

# Đa cộng tuyến

**Đa cộng tuyến (Multicollinearity)** là hiện tượng trong mô hình hồi quy (hoặc mô hình học máy tuyến tính) khi **hai hoặc nhiều biến độc lập (predictor variables)** có mối quan hệ tuyến tính chặt chẽ với nhau — nghĩa là **có thể dự đoán được một biến từ một hoặc nhiều biến khác**.

## Hiểu đơn giản:

“ Khi các biến giải thích (X) lại **tự giải thích lẫn nhau**, mô hình sẽ **khó xác định chính xác ảnh hưởng riêng của từng biến** đến biến phụ thuộc (Y).

## Hậu quả của đa cộng tuyến:

- Các hệ số hồi quy (coefficient) trở nên **không ổn định**, rất **nhạy cảm với dữ liệu**.
- Ý nghĩa thống kê (p-value)** của các biến có thể sai lệch → dễ **hiểu nhầm** rằng biến không quan trọng.
- Dễ dẫn đến mô hình **overfitting** hoặc **dự đoán kém** khi dùng dữ liệu mới.

## Ví dụ minh họa:

Giả sử bạn xây dựng một mô hình dự đoán **giá nhà** với các biến đầu vào sau:

Biến	Ý nghĩa
house_size	Diện tích ngôi nhà (m <sup>2</sup> )
num_bedrooms	Số lượng phòng ngủ
total_area	Tổng diện tích bao gồm sân vườn, gara

Trong thực tế:

- house\_size** và **total\_area** **gần như tuyến tính với nhau** (vì tổng diện tích = diện tích nhà + diện tích phụ).
- Điều này gây **đa cộng tuyến**.

→ Khi huấn luyện mô hình, thuật toán khó xác định chính xác:

- Liệu giá nhà tăng là do **house\_size** hay **total\_area**?

## Dấu hiệu nhận biết đa cộng tuyến:

- **Hệ số hồi quy có dấu hiệu "lạ"** (ví dụ: âm khi đáng ra phải dương).
- Giá trị **p-value cao bất thường** mặc dù biến đó có vẻ quan trọng.
- Dùng chỉ số thống kê như:
  - **VIF (Variance Inflation Factor)**: nếu  $VIF > 5$  hoặc  $10 \rightarrow$  có thể có đa cộng tuyến.

## Cách xử lý đa cộng tuyến:

1. **Loại bỏ một trong các biến có tương quan cao.**
2. **Tổng hợp các biến lại** (dùng PCA, hoặc tạo một biến trung gian).
3. **Chuẩn hóa dữ liệu** (scaling) có thể giúp phần nào.
4. **Sử dụng mô hình ít nhạy cảm với đa cộng tuyến**, như:
  - **Ridge Regression** (hồi quy co giãn)
  - **Tree-based models** (Random Forest, XGBoost...)

## Kết luận:

**Đa cộng tuyến là một "kẻ thù thầm lặng"** của các mô hình tuyến tính. Khi xây dựng mô hình hồi quy, bạn nên kiểm tra sự tương quan giữa các biến đầu vào để đảm bảo tính chính xác và ổn định cho mô hình.

Nếu bạn muốn, mình có thể giúp bạn tạo một mô hình hồi quy đơn giản bằng Python để **minh họa trực tiếp hiện tượng đa cộng tuyến**. Bạn có muốn thử không?

Tác giả: **Đỗ Ngọc Tú**  
Công Ty Phần Mềm **VHTSoft**

Phiên bản #1

Được tạo 15 tháng 5 2025 10:14:14 bởi Đỗ Ngọc Tú

Được cập nhật 15 tháng 5 2025 10:16:49 bởi Đỗ Ngọc Tú