

□□□□ HỌC THỐNG KÊ QUA TRUYỆN TRANH □□□□

Bước 4: Tiến hành kiểm tra phân tích phương sai (ANOVA).

Dưới đây là các giả định của chúng ta về các hệ số hồi quy từng phần. a_1 , a_2 , b là các hệ số của toàn bộ tổng thể.



Nếu phương trình hồi quy thu được là

$$y = a_1x_1 + a_2x_2 + b$$

- A_1 xấp xỉ bằng a_1 .
- A_2 xấp xỉ bằng a_2 .
- B xấp xỉ bằng b .

$$\sigma = \sqrt{\frac{S_e}{\text{cỡ mẫu} - \text{số lượng biến dự đoán} - 1}}$$

Cậu có thể áp dụng điều này cho dữ liệu của tiệm bánh kazami không?



Được.



Phương trình hồi quy đa biến là:
 $y = 41.5x_1 - 0.3x_2 + 65.3$,
Vì vậy...

Đây là những giả định của chúng ta



- A_1 xấp xỉ bằng 41,5.
- A_2 xấp xỉ bằng -0,3.
- B xấp xỉ bằng 65,3.
- $\sigma = \sqrt{\frac{4173.0}{10-2-1}} = 24.4$.

Tuyệt vời!





Một là kiểm định tất cả các hệ số hồi quy từng phần cùng nhau.

Giả thuyết không	$A_1 = 0$ và $A_2 = 0$
Giả thuyết thay thế	Không phải $A_1 = A_2 = 0$

Nói cách khác, một trong những giả thuyết sau là đúng:

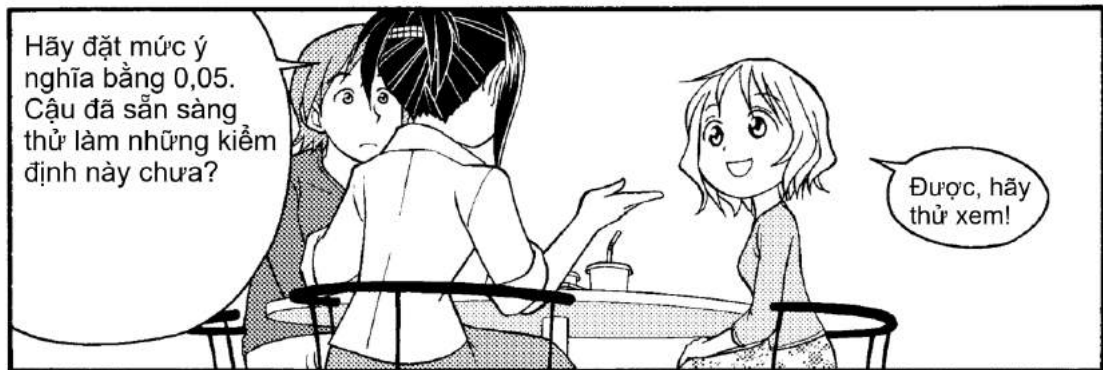
- $A_1 \neq 0$ $A_2 \neq 0$
- $A_1 \neq 0$ $A_2 = 0$
- $A_1 = 0$ $A_2 \neq 0$

Hai là kiểm định các hệ số hồi quy từng phần riêng lẻ.

Giả thuyết không	$A_i = 0$
Giả thuyết thay thế	$A_i \neq 0$

Vì vậy, chúng ta phải lập lại kiểm định này cho từng hệ số hồi quy từng phần?

Đúng!



Trước hết, chúng ta sẽ cùng nhau kiểm định tất cả các hệ số hồi quy từng phần.



