

В. И. РОМАНОВСКИЙ

ОБ ИНДУКТИВНЫХ ВЫВОДАХ В СТАТИСТИКЕ

(Представлено академиком А. Н. Колмогоровым 21 III 1940)

К числу основных проблем статистики относится вопрос о разыскании характеристик генеральных совокупностей по соответственным характеристикам случайных выборок из них. Разрешение этого вопроса чаще всего основывается на теореме Байеса или на теоремах, подобных ей, связанных с введением в рассуждения априорных вероятностей, от элементов произвола и необоснованности в которых невозможно освободиться. Но можно тот же вопрос разрешить иначе, совершенно не прибегая к априорным вероятностям.

На эту возможность указал впервые Р. Фишер. Он ввел понятие о доверительной (фидуциальной) вероятности и изложил свои идеи в применении к выборкам из нормальных совокупностей в ряде работ, начиная с 1930 г. (1-4). Эти идеи получили дальнейшее развитие у ряда других, главным образом английских, ученых: у Ю. Неймана, Э. Пирсона, Бартлетта, Вилкса и др. (5-9).

В настоящей заметке идеи Р. Фишера применяются к выборкам из произвольных совокупностей.

Будем для простоты рассматривать генеральную совокупность S с количественным аргументом x ; пусть x имеет в S конечное число значений

$$x_1, x_2, \dots, x_s \quad (1)$$

с долями

$$p_1, p_2, \dots, p_s. \quad (2)$$

S может быть конечной или бесконечной совокупностью. Предположим, что из S берется случайная выборка S' по схеме возвращенного шара; пусть в ней значения (1) имеют частоты

$$n_1, n_2, \dots, n_s$$

и пусть объем ее равен n .

Теперь можно доказать следующую, очень простую и очень важную теорему [(10), стр. 108].

Теорема 1. *Каковы бы ни были n, p_1, p_2, \dots, p_s и $\epsilon > 0$,*

$$P \left[\sum \left(\frac{n_i}{n} - p_i \right)^2 < \epsilon^2 \right] > 1 - \frac{1-s}{n\epsilon^2}. \quad (3)$$

Так как правая часть неравенства (3) не зависит от p_i , то из этой теоремы тотчас следует, что если доли p_i значений (1) совокупности S нам неизвестны, то мы можем их заменить относительными частотами $\frac{n_i}{n}$, наблюдаемыми в выборке S' , причем с вероятностью, как угодно близкой к достоверности, мы можем ожидать, что совершенные при этой замене ошибки будут сколь угодно малы, если n будет достаточно большим числом.

Мы видим, что это заключение совершенно свободно от априорных вероятностей*.

Теперь мы покажем, какую роль играет полученный нами вывод в индуктивных суждениях статистики.

Пусть θ представляет некоторую характеристику генеральной совокупности S , зависящую от значений (1) и неизвестных нам их долей (2); пусть затем θ' представляет соответствующую характеристику выборки S' , составленную из значений (1) и их частот $\frac{n_i}{n}$ так же, как θ составлена из (1) и (2). Тогда нетрудно видеть справедливость следующей теоремы, имеющей обширное применение в статистике и свободной от априорных вероятностей.

Теорема 2. Если функция $\theta(z_1, z_2, \dots, z_s)$ непрерывна для значений z_1, z_2, \dots, z_s , удовлетворяющих условиям

$$z_i \geq 0, \quad i = 1, 2, \dots, s; \quad \sum z_i \leq 1,$$

то вероятность неравенства

$$\left| \theta\left(\frac{n_1}{n}, \frac{n_2}{n}, \dots, \frac{n_s}{n}\right) - \theta(p_1, p_2, \dots, p_s) \right| > \varepsilon$$

стремится к нулю при $n \rightarrow \infty$ равномерно относительно

$$p_1, p_2, \dots, p_s,$$

каково бы ни было заданное $\varepsilon > 0$.

В заключение заметим, что очень просто можно доказать следующую теорему, подобную теореме 1.

