quad 1 \leq j \leq m, m \in N \right\}$ - набор по строчно независимых случайных величин распределенных как (Z, Y) , где Z принимает значения в конечном множестве Z , и Y - в $\{-1, 1\}$.

Предположим $\{f_m(z), m \in N, z \in Z\}$ - набор случайных величин со значениями в $\{-1, 1\}$.

Предположим существует $U \subset Z$ такое, что верны следующие условия:

1. $f_m(z) \rightarrow f(z)$ п.н. для любого $z \in U$, при $m \rightarrow \infty$, где неслучайная величина $f : Z \rightarrow \{-1, 1\}$.
2. $P(Y = 1 \mid Z = z) = P(Y = 1)$, если $z \in Z \setminus U$.
3. f - постоянная на $Z \setminus U$.

Тогда $m \rightarrow \infty$,

$$\frac{1}{2} \sum_{y \in \{-1, 1\}} \frac{\sum_{j=1}^m I \left\{ f_m \left(Z_j^{(m)} \right) \neq y, Y_j^{(m)} = y \right\}}{\sum_{j=1}^m I \left\{ Y_j^{(m)} = y \right\}} \rightarrow \text{Err}(f) \quad (12)$$

□. Авторы вводят $Q_m(y) = \sum_{j=1}^m I \{Y_j^{(m)} = y\}$ и определяют следующие события:

$$A_j^{(m)}(y) = \{f_m(Z_j^{(m)}) \neq y, Y_j^{(m)} = y\}$$

$$B_j^{(m)}(y) = \{f(Z_j^{(m)}) \neq y, Y_j^{(m)} = y\}$$

Тогда левая часть (12) равна

$$\sum_{y \in \{-1, 1\}} \frac{1}{2Q_m(y)} \sum_{j=1}^m I \{A_j^{(m)}(y)\}$$

Для $y \in \{-1, 1\}$ верно

$$\begin{aligned} & \frac{1}{Q_m(y)} \sum_{j=1}^m I \{A_j^{(m)}(y)\} \\ &= \frac{1}{mP\{Y = y\}} \sum_{j=1}^m I \{A_j^{(m)}(y)\} \\ &+ \frac{1}{m} \sum_{j=1}^m I \{A_j^{(m)}(y)\} \left(\frac{m}{Q_m(y)} - \frac{1}{P(Y = y)} \right) \end{aligned} \tag{13}$$

Абсолютная величина второго члена в правой части (13) не превышает $|m / (Q_m(y)) - 1/P(Y = y)|$ и стремится к 0 п.н. при $m \rightarrow \infty$. Это следует из усиленного закона больших чисел для массивов (УЗБЧМ). Заметим

$$\begin{aligned}
 & \frac{1}{m} \sum_{j=1}^m I \{A_j^{(m)}(y)\} \\
 &= \frac{1}{m} \sum_{j=1}^m I \{B_j^{(m)}(y)\} \\
 &+ \frac{1}{m} \sum_{j=1}^m I \{Z_j^{(m)} \in U\} \left(I \{A_j^{(m)}(y)\} - I \{B_j^{(m)}(y)\} \right) \\
 &+ \frac{1}{m} \sum_{j=1}^m I \{Z_j^{(m)} \notin U\} \left(I \{A_j^{(m)}(y)\} - I \{B_j^{(m)}(y)\} \right)
 \end{aligned} \tag{14}$$

Согласно УЗБЧМ первый член в правой части (14) стремится $P(f(Z) \neq y, Y = y)$ п.н. Предположим, второй член стремится к 0 п.н. Действительно, множество Z - конечно и функции f, f_m принимают два значения. Следовательно, из условия 1, для почти всех $\omega \in \Omega$, существует $N_1 = N_1(\omega)$ такой, что $f_m(z) = f(z)$ для всех $z \in U$ и $m > N_1$. Тогда второй член в правой части (14) равен 0 для всех $m > N_1$, что доказывает предположение. Таким образом остаётся оценить третий член.

Из условия 3, б.о.о. можем принять $f(z) = -1$ для $z \in Z \setminus U$. Тогда получается

$$\begin{aligned}
 V_m &:= \sum_{y \in \{-1, 1\}} \sum_{j=1}^m \frac{I\{Z_j^{(m)} \notin U\} \left(I\{A_j^{(m)}(y)\} - I\{B_j^{(m)}(y)\} \right)}{mP(Y = y)} \\
 &= \frac{1}{m} \sum_{j=1}^m I\{Z_j^{(m)} \notin U\} I\{f_m(Z_j^{(m)}) = 1\} R_j^{(m)} \\
 &= \frac{1}{m} \sum_{z \in Z \setminus U} I\{f_m(z) = 1\} \sum_{j=1}^m I\{Z_j^{(m)} = z\} R_j^{(m)}
 \end{aligned} \tag{15}$$

