

Bài giảng 4: A/B Testing

TS. Tô Đức Khánh

Khoa Toán-Tin Học, Trường Đại Học Khoa Học Tự Nhiên
Đại Học Quốc Gia Tp. HCM

Học phần: Xử lý Số liệu Thống kê

Nội dung buổi học

- 1 Kiểm định giả thuyết
- 2 A/B testing hai nhóm
- 3 A/B testing nhiều nhóm
- 4 Multi-Armed Bandit

1 Kiểm định giả thuyết

2 A/B testing hai nhóm

3 A/B testing nhiều nhóm

4 Multi-Armed Bandit

Khái niệm cơ bản

Kiểm định giả thuyết là một bài toán cơ bản trong thống kê lý thuyết cũng như thống kê ứng dụng.

Thông thường, kiểm định giả thuyết hướng tới trả lời một giả thuyết cho 1 đặc trưng của dữ liệu:

- Trung bình đối với biến định lượng;
- Tỷ lệ đối với biến định tính;
- Sự tương quan giữa hai biến.

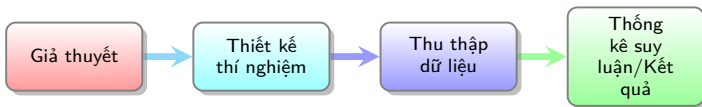
Khái niệm cơ bản

Kiểm định giả thuyết là một bài toán cơ bản trong thống kê lý thuyết cũng như thống kê ứng dụng.

Thông thường, kiểm định giả thuyết hướng tới trả lời một giả thuyết cho 1 đặc trưng của dữ liệu:

- Trung bình đối với biến định lượng;
- Tỷ lệ đối với biến định tính;
- Sự tương quan giữa hai biến.

Quy trình kiểm định giả thuyết thống kê:



Khái niệm cơ bản

Một giả thuyết thống kê gồm 2 thành phần:

Giả thuyết - Null hypothesis là một giả định bao hàm tính không khác biệt/không hiệu quả hoặc sự khác biệt chỉ là sự trùng hợp ngẫu nhiên.

Ví dụ: “Mức độ tương tác của trang web A và trang web B là không khác biệt.”

Đối thuyết - Alternative hypothesis là một giả định bao hàm sự ngược lại giả định đã được tuyên bố.

Ví dụ: “Mức độ tương tác của trang web A và trang web B là khác biệt.”

“Mức độ tương tác của trang web A là cao hơn của trang web B.”

“Mức độ tương tác của trang web A là thấp hơn của trang web B.”

Thông thường, ta mong muốn chứng minh rằng giả thuyết là **không đúng**, thông qua dữ liệu thu thập từ việc thiết kế thí nghiệm.

A/B testing

Một kiểm định A/B (A/B testing) là một thí nghiệm với hai nhóm nhằm xác định xem nhóm nào trong: hai phương pháp điều trị, hay hai sản phẩm, hay hai quy trình là tốt hơn cái còn lại.

Thông thường, hai nhóm gồm:

nhóm đối chứng - control group là nhóm gồm các đối tượng được tiếp xúc với điều kiện thông thường (không được áp dụng các biện pháp mới) - được ký hiệu là A;

nhóm điều trị - treatment group là nhóm gồm các đối tượng được tiếp xúc với 1 điều kiện điều trị cụ thể - được ký hiệu là B.

Giả thuyết và đối thuyết của A/B testing thường là

- Giả thuyết: Phương pháp điều trị có hiệu quả tương đồng với đối chứng.
- Đối thuyết: Phương pháp điều trị có hiệu quả tốt hơn đối chứng.

A/B testing

Ví dụ cho A/B testing:

- Thử nghiệm hai phương pháp xử lý đất để xác định phương pháp nào giúp hạt nảy mầm tốt hơn.
- Thử nghiệm hai liệu pháp để xác định liệu pháp nào ngăn chặn ung thư hiệu quả hơn.
- Kiểm tra hai mức giá để xác định mức giá nào mang lại nhiều lợi nhuận ròng hơn.
- Kiểm tra hai dòng tiêu đề trên web để xác định dòng tiêu đề nào tạo ra nhiều lượt nhấp chuột hơn.
- Thử nghiệm hai quảng cáo trên web để xác định quảng cáo nào tạo ra nhiều chuyển đổi hơn.

Đối tượng nghiên cứu trong một A/B test, có thể là:

- người;
- hạt mầm cây;
- người ghé thăm trang web.

A/B testing

Mỗi đối tượng sẽ được tiếp xúc với

- điều kiện thông thường (control); hoặc
- điều kiện điều trị (treatment).

Để có dữ liệu phục vụ cho việc phân tích, ta cần thiết kế thí nghiệm đảm bảo tính ngẫu nhiên (randomness).

↔ Các đối tượng được chỉ định một cách ngẫu nhiên (**randomization**) vào nhóm đối chứng hoặc nhóm điều trị.

↔ Khi đó, ta sẽ biết rằng bất kỳ sự khác biệt nào giữa các nhóm là do một trong hai điều sau:

- Hiệu quả của các phương pháp điều trị khác nhau.
- Sự may mắn trong kết quả do các đối tượng được chỉ định ngẫu nhiên vào các phương pháp điều trị (tức là, việc phân công ngẫu nhiên có thể dẫn đến việc các đối tượng có thành tích tốt hơn một cách tự nhiên tập trung vào A hoặc B).

Mở rộng A/B testing

A/B testing là khá phổ biến trong các lĩnh vực kinh doanh hoặc thương mại điện tử.

Tuy nhiên, vẫn có những loại kiểm định khác, là biến thể hoặc mở rộng của A/B testing:

- Số nhóm điều trị nhiều hơn hai.

Ví dụ: so sánh độ gắn kết khách hàng của 4 trang web.

- Sự liên hệ của hai biến định tính.

Ví dụ: kiểm tra sự tác động của 3 kiểu headline tới tương tác của người dùng (click hoặc no-click).

- Các đối tượng được quan sát lặp lại.

Mở rộng A/B testing

Trong thống kê, ta thường quan tâm tới các câu hỏi dạng:

“Sự khác biệt giữa phương pháp điều trị A và phương pháp điều trị B có ý nghĩa thống kê hay không?”

Tuy nhiên, các nhà khoa học dữ liệu quan tâm tới câu hỏi dạng:

“Trong số nhiều phương pháp điều trị khác nhau, cái nào là tốt nhất?”

tức là dạng câu hỏi trực diện hơn.

Với câu hỏi dạng này, một loại thiết kế thử nghiệm tương đối mới được sử dụng: the multi-armed bandit hay **Multi-Armed Bandit Algorithms**.

Sai lầm trong kiểm định

Bảng sau đây biểu diễn bốn hậu quả có thể của một quyết định trong bài toán kiểm định thống kê:

Quyết định	Sự thật của Giả thuyết	
	Đúng	Sai
Không bác bỏ	Chính xác	Sai lầm loại II
Bác bỏ	Sai lầm loại I	Chính xác

Sai lầm loại I - Type I error

Sai lầm loại I là sai lầm có thể mắc phải trong kiểm định giả thuyết thống kê, nó xảy ra khi ta quyết định **bác bỏ** Giả thuyết khi Giả thuyết thực sự đúng.

Sai lầm loại II - Type II error

Sai lầm loại II là sai lầm có thể mắc phải trong kiểm định giả thuyết thống kê, nó xảy ra khi ta quyết định **không bác bỏ** Giả thuyết khi Giả thuyết thực sự sai.

Sai lầm trong kiểm định

- Xác suất mắc sai lầm loại I thường được ký hiệu là α và nó là **mức ý nghĩa**.
- Xác suất mắc sai lầm loại II thường được ký hiệu là β .
- $1 - \beta$ được gọi là độ mạnh của kiểm định (power), tức là khả năng nhận ra được Giả thuyết là sai.

Trong khuôn khổ một kiểm định thống kê, chúng ta giả định rằng Giả thuyết là đúng. Vì vậy, chúng ta có thể nhận được:

- quyết định chính xác nếu ta không bác bỏ Giả thuyết; hoặc,
- sai lầm loại I nếu ta bác bỏ Giả thuyết.

