

Часть первая

Пусть произведено N опытов со случайной величиной Бернулли с вероятностью p выпадения единицы в одном отдельно взятом испытании. Получены N значений (нули и единицы): X_1, X_2, \dots, X_N . Пусть единиц из них K . Это число можно трактовать как сумму независимых случайных величин $K = \sum_{i=1}^N X_i$. Согласно центральной предельной теореме

сумма большого количества одинаково распределённых независимых случайных величин (имеющих математические ожидания и дисперсии) с ростом числа слагаемых стремится

к нормальной случайной величине. Количество успехов $K = \sum_{i=1}^N X_i$ как раз и есть такая

сумма. То есть при большом числе испытаний N случайная величина K приближённо может считаться нормальной. Для упрощения её можно «центрировать» (то есть вычесть из неё её же математическое ожидание, чтобы в итоге математическое ожидание оказалось равно нулю) и «нормировать» (то есть разделить на её же среднеквадратическое отклонение, чтобы в итоге дисперсия оказалась равна единице). Математическое ожидание $EK = Np$. Дисперсия $DK = Np(1-p)$. Среднеквадратическое отклонение $\sigma(K) = \sqrt{DK} = \sqrt{Np(1-p)}$. Итак, стандартное нормальное распределение будет иметь

$$\text{величина } Z = \frac{K - EK}{\sigma(K)} = \frac{K - Np}{\sqrt{Np(1-p)}}.$$

Построим критерий для оценки гипотезы $H_0: p = \theta$.

Альтернативной гипотезой выберем $H_1: p \neq \theta$.

Зададимся уровнем значимости α (это вероятность отвергнуть гипотезу H_0 , если на самом деле она верная). Тогда вероятность принять (точнее: не найти причины отвергнуть) гипотезу H_0 , если на самом деле она верная будет равна $(1-\alpha)$. При правильности гипотезы H_0 вероятность $p = \theta$, и тогда

$$Z = \frac{K - N\theta}{\sqrt{N\theta(1-\theta)}}$$

Запишем равенство из которого можно было бы найти промежуток значений K , для которых принимается гипотеза H_0 .

$$1 - \alpha = P\left(-c < \frac{K - N\theta}{\sqrt{N\theta(1-\theta)}} < c\right) = P(-c < Z < c) = \Phi(c) - \Phi(-c) = \Phi(c) - (1 - \Phi(c)) = 2\Phi(c) - 1$$

$$1 - \alpha = 2\Phi(c) - 1$$

$$2 - \alpha = 2\Phi(c)$$

$$\Phi(c) = 1 - \frac{\alpha}{2}$$

В программе Microsoft Office Excel встроена функция НОРМСТОБР(вероятность), с помощью которой можно находить значения функции, обратной к функции распределения $\Phi(x)$. Обозначим здесь эту функцию $G(x)$. Тогда

$$c = G\left(1 - \frac{\alpha}{2}\right)$$

Итак, c это известное число и можно относительно случайной величины K получить критерий

$$\begin{aligned}
-c &< \frac{K - N\theta}{\sqrt{N\theta(1-\theta)}} < c \\
-c\sqrt{N\theta(1-\theta)} &< K - N\theta < c\sqrt{N\theta(1-\theta)} \\
N\theta - c\sqrt{N\theta(1-\theta)} &< K < N\theta + c\sqrt{N\theta(1-\theta)}
\end{aligned}$$

Если это неравенство выполняется, то нет оснований отвергать гипотезу H_0 . Произведённые действия арифметически полностью аналогичны тем, которые были сделаны при нахождении доверительного интервала. Только логика обратная. Доверительный интервал находился так, чтобы он содержал истинное значение параметра. А критерий проверяет, входит ли некоторое значение θ в доверительный интервал. Создаётся впечатление, что методика нахождения доверительного интервала более гибкая для использования на практике. Однако если тестировать нужно выполнение какой либо конкретной гипотезы, то использование статистического критерия оказывается более простой задачей, чем построение доверительного интервала. Кроме того, используя тестируемую гипотезу как исходное данное, иногда удаётся получить более точные результаты.

Часть вторая

Пусть есть не одна случайная величина Бернулли, а две случайных величины: X и Y . Пусть у первой случайной величины Бернулли параметр p_X , а у второй p_Y . Пусть с первой из них произведено N_X опытов, а со второй N_Y опытов. Получены две последовательности (возможно разной длины) из нулей и единиц. Пусть в первом случае получены K_X единиц, а во втором случае K_Y единиц. Построим критерий, проверяющий равенство $p_X = p_Y$, при альтернативной гипотезе $p_X \neq p_Y$. Создадим случайную величину, которую можно принять за меру отклонения параметров p_X и p_Y друг от друга

$$U = \frac{K_X}{N_X} - \frac{K_Y}{N_Y}$$

Так разумно поступить, поскольку $EK_X = N_X p_X$, а $EK_Y = N_Y p_Y$ и поэтому математическое ожидание

$$\begin{aligned}
EU &= E\left(\frac{K_X}{N_X} - \frac{K_Y}{N_Y}\right) = E\left(\frac{K_X}{N_X}\right) - E\left(\frac{K_Y}{N_Y}\right) = E\left(\frac{1}{N_X} \cdot K_X\right) - E\left(\frac{1}{N_Y} \cdot K_Y\right) = \\
&= \frac{1}{N_X} \cdot E(K_X) - \frac{1}{N_Y} \cdot E(K_Y) = \frac{1}{N_X} \cdot N_X \cdot p_X - \frac{1}{N_Y} \cdot N_Y \cdot p_Y = p_X - p_Y
\end{aligned}$$

По свойству линейности случайная величина U тоже будет нормальной, как линейная комбинация нормальных случайных величин. Обратим внимание, что при построении критерия будем активно пользоваться тем, что $p_X = p_Y = p$. В частности $EK_X = N_X p_X = N_X p$, а $EK_Y = N_Y p_Y = N_Y p$. Тогда

$$EU = p_X - p_Y = p - p = 0$$

Случайная величина U уже оказалась центрированной.