Các bước của ANOVA

Bước 1	Xác định tổng thể.	Tổng thể là tất cả các cửa hàng Kazami Bakery.
Bước 2	Thiết lập một giả thuyết không và một giả thuyết thay thế.	Giả thuyết không là $A_1 = 0$ và $A_2 = 0$. Giả thuyết thay thế là A_1 hoặc A_2 hoặc cả hai $\neq 0$.
Bước 3	Lựa chọn kiểm định giả thuyết để thực hiện.	Chúng ta sẽ sử dụng kiểm định F.
Bước 4	Chọn mức ý nghĩa.	Chúng ta sẽ sử dụng mức ý nghĩa bằng 0,05
Bước 5	Tính toán kiểm định thống kê từ dữ liệu mẫu.	Kiểm định thống kê là: $\frac{\frac{S_{yy} - S_e}{\text{Số lượng biến dự đoán}}}{\frac{S_e}{\text{cỡ mẫu} - \text{số lượng biến dự đoán} - 1}} =$ $\frac{76199.6 - 4173.0}{2} \div \frac{4173.0}{10 - 2 - 1} = 60.4$ Kiểm định thống kê, 60,4, sẽ tuân theo phân phối F với bậc tự do thứ nhất là 2 (số biến dự đoán) và bậc tự do thứ hai là 7 (cỡ mẫu trừ đi số biến dự đoán trừ 1), nếu giả thuyết không là đúng.
Bước 6	Xác định xem giá trị p đối với kiểm định thống kê thu được ở Bước 5 có nhỏ hơn mức ý nghĩa hay không.	Ở mức ý nghĩa 0,05, với d là 2 và d là 7 (10 - 2 - 1), giá trị tới hạn là 4,7374. Kiểm định thống kê của chúng ta là 60,4.
Bước 7	Quyết định xem câu có thể bác bỏ giả thuyết không hay không.	Vì kiểm định thống kê của chúng ta lớn hơn giá trị tới hạn nên chúng ta bác bỏ giả thuyết không.

Tiếp theo, chúng ta sẽ kiểm định các hệ số hồi quy từng phần riêng lẻ. Tớ sẽ làm ví dụ với A.



Các bước của ANOVA

Bước 1	Xác định tổng thể.	Tổng thể là tất cả các cửa hàng Kazami Bakery.
Bước 2	Thiết lập một giả thuyết không và một giả thuyết thay thế.	Giả thuyết không là $A = 0$. Giả thuyết thay thế là $A \neq 0$.
Bước 3	Lựa chọn kiểm định giả thuyết để thực hiện.	Chúng ta sẽ sử dụng kiểm định F.
Bước 4	Chọn mức ý nghĩa.	Chúng ta sẽ sử dụng mức ý nghĩa bằng 0,05.
Bước 5	Tính toán kiểm định thống kê từ dữ liệu mẫu.	Kiểm định thống kê là: $\frac{\alpha_1^2}{S_{11}} \div \frac{S_e}{\text{cỡ mẫu} - \text{số lượng biến dự đoán} - 1} =$ $\frac{41.5^2}{0.0657} \div \frac{4173.0}{10 - 2 - 1} = 44$
Bước 6	Xác định xem giá trị p đối với kiểm định thống kê thu được ở Bước 5 có nhỏ hơn mức ý nghĩa hay không.	Kiểm định thống kê sẽ tuân theo phân phối F với bậc tự do thứ nhất là 1 và bậc tự do thứ hai là 7 (cỡ mẫu trừ đi số biến dự đoán trừ 1), nếu giả thuyết không là đúng. (Giá trị của S_{11} sẽ được giải thích ở trang tiếp theo) Ở mức ý nghĩa 0,05, với d là 1 và d là 7, giá trị tới hạn là 5,5914. Kiểm định thống kê của chúng ta là 44.
Bước 7	Quyết định xem cậu có thể bác bỏ giả thuyết không hay không.	Vì kiểm định thống kê của chúng ta lớn hơn giá trị tới hạn nên chúng ta bác bỏ giả thuyết không.