Để kiểm soát khả năng bị sai lầm loại I, ta cố định xác suất mắc sai lầm loại I (*mức ý nghĩa*), thường là 0.05 (hay 5%) và so sánh nó với một đại lượng được gọi là *p-value*, được tính từ dữ liệu.

Trong khi đó, sai lầm loại II được kiểm soát bằng việc xác định cỡ mẫu của dữ liệu.

p-value và quyết định

p -value là xác suất mà mô hình kiểm định tạo ra kết quả cực đoan hơn kết quả quan sát được, dưới giả định rằng Giả thuyết là đúng.

Ví dụ: p -value bằng 0.308 có nghĩa là ta có thể kỳ vọng rằng, ta sẽ thu được kết quả cực đoan như của dữ liệu được quan sát hoặc cực đoan hơn, một cách ngẫu nhiên trong hơn 30% số lần lấy mẫu, khi giả sử rằng Giả thuyết là đúng.

Nói cách khác, p -value có thể được dùng để đo sự tương thích của dữ liệu với Giả thuyết.

- nếu p -value là lớn, có nghĩa là dữ liệu tương thích với Giả thuyết
→ dữ liệu cung cấp bằng chứng ủng hộ Giả thuyết;
- nếu p -value là nhỏ, có nghĩa là dữ liệu không tương thích với Giả thuyết
→ dữ liệu cung cấp bằng chứng chống lại Giả thuyết
→ đưa ra quyết định bác bỏ Giả thuyết.

p-value và quyết định

Để xác định mức nhỏ của p -value, ta so sánh nó với mức ý nghĩa α :

- nếu $p\text{-value} \geq \alpha$, thì chưa đủ nhỏ;
- nếu $p\text{-value} < \alpha$, thì đủ nhỏ để bác bỏ Giả thuyết với sai lầm loại I thấp.

Ví dụ: với p -value bằng 0.308:

- so sánh với $\alpha = 0.05 \implies$ không bác bỏ Giả thuyết;
- so sánh với $\alpha = 0.1 \implies$ không bác bỏ Giả thuyết;
- so sánh với $\alpha = 0.35 \implies$ bác bỏ Giả thuyết, với mức sai lầm loại I là 35%;

1 Kiểm định giả thuyết

2 A/B testing hai nhóm

3 A/B testing nhiều nhóm

4 Multi-Armed Bandit

Bài toán

A/B testing là một thí nghiệm nhằm tìm ra phương pháp A hay phương pháp B là tốt hơn phương pháp còn lại.

↪ ta cần so sánh dữ liệu từ hai nhóm A và nhóm B.

- Khi dữ liệu ở dạng biến định lượng \Rightarrow so sánh giá trị **trung bình** của dữ liệu của hai nhóm \Rightarrow bài toán kiểm định hai trung bình.
- Khi dữ liệu ở dạng biến định tính \Rightarrow so sánh giá trị **tỷ lệ** của dữ liệu của hai nhóm \Rightarrow bài toán kiểm định hai tỷ lệ.

Ví dụ 1: Xét câu hỏi sau: “thời gian tương tác của trang web A và của trang web B có tương đồng nhau?”. Nếu:

- thời gian tương tác là không tương đồng
↪ một trong hai trang web (A hoặc B) giúp cải thiện (hoặc giảm) khả năng tương tác với khách hàng;
- thời gian tương tác là như nhau
↪ không có sự cải thiện (giảm) sự tương tác.

So sánh trung bình thời gian tương tác của web A và web B.

Bài toán

Ví dụ 2: Xét câu hỏi sau: “sự lựa chọn cánh gà chiên cay (hot wings) có phụ thuộc giới tính khách hàng (nam/nữ)?”. Nếu:

- có ảnh hưởng
↪ một trong hai giới (nam hoặc nữ) mua nhiều hơn giới còn lại
- không ảnh hưởng
↪ lượng mua cánh gà chiên cay là như nhau cho cả hai giới.

So sánh trung bình số lượng cánh gà chiên cay được mua bởi khách hàng nam và khách hàng nữ.

Ví dụ 3: Xét câu hỏi sau: “tỷ lệ chuyển đổi sản phẩm của một mặt hàng trên trang thương mại điện tử có khác biệt giữa hai mức giá?”. Nếu:

- có sự khác biệt
↪ một trong hai giá (A hoặc B) tạo ra sự hấp dẫn với khách hàng
- nếu không có sự khác biệt ↪ cả hai giá đều có sự hấp dẫn khách hàng như nhau.

So sánh tỷ lệ bán được sản phẩm của hai mức giá.

Kiểm định một phía, hai phía

Trong thống kê kiểm định nói chung và A/B testing nói riêng, ta luôn có:

- kiểm định hai phía (two sided)
- kiểm định một phía bên trái (left sided)
- kiểm định một phía bên phải (right sided)

Những kiểm định này được xác định theo tuyên bố của đối thuyết.

- $\mu_A \neq \mu_B \implies$ kiểm định hai phía;
- $\mu_A < \mu_B \implies$ kiểm định một phía bên trái;
- $\mu_A > \mu_B \implies$ kiểm định một phía bên phải.

Kiểm định hai phía được sử dụng khi ta không xác định được hướng so sánh từ ban đầu.

Nếu xác định được hướng/mục tiêu so sánh ngay từ đầu:

- trung bình thời gian tương tác của web A là thấp hơn so với web B;
- trung bình lượng cánh gà được mua bởi nam giới (nhóm A) là cao hơn bởi nữ giới (nhóm B);

thì kiểm định một phía sẽ được áp dụng.

Cách tiếp cận cổ điển

Như đã học ở trong Lý Thuyết Thống Kê hoặc Xác suất - Thống kê, ta sẽ sử dụng

■ thống kê

$$T = \frac{(\bar{Y}_A - \bar{Y}_B) - (\mu_A - \mu_B)}{s_p \sqrt{\frac{1}{n_A} + \frac{1}{n_B}}} \sim t_{n_A+n_B-2},$$

$s_p = \frac{(n_A - 1)s_A^2 + (n_B - 1)s_B^2}{n_A + n_B - 2}$, với giả định về phân phối chuẩn và đồng nhất phương sai của hai nhóm, khi kiểm định Giả thuyết “hai trung bình là bằng nhau”;

■ thống kê

$$Z = \frac{\hat{p}_A - \hat{p}_B}{\sqrt{\hat{p}(1 - \hat{p}) \left(\frac{1}{n_A} + \frac{1}{n_B} \right)}} \sim N(0, 1),$$

$\hat{p} = \frac{n_A \hat{p}_A + n_B \hat{p}_B}{n_A + n_B}$, khi kiểm định Giả thuyết “hai tỷ lệ là bằng nhau”.

Cách tiếp cận cổ điển

Nếu ta sử dụng cách tiếp cận cổ điển, ta cần:

- thông tin về phân phối của dữ liệu trong hai nhóm;
- giả định về sự đồng nhất phương sai giữa hai nhóm;
- cỡ mẫu tương đối lớn.

Tuy nhiên, trong thực tế, đối với khoa học dữ liệu, các điều kiện 1 và 2 thường:

- khó kiểm chứng, hoặc
- không thể đáp ứng.

Do đó, ta cần một phương pháp kiểm định khác, mà không yêu cầu các điều kiện 1 và 2.

