

MÔ HÌNH HỒI QUY

Tiết 1: Mô hình hồi quy tuyến tính

Ths Phạm Việt Anh

Viện Công nghệ HaUI
Trường Đại học Công nghiệp Hà Nội

Ngày 21 tháng 3 năm 2024



Nội dung

- ① Ôn tập kiến thức buổi 1
- ② Khái niệm cơ bản
- ③ Xây dựng mô hình hồi quy
- ④ Triển khai mô hình
- ⑤ Tổng kết



Nội dung

- 1 Ôn tập kiến thức buổi 1
- 2 Khái niệm cơ bản
- 3 Xây dựng mô hình hồi quy
- 4 Triển khai mô hình
- 5 Tổng kết



Kiểm tra kiến thức bài cũ

- Các em hãy cho biết học máy là gì?



Kiểm tra kiến thức bài cũ

■ Các em hãy cho biết học máy là gì?

Trả lời: Là phân nhánh của Khai phá dữ liệu quan tâm đến việc phát triển các thuật toán học cho máy tính nhằm giúp máy tính có thể trích rút ra thông tin và tri thức từ dữ liệu thông qua các phương pháp “học”.



Kiểm tra kiến thức bài cũ

- Các em hãy cho biết học máy là gì?

Trả lời: Là phân nhánh của Khai phá dữ liệu quan tâm đến việc phát triển các thuật toán học cho máy tính nhằm giúp máy tính có thể trích rút ra thông tin và tri thức từ dữ liệu thông qua các phương pháp “học”.

- Các bài toán cơ bản trong học máy là gì?



Kiểm tra kiến thức bài cũ

■ Các em hãy cho biết học máy là gì?

Trả lời: Là phân nhánh của Khai phá dữ liệu quan tâm đến việc phát triển các thuật toán học cho máy tính nhằm giúp máy tính có thể trích rút ra thông tin và tri thức từ dữ liệu thông qua các phương pháp “học”.

■ Các bài toán cơ bản trong học máy là gì?

Trả lời: Hai bài toán cơ bản trong học máy là học giám sát và học không giám sát.



Kiểm tra kiến thức bài cũ

- Các em hãy cho biết học máy là gì?

Trả lời: Là phân nhánh của Khai phá dữ liệu quan tâm đến việc phát triển các thuật toán học cho máy tính nhằm giúp máy tính có thể trích rút ra thông tin và tri thức từ dữ liệu thông qua các phương pháp “học”.

- Các bài toán cơ bản trong học máy là gì?

Trả lời: Hai bài toán cơ bản trong học máy là học giám sát và học không giám sát.

- Các em hãy truy cập và đường link sau và trả lời các câu hỏi dưới dạng trắc nghiệm



Kiểm tra kiến thức bài cũ

- Các em hãy cho biết học máy là gì?

Trả lời: Là phân nhánh của Khai phá dữ liệu quan tâm đến việc phát triển các thuật toán học cho máy tính nhằm giúp máy tính có thể trích rút ra thông tin và tri thức từ dữ liệu thông qua các phương pháp “học”.

- Các bài toán cơ bản trong học máy là gì?

Trả lời: Hai bài toán cơ bản trong học máy là học giám sát và học không giám sát.

- Các em hãy truy cập và đường link sau và trả lời các câu hỏi dưới dạng trắc nghiệm

Link: <https://bit.ly/ktktb1>



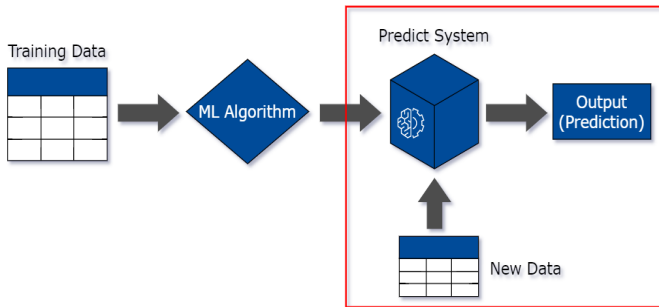
Nội dung

- ① Ôn tập kiến thức buổi 1
- ② Khái niệm cơ bản
- ③ Xây dựng mô hình hồi quy
- ④ Triển khai mô hình
- ⑤ Tổng kết



Ý nghĩa và mục đích

- Có rất nhiều ứng dụng phục vụ dự đoán
 - Ứng dụng phục vụ trong đời sống xã hội
 - Trong chiến lược của mỗi công ty, tổ chức



■ Dự đoán các thông tin trong tương lai

- Các thông tin chưa xảy ra
- Các thông tin xảy ra nhưng cần kiểm chứng
- Dự đoán giá nhà, giá vàng,...