где

$$R_j^{(m)} = \frac{I\{Y_j^{(m)} = -1\}}{P(Y = -1)} - \frac{I\{Y_j^{(m)} = 1\}}{P(Y = 1)}$$

УЗБЧН и условие 2 приводят к тому, что для $z \in Z \setminus U$ и $y \in \{-1, 1\}$:

$$\sum_{j=1}^m \frac{I\{Z_j^{(m)} = z\} I\{Y_j^{(m)} = y\}}{mP(Y = y)} \\ \rightarrow \frac{P(Z = z, Y = y)}{P(Y = y)} = P(Z = z)$$

почти наверное. Следовательно, для почти всех $\omega \in \Omega$, существует $N_2 = N_2(\omega)$ такое, что

$$\left| \frac{1}{m} \sum_{j=1}^m I\{Z_j^{(m)} = z\} R_j^{(m)} \right| < \varepsilon$$

для всех $z \in Z \setminus U$ и $m > N_2$. Используя последнюю оценку и (15), получается, что для $m > N_2$,

$$|V_m| \leq \sum_{z \in Z \setminus U} \varepsilon I\{f_m(z) = 1\} \leq \varepsilon \cdot \#Z$$

Тогда $V_m \rightarrow 0$ п.н. при $m \rightarrow \infty$. Собирая (12)-(16), получается необходимый результат. \square
 Теперь к доказательству Теоремы 1. Зафиксируем $1 \leq k \leq K$ и определим

$$f_m(z) := f_{PA} \left(x, \xi \left(\overline{S_k} \right) \right)$$

где $z \in Z := \{0, 1, 2\}^r$, $m := \#S_k$, S_k - определен в (11), и x - любой элемент из \mathcal{X} с $(x_{k_1}, \dots, x_{k_r}) = z$.
 Из условия 1 f_m корректно определён.

Применим лемму к набору

$$\left\{ \left(Z_j^{(m)}, Y_j^{(m)} \right), 1 \leq j \leq m \right\} := \left(\left(X_{k_1}^j, \dots, X_{k_r}^j \right), Y^j \right), j \in S_k \right\}$$

и $\{f_m(z), z \in Z\}$, получается, что п.н.

$$\frac{1}{2} \sum_{y \in \{-1, 1\}} \frac{\sum_{(k)} I \left\{ f_{PA} \left(X^j, \xi \left(\overline{S_k} \right) \right) \neq y, Y^j = y \right\}}{\sum_{(k)} I \{ Y^j = y \}} \rightarrow \text{Err}(f)$$

при $\#S_k \rightarrow \infty$.

Таким образом

$$\hat{\text{Err}}_K(f_{PA}(\cdot, \xi), \xi) \rightarrow \frac{1}{K} \sum_{k=1}^K \text{Err}(f) = \text{Err}(f)$$

при $N \rightarrow \infty$. \square

Как упомянуто ранее, отклик может зависеть лишь от ограниченного числа признаков. В статье вводится понятие набора значимых признаков - существует $l \in \mathbf{N}, l < n$, и вектор (k_1^*, \dots, k_l^*) , где $1 \leq k_1^* < \dots < k_l^* \leq n$ такой, что для каждого $x = (x_1, \dots, x_n) \in \mathcal{X}$ верно:

$$p(x) = P\left(Y = 1 \mid X_{k_1^*} = x_{k_1^*}, \dots, X_{k_l^*} = x_{k_l^*}\right) \quad (17)$$

Набор (k_1^*, \dots, k_l^*) имеющий наименьший размер l называется наиболее значимым.

Для $x \in \mathcal{X}$ и $\{k_1, \dots, k_r\} \subset \{1, \dots, n\}$, вводится

$$\begin{aligned} & f_{k_1, \dots, k_r}(x) \\ &= \begin{cases} 1, & P(Y = 1 \mid X_{k_1} = x_{k_1}, \dots, X_{k_r} = x_{k_r}) > P(Y = 1) \\ -1, & \text{иначе.} \end{cases} \end{aligned}$$

Аналогично вводится оценка $\hat{f}_{k_1, \dots, k_r}(x)$.

