

Новый критерий оценки качества классификационных моделей

Максим Гончаров
Декабрь 2013

В статье предлагается новый критерий оценки качества моделей классификации: увеличение вероятности совпадения значения целевой переменной со значением прогноза по сравнению с вероятностью их совпадения в случае независимости этих переменных. Получена формула для вычисления асимптотического доверительного интервала для этого критерия. Также получен статистический тест для проверки нулевой гипотезы, заключающейся в предположении, что вероятность совпадения целевой переменной и ее прогноза не превышает вероятности их совпадения если бы они были независимы. Получена асимптотическая значимость этого теста.

Пусть $(X_1, Y_1, Z_1), (X_2, Y_2, Z_2), \dots$ – последовательность одинаково распределенных независимых случайных величин, где множества значений Y_i и Z_i конечны и совпадают, т.е. $Y_i, Z_i : \Omega \rightarrow \{1, \dots, s\}$. Мы будем считать, что Y_i – это прогнозируемая целевая переменная, а Z_i – ее прогноз, зависящий только от набора входных переменных (предикторов) X_i , т.е. $Z_i = f(X_i)$. Далее, нас не будут интересовать предикторы X_i и метод получения прогноза f . Наша задача будет заключаться в оценке качества прогноза.

Введем обозначения:

$$\begin{aligned} p_k &:= P(Y_i = k) \\ q_k &:= P(Z_i = k) \\ r_{kl} &:= P(Y_i = k, Z_i = l) \end{aligned} \tag{1}$$

Определим случайные величины:

$$\begin{aligned} U_i^{(k)} &:= 1_{\{Y_i=k\}} \\ V_i^{(k)} &:= 1_{\{Z_i=k\}} \\ W_i^{(k)} &:= 1_{\{Y_i=k, Z_i=k\}} \end{aligned} \tag{2}$$

Из (2) и (1) получаем для математических ожиданий:

$$\begin{aligned} EU_i^{(k)} &= E1_{\{Y_i=k\}} = P(Y_i = k) = p_k \\ EV_i^{(k)} &= E1_{\{Z_i=k\}} = P(Z_i = k) = q_k \\ EW_i^{(k)} &= E1_{\{Y_i=k, Z_i=k\}} = P(Y_i = k, Z_i = k) = r_{kk} \end{aligned} \tag{3}$$

Далее получим выражения для дисперсий и ковариаций:

$$\begin{aligned}
\text{Var}\left(U_i^{(k)}\right) &= E1_{\{Y_i=k\}}^2 - \left(E1_{\{Y_i=k\}}\right)^2 = p_k - p_k^2 \\
\text{Cov}\left(U_i^{(k)}, U_i^{(l)}\right) &= E1_{\{Y_i=k\}}1_{\{Y_i=l\}} - E1_{\{Y_i=k\}}E1_{\{Y_i=l\}} = -p_k p_l, k \neq l \\
\text{Var}\left(V_i^{(k)}\right) &= E1_{\{Z_i=k\}}^2 - \left(E1_{\{Z_i=k\}}\right)^2 = q_k - q_k^2 \\
\text{Cov}\left(V_i^{(k)}, V_i^{(l)}\right) &= E1_{\{Z_i=k\}}1_{\{Z_i=l\}} - E1_{\{Z_i=k\}}E1_{\{Z_i=l\}} = -q_k q_l, k \neq l \\
\text{Var}\left(W_i^{(k)}\right) &= E1_{\{Y_i=k, Z_i=k\}}^2 - \left(E1_{\{Y_i=k, Z_i=k\}}\right)^2 = r_{kk} - r_{kk}^2 \\
\text{Cov}\left(W_i^{(k)}, W_i^{(l)}\right) &= E1_{\{Y_i=k, Z_i=k\}}1_{\{Y_i=l, Z_i=l\}} - E1_{\{Y_i=k, Z_i=k\}}E1_{\{Y_i=l, Z_i=l\}} = -r_{kk}r_{ll}, k \neq l \\
\text{Cov}\left(U_i^{(k)}, V_i^{(l)}\right) &= E1_{\{Y_i=k\}}1_{\{Z_i=l\}} - E1_{\{Y_i=k\}}E1_{\{Z_i=l\}} = r_{kl} - p_k q_l \\
\text{Cov}\left(U_i^{(k)}, W_i^{(l)}\right) &= E1_{\{Y_i=k\}}1_{\{Y_i=l, Z_i=l\}} - E1_{\{Y_i=k\}}E1_{\{Y_i=l, Z_i=l\}} = -p_k r_{ll}, k \neq l \\
\text{Cov}\left(U_i^{(k)}, W_i^{(k)}\right) &= E1_{\{Y_i=k\}}1_{\{Y_i=k, Z_i=k\}} - E1_{\{Y_i=k\}}E1_{\{Y_i=k, Z_i=k\}} = r_{kk} - p_k r_{kk} \\
\text{Cov}\left(V_i^{(k)}, W_i^{(l)}\right) &= E1_{\{Z_i=k\}}1_{\{Y_i=l, Z_i=l\}} - E1_{\{Z_i=k\}}E1_{\{Y_i=l, Z_i=l\}} = -q_k r_{ll}, k \neq l \\
\text{Cov}\left(V_i^{(k)}, W_i^{(k)}\right) &= E1_{\{Z_i=k\}}1_{\{Y_i=k, Z_i=k\}} - E1_{\{Z_i=k\}}E1_{\{Y_i=k, Z_i=k\}} = r_{kk} - q_k r_{kk}
\end{aligned}$$

