

Домашнее задание 4

Юрасов Никита Андреевич

Обновлено 2 декабря 2019 г.

Содержание

1	Проверка гипотез о виде распределения	2
1.1	Простые гипотезы	2
1.1.1	Критерий согласия Колмогорова-Смирнова	2
1.1.2	Критерий согласия χ^2	4
1.2	Сложные гипотезы	6
1.2.1	Критерий Колмогорова	6
1.2.2	Критерий χ^2	7

1 Проверка гипотез о виде распределения

1.1 Простые гипотезы

1.1.1 Критерий согласия Колмогорова-Смирнова

Введем статистику, которая представляет собой максимальное отклонение эмпирической функции распределения $\hat{F}(x)$, построенной по выборке X , от гипотетической функции распределения $F(x)$:

$$D_n = D_n(X) = \sup_{-\infty < x < \infty} |\hat{F}(x) - F(x)|$$

Пусть существует $X = (X_1, \dots, X_n)$ – выборка из $\mathcal{L}(\xi)$ с неизвестной функцией распределения $F_\xi(x)$, и пусть выдвинута гипотеза $H_0 : F_\xi(x) = F(x)$, где функция $F(x)$ полностью задана.

Для принятия или отвержения гипотезы H_0 необходимо по критерию Колмогорова сравнить $\sqrt{n}D_n$ с λ_α , которая определяется следующим равенством:

$$K(\lambda_\alpha) = 1 - \alpha,$$

где $K(x)$ – распределение Колмогорова.

На практике статистику D_n удобнее вычислять в следующем виде $D_n = \max(D_n^+, D_n^-)$, где

$$D_n^+ = \max_{1 \leq k \leq n} \left(\frac{k}{n} - F(X_{(k)}) \right), \quad D_n^- = \max_{1 \leq k \leq n} \left(F(X_{(k)}) - \frac{k-1}{n} \right)$$

Ответ на вопрос о виде распределения дает следующее сравнение:

- Если $\sqrt{n}D_n \geq \lambda_\alpha$, то гипотеза H_0 отвергается;
- Если $\sqrt{n}D_n \leq \lambda_\alpha$, то гипотеза H_0 принимается.

В неравенстве можно воспользоваться поправкой Большева о статистике $S(D_n)$, которая быстрее сходится к распределению Колмогорова:

$$S = \frac{6nD_n + 1}{6\sqrt{n}}$$

Преимущества

Критерий согласия Колмогорова начинает эффективно работать при выборке объемом $n \geq 20$, что допускает использование его при достаточно малых выборках данных.

Недостатки

Критерий Колмогорова-Смирнова применяется только для непрерывных распределений. Также вычисление статистики D_n предполагает достаточно большие аналитические вычисления, что затрудняет проверку.

Реализация для непрерывного распределения

В приложенном Jupyter Notebook написана функция `simple_kolmogorov_test`, которая по заданной выборке и уровне значимости проверяет критерий согласия Колмогорова-Смирнова. Далее будут представлены только результаты, а саму выборку размера 1000 можно будет в переменной `erlang_sample_for_KSTest`. Выборка генерировалась с параметрами $m = 2, \lambda = 0.2$

Результаты

`S_Bolshev = 1.064432117481404 < 1.2238478702170825` and K-S test accepts with `alpha=0.1`
`S_Bolshev = 1.064432117481404 < 1.3580986393225505` and K-S test accepts with `alpha=0.05`