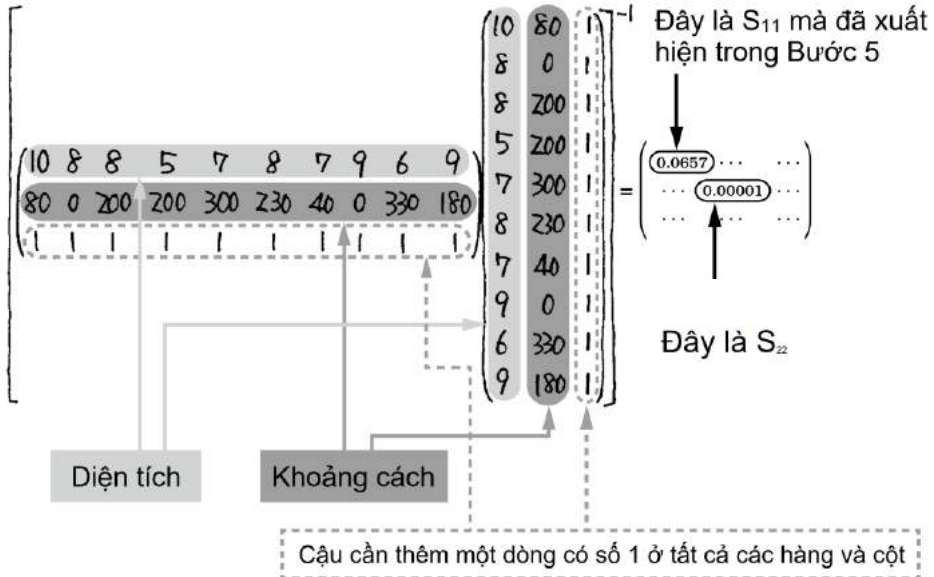
Bất kể kết quả của bước 7 là gì, nếu giá trị của kiểm định thống kê

$$\frac{\alpha_1^2}{S_{11}} \div \frac{S_e}{\text{cỡ mẫu} - \text{số lượng biến dự đoán} - 1}$$

bằng 2 trở lên, chúng ta vẫn coi biến dự đoán tương ứng với hệ số hồi quy từng phần đó là hữu ích cho việc dự đoán biến kết quả.



Tìm S_{11} và S_{22}

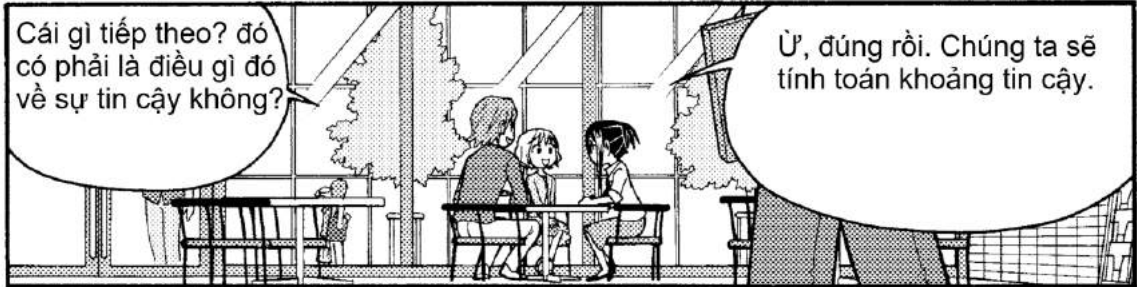


Chúng ta sử dụng ma trận để tìm S_{11} và S_{22} . Chúng ta cần S_{11} để tính toán kiểm định thống kê ở trang trước và chúng ta sử dụng S_{22} để kiểm định hệ số thứ hai một cách độc lập theo cách tương tự.*



*Một số người sử dụng phân phối t thay vì phân phối F khi giải thích "kiểm định hệ số hồi quy từng phần". Kết quả cuối cùng của cậu sẽ giống nhau cho dù cậu chọn phương pháp nào.