Теорема 3

Пусть (k_1^*, \dots, k_l^*) - наиболее значимый набор признаков. Тогда для любого фиксированного $\{k_1, \dots, k_r\} \subset \{1, \dots, n\}$ верно

1. $\text{Err}(f_{k_1^*, \dots, k_l^*}) \leq \text{Err}(f_{k_1, \dots, k_r})$;
2. $\hat{\text{Err}}_K(\hat{f}_{k_1, \dots, k_r})$ - сильная состоятельная асимптотически несмещенная оценка $\text{Err}(f_{k_1, \dots, k_r})$ при $N \rightarrow \infty$;
3. Для всех $\varepsilon, \delta > 0$ и для всех достаточно больших N

$$\mathbb{P}\left(\hat{\text{Err}}_K(\hat{f}_{k_1^*, \dots, k_l^*}) < \hat{\text{Err}}_K(\hat{f}_{k_1, \dots, k_r}) + \varepsilon\right) > 1 - \delta$$

□ 1) Из (17) следует, что $f_{k_1^*, \dots, k_l^*}$ совпадает с функцией f^* , которая является минимумом весовой ошибки прогнозирования.

2. Проверим условия теоремы 1 для $f_{PA}(x, \xi(S)) := \hat{f}_{k_1, \dots, k_r}(x, \xi(S))$. Условие 1 следует из определения $f_{k_1, \dots, k_r}(x, \xi(S))$. Далее

$$U := \{x \in X : P(Y = 1 \mid C_{k_1, \dots, k_r}(x)) \neq P(Y = 1)\} \quad (19)$$

где C_{k_1, \dots, k_r} введено раньше. Предположим, что для каждого $x \in U$ и любого $W_N \subset \{1, \dots, N\}$, что $\#W_N \rightarrow \infty$, верно:

$$\hat{f}_{k_1, \dots, k_r}(x, \xi(W_N)) \rightarrow f_{k_1, \dots, k_r}(x) \text{ п.н. при } N \rightarrow \infty$$

Действительно, предположим, что для некоторого $\varepsilon > 0$,

$$P(Y = 1 \mid X_{k_1} = x_{k_1}, \dots, X_{k_r} = x_{k_r}) - P(Y = 1) > \varepsilon$$

Из УЗБЧН следует, что

$$\hat{P}_{W_N}(Y = 1 | X \in C_{k_1, \dots, k_r}(x)) - \hat{P}_{W_N}(Y = 1)$$

сходится п.н. к

$$P(Y = 1 | X_{k_1} = x_{k_1}, \dots, X_{k_r} = x_{k_r}) - P(Y = 1)$$

$N \rightarrow \infty$. Тогда, для почти всех $\omega \in \Omega$ существует $N_0 = N_0(\omega)$ такой, что

$$\hat{P}_{W_N}(Y = 1 | X \in C_{k_1, \dots, k_r}(x)) - \hat{P}_{W_N}(Y = 1) > \varepsilon/2$$

для всех $N > N_0$. Следовательно, для всех $N > N_0$ имеем $\hat{f}_{k_1, \dots, k_r}(x, \xi(W_N)) = 1 = f_{k_1, \dots, k_r}(x)$, что доказывает предположение. Условие 2 Теоремы 1 удовлетворено.

Условия 3 и 4 следуют из (19) и определения $\hat{f}_{k_1, \dots, k_r}(x, \xi(S))$ и $f_{k_1, \dots, k_r}(x)$.

Следовательно, из Теоремы 1 следует, что $\widehat{Err}_K(\hat{f}_{k_1, \dots, k_r}) \rightarrow Err(f_{k_1, \dots, k_r})$ п.н. и в L_2 (в силу теоремы Лебега из ограниченности $\widehat{Err}_K(\hat{f}_{k_1, \dots, k_r})$ величиной 1).

3. Следует из 1) и 2).



Стоит отметить, что частично аналогичные результаты вместе с другими дополнительными результатами, но уже для небинарного отклика получены с статьях [4], [5], [6].

Рассмотрим процесс стратификации из [3]

$(X^1, Y^1), (X^2, Y^2), \dots$ - н.о.р как (X, Y) . Общий набор наблюдений разбивается на два набора:

$$\zeta_{N_1}^1 := \left\{ (X^{j_1^1}, 1), \dots, (X^{j_{N_1}^1}, 1) \right\}, \zeta_{N_{-1}}^{-1} := \left\{ (X^{j_{-1}^1}, -1), \dots, (X^{j_{-1}^{N_{-1}}}, -1) \right\}$$

где j_1^k и $j_{-1}^k, k \in \mathbb{N}$, упорядочены. $N := N_1 + N_{-1}$ для $\zeta_N := \zeta_{N_1}^1 \cup \zeta_{N_{-1}}^{-1}$. В отличие от $\xi_N := \{(X^1, Y^1), \dots, (X^N, Y^N)\}$, имеющий закон $\text{law}(X, Y)$, есть $\zeta_{N_1}^1$ и $\zeta_{N_{-1}}^{-1}$ с распределениями $X \mid Y = 1$ и $X \mid Y = -1$. Поэтому нельзя использовать оценки (например, of $P(Y = 1)$ или $P(X \in B, Y = 1)$, где $B \subset \mathbb{X}$) построенные на средних ζ_N . Пусть $N_1 = \max\{[aN], 1\}$ и $N_{-1} = N - N_1$