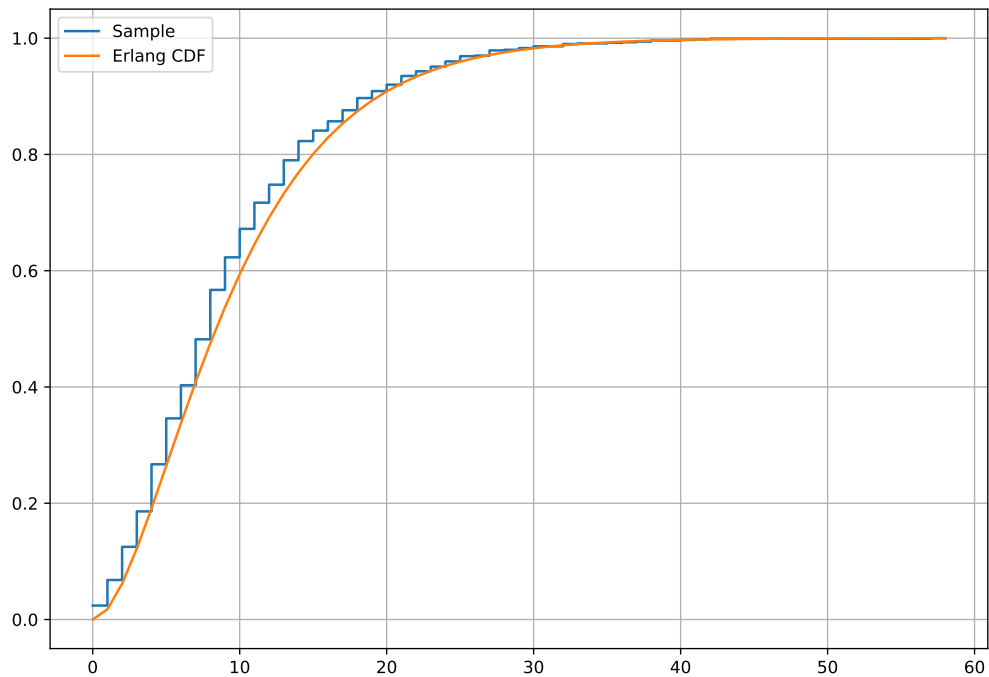


Рис. 1: Сравнение эмпирической функции и функции распределения Эрланга

1.1.2 Критерий согласия χ^2

Введем статистику \hat{X} , введенную К. Пирсоном, которая будет показывать отклонение эмпирических данных от гипотетических значений и которая называется мера *хи-квадрат*:

$$\hat{X}_n^2 = \hat{X}_n^2(\nu) = \sum_{i=1}^N \frac{(\nu_i - n\hat{p}_i)^2}{n\hat{p}_i}, \quad \text{где}$$

- N – количество принимаемых значений в эксперименте;
- $\nu = (\nu_1, \dots, \nu_N)$ – частоты появления каждого результата эксперимента;
- $n = \sum_{i=1}^N \nu_i$ – общий объем выборки;
- $\hat{p} = (\hat{p}_1, \dots, \hat{p}_N)$ – вероятность появления i -го события;

После подсчета *меры хи-квадрат* необходимо сравнить ее с критическим значением распределения хи-квадрат на уровне значимости α с $N - 1$ степенями свободы:

$$\chi_{1-\alpha, N-1}^2$$

Сравнение:

- Если $\hat{X}_n^2 > \chi_{1-\alpha, N-1}^2$, то говорят, что гипотеза H_0 отклоняется;
- Если $\hat{X}_n^2 \leq \chi_{1-\alpha, N-1}^2$, то говорят, что гипотеза H_0 принимается.

где гипотеза H_0 определена так же, как и в критерии согласия Колмогорова-Смирнова (см. страницу 2).

Преимущества

Критерий работает первоначально только с дискретными данными, но так как любые данные можно свести к дискретным методом группировки (см. правило Стёрджеса: [Википедия](#)). Также, этот критерий можно использовать для расчетов с хорошим приближением уже при $n \geq 50$.

Недостатки

Критерий χ^2 ошибается на выборках с низкочастотными (редкими) событиями. Решить эту проблему можно отбросив низкочастотные события, либо объединив их с другими событиями. Этот способ называется *коррекцией Йетса*.

