

## МАТЕМАТИЧЕСКАЯ СТАТИСТИКА

### Лекция 9

#### Повтор

#### 8.5. Проверка независимости бинарных признаков (таблица 2×2)

Пусть имеется  $n$  объектов, выбранных случайно из большой совокупности. Каждый объект характеризуется двумя признаками, которые могут или присутствовать, или отсутствовать в каждом из объектов. Возникает вопрос, имеется ли связь между этими признаками. Например, влияет ли **форма собственности** на **рентабельность** производства ( $A$  — частная собственность,  $\bar{A}$  — общественная собственность,  $B$  — рентабельно,  $\bar{B}$  — нерентабельно), влияет ли **курение на легочные заболевания** и т.д.?

Проверяется гипотеза  $H$  о независимости признаков  $A$  и  $B$ . Если обозначить  $p_A = P\{A\}$ ,  $p_B = P\{B\}$ , то гипотеза о независимости сводится к следующему:

$$H: P\{AB\} = p_A p_B,$$

(9)

т. е. вероятность встретить сочетание признаков  $A$  и  $B$  равна произведению вероятностей встретить  $A$  и встретить  $B$ . В результате анализа на присутствие признаков могут появиться события

$$AB, A\bar{B}, \bar{A}B, \bar{A}\bar{B}.$$

Пусть эти комбинации появились соответственно:

$$v_1, v_2, v_3, v_4$$

число раз.

**Гипотеза о независимости  $A$  и  $B$  сводится к гипотезе** о том, что вероятности этих событий

$$p_1^0 = p_A p_B, \quad p_2^0 = p_A(1 - p_B), \quad p_3^0 = (1 - p_A)p_B, \quad p_4^0 = (1 - p_A)(1 - p_B)$$

соответственно. Вероятности  $p_A$  и  $p_B$  — два неизвестных параметра.

Образует случайную величину  $X^2$ , зависящую от неизвестных параметров:

$$X^2(v_1, v_2, v_3, v_4, p_A, p_B) = \sum_{i=1}^4 \frac{(v_i - np_i^0(p_A, p_B))^2}{np_i^0(p_A, p_B)}$$

Решив задачу на минимум  $X^2$  по  $p_A$  и  $p_B$ , найдем оценки, которые получатся такими, как мы ожидаем:

$$\hat{p}_A = \frac{v_1 + v_2}{n}, \quad \hat{p}_B = \frac{v_1 + v_3}{n}$$

Подставив  $\hat{p}_A$  и  $\hat{p}_B$  в  $X^2$  получим:

$$\tilde{X}^2 = X^2(v_1, v_2, v_3, v_4, \hat{p}_A, \hat{p}_B).$$

При истинности  $H$  статистика  $\tilde{X}^2$  асимптотически распределена по  $\chi_{4-2-1}^2 = \chi_1^2$  — хи-квадрат с одной (4-1-2=1) степенью свободы. Пусть

**Примечание [MA1]:** Поправить в скобках  $v_1, v_2, v_3, v_4$

задано  $\alpha$ . Из условия  $\alpha = P\{\tilde{X}_1^2 > h\} \approx P\{\chi_1^2 > h\}$  находим порог  $h$ . Итак, если  $\tilde{X}^2 \geq h$ , то гипотезу  $H$  отклоняем.

Приведем окончательные **рабочие формулы**. Для этого представим данные в следующей таблице 2×2:

	<b>A</b>	<b><math>\bar{A}</math></b>	
<b>B</b>	$v_{11}$	$v_{12}$	$v_{1\bullet} = v_{11} + v_{12}$
<b><math>\bar{B}</math></b>	$v_{21}$	$v_{22}$	$v_{2\bullet} = v_{21} + v_{22}$
	$v_{\bullet 1} = v_{11} + v_{21}$	$v_{\bullet 2} = v_{12} + v_{22}$	$n$

Если провести необходимые выкладки, получим следующую формулу:

$$\tilde{X}^2 \equiv n \frac{(v_{11}v_{22} - v_{12}v_{21})^2}{v_{\bullet 1}v_{\bullet 2}v_{1\bullet}v_{2\bullet}}, \quad (10)$$

где в круглых скобках в числителе стоит квадрат определителя матрицы, а в знаменателе — произведение частных сумм (точка в обозначениях — суммирование по соответствующему индексу).