Предположим есть оценка \hat{P}_N^y от $P(Y = y)$:

$$\hat{P}_N^y \rightarrow P(Y = y) \text{ п.н., } N \rightarrow \infty, y \in \{-1, 1\} \quad (8)$$

Например, \hat{P}_N^y использует $\left\{ \left(X^k, Y^k \right), 1 \leq k \leq \max \left\{ j_1^{N_1}, j_{-1}^{N_{-1}} \right\} \right\}$. Тогда частотные оценки $P(Y = 1)$ или $P(X \in B, Y = 1), B \subset \mathbb{X}$, сильно состоятельны, $\max \left\{ j_1^{N_1}, j_{-1}^{N_{-1}} \right\} \rightarrow \infty$ п.н. $N \rightarrow \infty$.

$\hat{P}_N := \left(\hat{P}_N^{-1}, \hat{P}_N^1 \right)$. $P(Y = -1) = 0$ и $P(Y = 1) = 0$ не рассматривается.

Пусть $f_{PA} (x, \zeta_N, \hat{P}_N)$ - функция определяющая алгоритм предсказания, т.е. функция в $\{-1, 1\}$ зависящая от $x \in \mathbb{X}$, наблюдений ζ_N и \hat{P}_N . Аналогично определяется на поднаборах. $f_{PA} (x, \zeta_N(S), \hat{P}_N)$ for $\zeta_N(S) := \{(X^j, Y^j), j \in S\}$, $S \subset (\{j_1^1, \dots, j_1^{N_1}\} \cup \{j_{-1}^1, \dots, j_{-1}^{N_{-1}}\})$. Для каждой $f : \mathbb{X} \rightarrow \{-1, 1\}$ находится оценка f_{PA} , которая используется на оценки $\text{Err}(f)$. Для этого используется K -свёрточная кросс-валидация или по-другому поднаборный алгоритм усреднения. Для фиксированного $K \in \mathbb{N}$ и любого $y \in \{-1, 1\}$, вводится разбиение $\{j_y^1, \dots, j_y^{N_y}\}$ на K поднаборов $S_k^y(N_y, \omega)$, $k = 1, \dots, K$:

$$S_k^y(N_y, \omega) := \left\{ j_y^i(\omega) : i \in \left\{ (k-1) \left\lceil \frac{N_y}{K} \right\rceil + 1, \dots, k \left\lceil \frac{N_y}{K} \right\rceil \mathbb{I}\{k < K\} + N_y \mathbb{I}\{k = K\} \right\} \right\} \quad (9)$$

Вводится $S_k(N, \omega) = S_k^1(N_1, \omega) \cup S_k^{-1}(N_{-1}, \omega)$ и

$$\widehat{\text{Err}}_K(f_{PA}, \zeta_N, \hat{P}_N) := \frac{2}{K} \sum_{y \in \{-1, 1\}} \sum_{k=1}^K \sum_{j \in S_k^y(N_y)} \frac{\hat{\psi}\left(y, \zeta_N(\overline{S_k(N)}), \hat{P}_N\right) \mathbb{I}\left\{f_{PA}^j(N, k) \neq y\right\} \hat{P}_N^y}{\#S_k^y(N_y)} \quad (10)$$

где $f_{PA}^j(N, k) := f_{PA}(X^j, \zeta_N(\overline{S_k(N)}), \hat{P}_N)$ и $\hat{\psi}\left(y, \zeta_N(\overline{S_k(N)}), \hat{P}_N\right)$ - оценка величины $\psi(y)$, $y \in \{-1, 1\}$, построенная на наблюдениях $\zeta_N(\overline{S_k(N)})$ и \hat{P}_N ; $\overline{S_k(N)} = \{j_1^1, \dots, j_1^{N_1}\} \cup \{j_{-1}^1, \dots, j_{-1}^{N_{-1}}\} \setminus S_k(N)$, $\#$ - мощность множества.

Предполагается, что для любого $k = 1, \dots, N$,

$$\widehat{\psi} \left(y, \zeta_N \left(\overline{\mathcal{S}_k(N)} \right), \widehat{\mathbf{P}}_N \right) \rightarrow \psi(y) \text{ п.н., } N \rightarrow \infty, \quad y \in \{-1, 1\} \quad (11)$$

Определяется $L(x) = \psi(1)\mathbf{P}(X = x, Y = 1) - \psi(-1)\mathbf{P}(X = x, Y = -1), x \in \mathbb{X}$.