■ Dự đoán các thông tin trong tương lai

- Các thông tin chưa xảy ra
- Các thông tin xảy ra nhưng cần kiểm chứng
- Dự đoán giá nhà, giá vàng,...

■ Dựa trên các dữ liệu đã có

- Các thông tin đã được xử lý
- Các thông tin được giả định và có tính chắc chắn
- Dựa trên số lượng lớn các thông tin
- Dựa trên các thông tin xảy ra có tính quy luật



■ Dự đoán các thông tin trong tương lai

- Các thông tin chưa xảy ra
- Các thông tin xảy ra nhưng cần kiểm chứng
- Dự đoán giá nhà, giá vàng,...

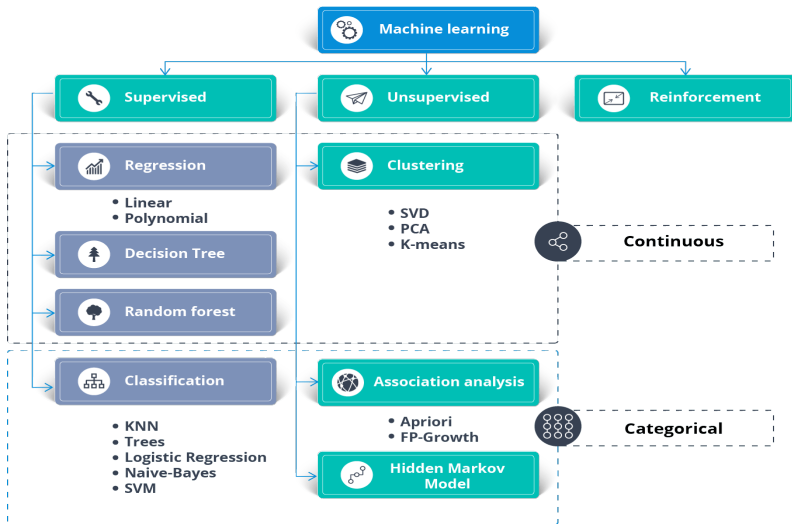
■ Dựa trên các dữ liệu đã có

- Các thông tin đã được xử lý
- Các thông tin được giả định và có tính chắc chắn
- Dựa trên số lượng lớn các thông tin
- Dựa trên các thông tin xảy ra có tính quy luật

Nhận xét: Các mô hình dự đoán có hiệu quả phụ thuộc rất nhiều vào các dữ liệu đã biết.



Lớp bài toán



■ Dữ liệu gồm các cặp đầu vào và nhãn tương ứng

- Các x_1, x_2, \dots, x_m được gọi là các đặc trưng/thuộc tính của dữ liệu đầu vào/bản ghi/đối tượng/quan sát \mathbf{x} .
- Trong thuật ngữ kinh tế lượng các x_1, x_2, \dots, x_m được gọi là biến độc lập và y gọi là biến phụ thuộc.
- y là dữ liệu đầu ra/biến mục tiêu.

	x_1	x_2	y
\mathbf{x}_1	6	87837	787
\mathbf{x}_2	7	78	5415
\mathbf{x}_3	545	778	7507
\mathbf{x}_4	545	18744	7560
\mathbf{x}_5	88	788	6344

- $Input = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- $Output : y = f(\mathbf{x}) \mid y_i \cong f(\mathbf{x}_i)$
- $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{im})$
- $\mathbf{x} = (x_1, x_2, \dots, x_m)$

Nhận xét: Việc xác định hàm $y = f(\mathbf{x})$ là việc xây dựng mô hình dự đoán.



Mô hình hồi quy

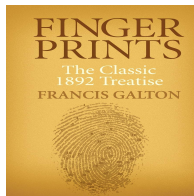
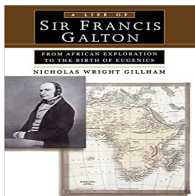
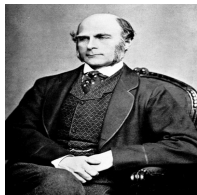
Một mô hình được gọi là hồi quy nếu mô hình đó mô tả được mối quan hệ giữa một hoặc nhiều biến độc lập với biến phụ thuộc.



