

Санкт-Петербургский политехнический университет Петра Великого
Институт прикладной математики и механики
Высшая школа прикладной математики и вычислительной физики

ОТЧЁТ ПО ЛАБОРАТОРНОЙ РАБОТЕ №7

по дисциплине
«Математическая статистика»

Выполнила студентка
группы 3630102/80401

Мамаева Анастасия Сергеевна

Проверил
Доцент, к.ф.-м.н.

Баженов Александр Николаевич

Санкт-Петербург
2021

СОДЕРЖАНИЕ

СПИСОК ТАБЛИЦ	3
1 Постановка задачи	4
1.1 Задание	4
2 Теория	4
2.1 Метод максимального правдоподобия	4
2.2 Проверка гипотезы о законе распределения генеральной совокупности. Метод хи-квадрат	5
3 Программная реализация	9
4 Результаты	9
4.1 Проверка гипотезы о законе распределения генеральной совокупности. Метод хи-квадрат	9
5 Обсуждение	10
5.1 Проверка гипотезы о законе распределения генеральной совокупности. Метод хи-квадрат	10
6 Приложение	11

СПИСОК ТАБЛИЦ

1	Вычисление χ_B^2 при проверке гипотезы H_0 о нормальном законе распределения $N(x, \hat{\mu}, \hat{\sigma})$	9
2	Вычисление χ_B^2 при проверке гипотезы H_0 о законе распределения $L(x, \hat{\mu}, \hat{\sigma})$, $n = 20$	10
3	Вычисление χ_B^2 при проверке гипотезы H_0 о законе распределения $U(x, \hat{\mu}, \hat{\sigma})$, $n = 20$	10

1 Постановка задачи

1.1 Задание

Сгенерировать выборку объёмом 100 элементов для нормального распределения $N(x, 0, 1)$. По сгенерированной выборке оценить параметры μ и σ нормального закона методом максимального правдоподобия. В качестве основной гипотезы H_0 будем считать, что сгенерированное распределение имеет вид $N(x, \hat{\mu}, \hat{\sigma})$. Проверить основную гипотезу, используя критерий согласия χ^2 . В качестве уровня значимости взять $\alpha = 0.05$. Привести таблицу вычислений χ^2 .

Исследовать точность (чувствительность) критерия χ^2 — сгенерировать выборки равномерного распределения и распределения Лапласа малого объема (например, 20 элементов). Проверить их на нормальность.

2 Теория

2.1 Метод максимального правдоподобия

Одним из универсальных методов оценивания является метод максимального правдоподобия, предложенный Р.Фишером (1921). Пусть x_1, \dots, x_n — случайная выборка из генеральной совокупности с плотностью вероятности $f(x, \theta)$; $L(x_1, \dots, x_n, \theta)$ — функция правдоподобия (ФП), представляющая собой совместную плотность вероятности независимых с.в. x_1, \dots, x_n и рассматриваемая как функция неизвестного параметра θ :

$$L(x_1, \dots, x_n, \theta) = f(x_1, \theta)f(x_2, \theta)\dots f(x_n, \theta) \quad (1)$$

Определение. Оценкой максимального правдоподобия (о.м.п) будем называть такое значение $\hat{\theta}_{\text{мп}}$ из множества допустимых значений параметра θ , для которого ФП принимает наибольшее значение при заданных x_1, \dots, x_n :

$$\hat{\theta}_{\text{мп}} = \arg \max_{\theta} L(x_1, \dots, x_n, \theta) \quad (2)$$

Если ФП дважды дифференцируема, то её стационарные значения даются корнями уравнения

$$\frac{\partial L(x_1, \dots, x_n, \theta)}{\partial \theta} = 0 \quad (3)$$

Достаточным условием того, чтобы некоторое стационарное значение $\tilde{\theta}$ было локальным максимумом, является неравенство

$$\frac{\partial^2 L}{\partial \theta^2}(x_1, \dots, x_n, \tilde{\theta}) < 0 \quad (4)$$

Определив точки локальных максимумов ФП (если их несколько), находят наибольший, который и даёт решение задачи (1). Часто проще искать максимум логарифма ФП, так как он имеет максимум в одной точке с ФП:

$$\frac{\partial \ln L}{\partial \theta} = \frac{1}{L} \frac{\partial L}{\partial \theta}, \text{ если } L > 0 \quad (5)$$

