

# Локализация сигнала в задаче интерпретации множества статистических тестов<sup>1</sup>

---

Анастасия Процветкина

14 мая, 2021

Санкт-Петербургский государственный электротехнический университет "ЛЭТИ"

Санкт-Петербург, Россия.

---

<sup>1</sup>Руководитель: Сергей Васильевич Малов

# Локализация в категориальных данных

---

# Постановка задачи

Имея таблицу сопряженности размера  $m \times n$ , можно проверить гипотезу независимости признаков с помощью хи-квадрат статистики.

Известно, что статистика критерия при нулевой гипотезе имеет асимптотическое распределение

$$\chi^2_{(m-1)(n-1)}$$

По определению, хи-квадрат величина  $X^2$  с  $k$  степенями свободы есть сумма квадратов  $k$  независимых одинаково распределенных нормальных случайных величин  $X_i$ , т.е.

$$X^2 = X_1^2 + X_2^2 + \dots + X_k^2$$

## Постановка задачи

С каждой таблицей  $m \times n$  мы теоретически имеем ассоциированный с ней вектор  $\mathbf{X} = (X_1, X_2, \dots, X_{(m-1)(n-1)})^T$ , где  $X_i$  асимптотически НОРСВ из  $\mathcal{N}(0, 1)$ .

Пусть имеется  $w$  таких таблиц, имеющих общий признак. Сформируем длинный вектор из всех асимптотически нормальных компонент:

$$\xi = (\underbrace{X_{11}, X_{12}, \dots, X_{1, (m-1)(n-1)}}_{\text{норм. компоненты из 1-й таб.}}, \dots, \underbrace{X_{w1}, X_{w2}, \dots, X_{w, (m-1)(n-1)}}_{\text{норм. компоненты из w-й таб.}})^T$$

Очевидно, что  $\xi \Rightarrow \mathcal{N}(0, \Sigma_q)$ ,  $q = w(m-1)(n-1)$ .

# Постановка задачи

## Теорема (Rao, Mitra)

$$\psi \sim \mathcal{N}(\mu, \Sigma), \quad r = \text{rk}(\Sigma),$$

$\Sigma^-$  – обобщенная обратная матрица для  $\Sigma$

1. Квадратичная форма  $\psi^T \Sigma^- \psi$  не зависит от выбора  $\Sigma^-$ .
2.  $\psi^T \Sigma^- \psi \sim \chi_{\delta, r}^2$  ( $\delta = \mu^T \Sigma^- \mu$  – параметр нецентральности).

## Комбинированная статистика

Для вектора  $\xi \Rightarrow \mathcal{N}(0, \Sigma)$  из всех асимптотически нормальных величин, соответствующих компонентам хи-квадрат статистики из  $w$  таблиц сопряженности, построим квадратичную форму

$$\xi^T \Sigma^+ \xi \stackrel{H_0}{\Rightarrow} \chi_r^2,$$

$H_0$ : “Ни в одной из  $W = \{1, \dots, w\}$  таблиц нет зависимости”,  
 $r = \text{rk}(\Sigma)$ ,  $\Sigma^+$  – обобщенная обратная матрица Мура-Пенроуза.

# Задача 1

Хи-квадрат статистика для проверки гипотезы независимости обычно записывается как

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{(O_{ij} - E_{ij})^2}{E_{ij}},$$

где  $O_{ij}$  – наблюдаемые частоты в ячейке  $(i, j)$ ,  $E_{ij}$  – соответствующие ожидаемые частоты при справедливости гипотезы независимости.

$\chi^2$  имеет  $(m - 1)(n - 1)$  степеней свободы, но алгебраически является суммой  $m \cdot n$  зависимых слагаемых, а нам необходимо представление вида:

$$\chi^2 = \chi_1^2 + \dots + \chi_{(m-1)(n-1)}^2$$

$\Rightarrow$  необходимо преобразование.

# Задача 1

## Решение

Ирвином и Ланкастером было предложено ортогональное преобразование, приводящее хи-квадрат статистику с  $m \cdot n$  слагаемыми в сумму квадратов  $(m - 1)(n - 1)$  асимптотических независимых нормальных величин.

## Задача 2

Вектор из всех асимп. нормальных компонент сформирован, т.е.

$$\xi = (\underbrace{X_{11}, X_{12}, \dots, X_{1, (m-1)(n-1)}}_{\text{норм. компоненты из 1-й таб.}}, \dots, \underbrace{X_{w1}, X_{w2}, \dots, X_{w, (m-1)(n-1)}}_{\text{норм. компоненты из w-й таб.}})^T$$

получен. Построение квадратичной формы  $\xi^T \Sigma^+ \xi$  требует оценивания ковариационной матрицы  $\Sigma$ , т.е. оценивания ковариации между нормальными величинами, соответствующими компонентам хи-квадрат статистики из разных таблиц.