Рассмотрим последовательность случайных величин при $i=1, 2, \dots$

$$T_i := \begin{pmatrix} U_i^{(1)} \\ \vdots \\ U_i^{(s)} \\ V_i^{(1)} \\ \vdots \\ V_i^{(s)} \\ W_i^{(1)} \\ \vdots \\ W_i^{(s)} \end{pmatrix} \tag{6}$$

Так как каждая T_i является измеримой функцией от (Y_i, Z_i) (см. формулу (2)), а последовательность $(Y_1, Z_1), (Y_2, Z_2), \dots$ представляет собой последовательность одинаково распределенных независимых случайных величин, то и последовательность T_1, T_2, \dots также является последовательностью одинаково распределенных независимых случайных величин с конечными вторыми моментами, определяемыми формулами (4) и (5). Таким образом мы можем применить центральную предельную теорему:

$$\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n T_i - ET_1 \right) \xrightarrow{d} N(0, \Sigma), \tag{7}$$

где $N(0, \Sigma)$ многомерно нормально распределенная случайная величина с нулевым вектором математического ожидания и матрицей ковариаций $\Sigma = \text{Var}(T_1)$.

Далее, с учетом (3):

$$ET_1 = \begin{pmatrix} EU_1^{(1)} \\ \vdots \\ EU_1^{(s)} \\ EV_1^{(1)} \\ \vdots \\ EV_1^{(s)} \\ EW_1^{(1)} \\ \vdots \\ EW_1^{(s)} \end{pmatrix} = \begin{pmatrix} p_1 \\ \vdots \\ p_s \\ q_1 \\ \vdots \\ q_s \\ r_{11} \\ \vdots \\ r_{ss} \end{pmatrix} \quad (8)$$

И с учетом (4) и (5):

