

ОПТИМАЛЬНОЕ ГРУППИРОВАНИЕ ПРИ ОБРАБОТКЕ
ЭКСПЕРИМЕНТАЛЬНЫХ ДАННЫХ

Статистическая обработка экспериментальных наблюдений в большинстве случаев сопровождается группированием данных. Это вызвано, во-первых, необходимостью сжатия больших объемов накапливаемой информации, вследствие чего упрощается математическая обработка результатов; во-вторых, довольно часто при испытаниях на надежность вообще отсутствует понятие индивидуального измерения, а регистрируется лишь число наблюдений, попавших в определенный интервал; в-третьих, группирование может быть вызвано применяемым методом статистической обработки данных.

В настоящей работе обсуждаются методы группирования регистрируемых данных, применяемые при исследовании законов распределения случайных величин. К группированию наблюдений прибегают при построении гистограмм, оценивании параметров распределений, проверке соответствия между гипотетическим и фактически наблюдаемым распределениями по критерию χ^2 Пирсона или его модификациям.

До последнего времени в практике статистических вычислений при группировании область определения случайной величины разбивают на интервалы, в основном, одним из двух способов [1]: на интервалы равной длины с последующим объединением смежных, если в них падает малое число наблюдений, и на интервалы равной вероятности. Во втором случае разбиение связано с видом распределения. Равновероятному группированию соответствует ситуация, когда выборку раз -

бывают на группы так, чтобы в каждой связалось одинаковое количество наблюдений.

Любое группирование влечет за собой потерю информации, понимаемой в самом общем смысле, по сравнению с исходной негруппированной выборкой. А это отрицательно сказывается на качестве результатов статистических выводов.

Оценки параметров распределений по группированным данным находят обычно методом максимального правдоподобия, методом минимума χ^2 или одним из родственных ему методов. При некоторых условиях [2-6] эти оценки существуют, единственны и обладают свойствами состоятельности и асимптотической эффективности. При этом асимптотическая дисперсионная матрица оценки вектора параметров θ распределения определяется соотношением

$$D(\hat{\theta}) = N^{-1} J^{-1}(\hat{\theta}),$$

где N - объем выборки;

$J(\theta)$ - информационная матрица Фишера по группированным данным. Так как $J(\theta)$ - положительно определенная матрица, то с ростом потерь информации, вызванным группированием, ухудшается дисперсионная матрица оценок параметров закона распределения.

Рассмотрим, как связана информационная матрица $J(\theta)$ с критерием χ^2 Пирсона. Статистика

$$\chi^2 = N \sum_{i=1}^k (n_i / N - p_i(\theta))^2 / p_i(\theta),$$

где k - число интервалов группирования;

n_i - число наблюдений, попавших в i -й интервал;

$p_i(\theta)$ - гипотетическая вероятность попадания наблюдения в i -й интервал,

подчиняется распределению χ^2 с $k - 1$ степенью свободы, если верна нулевая гипотеза, и подчиняется нецентральному распределе-

ний χ^2 с тем же числом степеней свободы и параметром нецентральности

$$\lambda = N \sum_{i=1}^k (P_i(\theta_i) - P_i(\theta))^2 / P_i(\theta),$$

если верна конкурирующая гипотеза и выборка соответствует распределению того же типа, но с параметром θ_1 . Пусть $\theta_1 = \theta + \delta\theta$. Разложим $P_i(\theta_1)$ в ряд Тейлора, пренебрегая членами высшего порядка

$$\begin{aligned} \lambda &\approx N \sum_{i=1}^k \frac{(P_i(\theta) + \nabla P_i^T(\theta) \delta\theta - P_i(\theta))^2}{P_i(\theta)} = N \sum_{i=1}^k \frac{\delta\theta^T \nabla P_i(\theta) \nabla P_i^T(\theta) \delta\theta}{P_i(\theta)} = \\ &= N \delta\theta^T \left(\sum_{i=1}^k \frac{\nabla P_i(\theta) \nabla P_i^T(\theta)}{P_i(\theta)} \right) \delta\theta = N \delta\theta^T J(\theta) \delta\theta. \end{aligned}$$

Отсюда ясно, что с ростом потерь информации при малых $\delta\theta$, то есть близких альтернативных гипотезах, уменьшается параметр нецентральности λ , а следовательно, и мощность критерия χ^2 .

Так как матрица $J(\theta)$ зависит от граничных точек интервалов, то можно подобрать их так, задаваясь определенным критерием, минимизировать потери информации от группирования и тем самым максимизировать мощность критерия χ^2 при близких альтернативных гипотезах и соответствующим образом минимизировать асимптотическую дисперсию оценок параметров.

Необходимо отметить, что в общем случае равномерное и равновероятное разбиения весьма далеки от оптимального.