Permutation test

Nhóm A

Mẫu gốc

1
2
3

HV #1

1
1
2

HV #2

3
2
2

HV #3

3
1
2

Nhóm B

Mẫu gốc

1
2
3

HV #1

2
3
3

HV #2

1
3
1

HV #3

2
3
1

Permutation test

Kiểm định hoán vị gồm các bước sau:

1. Từ dữ liệu gốc, tính trung bình mẫu \bar{y}_A và \bar{y}_B của hai nhóm, sau đó tính sự khác biệt giữa $\Delta = \bar{y}_A - \bar{y}_B$ (hoặc $\bar{y}_B - \bar{y}_A$, tùy theo đối thuyết).
2. Gộp chung các mẫu lại, tạo thành một mẫu lớn và xáo trộn chúng.
3. Lấy mẫu ngẫu nhiên không lặp cho nhóm A (với cỡ mẫu như ban đầu) từ mẫu chung.
4. Lấy mẫu ngẫu nhiên không lặp cho nhóm B, từ thành phần còn lại của mẫu chung.
5. Tính trung bình mẫu cho 2 mẫu mới tạo, \bar{y}_A^* và \bar{y}_B^* , sau đó, tính sự khác biệt giữa hai trung bình $\bar{y}_A^* - \bar{y}_B^*$ (hoặc $\bar{y}_B^* - \bar{y}_A^*$).
6. Lặp lại bước 3 tới 5 trong R lần, lưu trữ các kết quả của sự khác biệt giữa $\bar{y}_A^* - \bar{y}_B^*$ (hoặc $\bar{y}_B^* - \bar{y}_A^*$).
7. Tính p -value dựa vào mẫu gồm R giá trị $\Delta_i^* = \bar{y}_{A,i}^* - \bar{y}_{B,i}^*$:

- $p\text{-value} = \frac{1}{R} \sum_{i=1}^R I(|\Delta_i^*| > |\Delta|)$, kiểm định hai phía;
- $p\text{-value} = \frac{1}{R} \sum_{i=1}^R I(\Delta_i^* > \Delta)$, kiểm định một phía bên phải;
- $p\text{-value} = \frac{1}{R} \sum_{i=1}^R I(\Delta_i^* < \Delta)$, kiểm định một phía bên trái;

Ví dụ: Web stickiness

Bài toán

Một công ty bán một dịch vụ có giá trị cao, muốn kiểm tra xem trang web A hoặc B giúp công ty thực hiện công việc bán hàng tốt hơn.

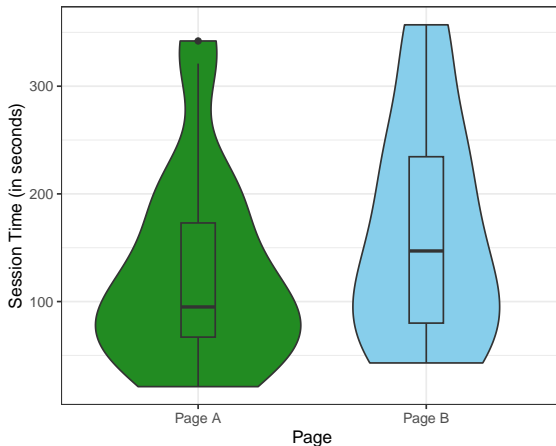
Cách tiếp cận

Một bản trình bày trên web thu hút sự chú ý của mọi người lâu hơn \Rightarrow doanh số bán hàng cao hơn.

\hookrightarrow dữ liệu là thời gian phiên trung bình (session time), đơn vị giây, so sánh trang web A với trang web B, được ghi nhận bởi Google Analytics.

	Page	Time
1	Page A	0.21
2	Page B	2.53
3	Page A	0.35
4	Page B	0.71
5	Page A	0.67
6	Page B	0.85

Ví dụ: Web stickiness



Có vẻ như thời gian phiên trung bình của trang web B là dài hơn so với trang web A.

Ví dụ: Web stickiness

Ta cần kiểm tra xem liệu thời gian phiên trung bình của trang web B có thực sự dài hơn so với trang web A, hay không.

↪ Ta phát biểu Giả thuyết và Đối thuyết như sau:

Giả thuyết: $\mu_B = \mu_A$

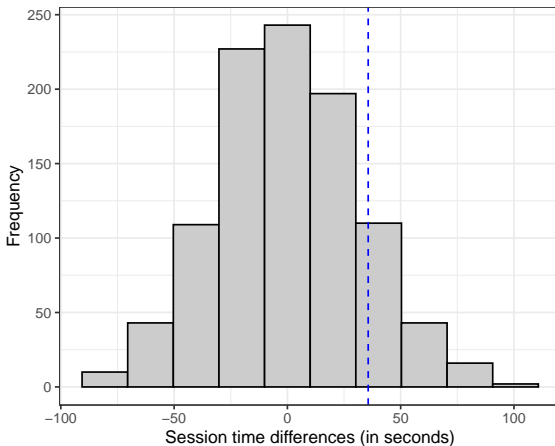
Đối thuyết: $\mu_B > \mu_A$

Từ dữ liệu ta có một số thông tin

- cỡ mẫu $n_A = 21$, $n_B = 15$;
- trung bình mẫu $\bar{y}_A = 126.33$ giây và $\bar{y}_B = 162$ giây;
- $\Delta = \bar{y}_B - \bar{y}_A = 35.67$.

Ta thực hiện kiểm định hoán vị với $R = 1000$ lần.

Ví dụ: Web stickiness



Ta tính được

$$p\text{-value} = \frac{1}{1000} \sum_{i=1}^{1000} \mathbb{I}(\Delta_i^* > 35.67) = 0.127$$

Ví dụ: Web stickiness

Chọn mức ý nghĩa $\alpha = 0.05$.

↪ $p\text{-value} > \alpha$

↪ ta không đủ cơ sở để bác bỏ Giả thuyết.

↪ Sự dài hơn của thời gian phiên trung bình của trang web B so với trang web A là không có ý nghĩa thống kê (**statistically significant**).

Kết luận: Có sự cải thiện trong thời phiên trung bình của trang web B so với trang web A. Nhưng sự cải thiện đó không có ý nghĩa thống kê, và có thể là kết quả của sự ngẫu nhiên.

Exhaustive permutation test

Kiểm định hoán vị toàn diện (Exhaustive permutation test) là một biến thể của kiểm định hoán vị.

- Thay vì chỉ xáo trộn và chia dữ liệu một cách ngẫu nhiên, kiểm định này xét tất cả các cách chia dữ liệu đó thành 2 nhóm A/B (hoặc nhiều hơn), sau đó xáo trộn và chia ngẫu nhiên dữ liệu vào các nhóm.
- Áp dụng cho các bài toán kiểm định với cỡ mẫu tương đối nhỏ.
- Đôi khi được gọi là kiểm định chính xác (**exact tests**), do đặc tính thống kê của chúng đảm bảo rằng nếu Giả định là đúng thì xác suất sai lầm loại I (α) sẽ luôn được duy trì ở mức ý nghĩa được mong muốn của kiểm định.
- Cách thức xác định p -value là tương tự như kiểm định hoán vị.

Ví dụ: với 2 nhóm dữ liệu, mỗi nhóm có lần lượt 9 và 8 quan sát. Lúc này,

- tổng số cách chia dữ liệu (chung) thành 2 nhóm là $2^{17} = 131072$ cách;
- số mẫu có 9 quan sát là $C_{17}^9 = 24310$;
- số mẫu có 8 quan sát là $C_{17}^8 = 24310$.

Sau đó, ta sẽ chọn ngẫu nhiên 9 quan sát (trong 17 quan sát) vào nhóm A; phần còn lại thì vào nhóm B.

Bootstrap permutation test

Kiểm định hoán vị bootstrap (bootstrap permutation test) là một biến thể của kiểm định hoán vị, trong đó, bước 3 và 4 được thực hiện với nguyên tắc lấy mẫu có hoàn lại.

Bằng cách này, ta có được

- chỉ định ngẫu nhiên cho đối tượng vào nhóm điều trị;
- lựa chọn ngẫu nhiên các đối tượng từ một quần thể.

Power Analysis

Effect size (cỡ hiệu ứng) là thước đo về sự chênh lệch giữa trung bình hoặc tỷ lệ giữa hai nhóm, và thông qua kiểm định thống kê, ta muốn phát hiện ra.

Ví dụ:

- trong kiểm định trung bình của hai nhóm, thì effect size là $\Delta = \mu_A - \mu_B$;
- trong kiểm định tỷ lệ của hai nhóm, thì effect size là sai số tương đối giữa hai tỷ lệ $\Delta = \frac{p_1 - p_2}{p_2}$.

↪ Giả thuyết tương ứng là $\Delta = 0$;

↪ Đối thuyết tương ứng là $\Delta \neq 0$.

Độ mạnh của kiểm định (Power) là khả năng nhận ra được Giả thuyết là sai, hay tương đương với nhận diện được effect size, tương ứng với một cỡ mẫu cố định.

Power Analysis



- effect size \iff kích cỡ cục bột cacao;
- cỡ mẫu \iff độ dày của lưới;
- power \iff khả năng lọc được cục bột của từng loại lưới.