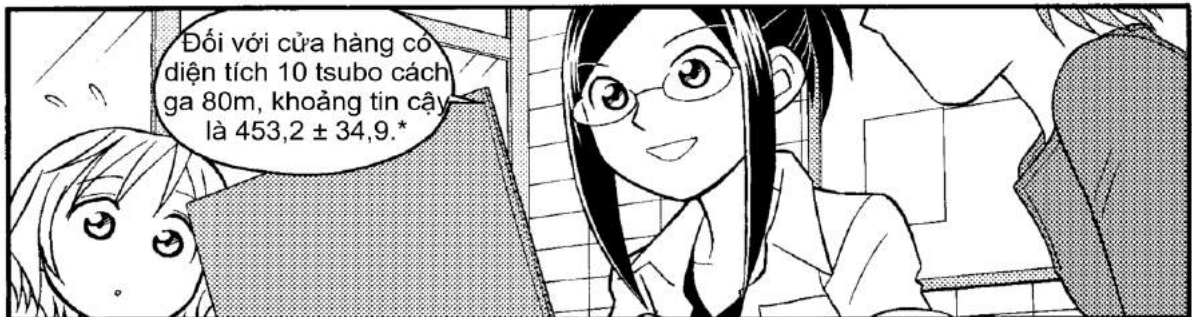
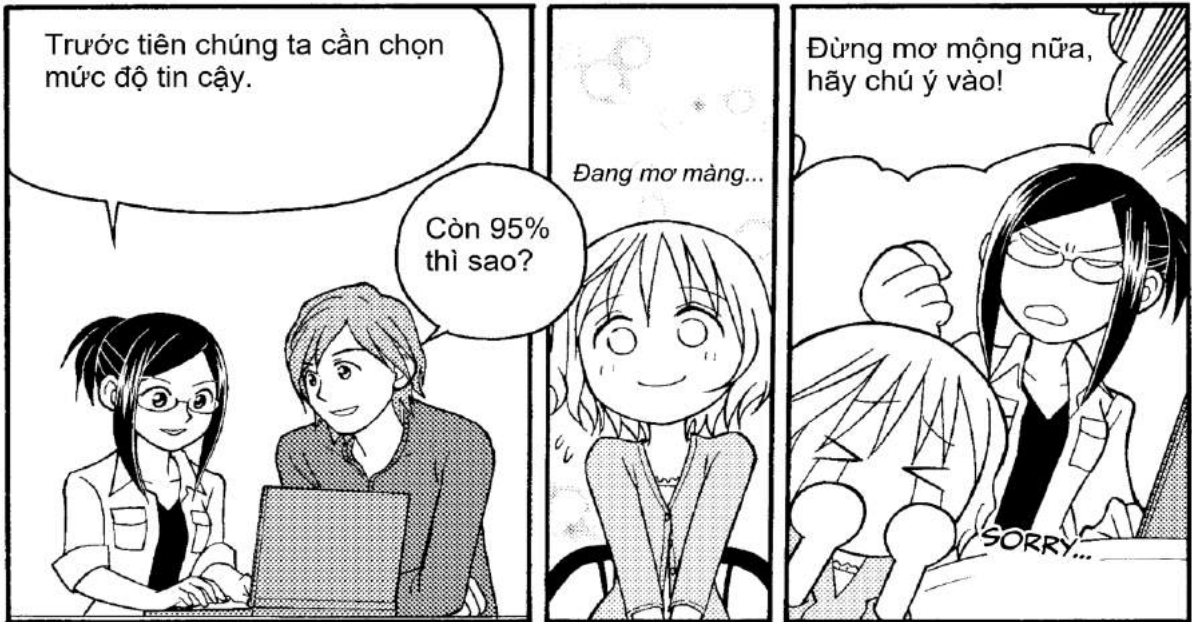
Bước 5: Tính khoảng tin cậy cho tổng thể.



* Nhà toán học P.C. Mahalanobis đã phát minh ra cách sử dụng khoảng cách đa biến để so sánh các tổng thể.

➤➤➤ HỌC THÔNG KÊ QUA TRUYỆN TRANH





* Cách tính này được giải thích chi tiết hơn ở phần sau.



➤➤➤ HỌC THỐNG KÊ QUA TRUYỆN TRANH

Bước 6: Thực hiện dự đoán

Đây là dữ liệu đối với một cửa hàng mới chúng ta dự định mở

	Diện tích cửa hàng (tsubo)	Khoảng cách đến ga gần nhất (m)
Isebashi Shop	10	110

Một cửa hàng ở Isebashi ư? Chỗ đó gần nhà tớ.

Cậu có thể dự đoán doanh thu không Mui?



Chính xác!

$$\begin{aligned}
 y &= 41.5x_1 - 0.3x_2 + 65.3 \\
 &= 41.5 \times 10 - 0.3 \times 110 + 65.3 \\
 &= \underline{447.3^*}
 \end{aligned}$$

¥4.473.000
một tháng

* Tính toán này được thực hiện bằng cách làm tròn số. Nếu không làm tròn thì kết quả sẽ là 442.96.

