

Проверка гипотез.

Часть 2

Лемма Неймана-Пирсона. Критерий отношения правдоподобия. Введение в А/В-тестирование. Критерий последовательного отношения правдоподобия. Непараметрические критерии.

ПМИ ФКН ВШЭ, 13 октября 2018 г.

Алексей Артемов^{1,2}

¹ Сколтех ² НИУ ВШЭ

Содержание лекции

- › Лемма Неймана-Пирсона.
- › Критерий отношения правдоподобия.
- › Функции штрафа и основы теории принятия статистических решений.
- › Введение в А/В-тестирование.
- › Критерий последовательного отношения правдоподобия.
- › Непараметрические критерии.

Лемма

Неймана-Пирсона

Показательный пример

- › Пусть $X_i \sim \mathcal{N}(\theta, \sigma^2)$, $i = 1, \dots, n$, а нулевая гипотеза заключается в том, что $\theta = \theta_0$.
- › Назовите достаточную статистику в этой задаче.

Показательный пример

- › Пусть $X_i \sim \mathcal{N}(\theta, \sigma^2)$, $i = 1, \dots, n$, а нулевая гипотеза заключается в том, что $\theta = \theta_0$.
- › Назовите достаточную статистику в этой задаче.
- › Достаточная статистика: $T(\mathbf{X}^\ell) = \overline{X}_n$.
- › Какое распределение $T(\mathbf{X}^\ell)$, когда верна \mathbb{H}_0 ?

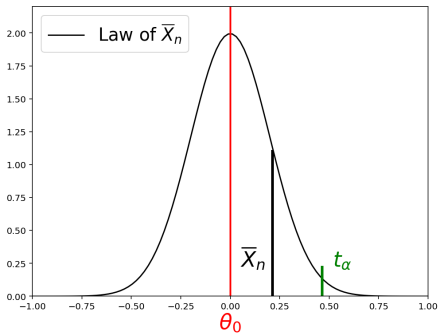
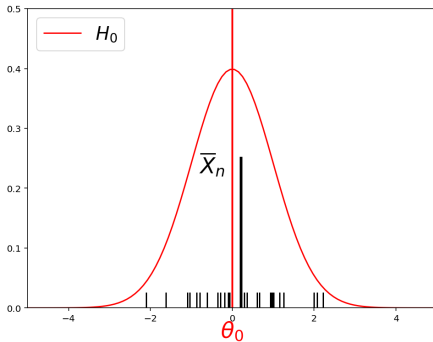
Показательный пример

- › Пусть $X_i \sim \mathcal{N}(\theta, \sigma^2)$, $i = 1, \dots, n$, а нулевая гипотеза заключается в том, что $\theta = \theta_0$.
- › Назовите достаточную статистику в этой задаче.
- › Достаточная статистика: $T(\mathbf{X}^\ell) = \bar{X}_n$.
- › Какое распределение $T(\mathbf{X}^\ell)$, когда верна \mathbb{H}_0 ?
- › $T(\mathbf{X}^\ell) \sim \mathcal{N}(\theta_0, \sigma^2/n)$
- › Какова критическая область? (Когда отклоняем \mathbb{H}_0 ?)

Показательный пример

- › Пусть $X_i \sim \mathcal{N}(\theta, \sigma^2)$, $i = 1, \dots, n$, а нулевая гипотеза заключается в том, что $\theta = \theta_0$.
- › Назовите достаточную статистику в этой задаче.
- › Достаточная статистика: $T(\mathbf{X}^\ell) = \bar{X}_n$.
- › Какое распределение $T(\mathbf{X}^\ell)$, когда верна \mathbb{H}_0 ?
- › $T(\mathbf{X}^\ell) \sim \mathcal{N}(\theta_0, \sigma^2/n)$
- › Какова критическая область? (Когда отклоняем \mathbb{H}_0 ?)
- › Критическая область $\mathcal{R}_\alpha = [t_\alpha, \infty)$, т.е. $T(\mathbf{X}^\ell) \geq t_\alpha$.

Показательный пример



Показательный пример

- › Подсчитайте вероятность ложной тревоги в этой задаче.

Показательный пример

› Подсчитайте вероятность ложной тревоги в этой задаче.

$$\begin{aligned}\alpha &= P_{\theta_0} \left(\frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma} \geq \frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right) = \\ &= 1 - \Phi \left(\frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right).\end{aligned}$$

› Как выбрать t_α , чтобы $\alpha \leq \alpha_0$?