Теорема 4

$\zeta_N, \psi, \hat{\psi}, f : \mathbb{X} \rightarrow \{-1, 1\}$ и f_{PA} - определены как ранее. Пусть существует не пустой $U \subset \mathbb{X}$ такой, что для каждого $x \in U$ и $k = 1, \dots, K$, верно:

$$f_{PA} \left(x, \zeta_N \left(\overline{S_k(N)} \right), \hat{P}_N \right) \rightarrow f(x) \text{ п.н., } N \rightarrow \infty \quad (12)$$

Тогда для каждого $a \in (0, 1)$ (с $N_1 = \max\{[aN], 1\}, N_{-1} = N - N_1$),

$$\widehat{\text{Err}}_K \left(f_{PA}, \zeta_N, \hat{P}_N \right) \rightarrow \text{Err}(f) \text{ a.s., } N \rightarrow \infty \quad (13)$$

равносильно

$$\sum_{k=1}^K \sum_{y \in \{-1, 1\}} \sum_{x \in \mathbb{X}_y} y \mathbb{I} \left\{ f_{PA} \left(x, \zeta_N \left(\overline{S_k(N)} \right), \hat{P}_N \right) = -y \right\} L(x) \rightarrow 0 \text{ a.s., } N \rightarrow \infty \quad (14)$$

где

$$\mathbb{X}_y = (\mathbb{X} \setminus U) \cap \{x \in \mathbb{X} : f(x) = y\}, \quad y \in \{-1, 1\} \quad (15)$$

Для доказательства сперва необходимо воспользоваться леммой

Пусть $(X, Y), (X^1, Y^1), (X^2, Y^2), \dots$ н.о.р. на вероятностном пространстве (Ω, \mathcal{F}, P) . Для каждого $\omega \in \Omega$ рассматривается $Y^1(\omega), Y^2(\omega), \dots$ и вводятся индексы $1 \leq j_{-1}^1(\omega) < j_{-1}^2(\omega) < \dots$, для которых $Y^{j_{-1}^k(\omega)}(\omega) = -1, k \in \mathbb{N}$. Аналогично для всех $Y^i(\omega)$ со значениями 1 как $\left\{ Y^{j_1^m(\omega)}(\omega) \right\}_{m \in \mathbb{N}}$, где $1 \leq j_1^1(\omega) < j_1^2(\omega) < \dots$. Используя распределение Бернулли с вероятностью успеха p , вводится отрицательная биномиальная величина $U_{r,p}$ - число успехов до r неудач где $r \in \mathbb{N}$ ($U_{r,p} \sim NB(r, p)$).

Тогда

$$P(U_{r,p} = k) = \binom{k+r-1}{k} p^k (1-p)^r, \quad k = 0, 1, \dots$$

$\{Y^i = 1\}$ и $\{Y^i = -1\}$ рассматриваются как успехи и неудачи (с вероятностью $p = P(Y = 1)$ успеха), тогда j_{-1}^r также распределён как $U_{r,p} + r$. Следовательно,

$$P(j_{-1}^r = m) = \begin{cases} \binom{m-1}{m-r} p^{m-r} (1-p)^r, & m = r, r+1, \dots \\ 0, & m = 1, \dots, r-1 \end{cases} \quad (4)$$

j_1^r распределён как $U_{r,1-p} + r$, где $U_{r,1-p} \sim NB(r, 1-p)$. j_1^r распределен также как $G_p^1 + \dots + G_p^r$, где G_p^1, \dots, G_p^r - независимые величины с геометрическими распределениями и с параметром p ($P(G_p^1 = k) = p(1-p)^{k-1}, k = 1, 2, \dots$). Тогда $Ej_1^r = \frac{r}{p}$ и $j_1^r < \infty$ п.н. для любого $r \in \mathbb{N}$ ($Ej_{-1}^r = \frac{r}{1-p}$ для $r \in \mathbb{N}$).

Определяются $Z^k := X_{j_1^k}^k$ для каждого $k \in \mathbb{N}$. \mathcal{B} - семейство всех подмножеств \mathbb{X} .

Лемма 2

Для каждого $m \in \mathbb{N}$, случайные величины Z^1, \dots, Z^m - независимы и распределены как X при $Y = 1$ ($X \mid Y = 1$), т.е., для каждого $B \in \mathcal{B}$ и $k = 1, \dots, m$,

$$\mathbb{P}(Z^k \in B) = \mathbb{P}(X \in B \mid Y = 1)$$

Теперь затронем тему логистической регрессии. Как упоминалось ранее:

$$f^*(x) = \begin{cases} 1, & p^*(x) > 1/2 \\ -1, & \text{otherwise} \end{cases}$$

Случаи $p^*(x) \in \{0, 1\}$ исключаются. Для оценки $p^*(x)$ используется

$$q^*(x) = \lambda(p^*(x)) \tag{24}$$

где $\lambda(z) = \log(z/(1-z))$, $z \in (0, 1)$, обратная логистическая функция. Логистическая функция - $\Lambda(t) = (1 + e^{-t})^{-1}$, $t \in \mathbf{R}$.

