

Bài giảng 5: Regression và Prediction - II

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

3 Mở rộng mô hình

1 Chuẩn đoán mô hình

2 Lựa chọn mô hình

3 Mở rộng mô hình

Chuẩn đoán mô hình

Chuẩn đoán mô hình là một trong những bước quan trọng của xây dựng và sử dụng mô hình hồi quy tuyến tính.

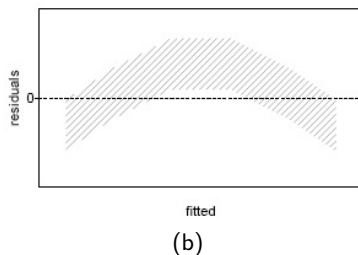
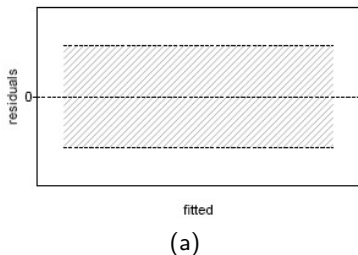
Chuẩn đoán mô hình bao gồm:

- Kiểm tra sự quan hệ tuyến tính giữa Y và các biến hồi quy X_1, X_2, \dots, X_p (giả định A1).
- Mỗi quan hệ tuyến tính giữa Y và một biến hồi quy X_j .
- Tính đồng nhất phương sai của thặng dư (giả định A4) - Homoskedasticity.
- Tính độc lập của thặng dư.
- Các giá trị ngoại lai xuất hiện trong (giả định A6).

↔ các chuẩn đoán mô hình này giúp ta biết được mô hình được xây dựng đã thực sự hợp lý với dữ liệu, và mô hình có vững.

Kiểm tra tính tuyến tính

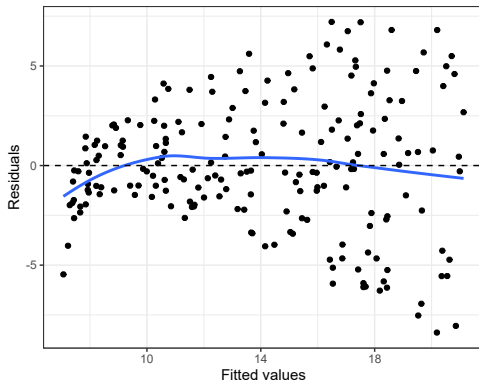
- Để kiểm tra sự quan hệ tuyến tính giữa Y và các biến hồi quy X_1, X_2, \dots, X_p (giả định A1), ta sử dụng biểu đồ **Residuals vs Fitted**.
- Nếu phần dư được phân bố ngẫu nhiên (khá đối xứng) xung quanh đường nằm ngang tương ứng với phần thặng dư $= 0$.
 ↪ giả định được thỏa mãn. (xem Hình a)
- Sự hiện diện của một xu hướng nào đó (một loại độ cong nào đó) có thể chỉ ra vấn đề với một số khía cạnh của mô hình tuyến tính.
 ↪ giả định không được thỏa mãn. (xem Hình b)



Kiểm tra tính tuyến tính

Xét biểu đồ **Residuals vs Fitted** của mô hình hồi quy tuyến tính:

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \varepsilon.$$



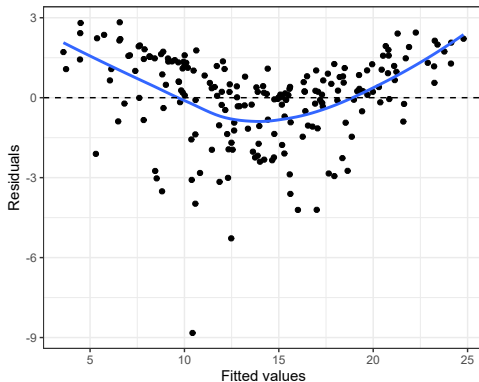
Hình vẽ không cho thấy xu hướng đường cong nào đáng kể.

↪ Giả định về tính tuyến tính của mô hình là phù hợp.

Kiểm tra tính tuyến tính

Xét biểu đồ **Residuals vs Fitted** của mô hình hồi quy tuyến tính:

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$



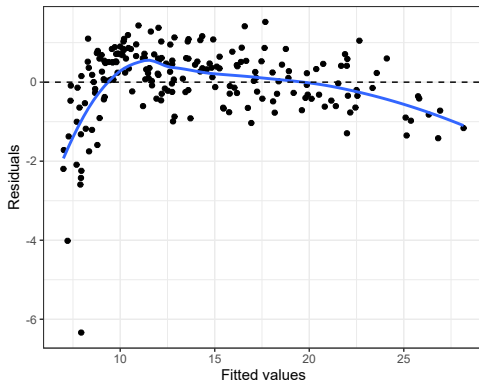
Hình vẽ cho thấy một xu hướng đường cong rõ ràng.

↪ Giả định về tính tuyến tính của mô hình là không phù hợp.

Kiểm tra tính tuyến tính

Xét biểu đồ **Residuals vs Fitted** của mô hình hồi quy tuyến tính:

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \beta_2 \text{radio} + \beta_3 \text{tv}:\text{radio} + \varepsilon.$$



Hình vẽ cho thấy một xu hướng đường cong tương đối.

↪ Giả định về tính tuyến tính của mô hình là không phù hợp.

Sự tuyến tính từng phần

Khi sự quan hệ tuyến tính giữa Y và các biến hồi quy X_1, X_2, \dots, X_p không hợp lý.

\hookrightarrow có lẽ có một vài biến không hỗ trợ quan hệ tuyến tính.

Để kiểm tra sự tuyến tính của Y đối với từng biến hồi quy X_j , ta sử dụng biểu đồ **thặng dư từng phần (partial residual plots)**:

$$\widehat{\varepsilon}_{ji} = \widehat{\varepsilon}_i + \widehat{\beta}_j X_{ji}.$$

Ý tưởng của biểu đồ thặng dư từng phần là để cô lập mối liên hệ giữa một biến hồi quy X_j và Y , nhưng tính đến mối liên hệ của tất cả các biến hồi quy còn lại.

Thông qua biểu đồ:

- nếu các điểm thặng dư từng phần $\widehat{\varepsilon}_{ji}$ có xu hướng theo một đường thẳng
 \hookrightarrow mối liên hệ tuyến tính giữa Y và X_j là hợp lý;
- nếu các điểm thặng dư từng phần $\widehat{\varepsilon}_{ji}$ có xu hướng theo một đường cong
 \hookrightarrow mối liên hệ tuyến tính giữa Y và X_j là không hợp lý.