Показательный пример

› Подсчитайте вероятность ложной тревоги в этой задаче.

$$\begin{aligned}\alpha &= P_{\theta_0} \left(\frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma} \geq \frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right) = \\ &= 1 - \Phi \left(\frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right).\end{aligned}$$

› Как выбрать t_α , чтобы $\alpha \leq \alpha_0$?

$$t_{\alpha_0} = \theta_0 + \sigma x_{1-\alpha_0} / \sqrt{n}$$

Показательный пример

- › Подсчитайте вероятность ложной тревоги в этой задаче.

$$\begin{aligned}\alpha &= P_{\theta_0} \left(\frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma} \geq \frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right) = \\ &= 1 - \Phi \left(\frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right).\end{aligned}$$

- › Как выбрать t_α , чтобы $\alpha \leq \alpha_0$?

$$t_{\alpha_0} = \theta_0 + \sigma x_{1-\alpha_0} / \sqrt{n}$$

- › Пусть на самом деле верна альтернатива $\mathbb{H}_1 : \theta = \theta_1$, причем $\theta_1 > \theta_0$. Какова вероятность ошибки 2-го рода?

Показательный пример

- › Подсчитайте вероятность ложной тревоги в этой задаче.

$$\begin{aligned}\alpha &= P_{\theta_0} \left(\frac{\sqrt{n}(\bar{X}_n - \theta_0)}{\sigma} \geq \frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right) = \\ &= 1 - \Phi \left(\frac{\sqrt{n}(t_\alpha - \theta_0)}{\sigma} \right).\end{aligned}$$

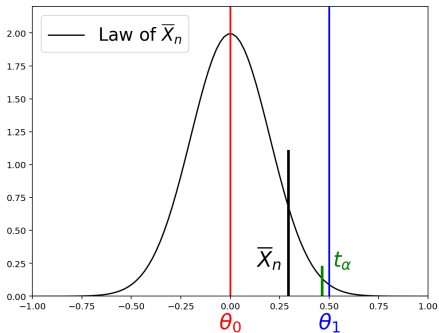
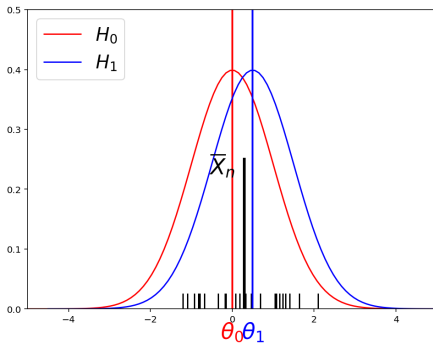
- › Как выбрать t_α , чтобы $\alpha \leq \alpha_0$?

$$t_{\alpha_0} = \theta_0 + \sigma x_{1-\alpha_0} / \sqrt{n}$$

- › Пусть на самом деле верна альтернатива $\mathbb{H}_1 : \theta = \theta_1$, причем $\theta_1 > \theta_0$. Какова вероятность ошибки 2-го рода?

$$\beta = P_{\theta_1} \left(\bar{X}_n < t_{\alpha_0} \right) = \Phi \left(x_{1-\alpha_0} - \frac{\sqrt{n}(\theta_1 - \theta_0)}{\sigma} \right).$$

Показательный пример

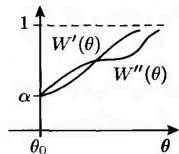


Сравнение двух критериев

- › Критерии на самом деле задаются критическими множествами
- › Пусть имеется два критерия, заданных множествами \mathcal{R}'_{α} и \mathcal{R}''_{α} .
Какой выбрать?

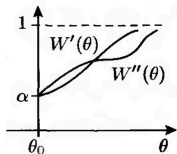
Сравнение двух критериев

- › Критерии на самом деле задаются критическими множествами
- › Пусть имеется два критерия, заданных множествами \mathcal{R}'_{α} и \mathcal{R}''_{α} .
Какой выбрать?
- › Сложная альтернатива: необходимо сравнивать функции мощности $W'(\theta)$ и $W''(\theta)$



Сравнение двух критериев

- › Критерии на самом деле задаются критическими множествами
- › Пусть имеется два критерия, заданных множествами \mathcal{R}'_α и \mathcal{R}''_α .
Какой выбрать?
- › Сложная альтернатива: необходимо сравнивать функции мощности $W'(\theta)$ и $W''(\theta)$
- › Простая альтернатива: существует наиболее мощный критерий (Неймана-Пирсона).
- › Идея: при заданной (достаточно малой) вероятности ошибки 1-го рода α постараться уменьшить вероятность ошибки 2-го рода β насколько возможно за счет подбора критического множества \mathcal{R}_α .