$p^*(x) \in (0, 1)$ для $x \in \mathcal{X}$, $q^*(x)$ может принимать любые действительные значения. Рассматривается класс G всех вещественно значных функций на тернарных значениях x_1, \dots, x_n . $M \subset G$ - называется моделью зависимости отклика и объясняющих переменных. Введем

$$\hat{\psi}(y, \xi(S)) = \frac{1}{4\hat{P}_S(Y = y)}, y \in \{-1, 1\}$$

$\hat{P}_S(Y = y)$ - оценка .

Введем функцию оценки регрессии:

$$L(h, \xi(S)) = \frac{1}{\#S} \sum_{j \in S} \varphi \left(-Y^j h \left(X^j \right) \right) \hat{\psi} \left(Y^j, \xi(S) \right) \quad (25)$$

$\varphi(t) = \log_2 (1 + e^t)$ для $t \in \mathbf{R}$, и $h \in M$. В отличие от других работ в этой статье авторы рассматривают нормировку, т.е. наблюдения с весами от отношения случаев в подвыборке $\xi(S)$. $\arg \min_{h \in M} L(h, \xi(S))$ равен $\arg \max_{h \in M}$ для

$$\frac{1}{\#S} \sum_{j \in S} \left(\frac{I \{ Y^j = 1 \}}{2\hat{P}_S(Y = 1)} \log \gamma_j + \frac{I \{ Y^j = -1 \}}{2\hat{P}_S(Y = -1)} \log (1 - \gamma_j) \right)$$

где $\gamma_j = \Lambda(h(X^j))$. Минимизация эквивалентна нормированной оценке максимума правдоподобия q^* .

Следующая теорема о строгой состоятельности оценки при корректной модели, т.е. $q^* \in M$. Введём $h(\cdot, \xi(S)) := \arg \min_{q \in M} L(q, \xi(S))$.

Теорема 5

Пусть $q^* \in M$, $h_0 \equiv 0$ принадлежит M и

$$\min_{(x,y) \in X \times \{-1,1\}} P(X = x \mid Y = y) > 0$$

Рассмотрим $W_N \subset \{1, \dots, N\}$ и множество

$$h_N(\cdot) = h(\cdot, \xi(W_N))$$

Тогда $h_N(x) \rightarrow q^*(x)$ п.н. для всех $x \in \mathcal{X}$, когда $\#W_N \rightarrow \infty$. Более того,

$$\text{Err}_K(f_{PA}(\cdot, \xi), \xi) \rightarrow \text{Err}(f^*) \text{ a.s., } N \rightarrow \infty$$

где $f_{PA}(\cdot, \xi) = 2I\{\Lambda(h(\cdot, \xi)) > 1/2\} - 1$.

□ Сперва показывается

$$h_N(x) < \varphi^{-1} \left(\frac{4}{\mathbb{P}(X = x \mid Y = -1)} \right) \text{ п.н.}$$

для всех $x \in \mathcal{X}$ и всех $N > N_1 = N_1(\omega)$. Положим $l_N = \#W_N$. По определению

$$\begin{aligned} \frac{\varphi(h_N(x))}{4l_N} \sum_{j \in W_N} \frac{I\{X^j = x, Y^j = -1\}}{\hat{\mathbb{P}}_{W_N}(Y = -1)} &\leq L(h_N, \xi(W_N)) \\ &\leq L(0, \xi(W_N)) = \frac{1}{4l_N} \sum_{(j,y) \in W_N \times \{-1,1\}} \frac{I\{Y^j = y\}}{\hat{\mathbb{P}}_N(Y = y)} \end{aligned}$$

Пользуясь УЗБЧН, получается сходимость п.н.

$$\begin{aligned} \max_{x \in \mathcal{X}} \left| \frac{1}{l_N} \sum_{j \in W_N} \frac{I\{X^j = x, Y^j = -1\}}{\hat{\mathbb{P}}_{W_N}(Y = -1)} \right. \\ \left. - \mathbb{P}(X = x \mid Y = -1) \right| \rightarrow 0 \end{aligned}$$

Очевидно

$$\frac{1}{l_N} \sum_{(j,y) \in W_N \times \{-1,1\}} \frac{I\{Y^j = y\}}{\hat{\mathbf{P}}_{W_N}(Y = y)} = 2.$$