и соответственно решать уравнение

$$\frac{\partial \ln L}{\partial \theta} = 0 \quad (6)$$

которое называют *уравнением правдоподобия*. В задаче оценивания векторного параметра $\theta = (\theta_1, \dots, \theta_m)$ аналогично (2) находится максимум ФП нескольких аргументов:

$$\hat{\theta}_{\text{МП}} = \arg \max_{\theta_1, \theta_2, \dots, \theta_m} L(x_1, x_2, \dots, x_n, \theta_1, \theta_2, \dots, \theta_m) \quad (7)$$

и в случае дифференцируемости ФП выписывается система уравнений правдоподобия

$$\frac{\partial L}{\partial \theta_k} = 0 \text{ или } \frac{\partial \ln L}{\partial \theta_k} = 0, k = 1, \dots, m \quad (8)$$

2.2 Проверка гипотезы о законе распределения генеральной совокупности. Метод хи-квадрат

Исчерпывающей характеристикой изучаемой случайной величины является её закон распределения. Поэтому естественно стремление исследователей построить этот закон приближённо на основе статистических данных.

Сначала выдвигается гипотеза о виде закона распределения.

После того как выбран вид закона, возникает задача оценивания его параметров и проверки (тестирования) закона в целом.

Для проверки гипотезы о законе распределения применяются критерии согласия. Таких критериев существует много. Мы рассмотрим наиболее обоснованный и наиболее часто используемый в практике — критерий χ^2 (хи-квадрат), введённый К.Пирсоном (1900 г.) для случая, когда параметры распределения известны. Этот критерий был существенно уточнён Р.Фишером (1924 г.), когда параметры распределения оцениваются по выборке, используемой для проверки.

Мы ограничимся рассмотрением случая одномерного распределения.

Итак, выдвинута гипотеза H_0 о генеральном законе распределения с функцией распределения $F(x)$.

Рассматриваем случай, когда гипотетическая функция распределения $F(x)$ не содержит неизвестных параметров.

Разобьём генеральную совокупность, т.е. множество значений изучаемой случайной величины X на k непересекающихся подмножеств $\Delta_1, \Delta_2, \dots, \Delta_k$.

Пусть $p_i = P(X \in \Delta_i), i = 1, \dots, k$.

Если генеральная совокупность — вся вещественная ось, то подмножества $\Delta_i = (a_{i-1}, a_i]$ — полуоткрытые промежутки ($i = 2, \dots, k-1$). Крайние промежутки будут полубесконечными: $\Delta_1 = (-\infty, a_1], \Delta_k = (a_{k-1}, +\infty)$. В этом случае $p_i = F(a_i) - F(a_{i-1})$; $a_0 = -\infty, a_k = +\infty$ ($i = 1, \dots, k$).

Отметим, что $\sum_{i=1}^k p_i = 1$. Будем предполагать, что все $p_i > 0$ ($i = 1, \dots, k$).

Пусть, далее, n_1, n_2, \dots, n_k — частоты попадания выборочных элементов в подмножества $\Delta_1, \Delta_2, \dots, \Delta_k$ соответственно.

В случае справедливости гипотезы H_0 относительные частоты n_i/n при большом n должны быть близки к вероятностям p_i ($i = 1, \dots, k$), поэтому за меру отклонения выборочного распределения от гипотетического с функцией $F(x)$ естественно выбрать величину

$$Z = \sum_{i=1}^k c_i \left(\frac{n_i}{n} - p_i \right)^2, \quad (9)$$

где c_i — какие-нибудь положительные числа (веса). К.Пирсоном в качестве весов выбраны числа $c_i = n/p_i$ ($i = 1, \dots, k$). Тогда получается статистика критерия хи-квадрат К.Пирсона

$$\chi^2 = \sum_{i=1}^k \frac{n}{p_i} \left(\frac{n_i}{n} - p_i \right)^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i}, \quad (10)$$

которая обозначена тем же символом, что и закон распределения хи-квадрат.

К.Пирсоном доказана теорема об асимптотическом поведении статистики χ^2 , указывающая путь её применения.

Теорема К.Пирсона. Статистика критерия χ^2 асимптотически распределена по закону χ^2 с $k-1$ степенями свободы.

Это означает, что независимо от вида проверяемого распределения, т.е. функции $F(x)$, выборочная функция распределения статистики χ^2 при $n \rightarrow \infty$ стремится к функции распределения случайной величины с плотностью вероятности

$$f_{k-1}(x) = \begin{cases} 0, & x \leq 0 \\ \frac{1}{2^{\frac{k-1}{2}} \Gamma(\frac{k-1}{2})} x^{\frac{k-3}{2}} e^{-\frac{x}{2}}, & x > 0 \end{cases} \quad (11)$$

Для прояснения сущности метода χ^2 сделаем ряд замечаний.