$$\Sigma = \text{Var}(T_1) = \begin{pmatrix} \text{Var}(U_1^{(1)}) & \cdots & \text{Cov}(U_1^{(1)}, U_1^{(s)}) & \text{Cov}(U_1^{(1)}, V_1^{(1)}) & \cdots & \text{Cov}(U_1^{(1)}, V_1^{(s)}) & \text{Cov}(U_1^{(1)}, W_1^{(1)}) & \cdots & \text{Cov}(U_1^{(1)}, W_1^{(s)}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(U_1^{(s)}, U_1^{(1)}) & \cdots & \text{Var}(U_1^{(s)}) & \text{Cov}(U_1^{(s)}, V_1^{(1)}) & \cdots & \text{Cov}(U_1^{(s)}, V_1^{(s)}) & \text{Cov}(U_1^{(s)}, W_1^{(1)}) & \cdots & \text{Cov}(U_1^{(s)}, W_1^{(s)}) \\ \text{Cov}(V_1^{(1)}, U_1^{(1)}) & \cdots & \text{Cov}(V_1^{(1)}, U_1^{(s)}) & \text{Var}(V_1^{(1)}) & \cdots & \text{Cov}(V_1^{(1)}, V_1^{(s)}) & \text{Cov}(V_1^{(1)}, W_1^{(1)}) & \cdots & \text{Cov}(V_1^{(1)}, W_1^{(s)}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(V_1^{(s)}, U_1^{(1)}) & \cdots & \text{Cov}(V_1^{(s)}, U_1^{(s)}) & \text{Cov}(V_1^{(s)}, V_1^{(1)}) & \cdots & \text{Var}(V_1^{(s)}) & \text{Cov}(V_1^{(s)}, W_1^{(1)}) & \cdots & \text{Cov}(V_1^{(s)}, W_1^{(s)}) \\ \text{Cov}(W_1^{(1)}, U_1^{(1)}) & \cdots & \text{Cov}(W_1^{(1)}, U_1^{(s)}) & \text{Cov}(W_1^{(1)}, V_1^{(1)}) & \cdots & \text{Cov}(W_1^{(1)}, V_1^{(s)}) & \text{Var}(W_1^{(1)}) & \cdots & \text{Cov}(W_1^{(1)}, W_1^{(s)}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(W_1^{(s)}, U_1^{(1)}) & \cdots & \text{Cov}(W_1^{(s)}, U_1^{(s)}) & \text{Cov}(W_1^{(s)}, V_1^{(1)}) & \cdots & \text{Cov}(W_1^{(s)}, V_1^{(s)}) & \text{Cov}(W_1^{(s)}, W_1^{(1)}) & \cdots & \text{Var}(W_1^{(s)}) \end{pmatrix} = \begin{pmatrix} p_1 - p_1^2 & \cdots & -p_1 p_s & r_{11} - p_1 q_1 & \cdots & r_{1s} - p_1 q_s & r_{11} - p_1 r_{11} & \cdots & -p_1 r_{ss} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ -p_s p_1 & \cdots & p_s - p_s^2 & r_{s1} - p_s q_1 & \cdots & r_{ss} - p_s q_s & -p_s r_{11} & \cdots & r_{ss} - p_s r_{ss} \\ r_{11} - p_1 q_1 & \cdots & r_{s1} - p_s q_1 & q_1 - q_1^2 & \cdots & -q_1 q_s & r_{11} - q_1 r_{11} & \cdots & -q_1 r_{ss} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ r_{1s} - p_1 q_s & \cdots & r_{ss} - p_s q_s & -q_s q_1 & \cdots & q_s - q_s^2 & -q_s r_{11} & \cdots & r_{ss} - q_s r_{ss} \\ r_{11} - p_1 r_{11} & \cdots & -p_s r_{11} & r_{11} - q_1 r_{11} & \cdots & -q_s r_{11} & r_{11} - r_{11}^2 & \cdots & -r_{11} r_{ss} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ -p_1 r_{ss} & \cdots & r_{ss} - p_s r_{ss} & -q_1 r_{ss} & \cdots & r_{ss} - q_s r_{ss} & -r_{ss} r_{11} & \cdots & r_{ss} - r_{ss}^2 \end{pmatrix} \quad (9)$$

Далее рассмотрим функцию:

$$\varphi: \mathbb{R}^{3s} \rightarrow \mathbb{R},$$

$$\varphi(u_1, \dots, u_s, v_1, \dots, v_s, w_1, \dots, w_s) := \sum_{k=1}^s w_k - \sum_{k=1}^s u_k v_k \quad (10)$$