Power Analysis

Về mặt ý tưởng, cỡ mẫu càng cao, thì độ mạnh của một kiểm định càng lớn.

↪ Câu hỏi:

“Cỡ mẫu bao nhiêu là đủ để đạt được độ mạnh chấp nhận được?”

Trong thống kê, độ mạnh bằng 0.8 (tức 80%) được coi là chấp nhận được.

Để tìm được cỡ mẫu cần thiết, ta cần có:

- effect size mong muốn được nhận ra bởi kiểm định;
- phân phối lý thuyết của dữ liệu (bài toán kiểm định trung bình);
- phương sai hoặc độ lệch chuẩn lý thuyết (bài toán kiểm định trung bình);
- độ mạnh mong muốn của kiểm định;
- kiểm định hai phía, hay kiểm định một phía.

Một quá trình ước tính độ mạnh theo từng kịch bản của cỡ mẫu và effect size được gọi là **phân tích độ mạnh (power analysis)**.

Phân tích độ mạnh cho kiểm định trung bình thường phức tạp hơn so với kiểm định tỷ lệ.

Power Analysis

Ví dụ: Xét kết quả về tỷ lệ chuyển đổi sản phẩm của một mặt hàng trên trang thương mại điện tử với hai mức giá A và B:

	Giá A	Giá B
Chuyển đổi	200	182
Không chuyển đổi	23539	22406

Ta tính được

■ $\hat{p}_A = 0.8425\%$

■ $\hat{p}_B = 0.8057\%$

- tỷ lệ chênh lệch tương đối là $\frac{\hat{p}_A - \hat{p}_B}{\hat{p}_B} \approx 4.57\%$, gần 5%, đủ lớn để có ý nghĩa trong một doanh nghiệp có khối lượng lớn.

Tuy nhiên, ta vẫn có thể hoài nghi về sự chênh lệch này có thực sự có nghĩa hay chỉ là kết quả của sự ngẫu nhiên.

Giả thuyết: “Nếu hai mức giá có cùng tỷ lệ chuyển đổi, liệu sự thay đổi ngẫu nhiên có thể tạo ra sự khác biệt lớn tới 5% không?”

Power Analysis

Nếu Giả thuyết là sai, thì với cỡ mẫu ta đang có (23739 và 22588), độ mạnh của kiểm định là bao nhiêu khi nhận diện effect size 5%?

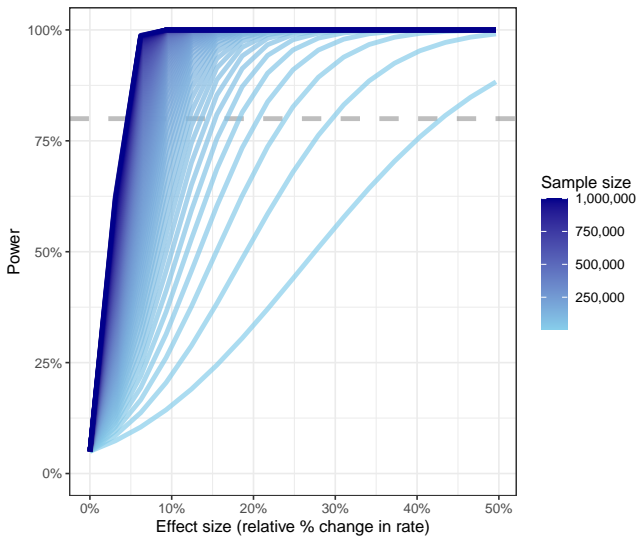
Áp dụng phân tích độ mạnh cho kiểm định tỷ lệ, độ mạnh \approx **11%**.

Hay có nghĩa xác suất mắc sai lầm loại II \approx **89%** (chấp nhận Giả thuyết khi nó thực sự sai).

Và để đạt được độ mạnh xấp xỉ 80%, ta cần cỡ mẫu \approx **746257** cho mỗi mức giá.

Tức gấp gần **32 lần** cỡ mẫu hiện tại.

Power Analysis



1 Kiểm định giả thuyết

2 A/B testing hai nhóm

3 A/B testing nhiều nhóm

4 Multi-Armed Bandit

Bài toán

Giả sử rằng, ta có một bộ dữ liệu về thời gian phiên trung bình của 4 trang web 1, 2, 3 và 4. Ta quan tâm tới câu hỏi:

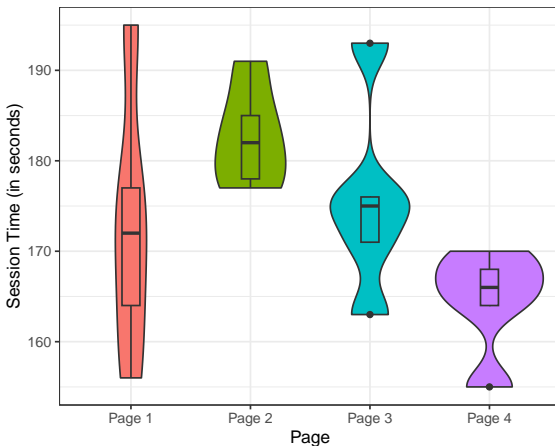
“Thời gian phiên trung bình của 4 trang web là như nhau?”

Lúc này, Giả thuyết cần kiểm chứng sẽ là

$$\mu_1 = \mu_2 = \mu_3 = \mu_4$$

tức là, ta cần so sánh 4 trung bình cùng một lúc thay vì 2.

Bài toán



Multiple testing

Một cách trực quan, ta sẽ tiến hành so sánh theo từng cặp:

$$\blacksquare \mu_1 = \mu_2$$

$$\blacksquare \mu_1 = \mu_4$$

$$\blacksquare \mu_2 = \mu_4$$

$$\blacksquare \mu_1 = \mu_3$$

$$\blacksquare \mu_2 = \mu_3$$

$$\blacksquare \mu_3 = \mu_4$$

↪ tổng cộng 6 cặp.

↪ ta cần lặp lại quá trình kiểm định hoán vị trong 6 lần!

Multiple testing

Một cách trực quan, ta sẽ tiến hành so sánh theo từng cặp:

$$\blacksquare \mu_1 = \mu_2$$

$$\blacksquare \mu_1 = \mu_4$$

$$\blacksquare \mu_2 = \mu_4$$

$$\blacksquare \mu_1 = \mu_3$$

$$\blacksquare \mu_2 = \mu_3$$

$$\blacksquare \mu_3 = \mu_4$$

↪ tổng cộng 6 cặp.

↪ ta cần lặp lại quá trình kiểm định hoán vị trong 6 lần!

Có hai vấn đề nảy sinh:

- thời gian tính toán;
- xác suất tạo ra sai lầm loại I tăng lên.

Multiple testing

Khi số nhóm cần so sánh là k thì ta cần kiểm định $k(k-1)/2$ cặp,

↪ Tức là lặp đi lặp lại $k(k-1)/2$ kiểm định hoán vị.

Ví dụ:

- với $k = 10$ thì cần kiểm định 45 Giả thuyết;
- với $k = 100$ thì cần kiểm định 450 Giả thuyết.

Khi ta phải kiểm định m Giả thuyết cùng một lúc, ta có thể có nhiều hơn một lần phạm sai lầm loại I.

Quyết định	Sự thật của Giả thuyết		Tổng
	Đúng	Sai	
Không bác bỏ	U	W	$m - R$
Bác bỏ	V	S	R
Tổng	m_0	$m - m_0$	m

Khi này, ta quan tâm tới xác suất $\Pr(V \geq 1)$.

Xác suất này được gọi là **family-wise error rate** (FWER).