**Конец повтора**

### 8.6. Обобщение. Проверка гипотезы о независимости признаков (таблица сопряженности признаков)

Предположим, имеется большая совокупность объектов, каждый из которых обладает двумя признаками  $A$  и  $B$ . Признак  $A$  имеет  $m$  уровней:  $A_1, A_2, \dots, A_m$ , а признак  $B$  —  $k$  уровней:  $B_1, B_2, \dots, B_k$ . Пусть уровень  $A_i$  встречается с вероятностью  $P(A_i)$ , а уровень  $B_j$  — с вероятностью  $P(B_j)$ . Независимость признаков  $A$  и  $B$  означает, что

$$P(A_i B_j) = P(A_i) \cdot P(B_j), \quad i = 1, 2, \dots, m, j = 1, 2, \dots, k,$$

т.е. вероятность встретить комбинацию  $A_i B_j$  равна произведению вероятностей. Пусть признаки определены на  $n$  объектах, случайно извлеченных из совокупности;  $v_{ij}$  — число объектов, имеющих комбинацию  $A_i B_j$ ,  $\sum_{i=1}^m \sum_{j=1}^k v_{ij} = n$ .

	<b>B<sub>1</sub></b>	<b>...</b>	<b>B<sub>j</sub></b>	<b>...</b>	<b>B<sub>k</sub></b>	<b><math>\Sigma</math></b>
<b>A<sub>1</sub></b>	$v_{11}$	$\dots$	$v_{1j}$	$\dots$	$v_{1k}$	$v_{1\bullet}$
$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
<b>A<sub>i</sub></b>	$v_{i1}$	$\dots$	$v_{ij}$	$\dots$	$v_{ik}$	$v_{i\bullet}$
$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
<b>A<sub>m</sub></b>	$v_{m1}$	$\dots$	$v_{mj}$	$\dots$	$v_{mk}$	$v_{m\bullet}$
<b><math>\Sigma</math></b>	$v_{\bullet 1}$	$\dots$	$v_{\bullet j}$	$\dots$	$v_{\bullet k}$	$n$

По совокупности наблюдений  $\{v_{ij}\}$  (таблица  $m \times k$ ) требуется проверить гипотезу  $H$  о независимости признаков  $A$  и  $B$ . Независимость означает, что вероятность встретить сочетание  $A_i$  и  $B_j$  равна произведению:

$$P(A_i B_j) = P(A_i) P(B_j).$$

Задача сводится к случаю с неизвестными параметрами, которыми являются вероятности

$$P(A_i), i = 1, 2 \dots m \text{ и } P(B_j), j = 1, 2 \dots k,$$

всего  $(m - 1) + (k - 1)$ . В выражение для статистики Пирсона входят неизвестные вероятности :

$$X^2 = \sum_{i=1}^m \sum_{j=1}^k \frac{v_{ij}^2}{nP(A_i)P(B_j)} - n \quad (8)$$

Оценки этих вероятностей по минимуму  $X^2$  приводят к вполне понятным значениям:

$$\hat{P}(A_i) = \frac{1}{n} \sum_{j=1}^k v_{ij} \equiv \frac{v_{i \cdot}}{n}, \quad \hat{P}(B_j) = \frac{1}{n} \sum_{i=1}^m v_{ij} \equiv \frac{v_{\cdot j}}{n}$$

