

МСПС | Критерии согласия и проверка нормальности

(1)

Задача определения типа распределения возникает во многих приложениях (Финансы, страхование и пр.). В страховании, например, важен распр. типа Парето, в т. частности - Вейбулла и его част. случаи нормальное распр. Критерии, основан. на нормальности данных, широко распр.ст. из-за 1) УПТ и 2) удобства вывода критериев. Сначала рассм. задачу проверки гипотезы о том, что ~~выбр.~~ ^{испыт.} распределение - какое-то конкретное (напр., $N(0,1)$), а не класс распр. (напр., все норм.). Пусть есть выборка X_1, \dots, X_n из неув. распр. с ф.р. F .

Рассм. гипотезу $H_0: F = F_0$ и альтер. $H_1: F \neq F_0$. Критерии проверки таких гипотез назыв. критериями согласия.

1) Критерий Колмогорова-Смирнова.

Рассм. $D_n = \sup_{x \in \mathbb{R}} |F_0(x) - \hat{F}_n(x)|$, где $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x)$ - эмпирич. ф.р.

Т. Колмогорова Если F_0 -непр., то при верной гипотезе H_0

$$P(\sqrt{n} D_n \leq x) \rightarrow K(x) = \sum_{k \in \mathbb{Z}} (-1)^k e^{-2k^2 x^2}$$

2) Критерий хи-квадрат Пирсона.

Разобьём \mathbb{R} на b -э интервалов $\{B_i\}_{i=1}^b$ так, чтобы $P_0(B_i) \cdot n \geq 5 \quad \forall i$ (размер выборки). Обозн. $p_i^0 = P_0(B_i)$; $\mu_i = \#\{j: X_j \in B_i\}$, тогда

Т. Карна Пирсона $\hat{\chi} = \sum_{i=1}^b \frac{(\mu_i - n p_i^0)^2}{n p_i^0} \xrightarrow{d} \chi_{b-1}^2$

Сам критерий (и в крит. К-С аналогично): если $\hat{\chi} > \chi_{1-\alpha, b-1}$, то отвергаем H_0 .

Здесь предпол. о непр. F_0 не нужно.

квантили распр. χ_{b-1}^2 уровня $1-\alpha$

Общее св-во, которыми обладают критерии согласия - их статистика не зависит от F_0 ! Это очень важно ^{для постро. критериев} и позволяет строить предельные распр. (иначе для каждого распр. оно было бы своим). Однако ф-ции мощностей на разных типах распр. будут разными. Ещё критерии согласия: Кроме того, все критерии согласия - состоятельные.

- Смирнова-Крамера-фон Мизеса, со стат. $\int (F_n(x) - F(x))^2 dF(x)$
- Андерсона-Дарлингта (ω^2 -крит.) $\int_{\mathbb{R}} \frac{(F_n(x) - F(x))^2}{F(x)(1-F(x))} dF(x)$
- Фридмана $\int |F_n(x) - F(x)| dF(x)$
- Ватсона $\int [F_n(x) - F(x) - \int (F_n(x) - F(x)) dF(x)]^2 dF(x)$

Иногда всё же исп. критерии, кот. зависят от вида распр.:

- Джамини: $\int |F_n(x) - F(x)| dx$
- Крамера-фон Мизеса: $\int_{\mathbb{R}} (F_n(x) - F(x))^2 dx$

Критерии согласия и проверка нормальности

(2)

Критерии согласия хороши ещё и тем, что они работают для произвольных распределений, однако на практике применяются критерии, которые имеют большую мощность по сравнению с критерием согласия, но ~~для~~ которые для других типов распределений не работают (напр., критерий Шапиро-Уилка), но чтобы их использовать, надо сначала примерно понять, что ~~я~~ распределение имеет (с помощью график. методов анализа).

• Есть общий метод проверки гипотезы $H_0: P \in \mathcal{P}_0 = \{P_\theta, \theta \in \Theta\}$ vs

$H_1: P \notin \mathcal{P}_0$:

Берём некие ~~с~~ состоят. оценки параметров $\theta = (\theta_1, \dots, \theta_d)$, тогда статистика критерия χ^2 -квдрат: $\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} \rightarrow \chi^2_{k-1-d}$

при не-усл. на дифф.

$P: (\theta_j)$

лучше так:

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \sum \frac{(\mu_i - np_i(\theta))^2}{np_i(\theta)}$$

считается с неким распределением P_θ , где $P_\theta \in \mathcal{P}_0$ и $\hat{\theta}$ - наша оценка

Проверка нормальности

• Так же можно поступить и с крит. Колмогорова - Смирнова

Всё Проверка гипотезы $H_0: P \in \{N(a, \sigma^2)\}$ против $H_1: P \notin \{N(a, \sigma^2)\}$

Нужно выбрать статистику $D_n = \sup_{x \in \mathbb{R}} |F_n^*(x) - \Phi(x)|$, где

F_n^* - эмпир. ф.р., $\Phi(x)$ - ф.р. $N(\bar{x}, s^2)$, где \bar{x} и s^2 - выбороч. сред. и дисперсия.