Sự tuyến tính từng phần

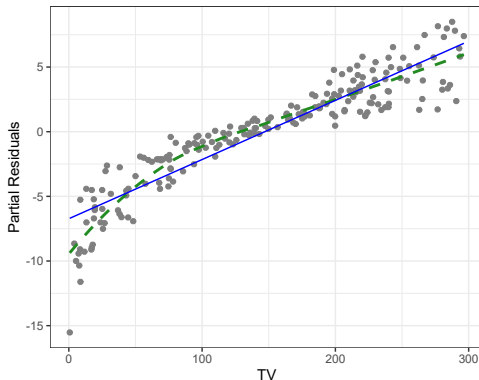
Xét mô hình hồi quy tuyến tính:

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$

Ta xét ba biểu đồ thặng dư từng phần cho:

- tv
- radio
- newspaper

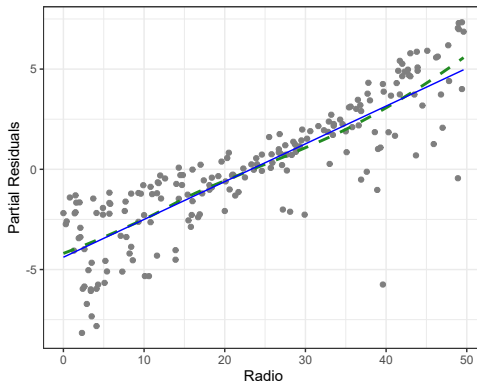
Sự tuyến tính từng phần



Kết quả cho thấy

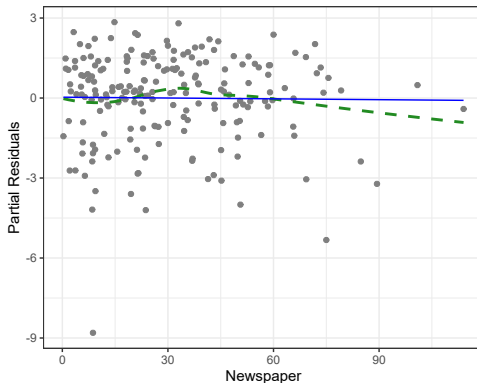
- đường thẳng tuyến tính (màu xanh dương) ước lượng không khớp với dữ liệu. Underestimate trong khoảng từ 80 tới 200;
- đường cong nét đứt (màu xanh lá) cho thấy xu hướng mối quan hệ phi tuyến tính của tv và sales.

Sự tuyến tính từng phần



Kết quả cho thấy đường thẳng tuyến tính (màu xanh dương) ước lượng tương đối khớp với dữ liệu.

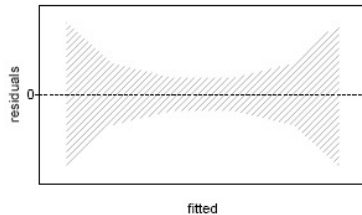
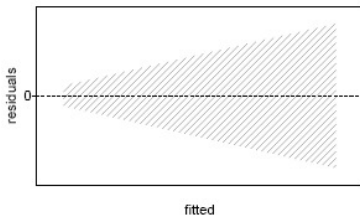
Sự tuyến tính từng phần



Kết quả cho thấy không có quan hệ tuyến tính giữa sales và newspaper trong mô hình.

Kiểm tra đồng nhất phương sai - Homoskedasticity

- Để kiểm tra giả định đồng nhất phương sai - Homoskedasticity, ta sử dụng biểu đồ **Residuals vs Fitted** hoặc **Scale-Location**.
- Đối với biểu đồ **Residuals vs Fitted**, nếu độ biến thiên của phần thặng dư thay đổi từ trái sang phải, như thế này:



thì điều kiện Homoskedasticity bị vi phạm, khi đó ta nói thặng dư của mô hình bị **heteroskedasticity**.

Kiểm tra đồng nhất phương sai - Homoskedasticity

Nhắc lại mô hình hồi quy tuyến tính

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon.$$

Ước lượng OLS $\hat{\beta}_n = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$. Do đó, giá trị fitted là

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_n = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} = \mathbf{H}\mathbf{Y}$$

Do đó, ta thăng dư của mô hình được viết thành

$$\hat{\varepsilon} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H}) \mathbf{Y}$$

Khi đó, ta tính được (dưới giả định đồng phương sai)

$$\text{Var}(\hat{\boldsymbol{\varepsilon}}) = \sigma^2 (\mathbf{I} - \mathbf{H})$$

Hay $\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - h_{ii})$, với h_{ii} là thành phần đường chéo của ma trận \mathbf{H} .

Ta định nghĩa standardized residuals là

$$r_i = \frac{\hat{\varepsilon}_i}{\sqrt{\text{Var}(\hat{\varepsilon}_i)}} = \frac{y_i - \hat{y}_i}{s\sqrt{1 - h_{ii}}},$$

s là ước lượng vững của σ .

Kiểm tra đồng nhất phương sai - Homoskedasticity

Một số tính chất thống kê của thặng dư chuẩn hóa (standardized residuals):

- r_i của độ lệch chuẩn bằng 1, nếu giả định về đồng nhất phương sai là đúng;
- các điểm r_i sẽ phân bố khá đồng đều theo chiều ngang khi so sánh với giá trị fitted.

Biểu đồ **Scale-Location** có trục

- x là giá trị fitted;
- trục y là $\sqrt{|\text{standardized residuals}|}$.

Nếu biểu đồ hiển thị đường nằm ngang, với các điểm trải đều bằng nhau thì đó là một dấu hiệu tốt về tính đồng nhất phương sai.

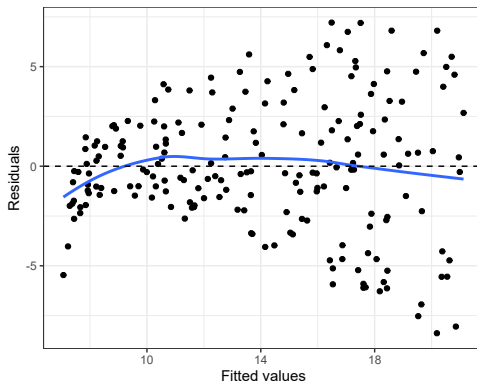
Chú ý:

Sự không đồng nhất phương sai trong mô hình (Heteroskedasticity) chỉ ra rằng các sai số dự đoán khác nhau đối với các phạm vi khác nhau của giá trị dự đoán và có thể gợi ý một mô hình không đầy đủ.

Kiểm tra đồng nhất phương sai - Homoskedasticity

Xét biểu đồ **Residuals vs Fitted** của mô hình

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \varepsilon.$$

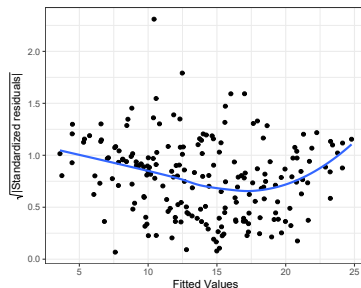
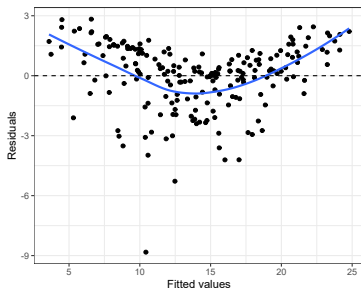


Các điểm thặng dư tăng theo giá trị ước lượng của biến đáp ứng, cho thấy phương sai không cố định trong các phần thặng dư (hoặc phương sai thay đổi - heteroscedasticity).

Kiểm tra đồng nhất phương sai - Homoskedasticity

Xét biểu đồ **Residuals vs Fitted** và **Scale-Location** của mô hình hồi quy tuyến tính:

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$



Nhìn vào hình thứ hai, có thể thấy xu hướng của dữ liệu không xấp xỉ đường thẳng 1, ám chỉ phương sai thặng dư là không đồng nhất.

Tính độc lập của thặng dư

Khi tính độc lập của thặng dư không được thỏa mãn, thì

- mô hình được xây dựng là không phù hợp với dữ liệu,
- ước lượng của hệ số không đảm bảo tính vững,
- các kết quả thống kê suy luận không còn chính xác.