Теорема 3. Если из генеральной совокупности S произведены две случайные выборки S' и S'' по схеме возвращенного шара, объемов n' и n'' и с частотами $\frac{n'_i}{n'}$ и $\frac{n''_i}{n''}$ для значений (1), то для любого $\varepsilon > 0$ справедливо неравенство

$$P \left[\sum \left(\frac{n'_i}{n'} - \frac{n''_i}{n''} \right)^2 < \varepsilon^2 \right] > 1 - \frac{1}{\varepsilon^2} \left(\frac{1}{n'} + \frac{1}{n''} \right). \quad (4)$$

Рассматривая выборку S' , как уже произведенную, а выборку S'' , как предстоящую и неизвестную, мы выводим из написанного неравенства, что с вероятностью, как угодно близкой к достоверности, все частоты $\frac{n''_i}{n''}$ будут отличаться от соответственных частот $\frac{n'_i}{n'}$ как угодно мало, если числа n' и n'' будут достаточно велики. И этот вывод свободен от априорных вероятностей.

Все изложенное можно без труда распространить на случайные выборки, произведенные по схеме возвращенного шара из генеральной

* Следует заметить, что в доказательстве неравенства (3) они также не принимают никакого участия.

совокупности с любым конечным числом аргументов, как угодно распределенных в ней*.

Поступило
21 III 1940

ЦИТИРОВАННАЯ ЛИТЕРАТУРА

¹ R. A. Fisher, Proc. Camb. Phil. Soc., 26, 528—535 (1930). ² R. A. Fisher, Proc. Roy. Soc., A, 139, 343—348 (1933). ³ R. A. Fisher, Proc. Roy. Soc., A, 144, 285—307 (1934). ⁴ R. A. Fisher, Ann. of Eugenics, VI, 391—398 (1935). ⁵ C. J. Clopper and E. S. Pearson, Biometrika, 26, 404—413 (1934). ⁶ G. Neyman, Phil. Trans., A, 236, 333—380 (1937). ⁷ M. S. Bartlett, Proc. Roy. Soc., A, 160, 268—282 (1937). ⁸ M. S. Bartlett, Ann. of Math. Stat., X, 129—138 (1939). ⁹ S. S. Wilks, Ann. of Math. Stat., IX, 272—280 (1938). ¹⁰ В. Романовский, Математическая статистика (1938).

* Некоторые формулировки сообщения В. И. Романовского могут дать повод к недоразумениям. Поэтому я считаю необходимым отметить здесь следующее:

В теореме 1 вероятность

$$P \left[\sum \left(\frac{n_i}{n} - p_i \right)^2 < \varepsilon^2 \right]$$

берется при заданных n, p_1, p_2, \dots, p_s и случайных n_1, n_2, \dots, n_s . Теорема утверждает, что эта вероятность удовлетворяет неравенству (3), каковы бы ни были n, p_1, p_2, \dots, p_s . В этой форме теорема, безусловно, правдива. Было бы, однако, ошибочно понимать дело так, что отсюда получается какая-либо оценка вероятности неравенства

$$\sum \left(\frac{n_i}{n} - p_i \right)^2 < \varepsilon^2$$

при заданных n_1, n_2, \dots, n_s и случайных p_1, p_2, \dots, p_s . Получить такую оценку без введения закона распределения для самих p_1, p_2, \dots, p_s нельзя.

Точно так же в теореме 3 вероятность

$$P \left[\sum \left(\frac{n'_i}{n'} - \frac{n''_i}{n''} \right)^2 < \varepsilon^2 \right]$$

берется при заданных $n', n'', p_1, p_2, \dots, p_s$ и случайных $n'_1, n'_2, \dots, n'_s, n''_1, n''_2, \dots, n''_s$. Теорема утверждает, что эта вероятность удовлетворяет неравенству (4), каковы бы ни были $n', n'', p_1, p_2, \dots, p_s$. Никаких заключений о вероятности неравенства

$$\sum \left(\frac{n'_i}{n'} - \frac{n''_i}{n''} \right)^2 < \varepsilon^2$$

при заданных n'_1, n'_2, \dots, n'_s и случайных $n''_1, n''_2, \dots, n''_s$ вывести из теоремы 3 нельзя. В ясной форме такого рода ошибочные заключения можно найти в цитированной В. И. Романовским работе R. A. Fisher'a [Annals of Eugenics, 6, стр. 391—392 (1935)]. Тем не менее я считаю, что теоремы типа теорем 1—3 и аналогичных теорем R. A. Fisher'a и его последователей при их правильном понимании могут служить достаточным обоснованием практических приемов математической статистики. Подробнее свою точку зрения по этому важному вопросу я изложу в особой статье.

Академик А. Н. Колмогоров