Всё D_n сходится к так наз. распредел. Лиллиефорса.

Недостатки: имеет низкую мощность; если хвосты отклоняются от нормальных, то не различает это.

• Критерий Шапиро-Уилка (самый лучший)

~~$H_0: P$~~ $H_0: X \sim N(\mu, \sigma^2)$; $H_1: X \not\sim N(\mu, \sigma^2)$

Статистика $W = \frac{(\sum_{i=1}^n a_i X_{(i)})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$, где $(a_1, \dots, a_n) = \frac{m^T V^{-1}}{\|m^T V^{-1}\|^2}$, где

$m = (m_1, \dots, m_n)^T$ - матрица порядковых статистик выборки из $N(0, 1)$, V - их ковар. матрица.

W при верной H_0 имеет таблич. распредел., значения (a_1, \dots, a_n) также табу- мированы.

Есть ещё критерий K^2 Д'Агостино, он тоже хороший, но хуже, чем Шапиро-Уилка. Он основан на статистиках skewness и kurtosis (выбороч. асимметрия и выбороч. эксцесс (поги)).

Есть более простой критерий, основ. на тех статистиках:

Критерии согласия и проверка нормальности

(3)

Критерий Харка-Бера (в рус.-яз. мире ещё ув. как Харке-Бера)

$$JB = \frac{n}{6} (\delta K^2 + \frac{1}{4} K^2), \text{ где } \delta K = \frac{m_3}{(m_2)^{3/2}}; K = \frac{m_4}{(m_2)^2} - 3, \text{ где}$$

$$m_j = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^j$$

При $n > 2000$ при верной H_0 : $JB \xrightarrow{d} \chi^2(2)$. При малых n следует искать квантили с пом. моделирования.

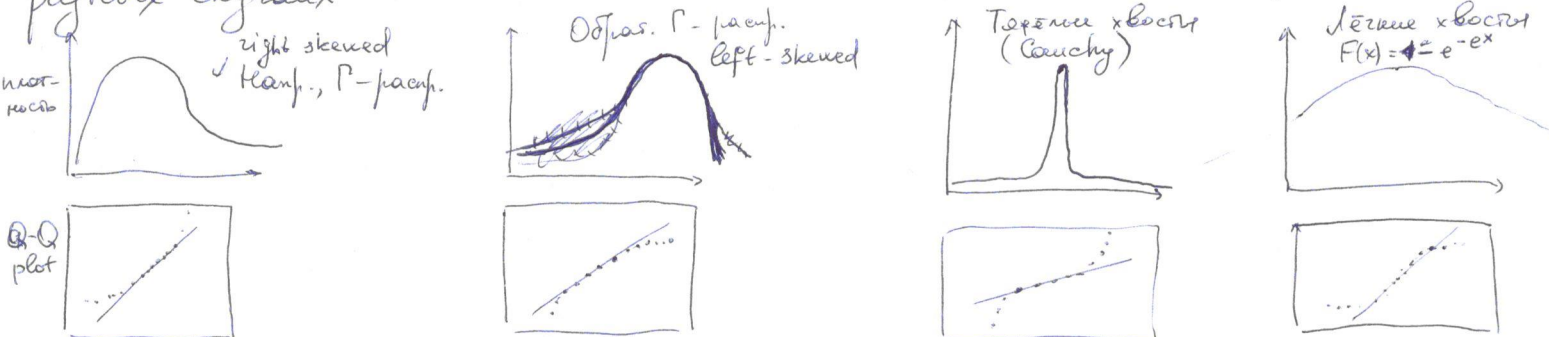
Сам крит., естеств., такой: отвергаю H_0 , если $JB > u_{1-\alpha}$.