Mô hình hồi quy

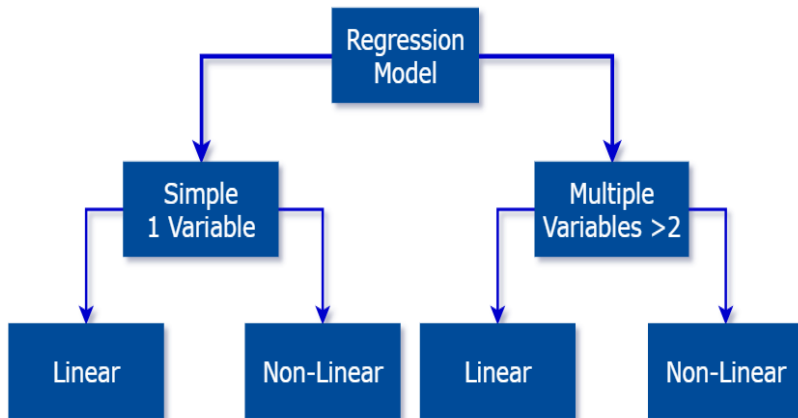
Một mô hình được gọi là hồi quy nếu mô hình đó mô tả được mối quan hệ giữa một hoặc nhiều biến độc lập với biến phụ thuộc.

- Được Francis Galton (1822-1911) đưa ra lần đầu tiên.
- Mô hình hồi quy có thể dự đoán hoặc ước lượng giá trị của một biến số từ các giá trị của một hay nhiều biến số khác.
- Phân tích hồi quy là cơ sở cho nhiều loại dự đoán và xác định sự tác động lên các biến mục tiêu



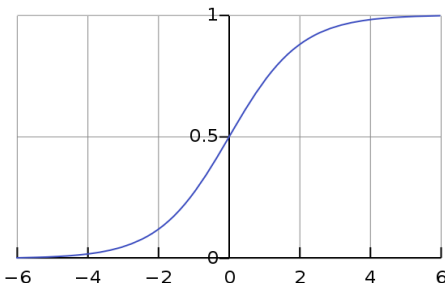
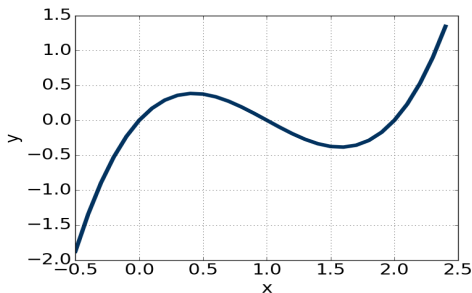
Các mô hình hồi quy

- Có rất nhiều cách để phân loại các mô hình hồi quy.



■ Có 6 mô hình hồi quy thường được sử dụng.

- Hồi quy tuyến tính
- Hồi quy Ridge
- Hồi quy Lasso
- Hồi quy Logistic
- Hồi quy đa thức
- Hồi quy bayesian



Nhận xét: Các mô hình hồi quy được xây dựng dựa trên sự phân bố dữ liệu đầu vào.



Mô hình hồi quy tuyến tính

Định nghĩa

Là mô hình mô tả được mối quan hệ giữa một hoặc nhiều biến độc lập với biến phụ thuộc **dựa trên một hàm tuyến tính**.

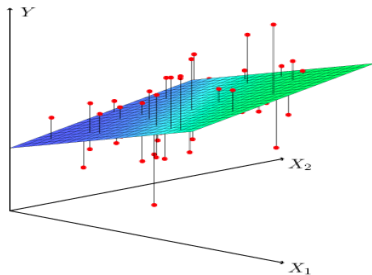
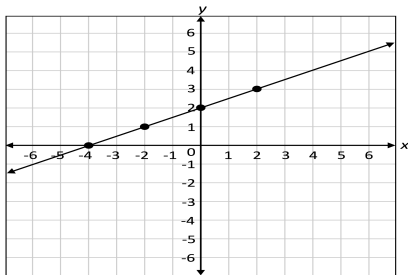


Mô hình hồi quy tuyến tính

Định nghĩa

Là mô hình mô tả được mối quan hệ giữa một hoặc nhiều biến độc lập với biến phụ thuộc **dựa trên một hàm tuyến tính**.