Очевидно, φ дифференцируемая функция, для которой справедливо:

$$\frac{\partial \varphi}{\partial u_j}(p_1, \dots, p_s, q_1, \dots, q_s, r_{11}, \dots, r_{ss}) = -q_j$$

$$\frac{\partial \varphi}{\partial v_j}(p_1, \dots, p_s, q_1, \dots, q_s, r_{11}, \dots, r_{ss}) = -p_j$$

$$\frac{\partial \varphi}{\partial w_j}(p_1, \dots, p_s, q_1, \dots, q_s, r_{11}, \dots, r_{ss}) = 1$$

Отсюда следует:

$$D\varphi(p_1, \dots, p_s, q_1, \dots, q_s, r_{11}, \dots, r_{ss}) = D\varphi(ET_1) = (-q_1, \dots, -q_s, -p_1, \dots, -p_s, 1, \dots, 1) \quad (11)$$

Используя дельта-метод, получаем из (7):

$$\sqrt{n} \left(\varphi \left(\frac{1}{n} \sum_{i=1}^n T_i \right) - \varphi(ET_1) \right) \xrightarrow{d} D\varphi(ET_1) \cdot N(0, \Sigma) = N \left(0, D\varphi(ET_1) \cdot \Sigma \cdot [D\varphi(ET_1)]' \right) \quad (12)$$

Далее,

$$\begin{aligned} \hat{T}_n &:= \varphi \left(\frac{1}{n} \sum_{i=1}^n T_i \right) = \varphi \left(\begin{array}{c} \frac{1}{n} \sum_{i=1}^n U_i^{(1)} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n U_i^{(s)} \\ \frac{1}{n} \sum_{i=1}^n V_i^{(1)} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n V_i^{(s)} \\ \frac{1}{n} \sum_{i=1}^n W_i^{(1)} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n W_i^{(s)} \end{array} \right) = \sum_{k=1}^s \frac{1}{n} \sum_{i=1}^n W_i^{(k)} - \sum_{k=1}^s \left(\frac{1}{n} \sum_{i=1}^n U_i^{(k)} \right) \left(\frac{1}{n} \sum_{j=1}^n V_j^{(k)} \right) = \\ &= \sum_{k=1}^s \sum_{i=1}^n \frac{1}{n} 1_{\{Y_i=k, Z_i=k\}} - \sum_{k=1}^s \left(\frac{1}{n} \sum_{i=1}^n 1_{\{Y_i=k\}} \right) \left(\frac{1}{n} \sum_{j=1}^n 1_{\{Z_j=k\}} \right) = \sum_{k=1}^s \frac{v_{kk}}{n} - \sum_{k=1}^s \frac{\mu_k}{n} \frac{\eta_k}{n} \end{aligned} \quad (13)$$

Здесь

$$v_{kk} := \sum_{i=1}^n 1_{\{Y_i=k, Z_i=k\}} \text{ — число записей, у которых } Y \text{ и } Z \text{ равны } k,$$

$$\mu_k := \sum_{i=1}^n 1_{\{Y_i=k\}} \text{ — число записей, у которых } Y \text{ равно } k,$$

$$\eta_k := \sum_{i=1}^n 1_{\{Z_i=k\}} \text{ — число записей, у которых } Z \text{ равно } k.$$

$$\theta := \varphi(ET_i) = \varphi \begin{pmatrix} p_1 \\ \vdots \\ p_s \\ q_1 \\ \vdots \\ q_s \\ r_{11} \\ \vdots \\ r_{ss} \end{pmatrix} = \sum_{k=1}^s r_{kk} - \sum_{k=1}^s p_k q_k =$$

$$= \sum_{k=1}^s \mathbb{P}(Y_i = k, Z_i = k) - \sum_{k=1}^s \mathbb{P}(Y_i = k) \mathbb{P}(Z_i = k) =$$

$$= \mathbb{P}(Y_i = Z_i) - \sum_{k=1}^s \mathbb{P}(Y_i = k) \mathbb{P}(Z_i = k) = \mathbb{P}(Y_i = Z_i) - \mathbb{P}_{\perp}(Y_i = Z_i)$$
(14)

$$\text{Здесь } \mathbb{P}_{\perp}(Y_i = Z_i) := \sum_{k=1}^s \mathbb{P}(Y_i = k) \mathbb{P}(Z_i = k).$$
(15)