(точка в обозначениях — суммирование по соответствующему индексу) и статистика (8) принимает вид:

$$\tilde{X}^2 = \sum_{i=1}^m \sum_{j=1}^k \frac{v_{ij}^2}{n\hat{P}(A_i)\hat{P}(B_j)} - n = n \left( \sum_{i=1}^m \sum_{j=1}^k \frac{v_{ij}^2}{v_{i \cdot} v_{\cdot j}} - 1 \right). \quad (11)$$

Если гипотеза  $H$  верна, то по теореме Фишера статистика  $\tilde{X}^2$  асимптотически распределена по закону хи-квадрат с числом степеней свободы

$$f = mk - 1 - (m - 1) - (k - 1) = (m - 1)(k - 1),$$

Мы можем определить вероятность получить наше значение  $\tilde{X}^2$  (меры различия), которое мы получили, если  $H$  верна:

$$P\{\text{получить} \geq \tilde{X}^2 | H\} \approx P\{\chi_f^2 \geq \tilde{X}^2\} \leq \alpha; \quad (12)$$

если она мала, то гипотезу о независимости признаков следует отклонить.

Ясно, что по (11), (12) можно проверять независимость двух случайных величин, разбив диапазоны их значений на  $m$  и  $k$  частей.

### 8.7. Проверка гипотезы об однородности выборок

Пусть имеется  $k$  выборок объемами  $n_i$ ,  $i = 1, 2 \dots m$ , извлеченных из различных совокупностей (например, имеем наблюдения доходов  $Z$  служащих в городе А и в городе В; можно ли считать, что распределения доходов одинаковы?).

Измеряемая величина в каждой из выборок может иметь  $k$  уровней  $B_j$ ,  $j = 1, 2 \dots k$ . Пусть  $p_{ij}$  — истинная (неизвестная) вероятность получить  $Z_j$  в  $i$ -й совокупности. Требуется проверить гипотезу  $H$  о том, что исходные совокупности распределены одинаково, т.е.

$$p_{ij} = q_j.$$

Обозначим  $v_{ij}$  — число наблюдений в  $i$ -й выборке, имеющих уровень  $Z_j$ , причем  $\sum_j v_{ij} \equiv v_{i\bullet} = n_i$ ,  $\sum_i n_i = n$ . Имеем таблицу наблюдений размером  $m \times k$ , аналогично предыдущему разделу 8.5.

	$B_1$	...	$B_j$	...	$B_k$	$\Sigma$
1-я совокупность	$v_{11}$	...	$v_{1j}$	...	$v_{1k}$	$v_{1\bullet} = n_1$
$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
$i$ -я совокупность	$v_{i1}$	...	$v_{ij}$	...	$v_{ik}$	$v_{i\bullet} = n_i$
$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$
$m$ -я совокупность	$v_{m1}$	...	$v_{mj}$	...	$v_{mk}$	$v_{m\bullet} = n_m$
$\Sigma$	$v_{\bullet 1}$	...	$v_{\bullet j}$	...	$v_{\bullet k}$	$n$

Если  $H$  верна, то имеем неизвестные параметры  $q_1, q_2 \dots q_k$ . Если бы они были известны, можно было бы определить теоретические частоты и составить статистику

$$X^2(q_1, \dots, q_m) = \sum_{i=1}^m \left( \sum_{j=1}^k \frac{(v_{ij} - n_i q_j)^2}{n_i q_j} \right) = \sum_{i=1}^m \left( \sum_{j=1}^k \frac{v_{ij}^2}{n_i q_j} - n_i \right), \quad (13)$$

причем  $\sum_{j=1}^k \frac{(v_{ij} - n_i q_j)^2}{n_i q_j}$  при истинности  $H$  приближено подчинялась бы распределению хи-квадрат с  $(k-1)$  степенями свободы, а вся сумма в (13) — распределению хи-квадрат с  $m(k-1)$  степенями свободы в силу независимости выборок.