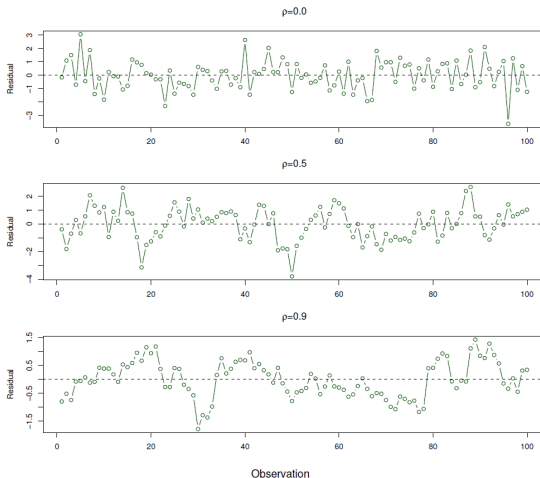
Điều kiện này thường sai khi dữ liệu được thu thập theo thời gian:

- thời gian máy bay khởi hành trễ theo các ngày trong năm;
- số lượng sinh viên theo năm học;
- doanh số bán hàng theo tháng.

Tính độc lập của thặng dư

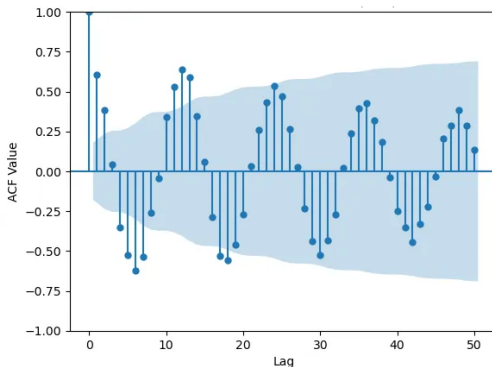
Sự khác biệt của thặng dư trong các tình huống

- độc lập ($\rho = 0$),
- có tương quan ($\rho = 0.5$ và $\rho = 0.8$).



Tính độc lập của thặng dư

Trong thực nghiệm, ta có thể dùng đồ thị Auto Correlation Function (ACF) của các thặng dư để kiểm tra sự tương quan giữa các thặng dư (có ý nghĩa hay không):

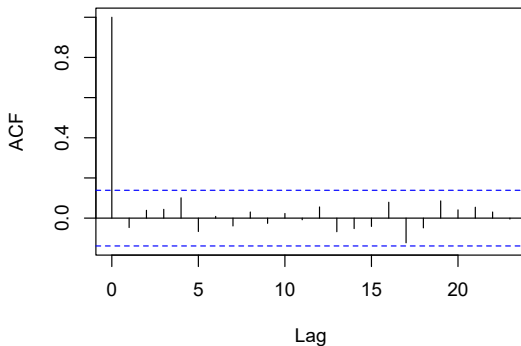


Đồ thị trên cho thấy sự tương quan giữa các thặng dư là có ý nghĩa.

Tính độc lập của thặng dư

Đồ thị Auto Correlation Function (ACF) cho thặng dư của mô hình

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$



Kết quả cho thấy sự độc lập của các thặng dư.

Ảnh hưởng của outliers lên ước lượng OLS

Dưới đây, ta mô phỏng sự ảnh hưởng của outliers tới ước lượng của đường hồi quy tuyến tính.

Phát hiện giá trị ngoại lai

Ta sử dụng biểu đồ **Residuals vs Leverage** để xác định các trường hợp có ảnh hưởng tới mô hình, đó là các giá trị ngoại lai có thể ảnh hưởng đến kết quả hồi quy khi được đưa vào hoặc loại trừ khỏi phân tích.

Leverage là giá trị h_{ii} trên đường chéo của ma trận **H**.

$$\forall i \text{ Var}(\hat{y}_i) = \sigma^2 h_{ii}, \text{ Var}(\hat{\varepsilon}_i) = \sigma^2(1 - h_{ii}), \text{ và } 0 \leq h_{ii} \leq 1, \text{ nên}$$

- giá trị h_{ii} lớn gần 1, thì $\text{Var}(\hat{y}_i) \approx \sigma^2 = \text{Var}(y_i)$ và $\text{Var}(\hat{\varepsilon}_i) \approx 0$
 $\hookrightarrow y_i$ có ảnh hưởng lớn tới \hat{y}_i .
- giá trị h_{ii} nhỏ gần 0, thì $\text{Var}(\hat{y}_i)$ là tương đối nhỏ
 $\hookrightarrow \hat{y}_i$ là được dựa trên đóng góp của nhiều quan sát, và y_i không có ảnh hưởng lớn tới \hat{y}_i .

Phát hiện giá trị ngoại lai

Một độ đo khác là khoảng cách Cook (Cook's distance), có thể dùng để xác định điểm ngoại lai trong mô hình tuyến tính:

$$D_i = r_i^2 \frac{h_{ii}}{p(1 - h_{ii})},$$

trong đó,

- p là số hệ số trong mô hình,
- r_i là thặng dư được chuẩn hóa,
- h_{ii} là giá trị leverage.

Những quan sát có giá trị khoảng cách Cook, D_i , lớn (thường là lớn hơn 0.5 - 1), tương ứng với trường hợp cả r_i và h_{ii} là lớn.

⇒ những quan sát đó là các điểm ngoại lai (có ảnh hưởng lớn) trong mô hình.

⇒ kết quả hồi quy sẽ bị thay đổi nếu chúng ta loại trừ những quan sát này.

Ta có thể tích hợp khoảng cách Cook vào trong biểu đồ **Residuals vs Leverage**, trong đó, khoảng cách Cook biểu thị kích cỡ của các điểm quan sát:

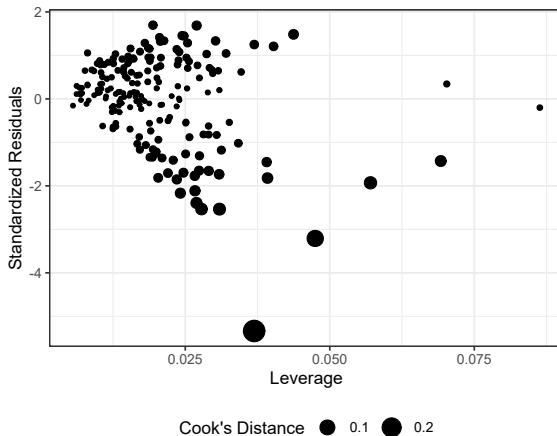
- điểm có kích cỡ nhỏ tương ứng khoảng cách Cook nhỏ;
- điểm có kích cỡ lớn tương ứng khoảng cách Cook lớn.

Phát hiện giá trị ngoại lai

Với mô hình

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$

biểu đồ Residuals vs Leverage:



Multicollinearity - Đa cộng tuyến

Multicollinearity - Đa cộng tuyến

Multicollinearity - Đa cộng tuyến là hiện tượng xảy ra trong mô hình hồi quy tuyến khi 1 biến hồi quy có tương quan (*correlated*) với 1 hoặc nhiều biến hồi quy khác có mặt trong mô hình.

Multicollinearity thường xuất hiện khi:

- 1 biến hồi quy được đưa vào nhiều lần, do lỗi trong xử lý;
- k giả biến (dummy variable), thay vì $k - 1$ giả biến, được tạo ra từ một biến nhân tố (biến định tính có k danh mục - categories);
- hai hoặc nhiều biến hồi quy có tương quan gần như hoàn hảo với nhau.