Указания для проверки непрерывных распределений

Так как вероятность попадания в одну конкретную точку (в случае непрерывных распределений) равна 0, воспользуемся отмеченным ранее правилом Стёрджеса для разбиения отрезка на k не пересекающихся интервалов:

$$k = 1 + \lfloor \log_2 N \rfloor$$

Также необходимо вместо вектора гипотетических вероятностей в каждой точке использовать вероятность попадания в каждый из полученных интервалов. Для этого нужно вычислить значение интеграла:

$$\int_{x_i}^{x_{i+1}} f(x) dx, \quad \text{где}$$

$f(x)$ – плотность распределения, а x_i – точки разбиения отрезка.

Реализация для дискретного распределения

Сгенерируем выборку (распределение Пуассона) с параметром $\lambda = 2$ и размером 1000, которую будем хранить в переменной `poisson_sample_for_Chi2Test`. Результаты для двух уровней значимости выглядят следующим образом:

Результаты

`S = 9.32952003034001 <= 12.591587243743977 and Chi2 test accepts with alpha=0.05`
`S = 9.32952003034001 <= 10.644640675668422 and Chi2 test accepts with alpha=0.1`

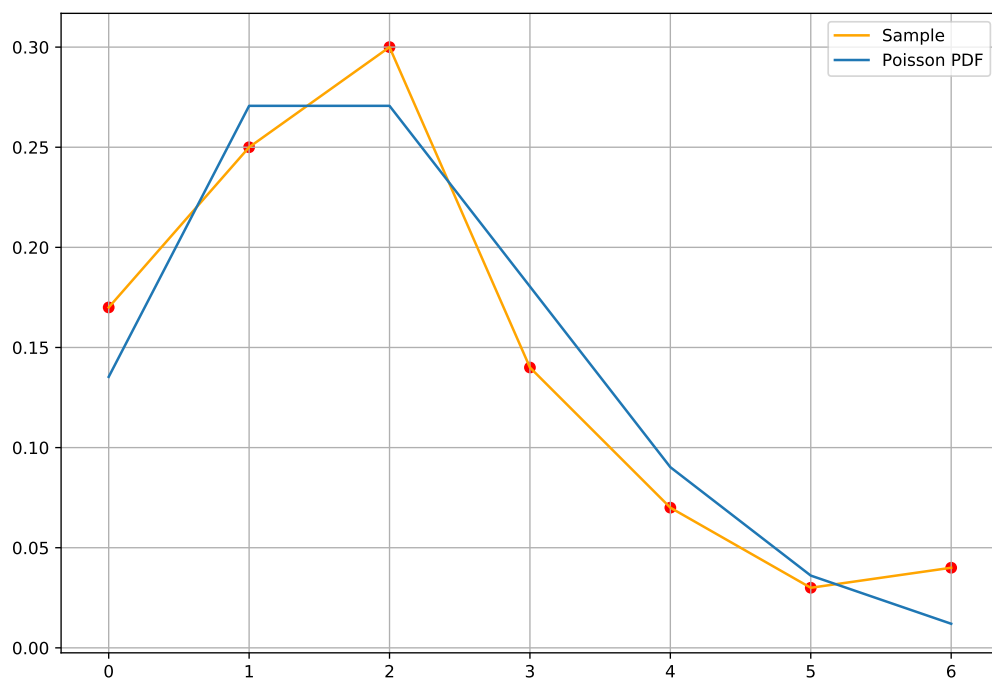


Рис. 2: Сравнение полигона частот и распределения Пуассона

Реализация для непрерывного распределения

Выборка, хранящаяся в переменной `erlang_sample_for_Chi2Test`, генерировалась с параметрами $m = 2$, $\lambda = 0.2$.