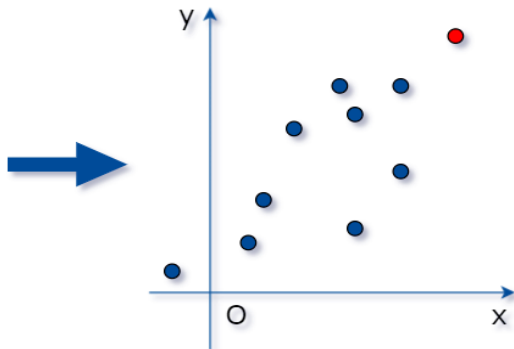
- Tuyến tính là các hàm số dạng bậc 1 (đường thẳng, mặt phẳng, siêu phẳng).



Mục đích

■ Xét một ví dụ sau:

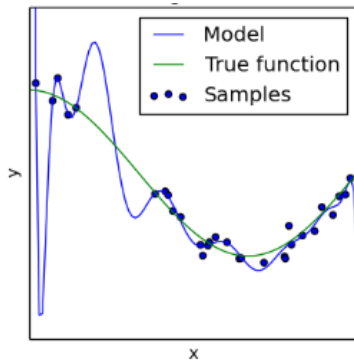
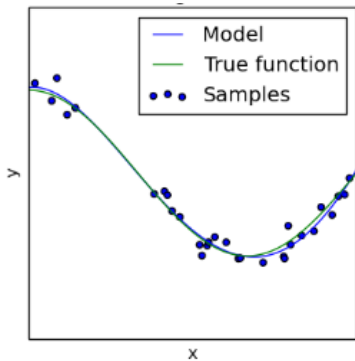
	x_1	y
x_1	0.13	-0.91
x_2	1.02	-0.17
x_3	3.17	1.61
x_4	-2.76	-3.31
x_5	1.44	0.18
x_6	5.28	3.36
x_7	-1.74	-2.46
x_8	7.93	5.56
...



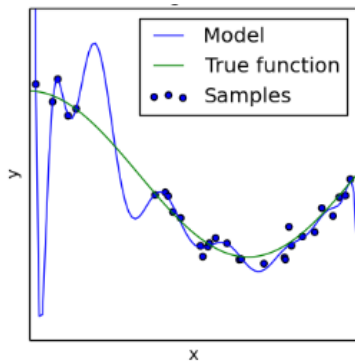
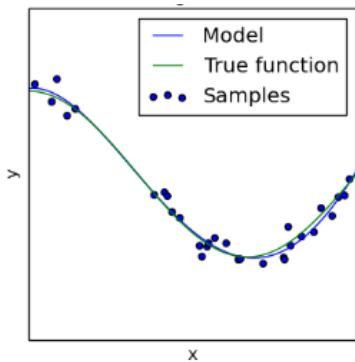
- Xác định một hàm $f(x)$ sao cho đi qua hết các điểm dữ liệu.
- Hàm $f(x)$ sử dụng để dự đoán cho các quan sát tiếp theo.



■ Các hàm có tính phức tạp (đa thức)



■ Các hàm có tính phức tạp (đa thức)



Nhận xét:

- Khó khăn trong việc xác định.
- Ảnh hưởng tốc độ tính toán với các dữ liệu tương lai.
- Không đảm bảo hiệu quả khi dự đoán dữ liệu tương lai.

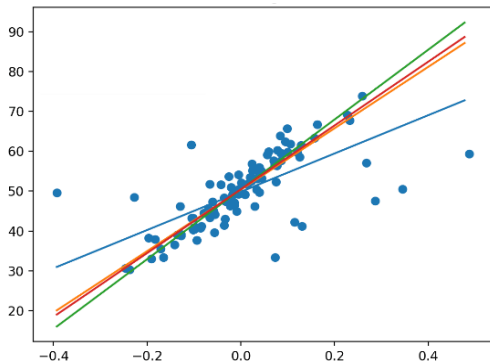


- Xác định một hàm tuyến tính đảm bảo sai số dự đoán là ít nhất, dễ dàng tính toán với các dữ liệu tương lai và dễ dàng xây dựng.