Multiple testing

Bằng việc giả sử rằng

- m Giả định là đúng;
- m kiểm định là độc lập;
- đặt mức sai lầm loại I cho mỗi kiểm định là α

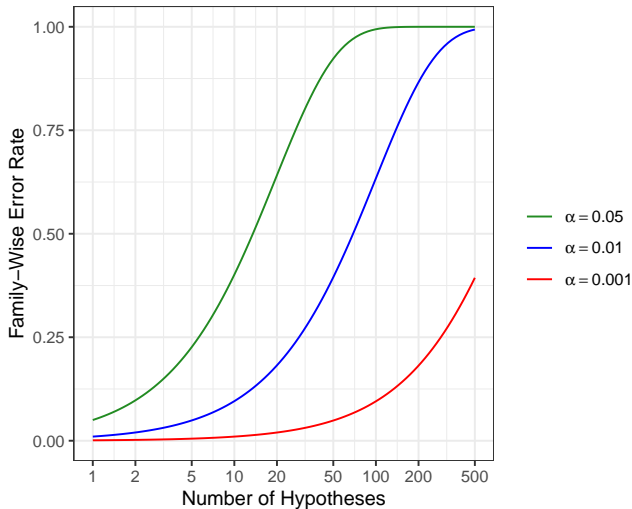
ta xác định được rằng

$$FWER(\alpha) = 1 - (1 - \alpha)^m.$$

Với $\alpha = 0.5$, nếu

- $m = 1$, thì $FWER(\alpha) = 0.05$, chính là mức kiểm soát sai lầm I;
- $m = 6$, thì $FWER(\alpha) = 0.27$;
- $m = 90$, thì $FWER(\alpha) = 0.99$, tức là chắc chắn sẽ có ít nhất 1 sai lầm loại I.

Multiple testing



Multiple testing

Để điều khiển FWER khi kiểm định m Giả thuyết, ta có một số phương pháp hiệu chỉnh:

- Phương pháp Bonferroni (**Bonferroni method**): phương pháp này đề xuất ngưỡng xác suất sai lầm loại I cho từng kiểm định là α/m ;
↪ khi đó, ta bác bỏ tất cả Giả thuyết có p -value $< \alpha/m$;
↪ làm như vậy

$$FWER(\alpha) \leq m \times \frac{\alpha}{m} = \alpha$$

Ví dụ, với $m = 100$ và $\alpha = 0.05$, khi đó, p -value sẽ so sánh với $0.05/100 = 0.0005$.

- Phương pháp Holm (**Holm's method**)
- Phương pháp Tukey (**Tukey's method**)
- Phương pháp Scheffé (**Scheffé's method**)

ANOVA

Lời giải tốt nhất cho việc kiểm định Giả thuyết

$$\mu_1 = \mu_2 = \mu_3 = \mu_4 = \dots = \mu_k$$

là sử dụng phương pháp **analysis of variance**, hay **ANOVA**.

Phương pháp này sẽ so sánh đồng thời trung bình của nhiều nhóm độc lập trong một phân tích duy nhất.

Phân tích phương sai có vẻ như là một cách gọi sai, do mục đích của chúng ta là so sánh *trung bình*, nhưng việc kiểm tra sự biến động giữa các nhóm cũng tương đương với việc hỏi liệu các trung bình có khác nhau hay không.

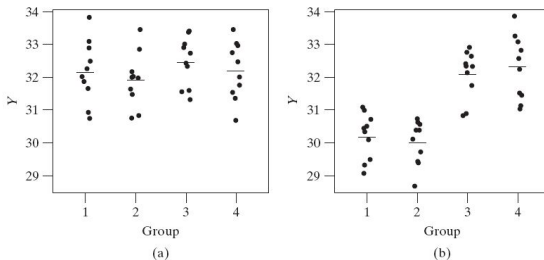
ANOVA được phát triển bởi nhà sinh vật học và nhà thống kê Sir. Ronald Aylmer Fisher¹ (1890 - 1962).

¹https://en.wikipedia.org/wiki/Ronald_Fisher

ANOVA

Giả sử, ta muốn kiểm định Giả thuyết: $\mu_1 = \mu_2 = \mu_3 = \mu_4$.

Biểu đồ dưới đây minh họa giá trị quan sát của biến Y tương ứng trong 4 nhóm.

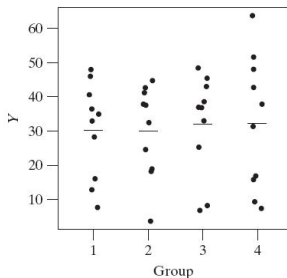
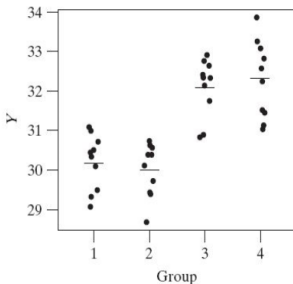


Các trung bình mẫu được hiển thị bởi đường nét liền trong từng nhóm.

- Hình (a), trực giác của chúng ta gợi ý rằng Giả thuyết có lẽ là đúng;
- Hình (b), rõ ràng rằng Giả thuyết là sai, ta có thể nhìn thấy sự khác biệt đáng kể giữa các trung bình của nhóm.

ANOVA

Tuy nhiên,



biểu đồ bên phải thể hiện một trường hợp kém rõ ràng khác.

- Thoạt đầu, ta nghĩ rằng Giả thuyết có lẽ đã đúng.
- Tuy nhiên, thực tế Giả thuyết là sai trong trường hợp - giá trị trung bình mẫu là giống như trong biểu đồ bên trái (tức là biểu đồ (b) ở slide trước đó).

Vấn đề ở đây là khi độ lệch chuẩn trong mỗi nhóm là tương đối lớn, nó sẽ khiến cho việc nhận định sự khác biệt của giá trị trung bình trở nên khó khăn.

ANOVA

Để tìm ra bằng chứng thuyết phục về sự khác biệt về trung bình, chúng ta cần tính đến:

- (1) sự biến động giữa các nhóm với nhau - variation between groups;
- (2) sự biến động vốn có trong các nhóm - variation within groups.

Thông qua việc so sánh độ lớn tương đối của hai loại biến động này - tức là “phân tích phương sai” - mà chúng ta có thể suy luận về trung bình.

Về mặt nguyên tắc, nếu

- sự biến động giữa các nhóm là khá giống với sự biến động trong các nhóm, điều đó có nghĩa là tất cả các nhóm được lấy mẫu từ cùng một quần thể, và do đó, tất cả các trung bình của nhóm đều giống nhau;
- ngược lại, có ít nhất một trung bình là khác với những cái còn lại.

ANOVA

Ta có:

one-way ANOVA so sánh trung bình của ba nhóm trở lên; thuật ngữ “one-way” đề cập đến thực tế là có một biến phân loại xác định các nhóm hoặc các phương pháp điều trị.

two-way ANOVA xem xét ảnh hưởng của hai biến định tính độc lập lên một biến định lượng.

k-way ANOVA xem xét ảnh hưởng của k biến định tính độc lập lên một biến định lượng.

One-way ANOVA: cách xử lý cổ điển

Giả thuyết thống kê

Đối với kiểm định one-way ANOVA, với k nhóm, ta quan tâm tới:

$$\begin{cases} \text{Giả thuyết : } \mu_1 = \mu_2 = \dots = \mu_k, \\ \text{Đối thuyết : Ít nhất có một trung bình là khác với những cái còn lại.} \end{cases}$$

Điều kiện cho one-way ANOVA

Để thực hiện được kiểm định one-way ANOVA, các điều kiện dưới đây phải được thỏa mãn:

1. Dữ liệu phải được lấy mẫu ngẫu nhiên.
2. Các nhóm phải độc lập với nhau.
3. Dữ liệu trong mỗi nhóm phải tuân theo phân phối chuẩn.
4. Phương sai trong các nhóm phải giống nhau (giả định về sự đồng nhất phương sai - the assumption of homogeneity of variance), tức là

$$\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2.$$

↪ ta cần kiểm chứng toàn bộ điều kiện này trước khi phân tích ANOVA.