Результаты

`S = 5.527258712276454 <= 15.50731305586545 and Chi2 test accepts with alpha=0.05`
`S = 5.527258712276454 <= 13.36156613651173 and Chi2 test accepts with alpha=0.1`

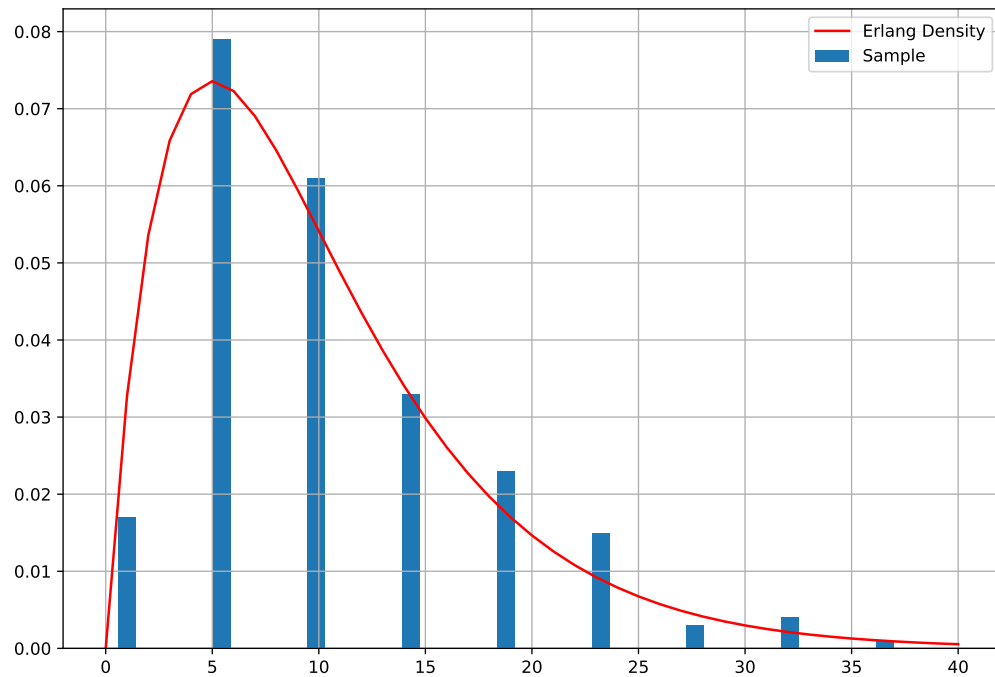


Рис. 3: Сравнение гистограммы частот и плотности распределения Эрланга

1.2 Сложные гипотезы

Для реализации сложных гипотез, будем использовать построение оценок параметров методом максимального правдоподобия. Каждую выборку, которую будем проверять по критериями хи-квадрат (дискретное и непрерывное распределение) и Колмогорова (только непрерывное распределение), разобьем предварительно на две части (50 на 50, но эту пропорцию можно варьировать): X_{est}, X_{test} . По первой части - X_{est} будем находить оценку максимального правдоподобия, а по X_{test} будем проверять критерий.

Напомним ОМП для распределений Пуассона и Эрланга:

- Пуассон: $\hat{\lambda} = \bar{X}$
- Эрланг: $\hat{\lambda} = \frac{m}{\bar{X}}$

1.2.1 Критерий Колмогорова

Критическая граница для критерия определяется следующим образом: $\tau_{1\alpha} = \{t \geq t_\alpha\}$, где $t_\alpha = \frac{\lambda_\alpha}{\sqrt{n}}$, а λ_α определяется из распределения Колмогорова $K(\lambda_\alpha) = 1 - \alpha$ при заданном уровне α .

Сгенерируем выборку `erlang_sample_for_complicated_KS_test` из 1000 элементов. Вот какие получились результаты:

S_Bolshev = 0.9483064237024541 < 1.3580986393225505 and K-S test accepts with alpha=0.05
 S_Bolshev = 0.9483064237024541 < 1.2238478702170825 and K-S test accepts with alpha=0.1

Значения критической границы указаны после знака <, соответствующие выбранному уровню α .