- Mô hình HQTТ tổng quát:

$$f(\mathbf{x}) = b + w_1x_1 + \dots + w_mx_m$$

- Xác định/ tìm kiếm $f(\mathbf{x})$?
- Tổng quát hóa là tốt nhất
- $f(\mathbf{x})$ tốt hơn so với $g(\mathbf{x})$?

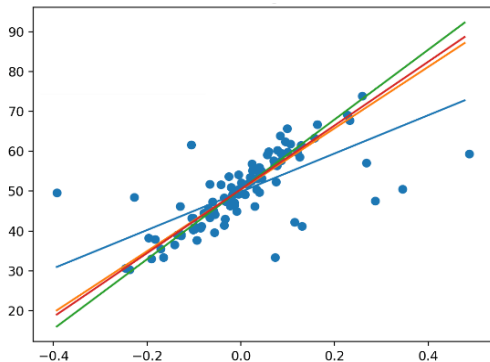


- Xác định một hàm tuyến tính đảm bảo sai số dự đoán là ít nhất, dễ dàng tính toán với các dữ liệu tương lai và dễ dàng xây dựng.

■ Mô hình HQT tổng quát:

$$f(\mathbf{x}) = b + w_1x_1 + \dots + w_mx_m$$

- Xác định/ tìm kiếm $f(\mathbf{x})$?
- Tổng quát hóa là tốt nhất
- $f(\mathbf{x})$ tốt hơn so với $g(\mathbf{x})$?



Nhận xét: Hệ số b được đưa vào mô hình để mang lại tính tổng quát hóa. Việc xác định/ tìm $f(\mathbf{x})$ là quá trình **xây dựng mô hình hồi quy tuyến tính**.



Nội dung

- ① Ôn tập kiến thức buổi 1
- ② Khái niệm cơ bản
- ③ Xây dựng mô hình hồi quy**
- ④ Triển khai mô hình
- ⑤ Tổng kết

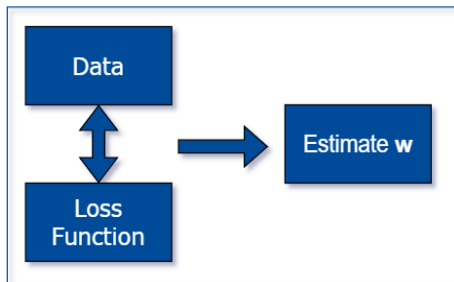


Xây dựng mô hình

Là việc xác định các tham số của mô hình dựa trên các dữ liệu đầu vào thông qua một độ đo tiêu chuẩn nào đó.

■ Độ đo:

- Dựa trên các dữ liệu đầu vào
- Đánh giá mức độ chênh lệch giữa giá trị đầu ra thực tế và giá trị đầu ra được dự đoán.

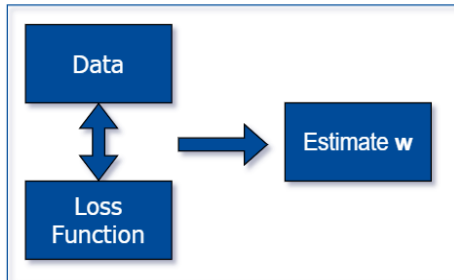


Xây dựng mô hình

Là việc xác định các tham số của mô hình dựa trên các dữ liệu đầu vào thông qua một độ đo tiêu chuẩn nào đó.

■ Độ đo:

- Dựa trên các dữ liệu đầu vào
- Đánh giá mức độ chênh lệch giữa giá trị đầu ra thực tế và giá trị đầu ra được dự đoán.



Nhận xét: Hàm mất mát/lỗi thường được sử dụng. Quá trình tối ưu hàm mất mát được gọi là **quá trình học** của mô hình học máy.



Hàm mất mát

- Xét mỗi điểm dữ liệu $\mathbf{x} = (x_1, x_2, \dots, x_m)$, ta cần xây dựng hàm $f(\mathbf{x}) = b + w_1x_1 + \dots + w_mx_m$ sao cho $y \cong f(\mathbf{x})$.



Hàm mất mát

- Xét mỗi điểm dữ liệu $\mathbf{x} = (x_1, x_2, \dots, x_m)$, ta cần xây dựng hàm $f(\mathbf{x}) = b + w_1x_1 + \dots + w_mx_m$ sao cho $y \cong f(\mathbf{x})$.
- Mỗi điểm dữ liệu bổ sung thêm thuộc tính $x_0 = 1$ và đặt $w_0 = bx_0$.