Эти выражения приводят к желаемой оценке $h_N(x)$. Аналогично

$$h_N(x) > -\varphi^{-1} \left(\frac{4}{\mathbf{P}(X = x \mid Y = 1)} \right)$$

для каждого $x \in \mathcal{X}$ и всех $N > N_2 = N_2(\omega)$. Следовательно, $h_N \in M_C := M \cap \{h : \|h\|_\infty \leq C\}$ для $N > \max(N_1, N_2)$, здесь $\|h\|_\infty = \max_{x \in X} |h(x)|$ и

$$C = \max_{(x,y) \in X \times \{-1,1\}} \varphi^{-1} \left(\frac{4}{\mathbf{P}(X = x \mid Y = y)} \right)$$

Если $h \in M_C$, тогда $|L(h, \xi(W_N)) - E\varphi(-Yh(X))\psi(Y)|$ меньше

$$\begin{aligned} & \frac{\varphi(\|h\|_\infty)}{2} \sum_{y \in \{-1, 1\}} \left(\left| \frac{1}{\hat{P}_{W_N}(Y=y)} - \frac{1}{P(Y=y)} \right| \right. \\ & \left. + \sum_{x \in X} \left| \sum_{j \in W_N} \frac{I\{X^j = x, Y^j = y\}}{l_N P(Y=y)} - \frac{P(X^j = x, Y^j = y)}{P(Y=y)} \right| \right) \end{aligned}$$

Из УЗБЧН следует

$$L(h, \xi(W_N)) - E\varphi(-Yh(X))\psi(Y) \rightarrow 0 \text{ a.s.} \quad (26)$$

равномерно по $\{h : \|h\|_\infty \leq C\}$.

Также

$$\begin{aligned} 2E\varphi(-Yh(X))\psi(Y) &= E\varphi(-Y^*h(X^*)) \\ &= -E \log_2 \left(\Lambda(h(X^*))^{I\{Y^*=1\}} (1 - \Lambda(h(X^*)))^{1-I\{Y^*=1\}} \right) \end{aligned}$$

Из (условного) информационного неравенства,

$$\mathbb{E} \log_2 \left(\Lambda(h(X^*))^{I\{Y^*=1\}} (1 - \Lambda(h(X^*)))^{1-I\{Y^*=1\}} \right)$$

достигает своего максимума по функциям h только на q^* , представленный в (24). В условиях теоремы, $q^* \in M$. Следовательно, по определению h_N и q^* получается

$$\begin{aligned} L(h_N, \xi(W_N)) &\leq L(q^*, \xi(W_N)) \\ \mathbb{E} \varphi(-Y^* h(X^*))|_{h=h_N} &\geq \mathbb{E} \varphi(-Y^* q^*(X^*)) \end{aligned}$$

Из (26) и УЗБЧН

$$\begin{aligned} L(h_N, \xi(W_N)) - \frac{1}{2} \mathbb{E} \varphi(-Y^* h(X^*)) \Big|_{h=h_N} &\rightarrow 0 \text{ a.s.}, \\ L(q^*, \xi(W_N)) - \frac{1}{2} \mathbb{E} \varphi(-Y^* q^*(X^*)) &\rightarrow 0 \text{ a.s.} \end{aligned}$$

Тогда

$$\mathbb{E} \varphi(-Y^* h(X^*))|_{h=h_N} \rightarrow \mathbb{E} \varphi(-Y^* q^*(X^*)) \text{ a.s.}$$

Это возможно только когда $h_N(x) \rightarrow q^*(x)$ п.н. для каждого $x \in X$. Действительно, для почти всех $\omega \in \Omega$, можно всегда выбрать подпоследовательность $h_{N_k} = h_{N_k(\omega)}(\cdot, \omega)$ сходящаяся к некоторой функции $\mu = \mu(\cdot, \omega)$ и

$$E\varphi(-Y^* \mu(X^*)) = E\varphi(-Y^* q^*(X^*))$$

Тогда из информационного неравенства $\mu(\cdot, \omega) = q(\cdot)$.

Для доказательства второй части Теоремы 3 нужно заметить, что $f_{PA}(x, \xi) = 2I\{\Lambda(h(x, \xi)) > 1/2\} - 1$ сходится п.н. к $f^*(x) = 2I\{\Lambda(q^*(x)) > 1/2\} - 1$ для почти всех $x \in U$, где

$$\begin{aligned} U &:= \{x \in X : p^*(x) \neq 1/2\} \\ &= \{x \in X : P(Y = 1 | X = x) \neq P(Y = 1)\} \end{aligned}$$

Остаётся только воспользоваться Теоремой 1 и замечанием к ней. \square

У меня есть некоторые вопросы про нефиксированное число признаков.