Сравнение двух критериев

- › Пусть задана выборка \mathbf{X}^ℓ
- › Гипотеза \mathbb{H}_0 и альтернатива \mathbb{H}_1 порождают в выборочном пространстве \mathbb{R}^ℓ меры P_0 и P_1
- › Таким образом, необходимо найти множество G такое, что $P_0(G) \leq \alpha$ и $P_1(G) \rightarrow \sup_{G: P_0(G) \leq \alpha} P_1(G)$
- › Рассмотрим систему вложенных множеств $G_c = \{\mathbf{x} \in \mathbb{R}^\ell : \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} \geq c\}$
- › Пусть $\varphi(c) = P_0(G_c)$, тогда $\varphi(c)$ убывает с ростом c

Сравнение двух критериев

- › На самом деле, $\varphi(c)$ убывает быстрее, чем $1/c$:

$$1 \geq P_1(G_c) = \int_{G_c} p_1(\mathbf{x}) d\mathbf{x} \geq c \int_{G_c} p_0(\mathbf{x}) d\mathbf{x} = c P_0(G_c) = c\varphi(c).$$

- › Далее еще потребуем, чтобы плотности $p_1(\mathbf{x})$ и $p_0(\mathbf{x})$ были всюду положительны
- › Дополнительно потребуем, чтобы
 $\forall \alpha \in (0, 1) \quad \exists c = c_\alpha : \quad \varphi(c_\alpha) = \alpha.$

Лемма Неймана-Пирсона

Лемма (Неймана-Пирсона)

Наиболее мощный критерий уровня α задается критическим множеством

$$G^* = G_{c_\alpha} = \left\{ \mathbf{x} \in \mathbb{R}^\ell : \frac{p_1(\mathbf{x})}{p_0(\mathbf{x})} \geq c_\alpha \right\}$$

- › Пусть G — критическое множество уровня α .
- › Тогда $P_0(G_c) \geq \alpha = P_0(G_{c_\alpha})$
- › Пусть $I(\mathbf{x})$ — индикатор G_c , $I^*(\mathbf{x})$ — индикатор G_{c_α}
- › Функция

$$f(\mathbf{x}) = (I^*(\mathbf{x}) - I(\mathbf{x}))(p_1(\mathbf{x}) - c_\alpha p_0(\mathbf{x}))$$

неотрицательна при всех $\mathbf{x} \in \mathbb{R}^\ell$

Лемма Неймана-Пирсона

› Функция

$$f(\mathbf{x}) = (I^*(\mathbf{x}) - I(\mathbf{x}))(p_1(\mathbf{x}) - c_\alpha p_0(\mathbf{x}))$$

неотрицательна при всех $\mathbf{x} \in \mathbb{R}^\ell$

› Поэтому

$$\begin{aligned} 0 &\leq \int_{\mathbf{x} \in \mathbb{R}^\ell} f(\mathbf{x}) d\mathbf{x} = \int_{\mathbf{x} \in \mathbb{R}^\ell} I^*(\mathbf{x}) p_1(\mathbf{x}) d\mathbf{x} - \int_{\mathbf{x} \in \mathbb{R}^\ell} I(\mathbf{x}) p_1(\mathbf{x}) d\mathbf{x} - \\ &- c_\alpha \left[\int_{\mathbf{x} \in \mathbb{R}^\ell} I^*(\mathbf{x}) p_0(\mathbf{x}) d\mathbf{x} - \int_{\mathbf{x} \in \mathbb{R}^\ell} I(\mathbf{x}) p_0(\mathbf{x}) d\mathbf{x} \right] = \\ &= P_1(G^*) - P_1(G) - c_\alpha \underbrace{[P_0(G^*) - P_0(G)]}_{\geq 0}. \end{aligned}$$

Критерий

Неймана-Пирсона:

2 простые гипотезы

Критерий Неймана-Пирсона

- › $\mathbb{H}_0 : \theta = \theta_0$ vs. $\mathbb{H}_1 : \theta = \theta_1$
- › Статистика Неймана-Пирсона:

$$T = \frac{\mathcal{L}(\theta_1)}{\mathcal{L}(\theta_0)} = \frac{\prod_{i=1}^n f(X_i; \theta_1)}{\prod_{i=1}^n f(X_i; \theta_0)}. \quad (1)$$

- › Допустим, что \mathbb{H}_0 отвергается при $T > k$. Выберем k так, что $P_{\theta_0}(T > k) = \alpha$.
- › Тогда, критерий Неймана-Пирсона (на основе статистики (2)) будет иметь наибольшую мощность $W(\theta_1)$ среди всех критериев размера α .

Пример

- › $X_i \sim \mathcal{N}(\mu, \sigma^2)$, причем дисперсия σ^2 известна
- › $\mathbb{H}_0 : \mu = \mu_0$ vs. $\mathbb{H}_1 : \mu = \mu_1$
- › Статистика Неймана-Пирсона:

$$T = \frac{\mathcal{L}(\theta_1)}{\mathcal{L}(\theta_0)} = \frac{\prod_{i=1}^n \mathcal{N}(X_i; \mu_1, \sigma^2)}{\prod_{i=1}^n \mathcal{N}(X_i; \mu_0, \sigma^2)}. \quad (2)$$

- › Подсчитайте статистику критерия (упростите)

Пример

- › $X_i \sim \mathcal{N}(\mu, \sigma^2)$, причем дисперсия σ^2 известна
- › $\mathbb{H}_0 : \mu = \mu_0$ vs. $\mathbb{H}_1 : \mu = \mu_1$
- › Статистика Неймана-Пирсона:

$$T = \frac{\mathcal{L}(\theta_1)}{\mathcal{L}(\theta_0)} = \frac{\prod_{i=1}^n \mathcal{N}(X_i; \mu_1, \sigma^2)}{\prod_{i=1}^n \mathcal{N}(X_i; \mu_0, \sigma^2)}. \quad (2)$$

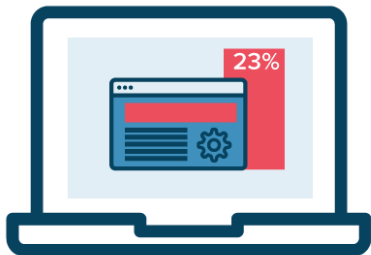
- › Подсчитайте статистику критерия (упростите)
- › Получается

$$\begin{aligned} T &= \exp \left[-\frac{1}{2\sigma^2} \sum_{i=1}^{\ell} [(X_i - \mu_1)^2 - (X_i - \mu_0)^2] \right] = \\ &= \exp \left[\frac{n}{\sigma^2} (\mu_1 - \mu_0) \underbrace{\left[\bar{X}_\ell - \frac{\mu_1 + \mu_0}{2} \right]}_{\text{важен знак!}} \right] \end{aligned}$$

Краткий экскурс в А/В тестирование

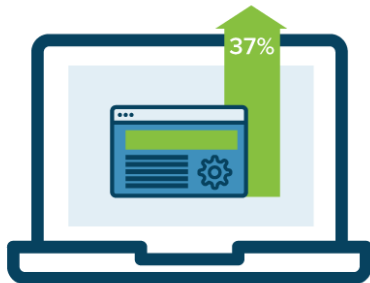
Идея А/В тестирования

A



CONTROL

B



VARIATION

Техника А/В тестирования

- › Варианты разбиения
 - › По пользователям (куки)
 - › По визитам (сессии)
 - › По действиям (запросы)
 - › ...
- › Типы экспериментов
 - › явные (интерфейсы, функциональность)
 - › неявные (ранжирование, персонализация)
 - › сервисы / части сервисов / кросс- сервисные
 - › улучшения / ухудшения / АА-тесты

Внедрения А/В тестирования

- › Масштабы экспериментов:

Яндекс @ 2017:

- › 4778 экспериментов,
- › ~400 экспериментов одновременно

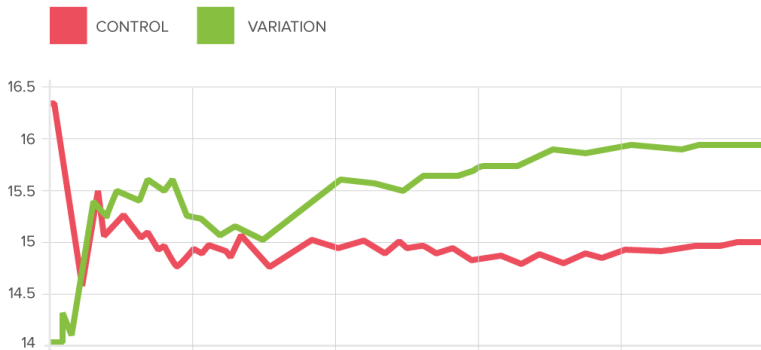
- › Результаты:

- › Яндекс, 2014: 21% принимаются, 2017: 28% принимаются
- › Бинг, 201?: 30% принимаются
- › Гугл, 201?: 10% принимаются
- › В среднем: 20% принимаются, 50% непонятных

Техника А/В тестирования

› Метрики

- › абсолютные (клики, временные)
- › относительные (ctr, доля некликнувших)



Критерий Стьюдента (t-test)

Определение

Случайная величина имеет распределение Стьюдента (t -распределение) с k степенями свободы, если:

$$f(t) = \frac{\Gamma(\frac{k+1}{2})}{\sqrt{k\pi}\Gamma(\frac{k}{2})(1 + \frac{t^2}{k})^{\frac{k+1}{2}}}$$

При $k \rightarrow \infty$ t -распределение стремится к стандартному нормальному распределению. При $k = 1$ t -распределение совпадает с распределением Коши.

- › t -критерий используют, когда распределение данных близко к нормальному, а размер выборки невелик.

Критерий Стьюдента (t-test)

Теорема (Критерий Стьюдента (t-test))

Пусть $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, где параметры (μ, σ^2) неизвестны.

$$\mathbb{H}_0 : \mu = \mu_0 \quad vs. \quad \mathbb{H}_1 : \mu \neq \mu_0$$

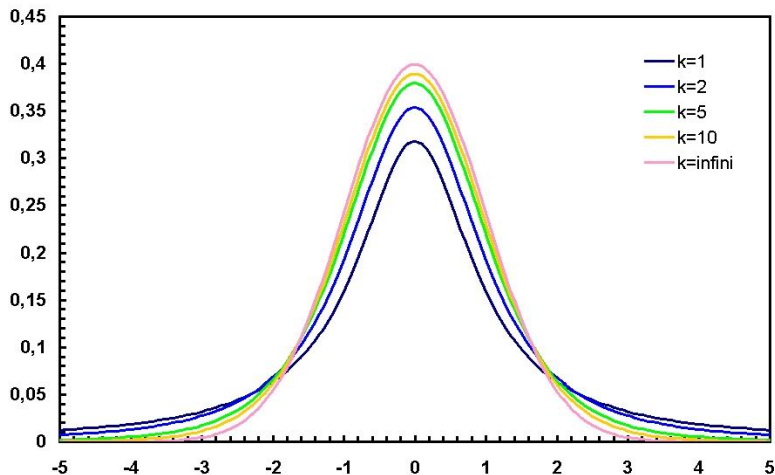
Обозначим через S_n^2 выборочную дисперсию. Тогда статистика критерия:

$$T = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{S_n}$$

Основная гипотеза отвергается, если $|T| > t_{n-1, \alpha/2}$, где $t_{n-1, \alpha/2}$ — квантиль распределения Стьюдента с $n - 1$ степенями свободы.

При больших n выполняется $T \sim \mathcal{N}(0, 1)$, то есть при больших n t-критерий эквивалентен критерию Вальда.

Критерий Стьюдента (t-test)



Критерий Стьюдента (t-test)