Эта величина равна вероятности совпадения значений Y и Z если бы они были независимы. Действительно, в случае независимости Y и Z получаем:

$$\sum_{k=1}^s \mathbb{P}(Y_i = k) \mathbb{P}(Z_i = k) = \sum_{k=1}^s \mathbb{P}(Y_i = k, Z_i = k) = \mathbb{P}(Y_i = Z_i)$$

Итак, согласно (14), $\varphi(ET_i)$ есть разность вероятности совпадения Y и Z и вероятности их совпадения, если бы они были независимы.

Получим теперь выражение дисперсии для предельного нормального распределения в (12):

$$S := D\varphi(ET_1) \cdot \Sigma \cdot [D\varphi(ET_1)]^t =$$

$$= (-q_1, \dots, -q_s, -p_1, \dots, -p_s, 1, \dots, 1) \cdot$$

$$\begin{pmatrix} p_1 - p_1^2 & \cdots & -p_1 p_s & r_{11} - p_1 q_1 & \cdots & r_{1s} - p_1 q_s & r_{11} - p_1 r_{11} & \cdots & -p_1 r_{ss} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ -p_s p_1 & \cdots & p_s - p_s^2 & r_{s1} - p_s q_1 & \cdots & r_{ss} - p_s q_s & -p_s r_{11} & \cdots & r_{ss} - p_s r_{ss} \\ r_{11} - p_1 q_1 & \cdots & r_{s1} - p_s q_1 & q_1 - q_1^2 & \cdots & -q_1 q_s & r_{11} - q_1 r_{11} & \cdots & -q_1 r_{ss} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ r_{1s} - p_1 q_s & \cdots & r_{ss} - p_s q_s & -q_s q_1 & \cdots & q_s - q_s^2 & -q_s r_{11} & \cdots & r_{ss} - q_s r_{ss} \\ r_{11} - p_1 r_{11} & \cdots & -p_s r_{11} & r_{11} - q_1 r_{11} & \cdots & -q_s r_{11} & r_{11} - r_{11}^2 & \cdots & -r_{11} r_{ss} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ -p_1 r_{ss} & \cdots & r_{ss} - p_s r_{ss} & -q_1 r_{ss} & \cdots & r_{ss} - q_s r_{ss} & -r_{ss} r_{11} & \cdots & r_{ss} - r_{ss}^2 \end{pmatrix} \cdot$$

$$\cdot (-q_1, \dots, -q_s, -p_1, \dots, -p_s, 1, \dots, 1)^t = \sum_{k=1}^s (p_k^2 q_k + p_k q_k^2 - 4p_k^2 q_k^2 + 6p_k q_k r_{kk} + r_{kk} - r_{kk}^2 - 2p_k r_{kk} - 2q_k r_{kk}) -$$

$$- \sum_{k=1}^s \sum_{l=k+1}^s (2r_{kk} r_{ll} + 8p_k p_l q_k q_l) + \sum_{k=1}^s 4p_k q_k \sum_{l=1, l \neq k}^s r_{ll} + \sum_{k=1}^s \sum_{l=1, l \neq k}^s 2p_k q_l r_{lk}$$

Итак,

$$\begin{aligned}
S &= D\varphi(ET_1) \cdot \Sigma \cdot [D\varphi(ET_1)]^t = \\
&= \sum_{k=1}^s \left(p_k^2 q_k + p_k q_k^2 - 4p_k^2 q_k^2 + 6p_k q_k r_{kk} + r_{kk} - r_{kk}^2 - 2p_k r_{kk} - 2q_k r_{kk} \right) - \\
&- \sum_{k=1}^s \sum_{l=k+1}^s (2r_{kk} r_{ll} + 8p_k p_l q_k q_l) + \sum_{k=1}^s 4p_k q_k \sum_{l=1, l \neq k}^s r_{ll} + \sum_{k=1}^s \sum_{l=1, l \neq k}^s 2p_k q_l r_{lk}
\end{aligned} \tag{16}$$