S - множество номеров всех признаков, m - количество фиксированных признаков, n - число наблюдений, S_m - семейство множеств номеров релевантных признаков размера m . Q_m — все подмножества S размера m . Цель:

$$\operatorname{argmax}_{L \in Q_m} I(X_L, Y)$$

У Алексея Кожевина при дополнительных ограничениях доказывается в статье 2020 года:

$$\hat{S}_{n,k}(w) = \operatorname{argmax}_{L \in Q_m} \hat{I}_{n,k,L}(w)$$

$$P(\hat{S}_{n,k} \subset S_m) \longrightarrow 1$$

$$n \rightarrow \infty$$

У меня есть вопрос. Пусть введём максимум по всем наборам:

$$\hat{S}_{n,k}(w) = \operatorname{argmax}_{L \subset S} \hat{I}_{n,k,L}(w)$$

$$u_{n,k} = \min_{M \in \hat{S}_{n,k}(w)} |M|$$

$$\hat{S}_{n,k}^{min}(w) = \{M \in \hat{S}_{n,k}(w), |M| = u_{n,k}\}$$

Пусть \mathbb{S} - семейство множеств релевантных признаков. Пусть \mathbb{S}^{min} - семейство наименьших по размеру множеств релевантных признаков. Следующие предположения выглядят правдоподобно при некоторых ограничениях:

$$P(\hat{S}_{n,k}(w) \subset \mathbb{S}) \longrightarrow 1$$

$$P(\hat{S}_{n,k}^{min}(w) \subset \mathbb{S}^{min}) \longrightarrow 1$$

$$n \rightarrow \infty$$

Насколько вообще такие предположения могут быть интересными? Я полагаю, что слишком большой перебор признаков, если оценку двигать по всем возможным комбинациям признаков.

Перейду к аналогичному вопросу, но уже с последовательным выбором признаков.
Главная цель:

$$\operatorname{argmax}_{L \in Q_m} I(X_L, Y)$$

Последовательный выбор признаков на $(k+1)$ -том шагу:

$$j_{k+1} = \operatorname{argmax}_{j \in S_k^c} (I(X_{S_k \cup \{j\}}, Y) - I(X_{S_k}, Y)) = \operatorname{argmax}_{j \in S_k^c} I(X_j, Y | X_{S_k})$$

$$S_{k+1} = S_k \cup \{j_{k+1}\}$$

$$t^{stop} := \operatorname{argmin}_{1 \leq k \leq p} (I(X_j, Y | S_k) = 0)$$

Оценочный stopping rule:

$$\hat{t}_n^{stop} = \operatorname{argmin}_{1 \leq k \leq p} (\hat{I}(X_j, Y | S_k) = 0)$$

Насколько сильны и неправдоподобны следующие предположения? А именно:

1. $S_{t^{stop}} \subset \mathbb{S}$
2. $S_{t^{stop}} \subset \mathbb{S}^{min}$
3. $P(S_{\hat{t}_n^{stop}} \subset \mathbb{S}) \longrightarrow 1, \quad n \rightarrow \infty$
4. $P(S_{\hat{t}_n^{stop}} \subset \mathbb{S}^{min}) \longrightarrow 1, \quad n \rightarrow \infty$

Аналогично для функции ошибки в статье для бинарного отклика:

$$\hat{Err}_K(f_{t_n^{stop}}) \longrightarrow Err(f_M)$$

где $M \in \mathbb{S}^{min}$



Ссылки

- 1 Stopping rules for mutual information-based feature selection, Jan Mielniczuk, Pawe l Teisseyre
<https://www.sciencedirect.com/science/article/abs/pii/S0925231219307544>
- 2 Statistical Methods of SNP Data Analysis and Applications, Alexander Bulinski, Oleg Butkovsky, Victor Sadovnichy, Alexey Shashkin, Pavel Yaskov, Alexander Balatskiy , Larisa Samokhodskaya, Vsevolod Tkachuk
<https://www.scirp.org/journal/paperinformation?paperid=16881>
- 3 New version of the MDR method for stratified samples, Alexander Bulinski, Alexey Kozhevin
https://www.mathnet.ru/php/archive.phtml?wshow=paper&jrnid=soic&paperid=1&option_lang=rus
- 4 Estimation of nonbinary random response, Bulinski A.V., Rakitko A.S.
<https://link.springer.com/article/10.1134/S1064562414020306>
- 5 MDR method for nonbinary response variable, Bulinski A., Rakitko A.
<https://doi.org/10.1016/j.jmva.2014.11.008>
- 6 Forward Selection of Relevant Factors by Means of MDR-EFE Method, Alexander Bulinski
<https://www.mdpi.com/2227-7390/12/6/831>