Cậu đúng là thiên tài, Miu! Tôi nên đặt tên cửa hàng theo tên cậu



Có lẽ cậu nên đặt tên nó cho Risa..

Nhưng làm sao chúng ta có thể biết chính xác doanh thu của một cửa hàng chưa được xây dựng? Chúng ta có nên tính khoảng dự đoán không?

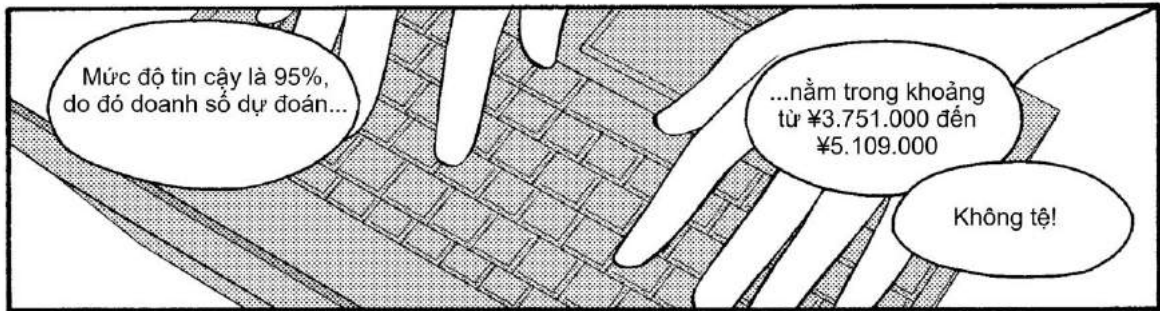


Đương nhiên!

Trong phân tích hồi quy đơn biến, phương pháp tìm cả khoảng tin cậy và khoảng dự đoán là tương tự nhau. Điều đó cũng đúng với phân tích hồi quy đa biến ư?



Đúng, nó cũng tương tự



Chọn sự kết hợp tốt nhất của các biến dự đoán



Phương trình hồi quy đa biến tốt nhất là cân bằng độ chính xác và độ phức tạp bằng cách chỉ bao gồm các biến dự đoán cần thiết để đưa ra dự đoán tốt nhất.

Khó $y = a_1x_1 + a_2x_2 + a_3x_3 + a_4x_4 + a_5x_5 + a_6x_6 + \dots + b$

Dễ $y = a_1x_1 + a_2x_2 + b$

Chính xác $y = a_1x_1 + a_2x_2 + b$ R^2 > R^2 Không chính xác $y = a_1x_1 + a_3x_3 + b$

Ngắn thì ngọt ngào.

Có một số cách để tìm ra phương trình mang lại cho cậu nhiều lợi ích nhất.

- Lựa chọn chuyển tiếp
- Loại bỏ ngược
- Lựa chọn từng bước tiến-lùi
- Hỏi chuyên gia về miền xem biến nào là quan trọng nhất.

Đây là một số cách phổ biến.

Phương pháp chúng ta sẽ sử dụng hôm nay đơn giản hơn bất kỳ phương pháp nào trong số đó. Nó được gọi là hồi quy tập hợp con tốt nhất hoặc đôi khi là phương pháp quay vòng (Round - Robin).

Quay vòng ư? Một con chim béo?

Cái quái gì vậy?

$x_1 \quad x_2 \quad x_3$

Tớ sẽ cho các cậu thấy. Giả sử x_1 , x_2 và x_3 là các biến dự báo tiềm năng.

Đầu tiên, chúng ta tính toán phương trình hồi quy đa biến cho mọi sự kết hợp của các biến dự đoán!

- x_1
- x_2
- x_3
- x_1 và x_2
- x_2 và x_3
- x_1 và x_3
- x_1 và x_2 và x_3

Haha. Điều này HAHA. chắc chắn THIS SURE là tròn trịa ROUND-ABOUT.

Biên dịch: Anh Tuấn (còn tiếp)