Hàm mất mát

- Xét mỗi điểm dữ liệu $\mathbf{x} = (x_1, x_2, \dots, x_m)$, ta cần xây dựng hàm $f(\mathbf{x}) = b + w_1x_1 + \dots + w_mx_m$ sao cho $y \cong f(\mathbf{x})$.
- Mỗi điểm dữ liệu bổ sung thêm thuộc tính $x_0 = 1$ và đặt $w_0 = bx_0$.
- Với mỗi điểm dữ liệu $\tilde{\mathbf{x}} = (x_0, x_1, x_2, \dots, x_m)$, ta cần xác định một giá trị dự đoán \hat{y} sao cho phương trình sau đây đạt giá trị nhỏ nhất:

$$\frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - w_0x_0 - w_1x_1 - \dots - w_mx_m)^2 \quad (1)$$



Hàm mất mát

- Xét mỗi điểm dữ liệu $\mathbf{x} = (x_1, x_2, \dots, x_m)$, ta cần xây dựng hàm $f(\mathbf{x}) = b + w_1x_1 + \dots + w_mx_m$ sao cho $y \cong f(\mathbf{x})$.
- Mỗi điểm dữ liệu bổ sung thêm thuộc tính $x_0 = 1$ và đặt $w_0 = bx_0$.
- Với mỗi điểm dữ liệu $\bar{\mathbf{x}} = (x_0, x_1, x_2, \dots, x_m)$, ta cần xác định một giá trị dự đoán \hat{y} sao cho phương trình sau đây đạt giá trị nhỏ nhất:

$$\frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - w_0x_0 - w_1x_1 - \dots - w_mx_m)^2 \quad (1)$$

- Nếu đặt $\mathbf{w} = (w_0, \dots, w_m)^T$ là một vector hệ số cần tối ưu (dạng cột), khi đó phương trình (1) trở thành:

$$\frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(y - \bar{\mathbf{x}}\mathbf{w})^2 \quad (2)$$



- Nếu xét trên các cặp dữ liệu đã biết (\mathbf{x}_i, y_i) , $i = 1, 2, \dots, n$, với n là số lượng dữ liệu quan sát được. Điều mong muốn là tổng sai số phải là nhỏ nhất. Điều này dẫn tới bài toán tìm \mathbf{w} sao cho phương trình sau là nhỏ nhất:

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n \frac{1}{2} (y_i - \bar{\mathbf{x}}_i \mathbf{w})^2 \quad (3)$$



- Nếu xét trên các cặp dữ liệu đã biết (\mathbf{x}_i, y_i) , $i = 1, 2, \dots, n$, với n là số lượng dữ liệu quan sát được. Điều mong muốn là tổng sai số phải là nhỏ nhất. Điều này dẫn tới bài toán tìm \mathbf{w} sao cho phương trình sau là nhỏ nhất:

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n \frac{1}{2} (y_i - \bar{\mathbf{x}}_i \mathbf{w})^2 \quad (3)$$

- $\mathcal{L}(\mathbf{w})$ được gọi là hàm mất mát của bài toán hồi quy tuyến tính và được sử dụng để xác định các tham số cho mô hình học máy.



- Nếu xét trên các cặp dữ liệu đã biết (\mathbf{x}_i, y_i) , $i = 1, 2, \dots, n$, với n là số lượng dữ liệu quan sát được. Điều mong muốn là tổng sai số phải là nhỏ nhất. Điều này dẫn tới bài toán tìm \mathbf{w} sao cho phương trình sau là nhỏ nhất:

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n \frac{1}{2} (y_i - \bar{\mathbf{x}}_i \mathbf{w})^2 \quad (3)$$

- $\mathcal{L}(\mathbf{w})$ được gọi là hàm mất mát của bài toán hồi quy tuyến tính và được sử dụng để xác định các tham số cho mô hình học máy.
- Đặt \mathbf{w}^* là giá trị của \mathbf{w} sao cho $\mathcal{L}(\mathbf{w})$ đạt giá trị nhỏ nhất. Khi đó:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$