Замечание 1. Выбор подмножеств $\Delta_1, \Delta_2, \dots, \Delta_k$ и их числа k в принципе ничем не регламентируется, так как $n \rightarrow \infty$. Но так как число n хотя и очень большое, но конечное, то k должно быть с ним согласовано. Обычно его берут таким же, как и для построения гистограммы, т.е. можно руководствоваться формулой

$$k \approx 1.72 \sqrt[3]{n} \quad (12)$$

или формулой Старджесса

$$k \approx 1 + 3.3 \lg n \quad (13)$$

При этом, если $\Delta_1, \Delta_2, \dots, \Delta_k$ — промежутки, то их длины удобно сделать равными, за исключением крайних — полубесконечных.

Замечание 2. (о числе степеней свободы). Числом степеней свободы функции (по старой терминологии) называется число её независимых аргументов. Аргументами статистики χ^2 являются частоты n_1, n_2, \dots, n_k . Эти частоты связаны одним равенством $n_1 + n_2 + \dots + n_k = n$, а в остальном независимы в силу независимости элементов выборки. Таким образом, функция χ^2 имеет $k - 1$ независимых аргументов: число частот минус одна связь. В силу теоремы Пирсона число степеней свободы статистики χ^2 отражается на виде асимптотической плотности $f_{k-1}(x)$.

На основе общей схемы проверки статистических гипотез сформулируем следующее правило.

Правило проверки гипотезы о законе распределения по методу χ^2 .

1. Выбираем уровень значимости α .
2. По таблице [3, с. 358] находим квантиль $\chi^2_{1-\alpha}(k-1)$ распределения хи-квадрат с $k-1$ степенями свободы порядка $1 - \alpha$.
3. С помощью гипотетической функции распределения $F(x)$ вычисляем вероятности $p_i = P(X \in$

Δ_i), $i = 1, \dots, k$.

4. Находим частоты n_i попадания элементов выборки в подмножества Δ_i , $i = 1, \dots, k$.

5. Вычисляем выборочное значение статистики критерия χ^2 :

$$\chi_B^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i}. \quad (14)$$

6. Сравниваем χ_B^2 и квантиль $\chi_{1-\alpha}^2(k-1)$.

а) Если $\chi_B^2 < \chi_{1-\alpha}^2(k-1)$, то гипотеза H_0 на данном этапе проверки принимается.

б) Если $\chi_B^2 \geq \chi_{1-\alpha}^2(k-1)$, то гипотеза H_0 отвергается, выбирается одно из альтернативных распределений, и процедура проверки повторяется.

Замечание 3. Из формулы (10) видим, что веса $c_i = n/p_i$ пропорциональны n , т.е. с ростом n увеличиваются. Отсюда следует, что если выдвинутая гипотеза неверна, то относительные частоты n_i/n не будут близки к вероятностям p_i , и с ростом n величина χ_B^2 будет увеличиваться. При фиксированном уровне значимости α будет фиксировано пороговое число - квантиль $\chi_{1-\alpha}^2(k-1)$, поэтому, увеличивая n , мы придём к неравенству $\chi_B^2 > \chi_{1-\alpha}^2(k-1)$, т.е. с увеличением объёма выборки неверная гипотеза будет отвергнута.

Отсюда следует, что при сомнительной ситуации, когда $\chi_B^2 \approx \chi_{1-\alpha}^2(k-1)$, можно попытаться увеличить объём выборки (например, в 2 раза), чтобы требуемое неравенство было более чётким.

Замечание 4. Теория и практика применения критерия χ^2 указывают, что если для каких-либо подмножеств Δ_i ($i = 1, \dots, k$) условие $np_i \geq 5$ не выполняется, то следует объединить соседние подмножества (промежутки).

Это условие выдвигается требованием близости величин $\frac{(n_i - np_i)}{\sqrt{np_i}}$, квадраты которых являются слагаемыми χ^2 к нормальным $N(0, 1)$. Тогда случайная величина в формуле (10) будет распределена по закону, близкому к хи-квадрат. Такая близость обеспечивается достаточной численностью элементов в подмножествах Δ_i [1, с. 481-485].

3 Программная реализация

Лабораторная работа выполнена на языке Python версии 3.7 в среде разработки JupyterLab. Использовались дополнительные библиотеки:

1. scipy
2. math
3. matplotlib
4. numpy

В приложении находится ссылка на GitHub репозиторий с исходным кодом.