Multicollinearity - Đa cộng tuyến

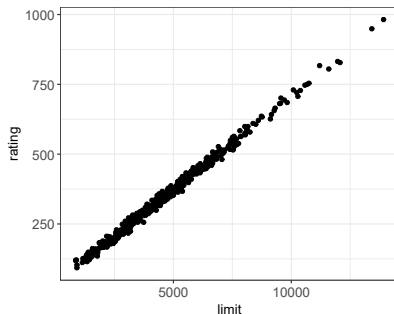
Ví dụ: Xét mô hình

$$\text{balance} = \beta_0 + \beta_1 \text{rating} + \beta_2 \text{limit} + \varepsilon,$$

trong đó,

- balance: số tiền cân bằng trong tài khoản;
- rating: xếp hạng tín dụng;
- limit: giới hạn tín dụng.

Biểu đồ phân tán của rating và limit cho thấy sự tương quan gần như hoàn hảo.



⇒ chắc chắn có hiện tượng multicollinearity trong mô hình.

Multicollinearity - Đa cộng tuyến

Tác động của multicollinearity

- làm giảm độ chính xác trong ước lượng hệ số của mô hình;
- làm sai số chuẩn của $\hat{\beta}_j$ tăng lên;
- làm giảm độ mạnh (power) của kiểm định giả thuyết $H_0 : \beta_j = 0$ (tức là, khả năng phát hiện đúng hệ số là khác không, bị giảm đi).

Ví dụ: Xét 2 mô hình

(MH1) : $\text{balance} = \beta_0 + \beta_1 \text{age} + \beta_2 \text{limit} + \varepsilon,$

(MH2) : $\text{balance} = \beta_0 + \beta_1 \text{rating} + \beta_2 \text{limit} + \varepsilon.$

Trong đó, (MH2) là mô hình có sự xuất hiện của multicollinearity (giữa limit và rating).

Multicollinearity - Đa cộng tuyến

Có kết quả như sau:

		Coefficient	Std. error	t-statistic	p-value
(MH1)	Intercept	-173.411	43.828	-3.957	< 0.0001
	age	-2.292	0.672	-3.407	0.0007
	limit	0.173	0.005	34.496	<0.0001
(MH1)	Intercept	-377.537	45.254	-8.343	<0.0001
	rating	2.202	0.952	2.312	0.0213
	limit	0.025	0.064	0.384	0.7012

⇒ sai số chuẩn của limit trong (MH2) tăng lên 12 lần so với trong (MH1).

Multicollinearity - Đa cộng tuyến

⇒ Cần phải giải quyết tình trạng đa công tuyến trong hồi quy.

⇒ Cần loại bỏ các biến cho đến khi tình trạng đa công tuyến biến mất.

Cách phát hiện đa công tuyến

- quan sát ma trận hệ số tương quan của các biến hồi quy trước khi thực hiện mô hình;
⇒ các giá trị lớn chỉ ra các cặp biến có sự tương quan cao ⇒ đa cộng tuyến có thể xuất hiện nếu sử dụng các cặp biến này;
- tính hệ số lạm phát phương sai (*variance inflation factor*) - VIF của từng biến hồi quy trong mô hình:

$$VIF_j = \frac{1}{1 - R_{X_j|X_{-j}}^2},$$

trong đó, $R^2_{X_j|X_{-j}}$ là hệ số R^2 của mô hình hồi quy của biến X_j theo các biến hồi quy còn lại trong mô hình;

$\Rightarrow VIF_i$ nhỏ nhất là 1 \Rightarrow không có đa cộng tuyến;

$\Rightarrow VIF_i$ vượt quá 5 hoặc 10 \Rightarrow vấn đề lớn với đa cộng tuyến;

Multicollinearity - Đa cộng tuyến

Ví dụ: Xét mô hình

$$\text{balance} = \beta_0 + \beta_1 \text{rating} + \beta_2 \text{limit} + \beta_3 \text{age} + \varepsilon$$

Khi đó, hệ số VIF được xác định như sau:

Biến	Mô hình	$R^2_{X_j X_{-j}}$	VIF _j
rating	$\text{rating} = \beta_0 + \beta_1 \text{limit} + \beta_2 \text{age} + \varepsilon$	0.9938	160.6683
limit	$\text{limit} = \beta_0 + \beta_1 \text{rating} + \beta_2 \text{age} + \varepsilon$	0.9937	160.5929
age	$\text{age} = \beta_0 + \beta_1 \text{limit} + \beta_2 \text{rating} + \varepsilon$	0.0113	1.0114

⇒ hai biến rating và limit là nguồn của đa công tuyến.

⇒ cần loại bỏ 1 trong hai biến này khỏi mô hình.

Multicollinearity - Đa cộng tuyến

Loại bỏ rating \Rightarrow mô hình mới

$$\text{balance} = \beta_0 + \beta_1 \text{limit} + \beta_2 \text{age} + \varepsilon$$

Khi đó, hệ số VIF được xác định như sau:

Biến	Mô hình	$R^2_{X_j X_{-j}}$	VIF _j
limit	$\text{limit} = \beta_0 + \beta_1 \text{age} + \varepsilon$	0.0102	1.0103
age	$\text{age} = \beta_0 + \beta_1 \text{limit} + \varepsilon$	0.0102	1.0103

\Rightarrow không còn đa cộng tuyến.

Loại bỏ limit \Rightarrow mô hình mới

$$\text{balance} = \beta_0 + \beta_1 \text{rating} + \beta_2 \text{age} + \varepsilon$$

Khi đó, hệ số VIF được xác định như sau:

Biến	Mô hình	$R^2_{X_j X_{-j}}$	VIF _j
rating	$\text{rating} = \beta_0 + \beta_1 \text{age} + \varepsilon$	0.0106	1.0108
age	$\text{age} = \beta_0 + \beta_1 \text{rating} + \varepsilon$	0.0106	1.0108

\Rightarrow không còn đa cộng tuyến.

1 Chuẩn đoán mô hình

2 Lựa chọn mô hình

3 Mở rộng mô hình

Lựa chọn mô hình

Trong xây dựng mô hình hồi quy đa biến

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon,$$

về mặt kỹ thuật, ta có thể có

- p mô hình hồi quy tuyến tính đơn;
- C_p^2 mô hình hồi quy tuyến tính hai biến;
- C_p^3 mô hình hồi quy tuyến tính ba biến;
- C_p^q mô hình hồi quy tuyến tính q biến, ($q \leq p$).

Khi này, một câu hỏi quan trọng nảy sinh:

“Mô hình nào là tốt nhất trong số các mô hình có thể có?”

→ ta cần có một phương pháp để chọn ra mô hình tốt nhất.

Lựa chọn mô hình

Trong xây dựng mô hình hồi quy đa biến

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon,$$

về mặt kỹ thuật, ta có thể có

- p mô hình hồi quy tuyến tính đơn;
- C_p^2 mô hình hồi quy tuyến tính hai biến;
- C_p^3 mô hình hồi quy tuyến tính ba biến;
- C_p^q mô hình hồi quy tuyến tính q biến, ($q \leq p$).

Khi này, một câu hỏi quan trọng nảy sinh:

“Mô hình nào là tốt nhất trong số các mô hình có thể có?”

→ ta cần có một phương pháp để chọn ra mô hình tốt nhất.

Hai vấn đề kỹ thuật:

- thể nào là một mô hình tốt?
- phương pháp lựa chọn?

Mô hình tốt?