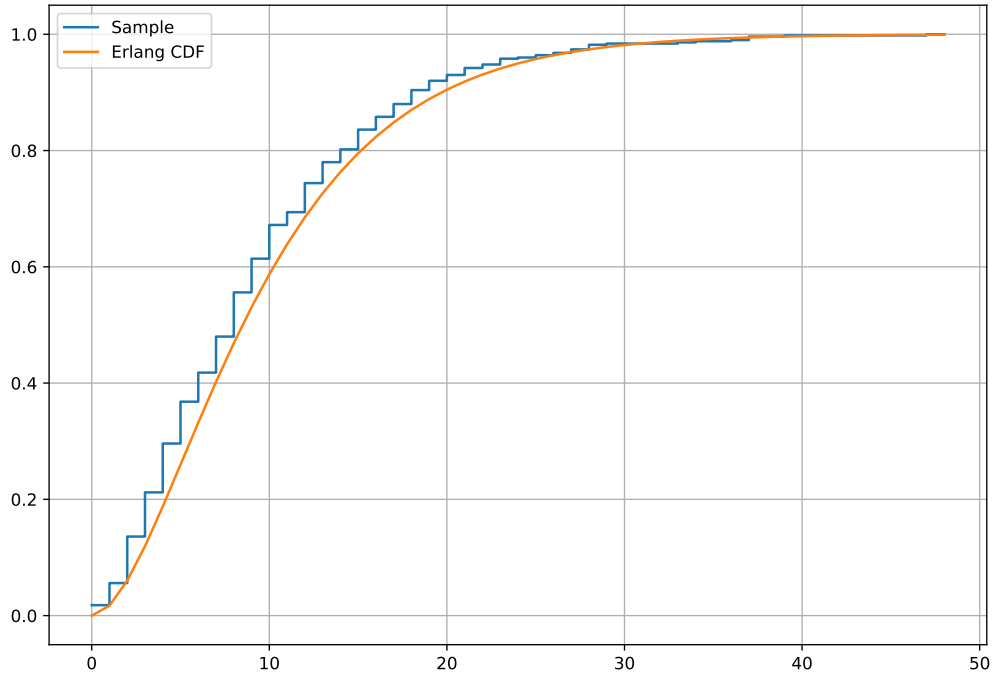


Рис. 4: Критерий согласия Колмогорова в случае сложной гипотезы

1.2.2 Критерий χ^2

Критическая граница $\tau_{1\alpha}$ рассчитывается как $\tau_{1\alpha} = \{\dot{X}_n^2 > t_\alpha\}$, где t_α вычисляется из распределения хи-квадрат с $N - 1 - r$ степенями свободы при заданном уровне α :

$$t_\alpha = \chi_{1-\alpha, N-1-r}^2$$

Реализация для непрерывного распределения

Сгенерируем выборку размера 1000 из распределения Эрланга, которая будет храниться в переменной `erlang_sample_for_complicated_Chi2_test`. Так как параметр m известен, то оцениваемый параметр один - λ .

S = 7.528579505480654 <= 14.067140449340169 and Chi2 test accepts with alpha=0.05
 S = 7.528579505480654 <= 12.017036623780532 and Chi2 test accepts with alpha=0.1