One-way ANOVA: cách xử lý cổ điển

Đặt các ký hiệu sau:

- Y_1, Y_2, \dots, Y_k là các mẫu của Y tương ứng với k nhóm;
- n_1, n_2, \dots, n_k là các cỡ mẫu;
- $N = \sum_{j=1}^k n_j$ tổng số quan sát trong dữ liệu;
- \bar{Y}_j là trung bình mẫu trong nhóm thứ j , với $j = 1, 2, \dots, k$:

$$\bar{Y}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} Y_{ji},$$

- s_j^2 là phương sai mẫu trong nhóm thứ j :

$$s_j^2 = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (Y_{ji} - \bar{Y}_j)^2,$$

- \bar{Y} là trung bình mẫu của tất cả N quan sát:

$$\bar{Y} = \frac{1}{N} \sum_{j=1}^k \sum_{i=1}^{n_j} Y_{ji}.$$

One-way ANOVA: cách xử lý cổ điển

Ta tính được, tổng bình phương trong các nhóm - sum of squares within groups (SSW) là:

$$SSW = \sum_{j=1}^k (n_j - 1) s_j^2,$$

và tổng bình phương giữa các nhóm - sum of squares between groups (SSB) là:

$$SSB = \sum_{j=1}^k n_j (\bar{Y}_j - \bar{Y})^2.$$

Khi đó, ta xác định được:

- độ đo cho sự biến động trong các nhóm - **mean square within** groups (MSW) là:

$$MSW = \frac{SSW}{N - k},$$

- độ đo cho sự biến động giữa các nhóm - **mean square between** groups (MSB) là:

$$MSB = \frac{SSB}{k - 1}.$$

One-way ANOVA: cách xử lý cổ điển

Giả sử rằng tất cả điều kiện (1) - (4) đều thỏa, khi đó, dưới giả định H_0 , tức là

$$\mu_1 = \mu_2 = \dots = \mu_k = \mu$$

ta có $Y_{ji} \sim \mathcal{N}(\mu, \sigma^2)$. Khi này, ta chứng minh được $SSW \sim \chi_{N-k}^2$ và $SSB \sim \chi_{k-1}^2$. Điều này, dẫn tới

$$F = \frac{MSB}{MSW} \sim F(k-1, N-k),$$

tuân theo phân phối¹ F với $k-1$ và $N-k$ bậc tự do

Ta ký hiệu giá trị F được tính từ dữ liệu quan sát là F_{obs} .

Ta định nghĩ p -value cho kiểm định one-way ANOVA bởi:

$$p\text{-value} = \Pr(F > F_{obs}),$$

tức là xác suất ở đuôi bên phải của phân phối $F(k-1, N-k)$.

¹<https://en.wikipedia.org/wiki/F-distribution>

One-way ANOVA: cách xử lý cổ điển

Thông thường, ta sẽ tổng hợp các kết quả thành một bảng, và gọi là bảng one-way ANOVA, như dưới đây:

Bảng: Bảng one-way ANOVA

Nguồn	df	Tổng bình phương	Trung bình bình phương	F_{obs}	p -value
Nhóm	$k - 1$	SSB	MSB	MSB / MSW	$\Pr(F > F_{obs})$
Sai số	$N - k$	SSW	MSW		

ở đây, df viết tắt của “bậc tự do”.

One-way ANOVA: ví dụ

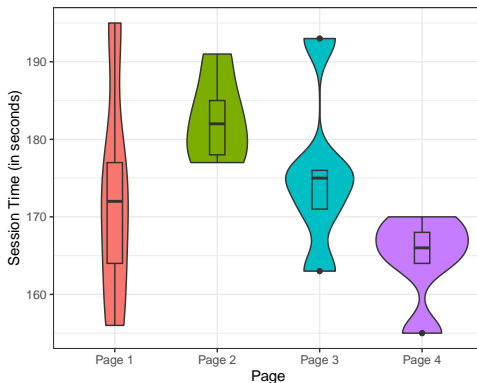
Xét lại bài toán về so sánh thời gian phiên trung bình của 4 trang web.

Ta quan tâm tới câu hỏi:

“Thời gian phiên trung bình của 4 trang web là như nhau?”

Lúc này, giả thuyết cần kiểm chứng sẽ là

$$\mu_1 = \mu_2 = \mu_3 = \mu_4$$



One-way ANOVA: ví dụ

Một số thông tin tổng hợp từ dữ liệu:

- số nhóm $k = 4$;
- tổng số quan sát $N = 20$.

Sau quá trình xử lý one-way ANOVA cổ điển, ta thu được kết quả:

Nguồn	df	Tổng bình phương	Trung bình bình phương	F_{obs}	p -value
Trang web	3	831.4	277.1	2.74	0.0776
Sai số	16	1618.4	101.2		

Cố định mức ý nghĩa $\alpha = 0.05$, ta thấy p -value > 0.05 .

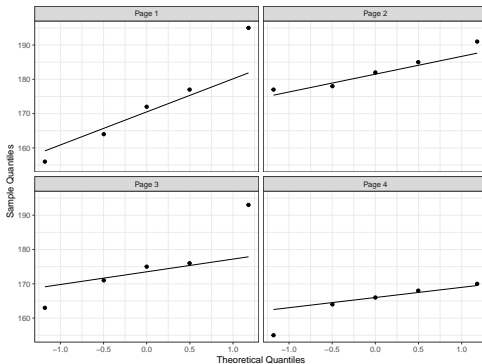
⇒ Không đủ bằng chứng để bác bỏ giả thuyết H_0 , tức là

Thời gian phiên trung bình của 4 trang web có thể là như nhau.

One-way ANOVA: ví dụ

Tuy nhiên, kết quả này liệu có đáng tin cậy?

Nhóm	Cỡ mẫu	Trung bình	Phương sai
Trang 1	5	172.8	217.7
Trang 2	5	182.6	32.3
Trang 3	5	175.6	120.8
Trang 4	5	164.6	33.8



One-way ANOVA: xử lý với Permutation ANOVA

Vấn đề của one-way ANOVA cổ điển

Khi có ít nhất 1 trong số các điều kiện:

- dữ liệu trong các nhóm là xấp xỉ một phân phối chuẩn;
- phương sai (lý thuyết) trong các nhóm là đồng nhất,

không thỏa mãn \Rightarrow xấp xỉ lý thuyết $F = \frac{MSB}{MSW} \sim F(k-1, N-k)$ không được chính xác.

Khắc phục

Để khắc phục, ta có thể thực hiện kiểm định độc lập bằng cách áp dụng ý tưởng của permutation test.

\Rightarrow tạo ra phân phối mẫu thực nghiệm cho F thay vì phân phối lý thuyết.

One-way ANOVA: xử lý với Permutation ANOVA

Cụ thể, ta có thể sử dụng thuật toán như sau:

1. Từ mẫu gốc, tính SSW , SSB , MSW , MSB và F .
2. Gộp chung các mẫu lại với nhau, tạo thành một mẫu lớn và xáo trộn chúng.
3. Lấy mẫu ngẫu nhiên không lặp cho các nhóm j , $j = 1, 2, \dots, k$ (với cỡ mẫu như ban đầu) từ mẫu chung.
4. Tính trung bình mẫu tổng thể và cho các nhóm mới tạo.
5. Tính phương sai mẫu tổng thể và cho các nhóm mới tạo.
6. Tính SSW^* , SSB^* , MSW^* , MSB^* và F^* .
7. Lặp lại các bước từ 3 tới 6, trong R lần.

Sau đó, xác định

$$p\text{-value} = \frac{1}{R} \sum_{i=1}^R \mathbf{I}(F_i^* > F),$$

với F_i^* là kết quả thu được trong bước 6 tương ứng lần lặp thứ i .

One-way ANOVA: xử lý với Permutation ANOVA

Xét lại bài toán về so sánh thời gian phiên trung bình của 4 trang web.

Giả thuyết cần kiểm chứng là

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4.$$

Áp dụng công thức tính SSW , SSB , MSW , MSB và F

Nguồn	df	Tổng bình phương	Trung bình bình phương	F_{obs}
Trang web	3	831.4	277.1	2.74
Sai số	16	1618.4	101.2	

Áp dụng cách xử lý bằng Permutation ANOVA, với 5000 lần lặp, ta thu được p -value là

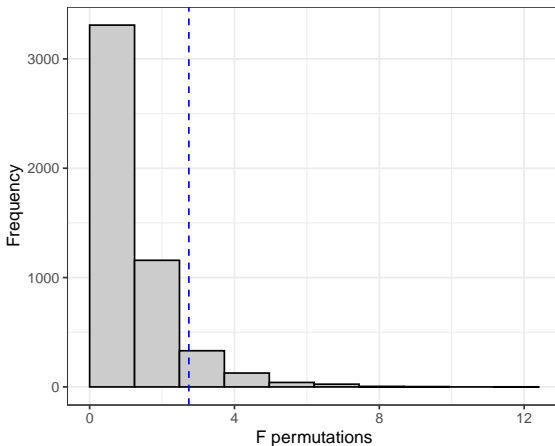
$$p\text{-value} = \frac{1}{5000} \sum_{i=1}^{5000} \mathbf{I}(F_i^* > 2.74) = 0.084$$

Cổ định mức ý nghĩa $\alpha = 0.05$, ta thấy $p\text{-value} > 0.05$.