Trong phần “Đánh giá mô hình”, ta đã biết rằng, một mô hình có

- RMSE nhỏ, hoặc
- RSE nhỏ, hoặc
- R^2 hoặc R_a^2 lớn, hoặc
- cross-validation error nhỏ,

sẽ là một mô hình tốt.

Tuy nhiên, khi ta thêm biến vào mô hình hồi quy

- RMSE và RSE sẽ luôn giảm (tối thiểu 0);
- R^2 luôn tăng (tối đa 1).

↪ không thực tế nếu dùng các chỉ số này để so sánh các mô hình có số lượng biến khác nhau.

Mô hình tốt?

Để so sánh các mô hình với số lượng biến khác nhau, ta dùng các chỉ số sau.

- AIC (Akaike's Information Criteria)¹ cho mô hình hồi quy tuyến tính:

$$AIC = 2p + n \log \left(\frac{RSS}{n} \right).$$

AIC được xây dựng dựa trên lý thuyết likelihood và giả định phân phối chuẩn.

- BIC (Bayesian Information Criteria) cho mô hình hồi quy tuyến tính:

$$BIC = p \log(n) + n \log \left(\frac{RSS}{n} \right).$$

BIC được xây dựng dựa trên lý thuyết likelihood, thống kê Bayes và giả định phân phối chuẩn.

- Mallows's C_p cho mô hình hồi quy tuyến tính

$$C_p = \frac{1}{n} (RSS + 2p\hat{\sigma}^2).$$

- Adjusted R-squares R_a^2 .

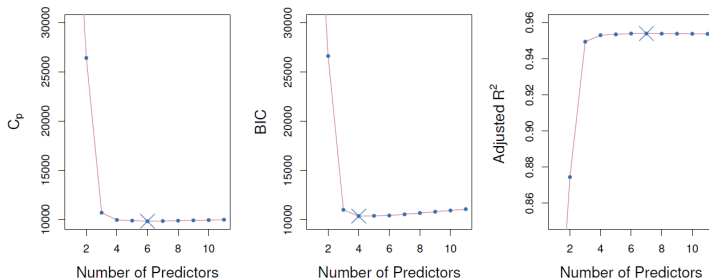
¹do nhà thống kê học Hirotugu Akaike, người Nhật Bản, đề xuất vào thập niên 70

Mô hình tốt?

Nhân xét:

- AIC, BIC và C_p sẽ có giá trị nhỏ tại các mô hình có sai số kiểm tra nhỏ.
→ ta cần tìm mô hình có giá trị AIC, BIC hoặc C_p nhỏ (tức là sai số kiểm tra nhỏ).
- Trong khuôn khổ mô hình hồi quy tuyến tính, AIC và BIC chỉ có nghĩa nếu giả định về phân phối chuẩn là hợp lý.
- Adjusted R-squares R_a^2 sẽ đạt giá trị lớn ở các mô hình có sai số kiểm tra nhỏ.
→ ta cần tìm mô hình có giá trị R_a^2 lớn.
- AIC sẽ tỷ lệ với C_p trong trường hợp giả định phân phối chuẩn là đúng.
- BIC sẽ cho kết quả mô hình đơn giản hơn (số biến nhỏ hơn) so với mô hình được tìm thấy bởi C_p .
- Adjusted R-squares R_a^2 thường ít được sử dụng trong thực tế hơn các chỉ số AIC, BIC và C_p .

Mô hình tốt?



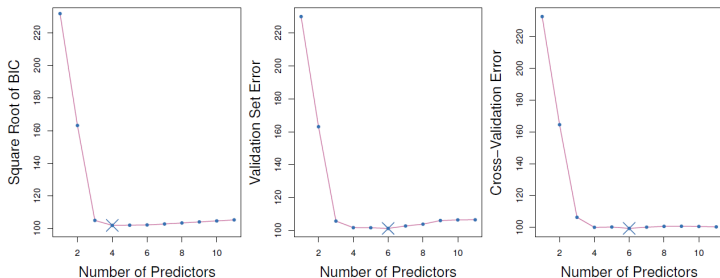
Hình trên, minh họa giá trị của C_p , BIC và R_a^2 trong việc lựa chọn mô hình hồi quy với số lượng biến hồi quy thay đổi (từ 1 tới 11).

- Tiêu chuẩn C_p cho kết quả là mô hình với 6 biến hồi quy.
- Tiêu chuẩn BIC cho kết quả là mô hình với 4 biến hồi quy.
- Tiêu chuẩn R_a^2 cho kết quả là mô hình với 7 biến hồi quy.

Quy tắc: Khi hiệu suất, sai số của các mô hình là tương đương nhau, thì mô hình đơn giản hơn là mô hình được ưu tiên.

Mô hình tốt?

Ngoài các tiêu chí trên, ta có thể sử dụng validation set hoặc cross-validation để tìm mô hình hồi quy với số biến q ($q < p$), sao cho sai số kiểm tra là nhỏ nhất.



Lựa chọn mô hình con tốt nhất

Xét mô hình hồi quy đa biến

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon,$$

ta cần áp dụng các chỉ số AIC, BIC, C_p , R_a^2 , cross-validation cho:

- p mô hình hồi quy tuyến tính đơn;
- C_p^2 mô hình hồi quy tuyến tính hai biến;
- C_p^3 mô hình hồi quy tuyến tính ba biến;
- C_p^q mô hình hồi quy tuyến tính q biến, ($q \leq p$).

Tổng quát, ta cần áp dụng các chỉ số đánh giá cho 2^p mô hình hồi quy có thể.

Việc này dẫn tới việc tính toán khá nặng.

Ví dụ: với $p = 200$ ta cần thực hiện $2^{20} = 1,048,576$ mô hình!

↪ ta sử dụng cách tiếp cận khác, có tên là hồi quy từng bước - stepwise regression.

Stepwise regression

Hồi quy tiến từng bước - Forward stepwise regression

1. Bắt đầu với mô hình \mathcal{M}_0 , chỉ chứa hệ số chặn: $Y = \beta_0 + \varepsilon$
2. Với $k = 0, 1, \dots, p - 1$:
 - (a) Xem xét tất cả $p - k$ mô hình làm tăng các yếu tố dự đoán trong \mathcal{M}_k bằng một yếu tố dự đoán bổ sung.
 - (b) Chọn mô hình tốt nhất trong số $p - k$ mô hình này và gọi nó là \mathcal{M}_{k+1} . Ở đây mô hình “tốt nhất” được xác định là mô hình có RSS nhỏ nhất hoặc R^2 cao nhất.
3. Chọn một mô hình tốt nhất trong số $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$ bằng cách sử dụng C_p , AIC, BIC, R_a^2 hoặc cross-validation.

Stepwise regression

Hồi quy lùi từng bước - Backward stepwise regression

1. Bắt đầu với mô hình \mathcal{M}_p , chứa tất cả các biến hồi quy:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

2. Với $k = p, p - 1, \dots, 1$:

(a) Xem xét tất cả k mô hình chứa tất cả ngoại trừ một trong các yếu tố dự đoán trong \mathcal{M}_k , để có tổng số $k - 1$ yếu tố dự đoán.

(b) Chọn mô hình tốt nhất trong số k mô hình này và gọi nó là \mathcal{M}_{k-1} . Ở đây mô hình “tốt nhất” được xác định là mô hình có RSS nhỏ nhất hoặc R^2 cao nhất.

3. Chọn một mô hình tốt nhất trong số $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$ bằng cách sử dụng C_p , AIC, BIC, R_a^2 hoặc cross-validation.