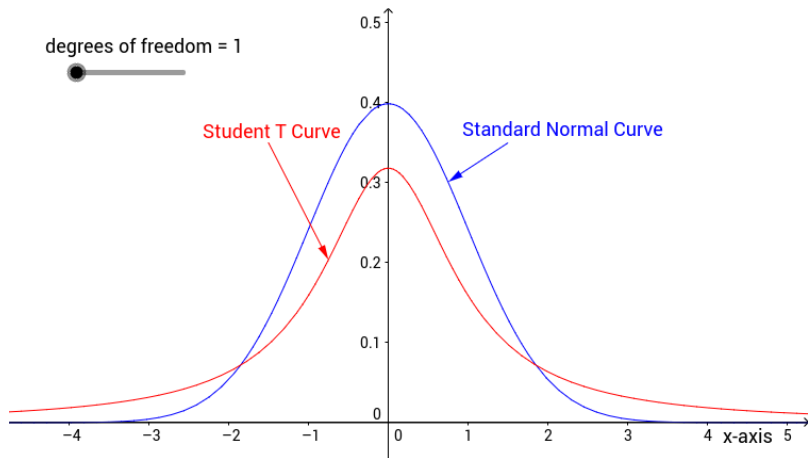


Рис.: <http://tananyag.geomatech.hu/m/53882>

Последовательный анализ

Мотивация

- › В классической теории математической статистики предполагается, что наблюдения (данные) заранее известны
- › Однако нужно ли заранее фиксировать размер выборки?
- › В действительности, размер выборки можно определять в зависимости от уже поступивших наблюдений!
- › Это приводит к последовательному тесту Вальда (последовательному критерию отношения правдоподобия, sequential probability ratio test, SPRT)
- › Основная мотивация: **экономить данные** (до 50% экономии — доказательство на семинаре!)

Последовательный критерий

- › Гипотеза \mathbb{H}_0 : наблюдения X_i имеют плотность $p_0(x)$, \mathbb{H}_1 : наблюдения X_i имеют плотность $p_1(x)$
- › Рассмотрим $Z_i = \log \frac{p_0(X_i)}{p_1(X_i)}$ и блуждание $S_k = Z_1 + \dots + Z_k$
- › Например, если $p_0(x) = \mathcal{N}(\mu_0, \sigma^2)$, $p_1(x) = \mathcal{N}(\mu_1, \sigma^2)$, то

$$Z_i = \frac{\mu_1 - \mu_0}{\sigma^2} \left[\bar{X}_i - \frac{\mu_1 + \mu_0}{2} \right]$$

- › Интуиция:

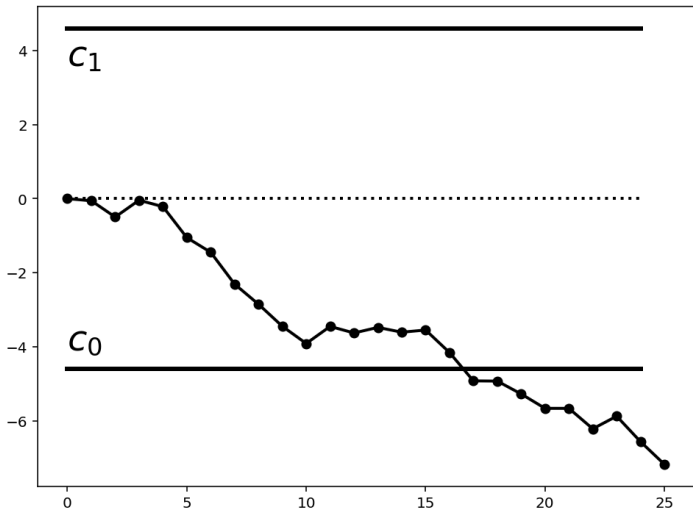
Последовательный критерий

- › Гипотеза \mathbb{H}_0 : наблюдения X_i имеют плотность $p_0(x)$, \mathbb{H}_1 : наблюдения X_i имеют плотность $p_1(x)$
- › Рассмотрим $Z_i = \log \frac{p_0(X_i)}{p_1(X_i)}$ и блуждание $S_k = Z_1 + \dots + Z_k$
- › Например, если $p_0(x) = \mathcal{N}(\mu_0, \sigma^2)$, $p_1(x) = \mathcal{N}(\mu_1, \sigma^2)$, то