Но вероятности  $q_1, q_2 \dots q_m$  неизвестны. Определяем оценки для них, полагая, что  $H$  верна:

$$\hat{q}_j = \frac{1}{n} \sum_{i=1}^m v_{ij}$$

и после подстановки оценок количество степеней свободы уменьшится на число  $(k-1)$  оцененных параметров и станет равным

$$f = m(k-1) - (k-1) = (k-1)(m-1).$$

Статистика  $\tilde{X}^2$  и вся процедура проверки гипотезы примет вид, аналогичный (11), (12) при проверке гипотезы о независимости. Решающая

$$\text{статистика: } \tilde{X}^2 = n \left( \sum_{i=1}^k \sum_{j=1}^m \frac{v_{ij}^2}{v_{i\bullet} v_{\bullet j}} - 1 \right).$$

Если  $P\{\chi_f^2 \geq \tilde{X}^2\} \leq \alpha$ , то гипотезу об однородности следует отклонить.

## § 9. Критерий согласия Колмогорова

Проверяется гипотеза  **$H$  о том**, что последовательность независимых наблюдений  $\xi_1, \xi_2 \dots \xi_n$  извлечена из совокупности с непрерывной функцией **распределения  $F(x)$** .

Эту задачу мы уже решали с помощью критерия хи-квадрат Там была одна мера различия между гипотезой и наблюдениями, здесь – другая. Там был общий метод (для произвольных гипотез) построения решающей статистики, здесь – частный метод (и более простой) для проверки  $H$  о заданной непрерывной ф.р.

Построим вариационный ряд  $\xi_{(1)} \leq \xi_{(2)} \leq \dots \leq \xi_{(n)}$  и функцию эмпирического распределения  $F_n^*(x)$ . Мерой различия (можно говорить качества согласования) **эмпирического и гипотетического распределения** примем максимальное отклонение функции эмпирического распределения  $F_n^*(x)$  от гипотетического  $F(x)$ , то есть верхняя грань модуля разности

$$D_n = \sup_x |F_n^*(x) - F(x)|.$$

Замечателен тот факт, что **распределение статистики  $D_n$  при истинности  $H$  не зависит от  $F(x)$**  [4, таблицы Бльшева и Смирнова], и поэтому критическое значение определяется.

Процедура проверки гипотезы  $H$ :

- вычисляется  $D_n$ ;
- выбирается порог  $\lambda_1$  (о выборе порога, см. ниже);
- если  $D_n \geq \lambda_1$ , то  $H$  отклоняется.

Критическое значение  $\lambda_1$  выбирается так, чтобы обеспечить заданную вероятность ошибки при истинности  $H$ :

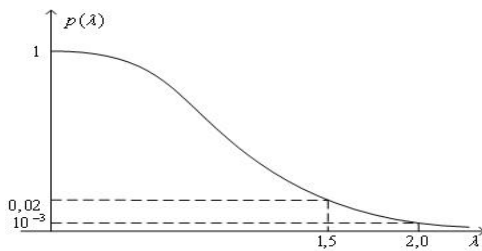
$$P\{\text{отклонить } H \mid H\} = P\{D_n \geq \lambda_1 \mid H\} = \alpha.$$

$\alpha$  - задаваемый уровень значимости. В таблицах [4] имеются значения  $\lambda_1(n, \alpha)$  для значений  $\alpha = 1, 2, 5, 10, 20 \%$ ,  $n = 1(1)100$ . Там же имеются асимптотические формулы.

Оказывается, что для любой непрерывной  $F(\cdot)$  при  $n \rightarrow \infty$  для распределения статистики  $D_n$  при истинности гипотезы  $H$  справедливо соотношение:

$$P\{D_n \sqrt{n} \geq \lambda \mid H\} \xrightarrow{n \rightarrow \infty} P(\lambda) \equiv 1 - K(\lambda),$$

где  $K(\lambda) = \sum_{k=-\infty}^{\infty} (-1)^k e^{-2k^2 \lambda^2}$  — функция распределения Колмогорова.