Для удобства её ещё нужно нормировать, то есть разделить на её среднеквадратическое отклонение. Найдём её дисперсию. При этом вспомним, что $DK_X = N_X p_X(1-p_X) = N_X p(1-p)$, а $DK_Y = N_Y p_Y(1-p_Y) = N_Y p(1-p)$.

$$DU = D\left(\frac{K_X}{N_X} - \frac{K_Y}{N_Y}\right) = D\left(\frac{K_X}{N_X}\right) + D\left(\frac{K_Y}{N_Y}\right) =$$

Обратите внимание, что дисперсия (независимых случайных величин) разности равна не разности, а сумме дисперсий. А постоянные множители выносятся за знак дисперсии в квадратах.

$$\begin{aligned}
&= D\left(\frac{1}{N_X} \cdot K_X\right) + D\left(\frac{1}{N_Y} \cdot K_Y\right) = \left(\frac{1}{N_X}\right)^2 \cdot D(K_X) + \left(\frac{1}{N_Y}\right)^2 \cdot D(K_Y) = \\
&= \frac{1}{(N_X)^2} \cdot D(K_X) + \frac{1}{(N_Y)^2} \cdot D(K_Y) = \frac{1}{(N_X)^2} \cdot N_X p(1-p) + \frac{1}{(N_Y)^2} \cdot N_Y p(1-p) = \\
&= \frac{p(1-p)}{N_X} + \frac{p(1-p)}{N_Y} = p(1-p) \left(\frac{1}{N_X} + \frac{1}{N_Y}\right)
\end{aligned}$$

Среднеквадратическое отклонение

$$\sigma(U) = \sqrt{DU} = \sqrt{p(1-p) \left(\frac{1}{N_X} + \frac{1}{N_Y}\right)}$$

Тогда случайная величина

$$Z = \frac{U}{\sigma(U)} = \frac{\frac{K_X}{N_X} - \frac{K_Y}{N_Y}}{\sqrt{p(1-p) \left(\frac{1}{N_X} + \frac{1}{N_Y}\right)}}$$

будет иметь стандартное нормальное распределение. Зададимся уровнем значимости α (это вероятность отвергнуть гипотезу H_0 , если на самом деле она верная). Тогда вероятность принять (точнее: не найти причины отвергнуть) гипотезу $H_0: p_X = p_Y$, если на самом деле она верная, будет равна $(1-\alpha)$. Это приводит к равенству

$$\begin{aligned}
1 - \alpha &= P\left(-c < \frac{\frac{K_X}{N_X} - \frac{K_Y}{N_Y}}{\sqrt{p(1-p) \left(\frac{1}{N_X} + \frac{1}{N_Y}\right)}} < c\right) = P(-c < Z < c) = \Phi(c) - \Phi(-c) = \\
&= \Phi(c) - (1 - \Phi(c)) = 2\Phi(c) - 1
\end{aligned}$$

Откуда следует знакомая арифметика

$$1 - \alpha = 2\Phi(c) - 1$$

$$2 - \alpha = 2\Phi(c)$$

$$\Phi(c) = 1 - \frac{\alpha}{2}$$

$$c = G\left(1 - \frac{\alpha}{2}\right)$$

Где $G(x)$ функция, обратная к функции распределения $\Phi(x)$.

Формально получен критерий

$$-c < \frac{\frac{K_X}{N_X} - \frac{K_Y}{N_Y}}{\sqrt{p(1-p) \left(\frac{1}{N_X} + \frac{1}{N_Y}\right)}} < c$$

Проблема его использования в том, что значение p неизвестно. Чтобы хоть как-то сделать вычисления возможными, заменим p на его точечную оценку, помня, что $p_X = p_Y = p$. Тогда X и Y это одна и та же случайная величина, с которой произвели $(N_X + N_Y)$ испытаний и в $(K_X + K_Y)$ из них была получена единица. То есть будем считать, что

$$p \approx \frac{K_X + K_Y}{N_X + N_Y}$$

Тогда критерий будет иметь вид

$$-c < \frac{\frac{K_X}{N_X} - \frac{K_Y}{N_Y}}{\sqrt{\frac{K_X + K_Y}{N_X + N_Y} \left(1 - \frac{K_X + K_Y}{N_X + N_Y}\right) \left(\frac{1}{N_X} + \frac{1}{N_Y}\right)}} < c$$

$$-c < \frac{\frac{K_X}{N_X} - \frac{K_Y}{N_Y}}{\sqrt{\frac{K_X + K_Y}{N_X + N_Y} \left(1 - \frac{K_X + K_Y}{N_X + N_Y}\right) \left(\frac{N_Y + N_X}{N_X N_Y}\right)}} < c$$

$$-c < \frac{\frac{K_X}{N_X} - \frac{K_Y}{N_Y}}{\sqrt{\frac{K_X + K_Y}{N_X N_Y} \left(1 - \frac{K_X + K_Y}{N_X + N_Y}\right)}} < c$$

Если это неравенство выполняется, то нет оснований отвергать гипотезу, что $p_X = p_Y$. Если неравенство не выполняется, то принимается гипотеза, что $p_X \neq p_Y$. Ещё раз обратите внимание на вид формулировки. Гипотезы не подтверждают, а говорят, что критерий не выявил оснований для отвержения гипотезы (возможно, что это не гипотеза правильная, а критерий не достаточно хороший).