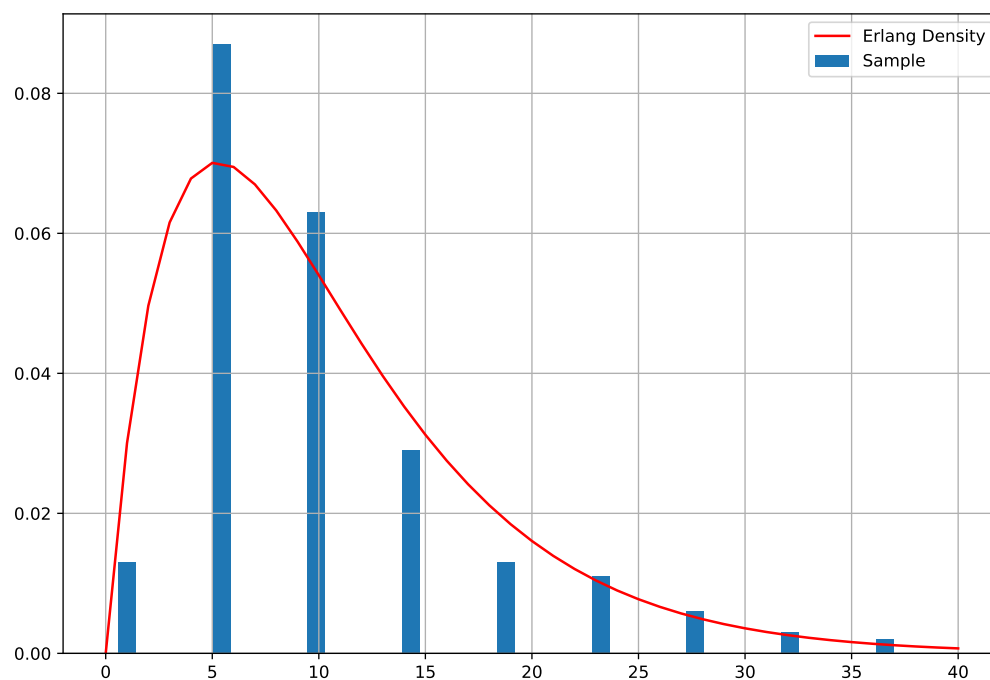


Рис. 5: Критерий хи-квадрат в случае сложной гипотезы для распределения Эрланга

Реализация для дискретного распределения

Создадим набор выборку данных из распределения Пуассона размера $n = 1000$ `poisson_sample_for_complicated_Chi2_test`. Функция плотности зависит всего от одного параметра, следовательно $r = 1$ - оцениваемый параметр.

Результаты

`S = 1.9708909736466178 <= 12.591587243743977 and Chi2 test accepts with alpha=0.05`
`S = 1.9708909736466178 <= 10.644640675668422 and Chi2 test accepts with alpha=0.1`

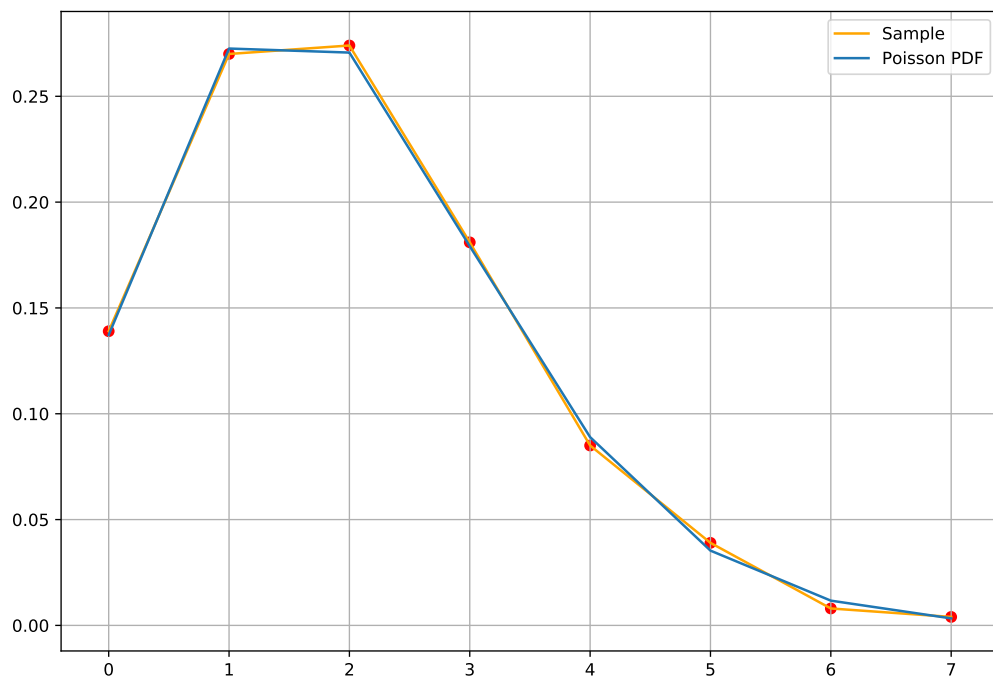


Рис. 6: Критерий хи-квадрат в случае сложной гипотезы для распределения Пуассона