- Nếu xét trên các cặp dữ liệu đã biết (\mathbf{x}_i, y_i) , $i = 1, 2, \dots, n$, với n là số lượng dữ liệu quan sát được. Điều mong muốn là tổng sai số phải là nhỏ nhất. Điều này dẫn tới bài toán tìm \mathbf{w} sao cho phương trình sau là nhỏ nhất:

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n \frac{1}{2} (y_i - \bar{\mathbf{x}}_i \mathbf{w})^2 \quad (3)$$

- $\mathcal{L}(\mathbf{w})$ được gọi là hàm mất mát của bài toán hồi quy tuyến tính và được sử dụng để xác định các tham số cho mô hình học máy.
- Đặt \mathbf{w}^* là giá trị của \mathbf{w} sao cho $\mathcal{L}(\mathbf{w})$ đạt giá trị nhỏ nhất. Khi đó:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}(\mathbf{w})$$

- \mathbf{w}^* được gọi là **điểm tối ưu** của mô hình.



Xác định tham số

■ Xác định vector hệ số \mathbf{w} .

- Có rất nhiều phương pháp để xác định \mathbf{w}^* . Phương pháp tổng bình phương nhỏ nhất (OLS-Ordinary Least Squares) thường được sử dụng.



Xác định tham số

■ Xác định vector hệ số \mathbf{w} .

- Có rất nhiều phương pháp để xác định \mathbf{w}^* . Phương pháp tổng bình phương nhỏ nhất (OLS-Ordinary Least Squares) thường được sử dụng.
- Đặt $\bar{\mathbf{X}} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n)$, $\mathbf{y} = (y_1, y_2, \dots, y_n)$ là ma trận chứa các vector hàng của dữ liệu đầu vào và vector cột chứa các dữ liệu đầu ra.



Xác định tham số

■ Xác định vector hệ số \mathbf{w} .

- Có rất nhiều phương pháp để xác định \mathbf{w}^* . Phương pháp tổng bình phương nhỏ nhất (OLS-Ordinary Least Squares) thường được sử dụng.
- Đặt $\bar{\mathbf{X}} = (\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2, \dots, \bar{\mathbf{x}}_n)$, $\mathbf{y} = (y_1, y_2, \dots, y_n)$ là ma trận chứa các vector hàng của dữ liệu đầu vào và vector cột chứa các dữ liệu đầu ra.
- Phương trình (3) trở thành:

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \bar{\mathbf{X}}\mathbf{w}\|_2^2 \quad (4)$$



Xác định tham số

■ Xác định vector hệ số \mathbf{w} .

- Có rất nhiều phương pháp để xác định \mathbf{w}^* . Phương pháp tổng bình phương nhỏ nhất (OLS-Ordinary Least Squares) thường được sử dụng.
- Đặt $\bar{\mathbf{X}} = (\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2, \dots, \bar{\mathbf{x}}_n)$, $\mathbf{y} = (y_1, y_2, \dots, y_n)$ là ma trận chứa các vector hàng của dữ liệu đầu vào và vector cột chứa các dữ liệu đầu ra.
- Phương trình (3) trở thành:

$$\mathcal{L}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \bar{\mathbf{X}}\mathbf{w}\|_2^2 \quad (4)$$

- $\|\mathbf{z}\|_2^2$ là tổng bình phương mỗi phần tử của \mathbf{z} .



- Lấy đạo hàm theo \mathbf{w} .

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \bar{\mathbf{X}} (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) \quad (5)$$



■ Lấy đạo hàm theo \mathbf{w} .

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \bar{\mathbf{X}} (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) \quad (5)$$

- $\bar{\mathbf{X}} = n \times m, \mathbf{w} = m \times 1, \mathbf{y} = n \times 1 \Rightarrow (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) = n \times 1$



- Lấy đạo hàm theo \mathbf{w} .

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \bar{\mathbf{X}} (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) \quad (5)$$

- $\bar{\mathbf{X}} = n \times m$, $\mathbf{w} = m \times 1$, $\mathbf{y} = n \times 1 \Rightarrow (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) = n \times 1$
- Để thỏa mãn phương trình đạo hàm cần chuyển thành:

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \bar{\mathbf{X}}^T (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) = 0 \Rightarrow \boxed{\mathbf{w} = (\bar{\mathbf{X}}^T \bar{\mathbf{X}})^{-1} \bar{\mathbf{X}}^T \mathbf{y}} \quad (6)$$



- Lấy đạo hàm theo \mathbf{w} .