Этот факт при  $n \geq 20$  позволяет простым способом приближенно обеспечить заданную вероятность ошибки первого рода (уровень значимости)

$$\alpha = P\{D_n \sqrt{n} \geq \lambda | H\} \approx P(\lambda).$$

**Рис. 12. График функции  $P(\lambda)$**

График функции  $P(\lambda)$  условно показан на рис. 12. Итак,  $\lambda$  выбирается по таблицам как функция  $\alpha : \lambda(\alpha)$ , т.е.

$$P\{D_n \sqrt{n} \geq \lambda(\alpha) | H\} \approx \alpha \Rightarrow \lambda_1 \approx \frac{\lambda(\alpha)}{\sqrt{n}},$$

и если  $D_n \geq \lambda_1 = \frac{\lambda(\alpha)}{\sqrt{n}}$ , то  $H$  отклоняется.

На рис.12 показаны две характерные точки функции  $P(\lambda)$ :

при  $\alpha = 0,02$   $\lambda(\alpha) \approx 1,5$ ,

при  $\alpha = 10^{-3}$   $\lambda(\alpha) \approx 2,0$ .

Для  $\lambda > 2,5$  хорошим приближением является  $P(\lambda) \approx 2e^{-2\lambda^2}$ , причем с ростом  $n$  точность улучшается.

## § 10. Различение двух простых гипотез

### 10.1. Фиксированный объем наблюдений

Пусть имеется совокупность наблюдений  $x = (x_1, x_2, \dots, x_n)$ , являющихся реализацией случайных величин  $\xi \equiv (\xi_1, \xi_2, \dots, \xi_n)$ , относительно которой имеется два предположения (гипотезы):

1.  $H_0$ :  $\xi$  распределена по закону  $p_0(x)$ ;

2.  $H_1$ :  $\xi$  распределена по закону  $p_1(x)$ .

(если  $\xi$  — непрерывна, то  $p_0(x)$ ,  $p_1(x)$  — плотности, если дискретна — вероятности). По наблюдениям  $x$  требуется принять одно из двух решений:

или «верна  $H_0$ » (это решение обозначим 0),

или «верна  $H_1$ » (решение 1).

Ясно, что дело сводится к определению решающей функции  $\delta(x)$ , имеющей два значения 0 и 1, т.е. к определению разбиения  $\Gamma = (\Gamma_0, \Gamma_1)$  пространства  $X$  всех возможных значений  $x$ :

$$\delta(x) = \begin{cases} 0, & \text{если } x \in \Gamma_0, \\ 1, & \text{если } x \in \Gamma_1. \end{cases}$$

При использовании любой решающей функции  $\delta(x)$  возможны ошибки двух типов:

1) ошибка первого рода — принятие  $H_1$  при истинности  $H_0$ ;

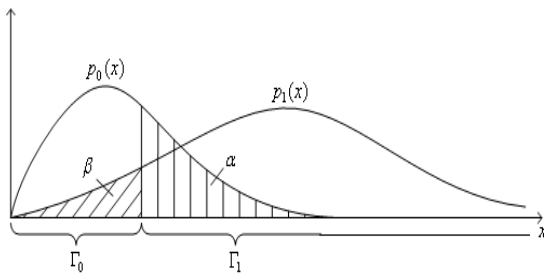
2) ошибка второго рода — принятие  $H_0$  при истинности  $H_1$ .