4 Результаты

4.1 Проверка гипотезы о законе распределения генеральной совокупности. Метод хи-квадрат

Метод максимального правдоподобия:

$$\hat{\mu} \approx -0.05, \hat{\sigma} \approx 0.95$$

Критерий согласия χ^2 :

Количество промежутков $k = 8$

Уровень значимости $\alpha = 0.05$

Тогда квантиль $\chi^2_{1-\alpha}(k-1) = \chi^2_{0.95}(7)$. Из таблицы [3, с. 358] $\chi^2_{0.95}(7) \approx 14.07$.

i	$limits$	n_i	p_i	np_i	$n_i - np_i$	$\frac{(n_i - np_i)^2}{np_i}$
1	[-inf, -1.1]	15	0.1357	13.57	1.43	0.15
2	[-1.1, -0.73]	9	0.096	9.6	-0.6	0.04
3	[-0.73, -0.37]	13	0.1253	12.53	0.47	0.02
4	[-0.37, 0.0]	20	0.1431	14.31	5.69	2.27
5	[0.0, 0.37]	9	0.1431	14.31	-5.31	1.97
6	[0.37, 0.73]	8	0.1253	12.53	-4.53	1.64
7	[0.73, 1.1]	14	0.096	9.6	4.4	2.02
8	[1.1, inf]	12	0.1357	13.57	-1.57	0.18
Σ	-	100	1	100	-0	8.27

Таблица 1: Вычисление χ^2_B при проверке гипотезы H_0 о нормальном законе распределения $N(x, \hat{\mu}, \hat{\sigma})$

Сравнивая $\chi_B^2 = 8.27$ и $\chi_{0.95}^2(7) \approx 14.07$, видим, что $\chi_B^2 < \chi_{0.95}^2(7)$.

i	$limits$	n_i	p_i	np_i	$n_i - np_i$	$\frac{(n_i - np_i)^2}{np_i}$
1	[-inf, -1.1]	3	0.1357	2.71	0.29	0.03
2	[-1.1, -0.37]	4	0.2213	4.43	-0.43	0.04
3	[-0.37, 0.37]	8	0.2861	5.72	2.28	0.91
4	[0.37, 1.1]	3	0.2213	4.43	-1.43	0.46
5	[1.1, inf]	2	0.1357	2.71	-0.71	0.19
Σ	-	20	1	20	-0	1.62

Таблица 2: Вычисление χ_B^2 при проверке гипотезы H_0 о законе распределения $L(x, \hat{\mu}, \hat{\sigma})$, $n = 20$

Сравнивая $\chi_B^2 = 1.62$ и $\chi_{0.95}^2(4) \approx 9.49$, видим, что $\chi_B^2 < \chi_{0.95}^2(4)$.

i	$limits$	n_i	p_i	np_i	$n_i - np_i$	$\frac{(n_i - np_i)^2}{np_i}$
1	[-inf, -1.1]	5	0.1357	2.71	2.29	1.93
2	[-1.1, -0.37]	5	0.2213	4.43	0.57	0.07
3	[-0.37, 0.37]	4	0.2861	5.72	-1.72	0.52
4	[0.37, 1.1]	4	0.2213	4.43	-0.43	0.04
5	[1.1, inf]	2	0.1357	2.71	-0.71	0.19
Σ	-	20	1	20	-0	2.75

Таблица 3: Вычисление χ_B^2 при проверке гипотезы H_0 о законе распределения $U(x, \hat{\mu}, \hat{\sigma})$, $n = 20$

Сравнивая $\chi_B^2 = 2.75$ и $\chi_{0.95}^2(4) \approx 9.49$, видим, что $\chi_B^2 < \chi_{0.95}^2(4)$.

5 Обсуждение

5.1 Проверка гипотезы о законе распределения генеральной совокупности. Метод хи-квадрат

Закключаем, что гипотеза H_0 о нормальном законе распределения $N(x, \hat{\mu}, \hat{\sigma})$ на уровне значимости $\alpha = 0.05$ согласуется с выборкой для нормального распределения $N(x, 0, 1)$.

Также видно, что для выборок сгенерированных по равномерному закону и закону Лапласа гипотеза H_0 оказалась принята.

6 Приложение

Код программы GitHub URL:

<https://github.com/Brightest-Sunshine/Math-Statistic-2021/blob/main/Lab7/Lab7.ipynb>