В таблицах представлены асимптотически оптимальные граничные точки, максимизирующие определитель информационной матрицы Фишера по группированным данным. В последних колонках таблиц показаны значения отношения определителя информационной матрицы Фишера по группированным наблюдениям к определителю информационной матрицы по негруппированной выборке. Для сравнения

заметим, что при равновероятном группировании значение этого отношения для экспоненциального распределения и числе интервалов $K = 10$ составляет всего 0,8928.

В табл. I приведены оптимальные граничные точки для экспоненциального распределения с плотностью распределения $f(x) = \theta \exp(-\theta x)$ распределения Рэлея - $f(x) = (x/\theta^2) \exp(-x^2/2\theta^2)$, распределения Парето - $f(x) = \theta \alpha^\theta x^{-(\theta+1)}$, в табл. 2 оптимальные граничные точки для распределения Максвелла с плотностью $f(x) = (2x^2/\theta^3 \sqrt{2\pi}) \exp(-x^2/2\theta^2)$, в табл. 3 оптимальные граничные точки для нормального распределения, в табл. 4 для распределения Вейбулла с плотностью $f(x) = \frac{\theta}{\sigma} (\frac{x}{\sigma})^{\theta-1} \exp(-(\frac{x}{\sigma})^\theta)$ и для распределения экстремального значения с плотностью $f(x) = \frac{1}{\sigma} \exp\left\{t \frac{M-x}{\sigma} - \exp\left[-\frac{M-x}{\sigma}\right]\right\}$, наибольшего значения со знаком (+), наименьшего значения со знаком (-). Для распределения Коши с плотностью $f(x) = \theta/\pi (\theta^2 + (x-\theta)^2)$ оптимальное группирование совпало с разбиением на интервалы равной вероятности. Поэтому таблица оптимальных граничных точек для него не приводится.

На рисунке в качестве иллюстрации преимуществ оптимального группирования при проверке гипотез о согласии даны функции мощности критерия χ^2 для экспоненциального распределения при 10 интервалах группирования, объема выборки $N = 1000$ и уровня значимости $\alpha = 0.05$. Оптимальное группирование на несколько процентов увеличивает мощность критерия.

Если необходимо проверить согласие выборки с некоторым распределением, то следует непосредственно пользоваться таблицами I-4. Для оценивания параметров распределений процедура использования таблиц оптимального группирования несколько усложняется. Для этого, во-первых, можно воспользоваться априорными сведениями о параметрах и выбрать граничные точки, опираясь на них, а, во-вторых, лучше и проще всего сгруппировать выборку так, чтобы

Таблица I

Оптимальные граничные точки для экспоненциального распределения в виде $t_i = \theta \cdot X_i$, для распределения Релея в виде $t_i = (x_i/\theta)^2/2$, для распределения Парето в виде $t_i = \theta^2 h(x_i/\omega)$ и соответствующие значения относительной асимптотической информации A

| K | t_1 | t_2 | t_3 | t_4 | t_5 | t_6 | t_7 | t_8 | t_9 | t_{10} | A |
|----|--------|--------|--------|--------|--------|--------|--------|--------|--------|----------|--------|
| 2 | 1.5936 | | | | | | | | | | 0.6476 |
| 3 | 1.0176 | 2.6112 | | | | | | | | | 0.8203 |
| 4 | 0.7541 | 1.7716 | 3.3652 | | | | | | | | 0.8910 |
| 5 | 0.6004 | 1.3545 | 2.3720 | 3.9557 | | | | | | | 0.9269 |
| 6 | 0.4993 | 1.0997 | 1.8538 | 2.8714 | 4.4650 | | | | | | 0.9476 |
| 7 | 0.4276 | 0.9269 | 1.5273 | 2.2813 | 3.2989 | 4.8925 | | | | | 0.9606 |
| 8 | 0.3739 | 0.8015 | 1.3008 | 1.9012 | 2.6553 | 3.6729 | 5.2665 | | | | 0.9693 |
| 9 | 0.3323 | 0.7063 | 1.1338 | 1.6331 | 2.2336 | 2.9876 | 4.0052 | 5.5988 | | | 0.9754 |
| 10 | 0.2990 | 0.6314 | 1.0053 | 1.4329 | 1.9322 | 2.5326 | 3.2866 | 4.3042 | 5.8979 | | 0.9798 |
| 11 | 0.2716 | 0.5695 | 0.9014 | 1.2746 | 1.7015 | 2.1989 | 2.7955 | 3.5429 | 4.5480 | 6.1176 | 0.9832 |