Из (12) с учетом (13), (14) и (16) получаем:

$$\begin{aligned}
\sqrt{n}(\hat{T}_n - \theta) &\xrightarrow{d} N(0, S), \text{ и, следовательно} \\
\frac{1}{\sqrt{S}} \sqrt{n}(\hat{T}_n - \theta) &\xrightarrow{d} N(0, 1)
\end{aligned} \tag{17}$$

Определим

$$\begin{aligned}
\hat{p}_0^{(n)} &:= \frac{1}{n} \sum_{i=1}^n 1_{\{Y_i=0\}}, \hat{p}_1^{(n)} := \frac{1}{n} \sum_{i=1}^n 1_{\{Y_i=1\}}, \\
\hat{q}_0^{(n)} &:= \frac{1}{n} \sum_{i=1}^n 1_{\{Z_i=0\}}, \hat{q}_1^{(n)} := \frac{1}{n} \sum_{i=1}^n 1_{\{Z_i=1\}}, \\
\hat{r}_{00}^{(n)} &:= \frac{1}{n} \sum_{i=1}^n 1_{\{Y_i=0, Z_i=0\}}, \hat{r}_{01}^{(n)} := \frac{1}{n} \sum_{i=1}^n 1_{\{Y_i=0, Z_i=1\}}, \hat{r}_{10}^{(n)} := \frac{1}{n} \sum_{i=1}^n 1_{\{Y_i=1, Z_i=0\}}, \hat{r}_{11}^{(n)} := \frac{1}{n} \sum_{i=1}^n 1_{\{Y_i=1, Z_i=1\}}
\end{aligned}$$

Из усиленного закона больших чисел следует при $n \rightarrow \infty$ почти наверняка:

$$\begin{aligned}
\hat{p}_0^{(n)} &\rightarrow p_0, \hat{p}_1^{(n)} \rightarrow p_1, \hat{q}_0^{(n)} \rightarrow q_0, \hat{q}_1^{(n)} \rightarrow q_1, \\
\hat{r}_{00}^{(n)} &\rightarrow r_{00}, \hat{r}_{01}^{(n)} \rightarrow r_{01}, \hat{r}_{10}^{(n)} \rightarrow r_{10}, \hat{r}_{11}^{(n)} \rightarrow r_{11}
\end{aligned}$$

Следовательно, с вероятностью 1:

$$\begin{aligned}
\hat{S}_n &:= \\
&= \sum_{k=1}^s \left(\left(\hat{p}_k^{(n)} \right)^2 \hat{q}_k^{(n)} + \hat{p}_k^{(n)} \left(\hat{q}_k^{(n)} \right)^2 - 4 \left(\hat{p}_k^{(n)} \right)^2 \left(\hat{q}_k^{(n)} \right)^2 + 6 \hat{p}_k^{(n)} \hat{q}_k^{(n)} \hat{r}_{kk}^{(n)} + \hat{r}_{kk}^{(n)} - \left(\hat{r}_{kk}^{(n)} \right)^2 - 2 \hat{p}_k^{(n)} \hat{r}_{kk}^{(n)} - 2 \hat{q}_k^{(n)} \hat{r}_{kk}^{(n)} \right) - \\
&- \sum_{k=1}^s \sum_{l=k+1}^s \left(2 \hat{r}_{kk}^{(n)} \hat{r}_{ll}^{(n)} + 8 \hat{p}_k^{(n)} \hat{p}_l^{(n)} \hat{q}_k^{(n)} \hat{q}_l^{(n)} \right) + \sum_{k=1}^s 4 \hat{p}_k^{(n)} \hat{q}_k^{(n)} \sum_{l=1, l \neq k}^s \hat{r}_{ll}^{(n)} + \sum_{k=1}^s \sum_{l=1, l \neq k}^s 2 \hat{p}_k^{(n)} \hat{q}_l^{(n)} \hat{r}_{lk}^{(n)} \rightarrow S, (n \rightarrow \infty)
\end{aligned} \tag{18}$$