Есть ещё QQ-plots и PP-plots (последняя реже используется)

Проверка гипотезы $H_0: F = F_0(\frac{x-\mu}{\sigma})$ (т.е. $F \in$ семейству распредел. с пар. разм. сдвига и масштаба, например, $N(\mu, \sigma^2)$)

Делается след.: на график наносится точки $(X_{(i)}, F_0^{-1}(\frac{i-0.5}{n}))$, $i=1 \dots n$. ~~Если $F_0(\frac{X_{(i)}-\mu}{\sigma})$ будет иметь распредел. $N(0,1)$ (известно заранее), то $F_0(X_{(i)})$ будет иметь распр. $N(\mu, \sigma^2)$ (т.е. $F_0(X_{(i)}) = \mu + \sigma \cdot N(0,1)$). Распр. $N(\mu, \sigma^2)$ и $N(0,1)$ будут иметь на прямой $A \cdot E X_{(i)} = u_{\frac{i}{n}}$, где $u_{\frac{i}{n}}$ - квантили $N(\mu, \sigma^2)$.~~ ~~Потому именно так, см. стр. 4~~

На примере норм. закона посмотрим, как ведут себя QQ-plots в разных случаях



Общие выводы о проверке нормальности:

- 1) Если выборка мала, то ничего нельзя сказать, и по графикам тоже.
- 2) Если выборка большая, то \forall крит. может выявлять небольшие отклонения, кот. не вл. значимыми для практики.
- 3) Критерий Лиллиефорса предст. только историч. ценность. Лучшее Шампир-Уилка
- 4) Критерий хи-квадрат тоже не самый мощный и довольно общий, но статистиками используется.

Вопрос: а что делать с дискрет. распределением?

- 1) Крит. хи-квадрат; 2) сделать свёртку с каким-нибудь ~~непр.~~ адром (см. оценку плотности) - и тогда будет непр. распредел.

Для проверки гипотезы $H_0: p = p_0$ в распр. Бернулли можно брать статистику

$$\frac{\bar{X} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \xrightarrow{d} N(0,1) \text{ по ЦПТ.}$$

Теорема Пирса о χ^2 -критерии

Итак, проверяем стат. гипотезу $H_0: P \in \mathcal{P} = \{P_\theta, \theta \in \Theta\}$, $\dim \Theta = d$
 против альтернативы $H_1: P \notin \mathcal{P}$.
 Как и ранее, разбиваем чис. ось на интервалы $\{B_i\}_{i=1}^k$, $\bigcup B_i = \mathbb{R}$.
 $\mu_i = \# \{X_j: X_j \in B_i\}$; $p_i(\theta) = P_\theta(B_i)$, где $k > d$

Теорема Пирса

Пусть Θ - отк. мн-во в \mathbb{R}^d . Пусть вып. условие

- 1) $\forall \theta \in \Theta$ ~~$p_i(\theta) > 0$~~ $p_i(\theta) > c > 0 \quad \forall i$
- 2) $\forall \theta \in \Theta$ $p_i(\theta)$ дважды непрерывно дифф. $\forall i$
- 3) $\forall \theta \in \Theta$ матрица $\left\| \frac{\partial p_i(\theta)}{\partial \theta_j} \right\|_{\substack{i,j=1 \dots k \\ j=1 \dots d}}$ имеет ранг d .

Пусть $\hat{\theta}$ - оценка, найд. методом макс. прав. по выборке μ_1, \dots, μ_k ,
 т.е. $\hat{\theta} = \arg \max_{\theta \in \Theta} L(\{\mu_i\}, \theta)$, где $L(\{\mu_i\}, \theta) = \frac{n!}{\mu_1! \dots \mu_k!} \prod_{i=1}^k p_i^{\mu_i}(\theta)$.

или $\hat{\theta}$ - оценка по минимуму хи-квадрат:

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \sum_{i=1}^k \frac{(\mu_i - n p_i(\theta))^2}{n p_i(\theta)}$$

Тогда если H_0 верна, то

$$\chi^2(\hat{\theta}) = \sum_{i=1}^k \frac{(\mu_i - n p_i(\hat{\theta}))^2}{n p_i(\hat{\theta})} \xrightarrow{n \rightarrow \infty} \chi^2_{k-d-1}$$

Объяснение, почему работает Q-Q plot.

Вроде как точки $(X_{(i)}, F_0^{-1}(\frac{i-0.5}{n}))$ должны лежать на одной прямой, если X_i - выборка из $F_0(\frac{x-\mu}{\sigma})$.

Итак, пусть Y_1, \dots, Y_n - выборка из распредел. с ф.р. F_0 , тогда

$\frac{Y_{(i)}}{p} \xrightarrow{p \rightarrow 0} F_0^{-1}(\frac{i}{n})$ по т. о выбороч. квантили, т.е. можно считать, что $Y_{(i)} \approx F_0^{-1}(\frac{i-0.5}{n})$

Но $\frac{X_{(i)} - \mu}{\sigma} \stackrel{d}{=} Y_{(i)}$. Отсюда следует, что $\frac{X_{(i)} - \mu}{\sigma} \stackrel{d}{=} Y_{(i)}$,

т.е. точки $(X_{(i)}, F_0^{-1}(\frac{i-0.5}{n}))$ должны лежать недалеко от прямой $\frac{x-\mu}{\sigma} = y$.