### Замечание

Каждая величина  $X_{ij}$  представляет одну степень свободы статистики хи-квадрат, поэтому  $\text{var}(X_{ij}) = 1$ . В дальнейшем матрицу  $\Sigma$  будем называть корреляционной, нежели просто ковариационной.



## Задача 2

### Решение

Оценки корреляционной матрицы  $\Sigma$  были получены в частных случаях  $2 \times 2$  и  $3 \times 2$ . Формулы громоздки, поэтому не приводятся.

## Задача 3

Теоретическая и неизвестная матрица  $\Sigma$  обладает неустойчивым свойством положительной полу-определенности. Оценивание матрицы может привести к его потере, поэтому необходимо учитывать такую потенциальную проблему и иметь решение в случае возникновения.

### Решение

Согласно Nigham, ближайшую положительно-определенную матрицу для  $\Sigma$  можно построить, обнулив отрицательные собственные числа  $\Sigma$ .

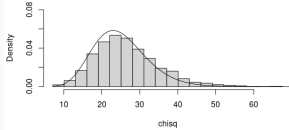
# Результаты

---

# Сравнение распределений: $2 \times 2$ случай

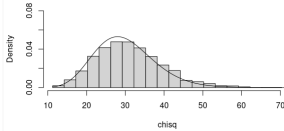
2a.

Hist of 1969 observed chisqs  
with 25 df, sample size: 1000



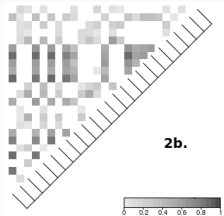
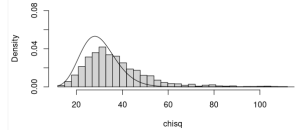
3a.

Hist of 2000 observed chisqs  
with 30 df, sample size: 1000

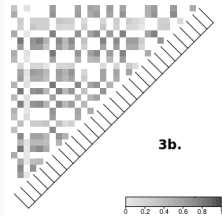


4a.

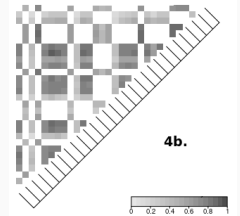
Hist of 1692 observed chisqs  
with 30 df, sample size: 1000



2b.



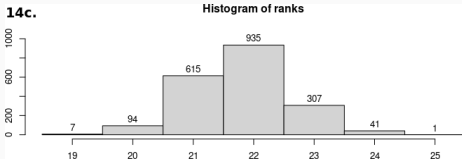
3b.



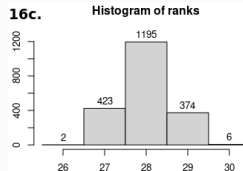
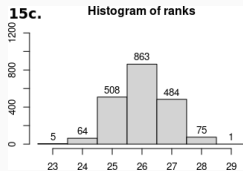
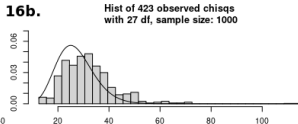
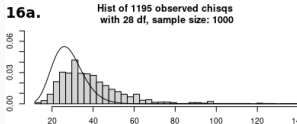
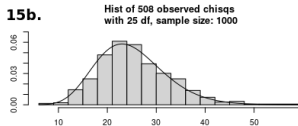
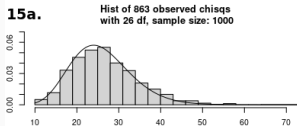
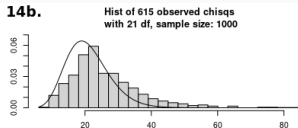
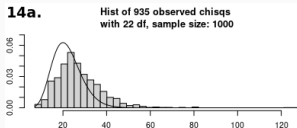
4b.

2a,b:  $w = 25$ , 3a,b:  $w = 30$ , 4a,b:  $w = 30$   
 $N = 1000$ , маргинальные вероятности по 0.5.

# Сравнение распределений: $3 \times 2$ случай



Распределение для  
2 наиболее частых  
рангов.  $w = 15$ ,  
 $N = 1000$



# Заключение

- Разработан новый метод, позволяющий объединять отдельные слабые сигналы в единый более мощный сигнал
- Получены оценки корреляционной матрицы для асимптотически нормального базисного вектора
- Потенциальные проблемы вырождаемости матрицы  $\Sigma$  и потеря положительной определенности были предвидены и разрешены
- Разработанный метод протестирован на симулированных зависимых данных
- Отмечена сходимость распределений, но для окон малого размера
- Матрица  $\Sigma$  вырождается с увеличением размера окна и зависимости между маркерами