$$\text{Отсюда } \frac{S}{\hat{S}_n} \xrightarrow{P} 1, (n \rightarrow \infty)$$

Используя теорему Слуцкого, получаем из (17) и (18):

$$\frac{1}{\sqrt{\hat{S}_n}} \sqrt{n}(\hat{T}_n - \theta) = \underbrace{\sqrt{\frac{S}{\hat{S}_n}}}_{\xrightarrow{P} 1} \underbrace{\frac{1}{\sqrt{S}} \sqrt{n}(\hat{T}_n - \theta)}_{\xrightarrow{d} N(0,1)} \xrightarrow{d} N(0, 1) \tag{19}$$

Мы получили, что случайная величина $\frac{1}{\sqrt{\hat{S}_n}} \sqrt{n}(\hat{T}_n - \theta)$ имеет асимптотическое стандартное

нормальное распределение. Это позволяет получить асимптотический доверительный интервал для $\theta = P(Y_i = Z_i) - P_{\perp}(Y_i = Z_i)$:

$$\begin{aligned}
1 - \alpha &= \mathbb{P} \left(y_{\frac{\alpha}{2}} \leq N(0,1) \leq y_{1-\frac{\alpha}{2}} \right) \approx \mathbb{P} \left(y_{\frac{\alpha}{2}} \leq \frac{1}{\sqrt{\hat{S}_n}} \sqrt{n} (\theta - \hat{T}_n) \leq y_{1-\frac{\alpha}{2}} \right) = \\
&= \mathbb{P} \left(\hat{T}_n + y_{\frac{\alpha}{2}} \sqrt{\frac{\hat{S}_n}{n}} \leq \theta \leq \hat{T}_n + y_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{S}_n}{n}} \right)
\end{aligned} \tag{20}$$

Здесь y_{β} – это β -ый квантиль стандартного нормального распределения, т.е. $\mathbb{P}(N(0,1) \leq y_{\beta}) = \beta$.

Далее рассмотрим следующие гипотезы:

H_0 : Вероятность совпадения Y и Z не превышает вероятности их совпадения в случае независимости, т.е. $\theta = \mathbb{P}(Y_i = Z_i) - \mathbb{P}_{\perp}(Y_i = Z_i) \leq 0$

H_1 : Вероятность совпадения Y и Z больше вероятности их совпадения в случае независимости, т.е.

$\theta = \mathbb{P}(Y_i = Z_i) - \mathbb{P}_{\perp}(Y_i = Z_i) > 0$

Пусть задано (малое) число $0 < \alpha < 1$.

Определим следующий статистический тест: если $\sqrt{\frac{n}{\hat{S}_n}} \hat{T}_n \geq y_{1-\alpha}$, то отвергаем гипотезу H_0 и

принимаям альтернативную гипотезу H_1 . Этот тест имеет асимптотический уровень значимости $1 - \alpha$.

Действительно: пусть справедлива гипотеза H_0 , следовательно $\theta \leq 0$, и при этом $\sqrt{\frac{n}{\hat{S}_n}} \hat{T}_n \geq y_{1-\alpha}$. Тогда

получаем $\sqrt{\frac{n}{\hat{S}_n}} \left(\hat{T}_n - \theta \right) \geq \sqrt{\frac{n}{\hat{S}_n}} \hat{T}_n \geq y_{1-\alpha}$, следовательно, мы получили событие $\sqrt{\frac{n}{\hat{S}_n}} \left(\hat{T}_n - \theta \right) \geq y_{1-\alpha}$.

Но согласно (19) вероятность этого события в пределе равна

$\mathbb{P} \left(\sqrt{\frac{n}{\hat{S}_n}} \left(\hat{T}_n - \theta \right) \geq y_{1-\alpha} \right) \rightarrow \mathbb{P} \left(N(0,1) \geq y_{1-\alpha} \right) = \alpha$. Т.е. в случае справедливости гипотезы H_0 , мы

ошибочно отвергаем ее с вероятностью, асимптотически равной α , что доказывает асимптотический уровень значимости теста $1 - \alpha$.