⇒ Không đủ bằng chứng để bác bỏ giả thuyết H_0 , tức là

Thời gian phiên trung bình của 4 trang web có thể là như nhau.

One-way ANOVA: xử lý với Permutation ANOVA



Bài toán A/B testing nhiều nhóm định tính

Giả sử một công ty đang thử nghiệm ba tiêu đề web khác nhau: A, B và C, và công ty chạy mỗi tiêu đề trên 1.000 khách truy cập, với kết quả được hiển thị trong Bảng sau:

		Tiêu đề web		
		Tiêu đề A	Tiêu đề B	Tiêu đề C
Hành động	Click	14	8	12
	No-click	986	992	988

Dựa vào kết quả này, công ty muốn kiểm tra giả thiết

“Tỷ lệ click là **phụ thuộc** vào tiêu đề trang web.”

Bài toán A/B testing nhiều nhóm định tính

Do hai biến

- “Tiêu đề trang web”
- “Hành động”

là hai biến định tính (định danh).

⇒ Giả thuyết tương đương với kiểm tra **sự độc lập** của hai biến định tính.

Kiểm định độc lập: xử lý cổ điển

Kiểm định độc lập

Kiểm định độc lập (Chi-square contingency test) (hay còn được gọi là **Chi-square independent test**) là kiểm định được sử dụng phổ biến nhất trong việc kiểm tra mối liên hệ giữa hai biến định tính.

Đối với kiểm định độc lập, các phát biểu về giả thuyết H_0 và đối thuyết H_1 có thể khác nhau tùy theo ngữ cảnh, nhưng chúng luôn có ý nghĩa sau:

H_0 : Hai biến định tính là độc lập (nghĩa là không có mối liên hệ nào giữa hai biến định tính).

H_1 : Hai biến định tính là **không** độc lập (nghĩa là có mối liên hệ giữa hai biến định tính là có ý nghĩa thống kê).

Đối với bài toán đang xét, dữ liệu tiêu đề web, ta có các giả thuyết như sau:

H_0 : Hành động click không có xu hướng đối với loại tiêu đề.

H_1 : Tồn tại một xu hướng nhất định trong hành động click của khách hàng và tiêu đề web.

Kiểm định độc lập: xử lý cổ điển

Nhắc lại định nghĩa hai biến ngẫu nhiên độc lập

Ta biết rằng, hai biến cố A và B là độc lập khi và chỉ khi

$$\Pr(A \cap B) = \Pr(A) \Pr(B)$$

Hai biến ngẫu nhiên rời rạc X và Y được gọi là độc lập nếu và chỉ nếu

$$\Pr(X = i, Y = j) = \Pr(X = i) \Pr(Y = j),$$

với mọi i và j trong tập giá trị của X và Y , tương ứng.

Nhắc lại rằng hai biến định tính tương ứng là 2 biến đa thức thông qua phép đặt giả biến.

Kiểm định độc lập: xử lý cổ điển

Giả sử rằng biến định tính X và Y có lần lượt k và m tính chất. Ta biểu diễn trong bảng dưới đây:

		Biến định tính X			
		Cột 1	Cột 2	...	Cột k
Biến định tính Y	dòng 1	n_{11}	n_{12}	...	n_{1k}
	dòng 2	n_{21}	n_{22}	...	n_{2k}

	dòng m	n_{m1}	n_{m2}	...	n_{mk}

} tổng dòng r_i

{
tổng cột c_j

Ta đặt

- n_{ij} là số lượng quan sát có tính chất i của X và tính chất j của Y ;
- r_i là tổng các quan sát trên dòng i ;
- c_j là tổng các quan sát trên cột j ;
- N là tổng số lượng quan sát.

Kiểm định độc lập: xử lý cổ điển

Theo bảng ta có

- trên mỗi cột, $\Pr(X = j) = \frac{c_j}{N}$,
- trên mỗi dòng, $\Pr(Y = i) = \frac{r_i}{N}$,
- $\Pr(X = j, Y = i) = \frac{n_{ij}}{N}$

Do đó, nếu H_0 là đúng thì:

$$\frac{n_{ij}}{N} \approx \frac{c_j}{N} \times \frac{r_i}{N}$$

Đặt

$$E_{ij} = \frac{r_i \times c_j}{N},$$

là tần số kỳ vọng của ô tương ứng dòng thứ i và cột thứ j trong bảng tổng hợp. Khi đó, nếu H_0 là đúng thì n_{ij} sẽ gần bằng E_{ij} .

Kiểm định độc lập: xử lý cổ điển

Xét lại ví dụ tiêu đề trang, ta có tổng dòng và tổng cột là:

		Tiêu đề web			
		Tiêu đề A	Tiêu đề B	Tiêu đề C	Tổng dòng
Hành động	Click	14	8	12	34
	No-click	986	992	988	2966
Tổng cột		1000	1000	1000	

Từ đây, ta tính được, ví dụ

$$E_{11} = \frac{r_1 \times c_1}{N} = \frac{34 \times 1000}{3000} = 11.33$$

Kiểm định độc lập: xử lý cổ điển

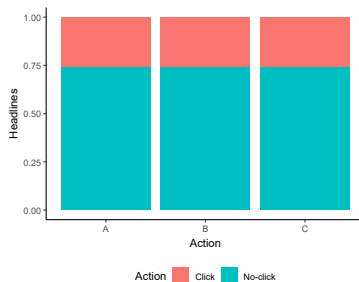
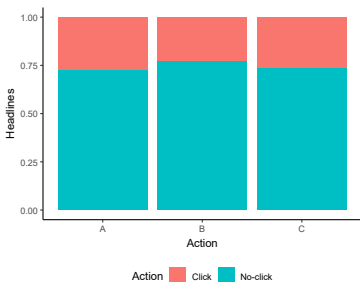
Và ta thu được bảng tần số kỳ vọng như sau:

		Tiêu đề web		
		Tiêu đề A	Tiêu đề B	Tiêu đề C
Hành động	Click	11.33	11.33	11.33
	No-click	988.67	988.67	988.67
Tổng cột		1000	1000	1000

So sánh với bảng tần số (quan sát), ta dễ dàng nhận ra sự không quá khác biệt.

Kiểm định độc lập: xử lý cổ điển

Để hiểu hơn vì giả thuyết H_0 , ta so sánh biểu đồ của bảng quan sát và tần số kỳ vọng (dưới giả định H_0 là đúng).



Tỷ lệ hành động là không đổi trong các cột loại tiêu đề.

Kiểm định độc lập: xử lý cổ điển

Với ý tưởng so sánh tần số quan sát với tần số kỳ vọng, đưa ta tới thống kê

$$Q_n = \sum_{i=1}^m \sum_{j=1}^k \frac{(n_{ij} - E_{ij})^2}{E_{ij}}$$

Nếu H_0 là đúng và tất cả các $E_{ij} \geq 5$, thì $Q_n \sim \chi^2_{(m-1) \times (k-1)}$.

Với một dữ liệu quan sát, ta tính được giá trị quan sát w của thống kê Q_n , ta định nghĩa

$$p\text{-value} = \Pr(Q_n > w).$$

tức là xác suất ở đuôi bên phải của phân phối $\chi^2_{(m-1) \times (k-1)}$.

Kiểm định độc lập: xử lý cổ điển

Quay trở lại với ví dụ tiêu đề trang web.

- Ta tính được giá trị quan sát

$$w = \frac{(12 - 11.33)^2}{11.33} + \frac{(8 - 11.33)^2}{11.33} + \dots + \frac{(988.67 - 988)^2}{988} = 1.666$$

- Thống kê Q_n tuân theo phân phối χ_2^2 (do $m = 2$ và $k = 3$).
- Giá trị p -value là 0.4348 (lớn hơn 0.05) \Rightarrow ta không đủ cơ sở để bác bỏ H_0
 \Rightarrow Hành động click không có liên quan tới dạng tiêu đề web.