Stepwise regression

- Các phương pháp hồi quy tiến từng bước và lùi từng bước tốt nhất thường cho các mô hình tương tự nhưng không giống nhau.
- Để cho thuận tiện cho việc áp dụng và diễn giải, ta có thể sử dụng phương pháp hỗn hợp (kết hợp cả tiến và lùi).

Hồi quy từng bước hỗn hợp - Hybrid stepwise regression

1. Bắt đầu với mô hình \mathcal{M}_0 , chỉ chứa hệ số chặn: $Y = \beta_0 + \varepsilon$
2. Với $k = 0, 1, \dots, p - 1$:
 - (a) Xem xét tất cả $p - k$ mô hình làm tăng các yếu tố dự đoán trong \mathcal{M}_k bằng một yếu tố dự đoán bổ sung.
 - (b) Loại bỏ bất kỳ một biến dự đoán (hồi quy) trong mô hình \mathcal{M}_k .
 - (c) Chọn mô hình tốt nhất trong số p mô hình này và gọi nó là \mathcal{M}_{k+1} . Ở đây mô hình “tốt nhất” được xác định là mô hình có RSS nhỏ nhất hoặc R^2 cao nhất.
3. Chọn một mô hình tốt nhất trong số $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$ bằng cách sử dụng C_p , AIC, BIC, R_a^2 hoặc cross-validation.

Phương pháp co hệ số - Shrinkage methods

Xét mô hình hồi quy:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

Một cách tiếp cận khác cho việc lựa chọn biến hồi quy, là ước lượng hệ số của mô hình với một ràng buộc sao cho các hệ số không cần thiết sẽ “co lại” về 0.

- Ridge regression

$$\hat{\beta}_R = \arg \min \sum_{i=1}^n \left(Y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ji} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

- Lasso regression

$$\hat{\beta}_L = \arg \min \sum_{i=1}^n \left(Y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ji} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

Phương pháp co hệ số - Shrinkage methods

- SCAD (Smoothly Clipped Absolute Deviation)

$$\hat{\beta}_S = \arg \min \sum_{i=1}^n \left(Y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ji} \right)^2 + n \sum_{j=1}^p p_\lambda(|\beta_j|)$$

với

$$p_\lambda(|\theta|) = \lambda^2 - (|\theta| - \lambda)^2 \mathbf{I}(|\theta| < \lambda).$$

Tham số λ được gọi là turning parameter.

- $\lambda = 0$, ước lượng tương đương với ước lượng OLS.
- $\lambda \rightarrow \infty$, các hệ số ước lượng sẽ xấp xỉ 0.

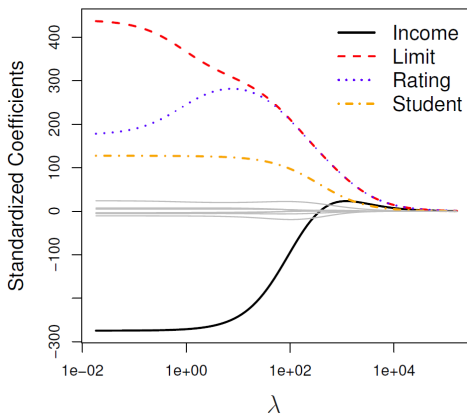
↪ một trong những vấn đề ứng dụng là xác định được λ hợp lý.

→ cross-validation được áp dụng để tìm λ . Ví dụ:

$$\text{CV}_{(k)}(\lambda) = \frac{1}{k} \sum_{j=1}^k \sum_{i=1}^n \left(Y_i - \hat{Y}_i(\lambda) \right)^2,$$

tìm λ sao cho cực tiểu hóa $CV_{(k)}(\lambda)$.

Phương pháp co hệ số - Shrinkage methods



Sự thay đổi của hệ số ước lượng tương ứng với giá trị λ của ridge regression.

1 Chuẩn đoán mô hình

2 Lựa chọn mô hình

3 Mở rộng mô hình

Weighted Regression

Hồi quy có trọng số - Weighted regression, là một biến thể của mô hình hồi quy tuyến tính:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{V}^{1/2}\boldsymbol{\varepsilon},$$

trong đó, $\mathbf{V} = \mathbf{W}\mathbf{W}^\top$, với $\mathbf{W} = (W_1, \dots, W_n)$ vectơ trọng số của n quan sát.

Ước lượng trọng số OLS (weighted OLS):

$$\hat{\boldsymbol{\beta}}_n = (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1} \mathbf{Y}.$$

Hồi quy có trọng số được dùng trong các trường hợp:

- Trọng số phương sai nghịch đảo khi các quan sát khác nhau được đo với độ chính xác khác nhau.
- Phân tích dữ liệu ở dạng tổng hợp sao cho biến trọng số mã hóa số lượng quan sát ban đầu mà mỗi hàng trong dữ liệu tổng hợp đại diện.

Đặc biệt, khi mô hình hồi quy thông thường bị heteroskedasticity (không đồng nhất phương sai) thì hồi quy có trọng số là một giải pháp thay thế.

Hồi quy đa thức

Trong các phần trước, ta đã giới thiệu mô hình hồi quy tuyến tính:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon.$$

Thuật ngữ “tuyến tính” ám chỉ sự tuyến tính theo hệ số β_j .

Do đó, X_j có thể được sử dụng dưới dạng một biến đổi: bình phương, log, căn bậc hai.

Đặc biệt, nếu ta sử dụng một đa thức bậc k của X_j , thì khi đó, ta có mô hình hồi quy đa thức (polynomial regression).

Ví dụ: với $k = 2$,

$$Y = \beta_0 + \alpha_1 X_1 + \alpha_2 X_1^2 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon.$$

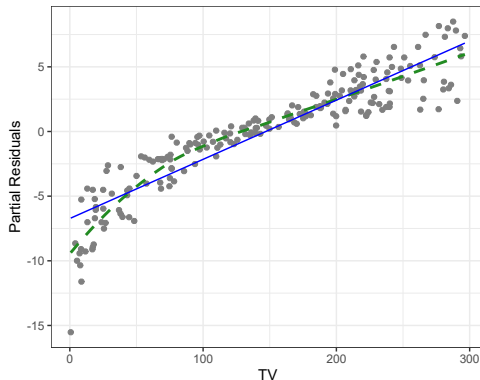
Hồi quy đa thức thường được dùng trong các trường hợp mà mối quan hệ giữa Y và X_j là phi tuyến.

Hồi quy đa thức

Ví dụ: Trong mô hình tuyến tính:

$$\text{sales} = \beta_0 + \beta_1 \text{tv} + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$

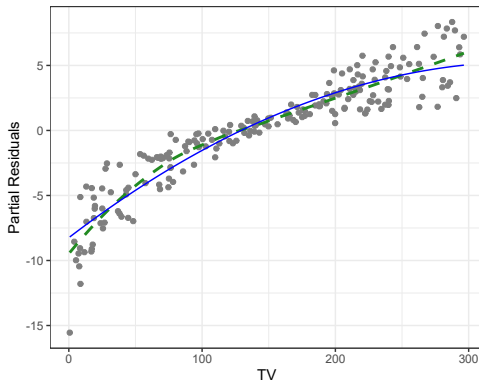
thông qua biểu đồ thặng dư từng phần cho biến tv , ước lượng tuyến tính là không phù hợp với dữ liệu



Hồi quy đa thức

Thay vì sử dụng tv , ta sẽ áp dụng dạng đa thức bậc k cho tv . Ta chọn $k = 2$

$$\text{sales} = \beta_0 + \alpha_1 tv + \alpha_2 tv^2 + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$

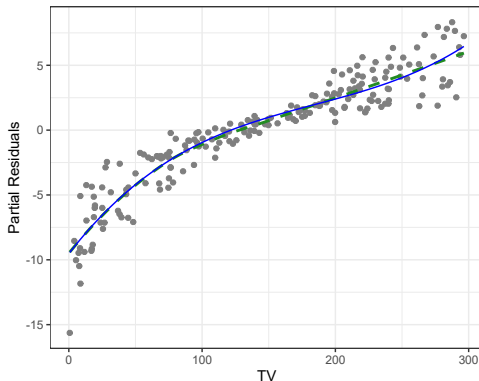


Biểu đồ thặng dư từng phần cho thấy mô hình đã được cải thiện đáng kể.