Оптимальные граничные точки для распределения Максвелла
 в виде $t_i = \alpha_i / \theta$ и соответствующие значения
 относительной асимптотической информации A

| K | t_1 | t_2 | t_3 | t_4 | t_5 | t_6 | t_7 | t_8 | t_9 | A |
|-----|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
| 2 | 2.0451 | | | | | | | | | 0.6451 |
| 3 | 1.6762 | 2.5366 | | | | | | | | 0.8179 |
| 4 | 1.4689 | 2.1292 | 2.8402 | | | | | | | 0.8692 |
| 5 | 1.3292 | 1.8879 | 2.4221 | 3.0583 | | | | | | 0.9254 |
| 6 | 1.2261 | 1.7205 | 2.1667 | 2.6379 | 3.2274 | | | | | 0.9464 |
| 7 | 1.1458 | 1.5947 | 1.9859 | 2.3759 | 2.8081 | 3.3649 | | | | 0.9596 |
| 8 | 1.0807 | 1.4952 | 1.8481 | 2.1879 | 2.5431 | 2.9480 | 3.4803 | | | 0.9685 |
| 9 | 1.0267 | 1.4138 | 1.7377 | 2.0423 | 2.3499 | 2.6803 | 3.0652 | 3.5789 | | 0.9747 |
| 10 | 0.9798 | 1.3447 | 1.6460 | 1.9252 | 2.2003 | 2.4857 | 2.7984 | 3.1668 | 3.6625 | 0.9792 |

Оптимальные граничные точки для нормального распределения
в виде $t_i (x_i - \mu) / \sigma$ и соответствующие значения
относительной асимптотической информации A

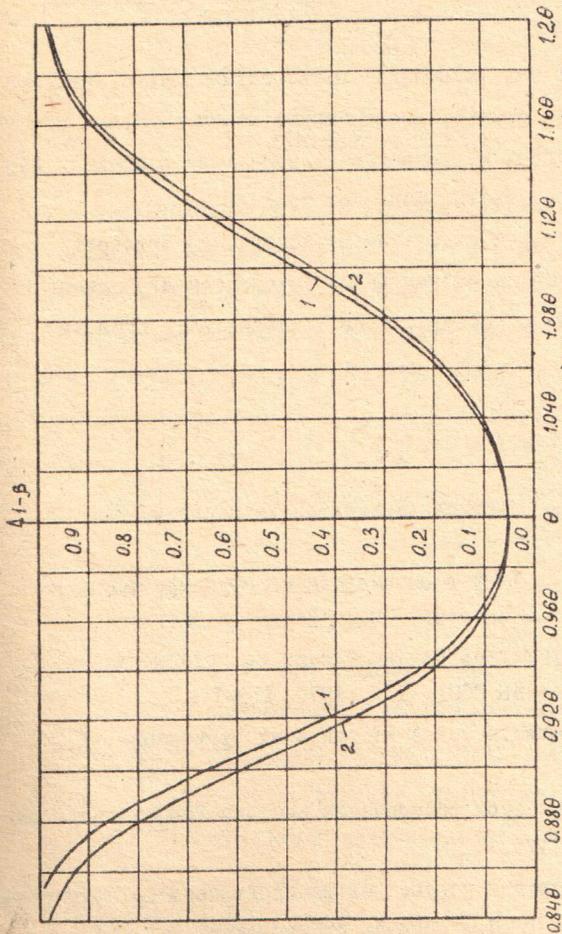
| k | t_1 | t_2 | t_3 | t_4 | t_5 | t_6 | t_7 | t_8 | t_9 | t_{10} | t_{11} | t_{12} | A |
|-----|--------|--------|--------|--------|--------|--------|-------|-------|-------|----------|----------|----------|--------|
| 3 | -1.111 | 1.111 | | | | | | | | | | | 0.4065 |
| 4 | -1.383 | 0.0 | 1.383 | | | | | | | | | | 0.5527 |
| 5 | -1.696 | -0.689 | 0.689 | 1.696 | | | | | | | | | 0.6827 |
| 6 | -1.882 | -0.997 | 0.0 | 0.997 | 1.882 | | | | | | | | 0.7557 |
| 7 | -2.060 | -1.265 | -0.492 | 0.492 | 1.265 | 2.060 | | | | | | | 0.8103 |
| 8 | -2.195 | -1.455 | -0.786 | 0.0 | 0.786 | 1.455 | 2.195 | | | | | | 0.8474 |
| 9 | -2.319 | -1.622 | -1.022 | -0.383 | 0.383 | 1.022 | 1.622 | 2.319 | | | | | 0.8753 |
| 10 | -2.423 | -1.758 | -1.205 | -0.650 | 0.0 | 0.650 | 1.205 | 1.758 | 2.423 | | | | 0.8960 |
| 11 | -2.517 | -1.878 | -1.360 | -0.862 | -0.314 | 0.314 | 0.862 | 1.360 | 1.878 | 2.517 | | | 0.9121 |
| 12 | -2.599 | -1.983 | -1.491 | -1.033 | -0.553 | 0.0 | 0.553 | 1.033 | 1.491 | 1.983 | 2.599 | | 0.9247 |
| 13 | -2.675 | -2.076 | -1.607 | -1.178 | -0.746 | -0.267 | 0.267 | 0.746 | 1.178 | 1.607 | 2.076 | 2.675 | 0.9348 |