Kiểm định độc lập: xử lý với Resampling method

Vấn đề của kiểm định độc lập

Khi có ít nhất 1 tần số kỳ vọng $E_{ij} < 5$, thì xấp xỉ lý thuyết $Q_n \sim \chi^2_{(m-1) \times (k-1)}$ không được chính xác.

Khắc phục

Để khắc phục, ta có thể thực hiện kiểm định độc lập bằng cách áp dụng phương pháp lấy lại mẫu (resampling method), tương tự như bootstrap hoặc permutation.

⇒ tạo ra phân phối mẫu thực nghiệm cho Q_n thay vì phân phối lý thuyết.

Kiểm định độc lập: xử lý với Resampling method

Cụ thể, ta có thể sử dụng thuật toán lấy mẫu lại này:

1. Từ mẫu gốc tính r_i , c_j , E_{ij} và w .
2. Tạo một mẫu có r_i số lượng của nhóm thứ i , tương ứng với $i = 1, \dots, m$.
3. Trộn, lấy k mẫu riêng biệt gồm c_j số lượng của nhóm thứ j , tương ứng với $j = 1, \dots, k$.
4. Dựa trên mẫu mới tạo ở bước 2, tính các E_{ij}^* và w^* .
5. Lặp lại các bước 3 và 4 trong R lần.

Sau đó, xác định

$$p\text{-value} = \frac{1}{R} \sum_{i=1}^R \mathbf{I}(w_i^* > w),$$

với w_i^* là kết quả thu được trong bước 5 tương ứng lần lặp thứ i .

Kiểm định độc lập: xử lý với Resampling method

Áp dụng vào trong ví dụ tiêu đề trang web:

1. Từ mẫu gốc tính được $w = 1.666$.
2. Tạo một mẫu có 34 số 1 cho hành động click và 2966 số 0 cho hành động no-click.
3. Trộn, lấy 3 mẫu riêng biệt gồm 1000 số lượng của 3 tiêu đề trang web.
4. Dựa trên mẫu mới tạo ở bước 2, tính các E_{ij}^* và w^* .
5. Lặp lại các bước 3 và 4 trong 5000 lần.

Sau đó, xác định

$$p\text{-value} = \frac{1}{5000} \sum_{i=1}^{5000} I(w_i^* > 1.666).$$

Kết quả thu được $p\text{-value} = 0.4176$ (lớn hơn 0.05) \Rightarrow ta không đủ cơ sở để bác bỏ H_0 .

\Rightarrow Hành động click không có liên quan tới dạng tiêu đề web.

1 Kiểm định giả thuyết

2 A/B testing hai nhóm

3 A/B testing nhiều nhóm

4 Multi-Armed Bandit

Giới thiệu

Ví dụ: Ta có ba máy chơi game, với tỷ lệ thắng cược lần lượt là

- Máy A: 10 chiến thắng trong 50 lượt chơi;
- Máy B: 2 chiến thắng trong 50 lượt chơi;
- Máy C: 4 chiến thắng trong 50 lượt chơi.

Trực quan, ta dễ dàng nói rằng nên tiếp tục chơi máy A.

Tuy nhiên, quyết định như vậy có vội vàng, khi mà ta mới thực hiện 50 lượt chơi?

Điều gì sẽ xảy ra nếu như ta tiếp chơi máy A trong 100 lượt tiếp?

Nếu ta tiếp tục chơi máy A, và số ván thua bắt đầu nhiều lên, liệu ta có tiếc nuôi máy C và B?

Có nên thử với máy C?

→ điều này dẫn tới nhu cầu về một thuật toán mô phỏng tiến trình chơi của 3 máy, để nhìn thấy cơ hội ở máy nào là rõ rệt hơn.

↪ Multi-Armed Bandit Algorithm.

Khai thác và khám phá

Multi-Armed Bandit Algorithm liên quan tới hai khái niệm

- khai thác (exploitation): chọn những lựa chọn có vẻ tốt nhất dựa trên kết quả trong quá khứ
- khám phá (exploration): chọn các lựa chọn chưa được thử (hoặc chưa thử đủ)

Khai thác có liên quan những khái niệm:

- “tham lam” (greedy) và “thiếu cận” (short-sighted).
- khai thác quá nhiều \iff tiếc nuối vì đã bỏ lỡ những cơ hội ngon ăn khác chưa được khám phá.

Khám phá liên quan tới

- “thu thập thông tin” và “nhìn xa trông rộng” (long-sighted).
- khám phá quá nhiều \iff hối tiếc vì đã lãng phí thời gian vào những việc “ngớ ngẩn”.

\hookrightarrow Làm thế nào để cân bằng giữa Khai thác và Khám phá để kết hợp giữa lợi ích thông tin và lợi nhuận một cách tối ưu nhất.

Khai thác và khám phá

Ví dụ:

- Lựa chọn nhà hàng
 - **Khai thác:** Đi tới nhà hàng yêu thích
 - **Khám phá:** Tới một cửa hàng mới
- Biểu ngữ quảng cáo online
 - **Khai thác:** Hiển thị dòng quảng cáo thành công nhất
 - **Khám phá:** Hiển thị một dòng quảng cáo mới
- Lựa chọn môn học
 - **Khai thác:** Lựa chọn những môn được cho là dễ học và điểm cao
 - **Khám phá:** Lựa chọn môn mới/giảng viên mới
- Lựa chọn phim
 - **Khai thác:** Xem những phim/dòng phim mà mình thích
 - **Khám phá:** Thử những phim/dòng phim mới

Các thuật toán Multi-Armed Bandit

Có nhiều thuật toán khác nhau đã được phát triển cho Multi-Armed Bandit:

- epsilon-greedy algorithm
- the softmax algorithm
- the upper confidence bound algorithm
- Multi-Armed Bandit with Thompson Sampling
- Epsilon-First

Mỗi thuật toán có ưu nhược điểm khác nhau.

Ví dụ: lựa chọn quảng cáo

Một khách muốn có lời khuyên nên sử dụng quảng cáo nào trong 5 biến thể.

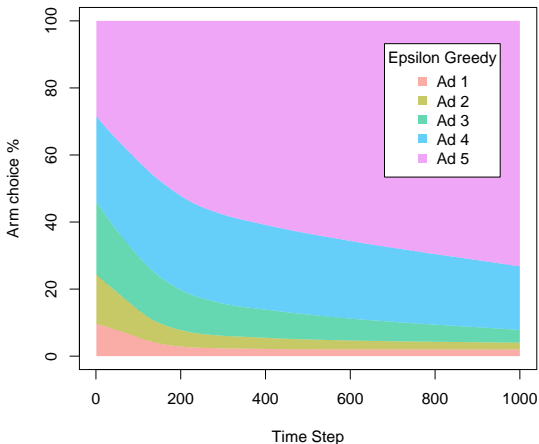
Biết rằng, tỷ lệ chuyển đổi sản phẩm của 5 quảng cáo:

- Quảng cáo 1: 5%
- Quảng cáo 2: 10%
- Quảng cáo 3: 15%
- Quảng cáo 4: 20%
- Quảng cáo 5: 25%

Quảng cáo đóng vai trò là các arms.

Để đưa ra lựa chọn quảng cáo nào là tốt, cũng như không phải hối tiếc những quảng cáo khác, ta áp dụng epsilon-greedy algorithm để mô phỏng % chọn lựa của 5 loại quảng cáo trong 1000 lần thử.

Ví dụ: lựa chọn quảng cáo



Thông qua hình minh họa kết quả, ta thấy rằng Quảng cáo số 5 được lựa chọn nhiều nhất.

Đề tài cuối kỳ

Ta có một số đề tài cuối cùng liên quan tới chủ đề “Multi-Armed Bandit Algorithms”.

1. Tìm hiểu thuật toán Epsilon-greedy algorithm và các ứng dụng trong thực tế.
2. Tìm hiểu thuật toán Softmax algorithm và các ứng dụng trong thực tế.
3. Tìm hiểu thuật toán Upper confidence bound algorithm và các ứng dụng trong thực tế.
4. Tìm hiểu thuật toán Multi-Armed Bandit with Thompson Sampling và các ứng dụng trong thực tế.