Hồi quy đa thức

Ta chọn $k = 3$

$$\text{sales} = \beta_0 + \alpha_1 \text{tv} + \alpha_2 \text{tv}^2 + \alpha_3 \text{tv}^3 + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$



Biểu đồ thặng dư từng phần cho thấy mô hình đã được cải thiện đáng kể hơn $k = 2$.

Hồi quy đa thức

Như vậy, để áp dụng hồi quy đa thức:

- Xây dựng mô hình hồi quy tuyến tính (không đa thức).
- Kiểm tra mối quan hệ giữa Y và X_j thông qua biểu đồ thặng dư từng phần.
- Nếu mối quan hệ giữa Y và X_j (cố định j) không phải là tuyến tính thì ta có sử dụng dạng đa thức bậc k của X_j để cải thiện mô hình.
- Để tìm k tối ưu, ta có thể áp dụng thuật toán cross-validation.
- Điểm hạn chế của hồi quy đa thức, đó là ta không thể làm một cách tự động.
- Mô hình có thể trở nên chồng kênh, với k lớn.

Hồi quy đa thức từng phần

Trong phần trước, ta đã xây dựng mô hình hồi quy đa thức trên toàn miền giá trị của X .

Tuy nhiên, trong thực tế, ta có thể xây dựng và ước lượng các dạng đa thức khác nhau trên từng đoạn giá trị của X :

$$Y = \begin{cases} \beta_{01} + \beta_{11}X + \beta_{21}X^2 + \dots + \beta_{k1}X^k + \varepsilon, & \text{nếu } X \leq c \\ \beta_{02} + \beta_{12}X + \beta_{22}X^2 + \dots + \beta_{k2}X^k + \varepsilon, & \text{nếu } X > c \end{cases}$$

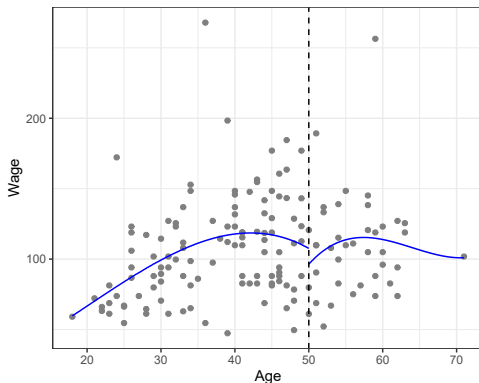
với c là một điểm nút (knot) được xác định trước.

Trên thực tế, ta có thể có nhiều hơn một điểm nút, tùy thuộc vào sự thay đổi của dữ liệu theo vùng giá trị của X .

Hồi quy đa thức từng phần

Ví dụ: Xét mô hình dự đoán tiền công của công nhân làm việc ở khu vực Mid-Atlantic, theo độ tuổi.

$$\text{wage} = \begin{cases} \beta_{01} + \beta_{11}\text{age} + \beta_{21}\text{age}^2 + \beta_{31}\text{age}^3 + \varepsilon, & \text{nếu } \text{age} \leq 50 \\ \beta_{02} + \beta_{12}\text{age} + \beta_{22}\text{age}^2 + \beta_{32}\text{age}^3 + \varepsilon, & \text{nếu } \text{age} > 50 \end{cases}$$



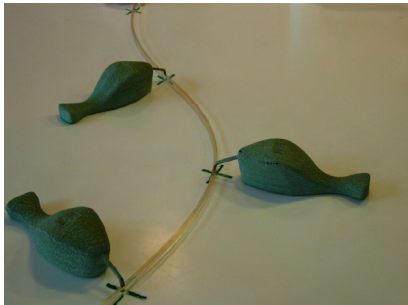
Hồi quy đa thức từng phần

Một số nhận xét:

- Có thể tương thích tốt hơn với sự thay đổi của Y trên các miền của X .
- Tuy nhiên, đường ước lượng là không liên tục tại điểm nút.
- Phải xác định bằng tay (manually) các nút và khoảng của X tương ứng với sự thay đổi đáng kể của dữ liệu.
- Phải xác định dạng cụ thể của mô hình trên mỗi vùng. Không nhất thiết là phải dạng đa thức.
- Khi có nhiều vùng, mô hình trở nên cồng kềnh.

Hồi quy splines

Splines là những dải linh hoạt (gỗ hoặc cao su) được sử dụng để vẽ các đường cong.



Hồi quy splines

Về mặt toán học, **Splines** là các hàm đa thức từng phần, được kết nối trơn tru.

Một cách chính thức hơn, một đường spline có bậc r với k điểm nút (knots) $\xi_1, \xi_2, \dots, \xi_k$ là hàm bất kỳ:

$$s(x) = \sum_{j=0}^{r+k} \alpha_j b_j(x),$$

với $b_j(x)$ là các hàm cơ sở của x .

- $r = 3$, ta có cubic spline
- nếu ta thêm điều kiện rằng $s(x)$ là hàm tuyến tính ngoài khoảng giá trị được quan sát của X , thì ta có natural spline.

Hồi quy splines

Ta có các cách chọn hàm cơ sở cơ bản:

- hàm cơ sở mũ cắt cụt (truncated power basis):

$$s(x) = \sum_{j=0}^r \alpha_j x^j + \sum_{j=1}^k \delta_j (x - \xi_j)_+^r,$$

với

$$(x - \xi_j)_+^r = \begin{cases} (x - \xi_j)^r & \text{nếu } x > \xi_j, \\ 0 & \text{nếu } x \leq \xi_j, \end{cases}$$

với $r = 3$ và sử dụng natural spline, ta có natural cubic spline;

Hồi quy splines

- hàm cơ sở B-spline:

$$s(x) = \sum_{j=0}^{r+k} \alpha_j B_{j,r}(x),$$

với $B_{j,r}(x)$ là hàm cơ sở xác định bởi:

(+) $r = 1$, các hàm cơ sở $B_{j,1}(x)$, $j = 0, \dots, r, r+1, \dots, k$:

$$B_{j,1}(x) = I(\xi_j \leq x < \xi_{j+1})$$

(+) $r > 1$,

$$B_{j,r}(x) = \frac{x - \xi_j}{\xi_{j+r} - \xi_j} B_{j,r-1}(x) + \frac{\xi_{j+r+1} - x}{\xi_{j+r+1} - \xi_{j+1}} B_{j+1,r-1}(x),$$

với $0/0 = 0$ và $\text{Inf}/0 = 0$, và các điểm nút cố định

$$a = \xi_0 < \xi_1 < \xi_2 < \dots < \xi_k < b = \xi_{k+1}$$

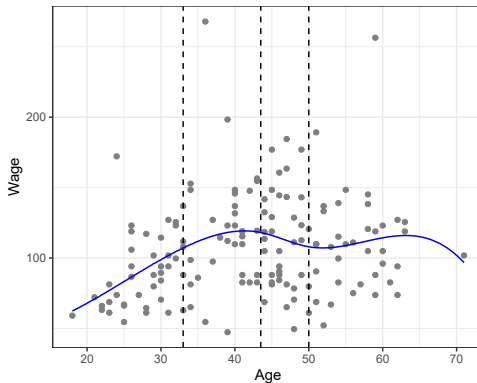
Phương pháp B-spline thường được dùng phổ biến hơn natural spline.