Для любой решающей функции имеем две условные вероятности:

$$\alpha = P(\text{принять } H_1 | H_0) = \int_{\Gamma_1} p_0(x) dx, \quad (1)$$

$$\beta = P(\text{принять } H_0 | H_1) = \int_{\Gamma_0} p_1(x) dx,$$

которые называются вероятностями ошибок первого и второго рода



соответственно (рис. 13). Хотелось бы иметь  $\alpha$  и  $\beta$  близкими к нулю, но из (1) ясно, что если одна из них уменьшается, например,  $\alpha$  (за счет уменьшения  $\Gamma_1$ ), то другая,  $\beta$ , увеличивается (за счет увеличения  $\Gamma_0$ ;  $\Gamma_0 \cup \Gamma_1 = X$ ,  $\Gamma_0 \setminus \Gamma_1 = \emptyset$ ). Существуют различные подходы к

определению оптимального правила.

**Рис. 13.** Вероятности ошибок  $\alpha$  и  $\beta$

#### Байесовский подход

Будем считать, что многократно сталкиваемся с проблемой выбора между  $H_0$  и  $H_1$ . В этом случае можно говорить о частоте, с которой истинна  $H_0$  (или  $H_1$ ), т.е. о том, что истинность  $H_0$  (или  $H_1$ ) — событие случайное, причем вероятность события, когда верна  $H_0$ , а когда верна  $H_1$ , известны:

$$P(H_0) = q_0, \quad P(H_1) = q_1, \quad q_0 + q_1 = 1.$$

Кроме того, будем считать, что при каждой ошибке первого рода несем потери  $W_0$  (вероятность этого события  $P(H_0) \cdot P(\text{принять } H_1 | H_0)$ ), а при ошибке второго рода — потери  $W_1$  (с вероятностью  $P(H_1) \cdot P(\text{принять } H_0 | H_1)$ ).

Если пользуемся правилом  $\delta$  (с разбиением  $\Gamma$ ), то средний штраф от однократного использования:

$$R(\Gamma) = q_0 \cdot \alpha(\Gamma) \cdot W_0 + q_1 \cdot \beta(\Gamma) \cdot W_1. \quad (1a)$$

Назовем правило  $\delta$  (соответственно, разбиение  $\Gamma \equiv (\Gamma_0, \Gamma_1)$ ) в байесовском смысле оптимальным, если

$$R(\Gamma) = \min_{\Gamma'} R(\Gamma').$$

Оказывается справедливой следующая теорема.

**Теорема.** Оптимальным является правило, для которого область  $\Gamma_1$  принятия гипотезы  $H_1$  определяется соотношением:

$$\Gamma_1 = \left\{ x: \frac{p_1(x)}{p_0(x)} \geq h \equiv \frac{q_0 W_0}{q_1 W_1} \right\}. \quad (2)$$

Действительно, пусть  $T = (T_0, T_1)$  — произвольное разбиение; для него средние потери

$$R(T) = q_0 \cdot \alpha(T) \cdot W_0 + q_1 \cdot \beta(T) \cdot W_1 = q_0 W_0 \int_{T_1} p_0(x) dx + q_1 W_1 \int_{T_0} p_1(x) dx - \\ - q_1 W_1 \int_{T_1} p_1(x) dx + q_1 W_1 \int_{T_1} p_1(x) dx.$$

Здесь во второй строке добавлено и вычтено одно и то же слагаемое. После объединения интегралов получаем

$$R(T) = q_0 W_0 \int_{T_1} [q_0 W_0 p_0(x) - q_1 W_1 p_1(x)] dx + q_1 W_1.$$

Очевидно, для того чтобы получить минимальное значение средних потерь, в область интегрирования  $T_1$  нужно включить те точки, в которых подынтегральная функция в квадратных скобках отрицательна, т.е.

$$q_0 W_0 p_0(x) - q_1 W_1 p_1(x) < 0,$$

что дает в эквивалентной записи оптимальную область  $\Gamma_1$  в (2).

**Замечание.** В частном случае, если  $W_0 = W_1 = 1$ , то  $R(\Gamma)$  в (1а) имеет смысл безусловной вероятности ошибки, а соответствующее оптимальное правило называется правилом “идеального наблюдателя” или правилом Зигерта- Котельникова.