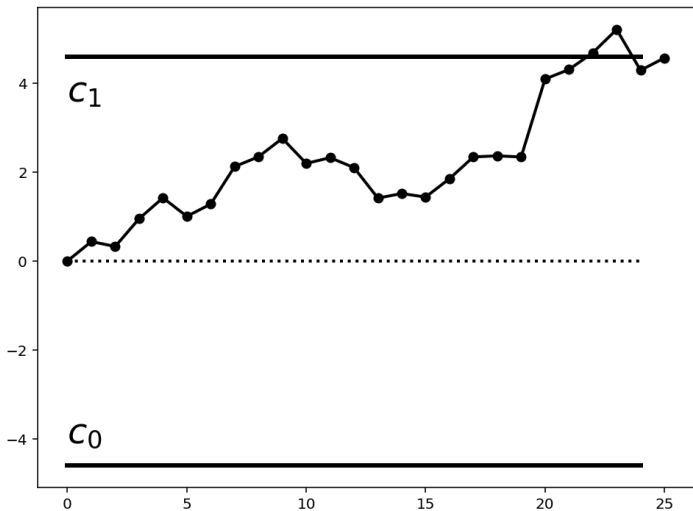
$$Z_i = \frac{\mu_1 - \mu_0}{\sigma^2} \left[\bar{X}_i - \frac{\mu_1 + \mu_0}{2} \right]$$

- › Интуиция:
 - › если верна \mathbb{H}_0 (сигнала нет), то в среднем $p_0(X_i) \geq p_1(X_i)$ (и тогда $Z_i \leq 0$) $\implies S_k$ убывает
 - › если верна \mathbb{H}_1 (сигнал!), то в среднем $p_0(X_i) \leq p_1(X_i)$ (и тогда $Z_i \geq 0$) $\implies S_k$ растёт

Пример: верна гипотеза



Пример: верна альтернатива



Последовательный критерий

› Последовательный тест:

1. Наблюдаем X_1, X_2, \dots последовательно
2. Вычисляем значения $S_k, k = 1, 2, \dots$
3. $S_k \geq c_1$: остановимся и отклоним \mathbb{H}_0
4. $S_k \leq c_0$: остановимся и не отклоним \mathbb{H}_0
5. $c_0 < S_k < c_1$: продолжим наблюдения

› Остается задать пределы роста S_k :

$$c_0 = \ln \frac{\beta'}{1 - \alpha'}, \quad c_1 = \ln \frac{1 - \beta'}{\alpha'}, \quad \alpha' + \beta' < 1$$

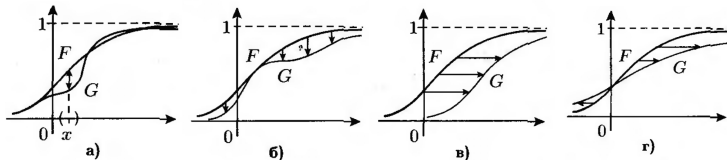
При этом для вероятностей α и β :

$$\alpha \leq \frac{\alpha'}{1 - \beta'}, \quad \beta \leq \frac{\beta'}{1 - \alpha'}, \quad \alpha + \beta < \alpha' + \beta'.$$

Непараметрические критерии

Альтернативы однородности

- › Имеем две выборки $X^n \sim F(x)$ и $Y^m \sim G(x)$
- › Гипотеза однородности $\mathbb{H}_0 : F(x) = G(x), x \in \mathbb{R}$



- › Бывает важно уловить отклонения от \mathbb{H}_0 только определенного типа (наличие прироста Y_j по сравнению с X_i)
- › Сужение типов альтернатив \implies более эффективные критерии
- › Проверяем гипотезу однородности против альтернативы доминирования \mathbb{H}_1 (варианты б) и в))

Критерий ранговых сумм MWW

- › Построим вариационный ряд из объединенной выборки $(X_1, \dots, X_n, Y_1, \dots, Y_m)$
 - › Верна $H_0 \implies$ значения Y_j рассеяны по всему ряду
 - › Иначе средний ранг значений Y_j относительно большой
- › Обозначим S_j ранг порядковой статистики $Y_{(j)}$ в этом ряду
- › Положим $V = S_1 + \dots + S_m$
- › Критическая область: $V \geq c$, где $c = \text{const}$
- › **Большие выборки:** (Mann-Whitney-Wilcoxon, MWW)

$$U = \sum_{i=1}^n \sum_{j=1}^m I_{X_i < Y_j} \rightarrow \mathcal{N}\left(\frac{nm}{2}, \frac{nm(n+m+1)}{12}\right)$$

Резюме лекции

- › Критерий отношения правдоподобия: оптимален, если у вас простые гипотезы.
- › А/В-тестирование — проверка гипотез в продакшене.
- › Критерий последовательного отношения правдоподобия: снижение затрат на тестирование.
- › Непараметрический ранговый критерий: рабочая лошадка АВТ (наряду с t -тестом!).