Một cách mặc định, $r = 3$ thường được sử dụng, bởi vì hàm ước lượng sẽ “đủ” trơn khi nhìn bằng mắt.

Hồi quy splines

Ta áp dụng B-spline cho mô hình dự đoán tiền công của công nhân dựa vào tuổi của họ:

- $r = 3$
- $k = 3$
- ξ_j là các điểm tương ứng với phân vị mức 0.25, 0.5 và 0.75 của age.

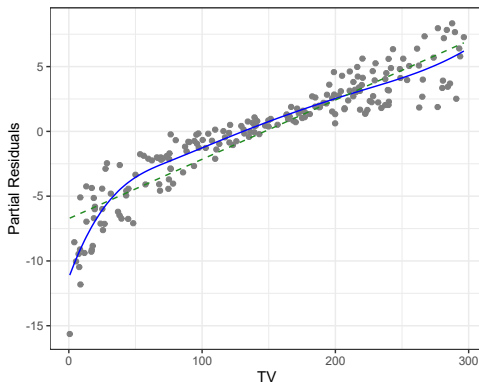


Hồi quy splines

Ví dụ: xét lại mô hình dự đoán doanh số bán hàng, với thành phần B-spline:

$$\text{sales} = \beta_0 + \text{bs}(\text{tv}, 3) + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$

với $\text{bs}(\text{tv}, 3)$ là thành phần được xác định bởi B-spline bậc 3.



Đường màu xanh dương tương ứng với ước lượng dựa trên B-spline bậc 3, đường nét đứt là đường ước lượng tuyến tính (mô hình ban đầu).

Hồi quy splines

So sánh với các mô hình tiên đoán sale đã ước lượng trước đây:

Mô hình	RSE	AIC	R_a^2
Tuyến tính	1.686	782.362	0.896
Đa thức bậc 2 cho tv	1.521	742.288	0.915
Đa thức bậc 3 cho tv	1.431	718.823	0.925
B-spline bậc 3 cho tv	1.399	712.626	0.928

- Mô hình với B-spline bậc 3 là mô hình tốt nhất.
- So với mô hình đa thức bậc 3, thì mô hình với B-spline bậc 3 là linh hoạt hơn, do ta không cần phải xác định số bậc của đa thức.

Trong áp dụng, ta cần quan tâm:

- vị trí đặt có điểm nút: có thể đặt theo điểm phân vị của dữ liệu;
- số lượng điểm nút: có thể sử dụng cross-validation để xác định;
- bậc của B-spline: có thể sử dụng cross-validation để xác định.

Chú ý: Ta có thể áp dụng B-spline cho nhiều hơn một biến hồi quy trong mô hình.

Generalized Additive Models

Một cách tự nhiên để mở rộng mô hình hồi quy tuyến tính

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon,$$

để cho phép biểu diễn linh hoạt quan hệ

- tuyến tính
- phi tuyến tính

giữa Y và X_j , là thay thế $\beta_j X_j$ bằng một hàm trơn $f_j(X_j)$ chưa biết của X_j :

$$Y = \beta_0 + f_1(X_1) + f_2(X_2) + \dots + f_p(X_p) + \varepsilon.$$

Mô hình này được gọi là Mô hình cộng tính tổng quát (Generalized Additive Models), gọi tắt là GAM.

Ta gọi là “cộng tính” vì ta ước lượng riêng biệt các f_j cho mỗi X_j , sau đó, cộng chúng lại với nhau, để ước đoán biến phản hồi.

Generalized Additive Models

Trong áp dụng, ta có thể:

- sử dụng GAM cho toàn bộ biến trong mô hình tuyến tính;
- sử dụng GAM cho một hoặc một vài biến trong mô hình tuyến tính.

Đối với cách sử dụng thứ 2, sẽ hữu dụng khi ta phát hiện ra mối quan hệ giữa Y và một X_j là không tuyến tính (thông qua biểu đồ của thặng dư từng phần).

Cách sử dụng thứ 1, là cách ước lượng mô hình tự động, hữu ích cho trường hợp

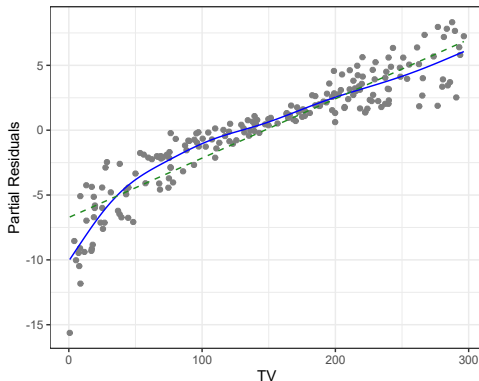
- có nhiều biến hồi quy,
- ta không thể kiểm tra quan hệ tuyến tính cho tất cả các biến.

Nếu mối quan hệ giữa Y và X_j thực sự là tuyến tính, f_j sẽ có dạng tuyến tính, và ngược lại.

Generalized Additive Models

Ví dụ: xét lại mô hình dự đoán doanh số bán hàng, với thành phần $f(\text{tv})$:

$$\text{sales} = \beta_0 + f(\text{tv}) + \beta_2 \text{radio} + \beta_3 \text{newspaper} + \varepsilon.$$

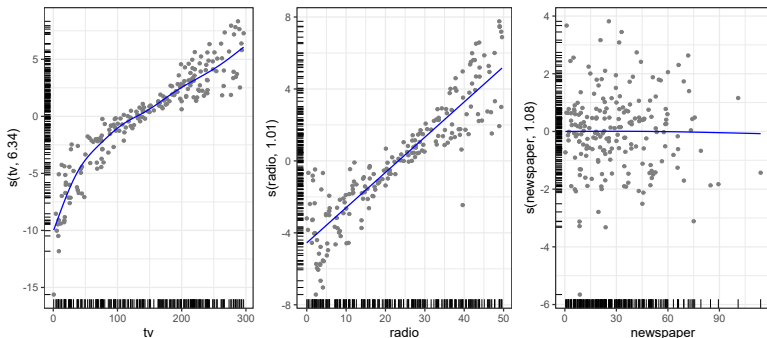


Đường màu xanh dương tương ứng với ước lượng dựa trên $f(\text{tv})$, đường nét đứt là đường ước lượng tuyến tính (mô hình ban đầu).

Generalized Additive Models

Ví dụ: xét mô hình GAM dự đoán doanh số bán hàng:

$$\text{sales} = \beta_0 + f_1(\text{tv}) + f_2(\text{radio}) + f_3(\text{newspaper}) + \varepsilon.$$



Chú ý: ta có thể thêm thành phần tương tác vào mô hình GAM, ví dụ

$$\text{sales} = \beta_0 + f_1(\text{tv}) + f_2(\text{radio}) + f_3(\text{tv}, \text{radio}) + \varepsilon.$$