Таблица 4

Оптимальные граничные точки для распределения Вейбулла в виде $t_i = (x_i/\theta)^{\theta}$,
 для распределения наибольшего экстремального значения в виде $t_i = \exp((1-x_i)/\sigma)$,
 для распределения наименьшего экстремального значения в виде $t_i = \exp(-(1-x_i)/\sigma)$
 и соответствующие значения относительной асимптотической информации A

| K | t_1 | t_2 | t_3 | t_4 | t_5 | t_6 | t_7 | t_8 | t_9 | t_{10} | t_{11} | t_{12} | A |
|----|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|--------|
| 3 | 0.273 | 2.607 | | | | | | | | | | | 0.4079 |
| 4 | 0.211 | 1.398 | 3.414 | | | | | | | | | | 0.5572 |
| 5 | 0.104 | 0.512 | 1.959 | 3.861 | | | | | | | | | 0.6836 |
| 6 | 0.077 | 0.365 | 1.227 | 2.573 | 4.410 | | | | | | | | 0.7571 |
| 7 | 0.050 | 0.232 | 0.676 | 1.719 | 2.992 | 4.795 | | | | | | | 0.8109 |
| 8 | 0.038 | 0.174 | 0.484 | 1.190 | 2.204 | 3.429 | 5.205 | | | | | | 0.8480 |
| 9 | 0.028 | 0.127 | 0.343 | 0.783 | 1.603 | 2.571 | 3.767 | 5.527 | | | | | 0.8756 |
| 10 | 0.021 | 0.099 | 0.264 | 0.577 | 1.181 | 1.993 | 2.927 | 4.102 | 5.848 | | | | 0.8963 |
| 11 | 0.017 | 0.077 | 0.205 | 0.436 | 0.856 | 1.534 | 2.319 | 3.232 | 4.593 | 6.127 | | | 0.9123 |
| 12 | 0.013 | 0.062 | 0.164 | 0.343 | 0.652 | 1.179 | 1.857 | 2.616 | 3.510 | 4.659 | 6.385 | | 0.9248 |
| 13 | 0.011 | 0.050 | 0.133 | 0.275 | 0.511 | 0.903 | 1.481 | 2.140 | 2.881 | 3.762 | 4.902 | 6.621 | 0.9349 |

| | | | | | | | | | | | | | |
|----|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|--------|
| 13 | 0.011 | 0.050 | 0.133 | 0.275 | 0.511 | 0.903 | 1.481 | 2.140 | 2.881 | 3.762 | 4.902 | 6.621 | 0.9349 |
|----|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|--------|



Функции модности критерия χ^2 при уровне значимости $\alpha = 0.05$,
 объёме выборки $N = 1000$, числе интервалов группирования $K = 10$:

1 - при оптимальном группировании;

2 - при равновероятных интервалах

частоты были пропорциональны вероятностям, соответствующим оптимальному группированию. Последнюю процедуру группирования целесообразно использовать при выравнивании эмпирических распределений по определенным законам.

Число интервалов в случае использования оптимального группирования определяется очень просто. Его выбирают таким образом, чтобы произведение объема выборки на вероятность попадания в крайний интервал с наименьшей вероятностью было больше 1.

Для проверки гипотез об отдельных параметрах по критерию χ^2 полезно использовать таблицы оптимального группирования, минимизирующие потери информации о соответствующих параметрах, приводимые в [7].

Л и т е р а т у р а

1. Кендалл М., Стьюарт А. Статистические выводы и связи. М., "Наука", 1973.
2. Куллдорф Г. Введение в теорию оценивания по группированным и частично группированным выборкам. М., "Наука", 1966.
3. Бодин Н.А. Оценка параметров распределения по группированным выборкам. - Тр. Матем. инст. АН СССР, III, 1970, 110-154.
4. Рао С.Р. Линейные статистические методы и их применение. М., "Наука", 1968.
5. Денисов В.И. Математическое обеспечение системы ЭВМ - экспериментатор. М., "Наука", 1977.
6. Лемешко Б.Д. Об оценивании параметров распределений по группированным наблюдениям. - В сб.: Вопросы кибернетики. М., 1977, вып. 30.
7. Губинский А.И., Денисов В.И., Гречко Ю.П., Лемешко Б.Д., Цой Е.Б. Методические рекомендации по планированию экспериментов и обработке экспериментальных данных при исследовании надежности и качества функционирования СЧТ. Л., ЛЭТИ, 1978.