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \bar{\mathbf{X}} (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) \quad (5)$$

- $\bar{\mathbf{X}} = n \times m$, $\mathbf{w} = m \times 1$, $\mathbf{y} = n \times 1 \Rightarrow (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) = n \times 1$
- Để thỏa mãn phương trình đạo hàm cần chuyển thành:

$$\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \bar{\mathbf{X}}^T (\bar{\mathbf{X}} \mathbf{w} - \mathbf{y}) = 0 \Rightarrow \boxed{\mathbf{w} = (\bar{\mathbf{X}}^T \bar{\mathbf{X}})^{-1} \bar{\mathbf{X}}^T \mathbf{y}} \quad (6)$$

Nhận xét:

- $\bar{\mathbf{X}}^T \bar{\mathbf{X}}$ phải là một **ma trận khả nghịch**.
- Trường hợp $\bar{\mathbf{X}}^T \bar{\mathbf{X}}$ không khả nghịch thì phương pháp OLS không thể sử dụng.



Nội dung

- ① Ôn tập kiến thức buổi 1
- ② Khái niệm cơ bản
- ③ Xây dựng mô hình hồi quy
- ④ Triển khai mô hình
- ⑤ Tổng kết



Ví dụ minh họa

Với 13 quan sát cho biết diện tích và giá của 13 căn nhà. Xây dựng một mô hình hồi quy tuyến tính thể hiện mối quan hệ giữa diện tích và giá nhà.

O	S (m2)	p (Tỷ VNĐ)
x_1	73.5	1.49
x_2	75.0	1.50
x_3	76.5	1.51
x_4	79.0	1.54
x_5	81.5	1.58
x_6	82.5	1.59
x_7	84	1.60
x_8	85	1.62
x_9	86.5	1.63
x_{10}	87.5	1.64
x_{11}	89	1.66
x_{12}	90	1.67
x_{13}	91.5	1.68

Công cụ triển khai

Ngôn ngữ	Python	3.7
Môi trường	Anaconda	Latest
Thư viện	Numpy, matplotlib	Latest
Soạn thảo	Jupyter Notebook	Latest

Khuyến nghị: Sử dụng Google Colab

■ Truy cập vào link sau:

<https://colab.research.google.com/>



Triển khai mô hình

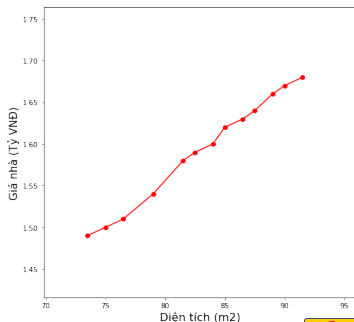
■ Thêm thư viện, chuẩn bị dữ liệu đầu vào và xây dựng hàm trực quan hóa dữ liệu

```
✓ [1] import matplotlib.pyplot as plt  
0s import numpy as np
```

```
[2] # Tiền xử lý dữ liệu  
# X (nxm) , y (nx1)  
X = np.array([[73.5,75.,76.5,79.,81.5,82.5,84.,85.,86.5,87.5,89.,90.,91.5]]).T  
y = np.array([[1.49,1.50,1.51,1.54,1.58,1.59,1.60,1.62,1.63,1.64,1.66,1.67,1.68]]).T
```

```
✓ [3] """  
0s Xây dựng hàm trực quan hóa dữ liệu  
"""  
def plotData(X, y, title="", xlabel="", ylabel=""):  
    # Tạo khung của đồ thị  
    plt.figure(figsize=(8,8))  
    plt.plot(X, y, 'r-o', label="price")  
  
    # Xác định các giá trị max và min của dữ liệu  
    X_min = np.min(X)  
    X_max = np.max(X)  
    y_min = np.min(y)  
    y_max = np.max(y)  
  
    plt.axis([X_min*0.95, X_max*1.05, y_min*0.95, y_max*1.05])  
    plt.xlabel(xlabel, fontsize = 16)  
    plt.ylabel(ylabel, fontsize = 16)  
    plt.show()
```

```
✓ [4] plotData(X, y, title="Giá nhà theo diện tích", xlabel="Diện tích (m2)", ylabel="Giá nhà (Tỷ VNĐ)")  
0s
```

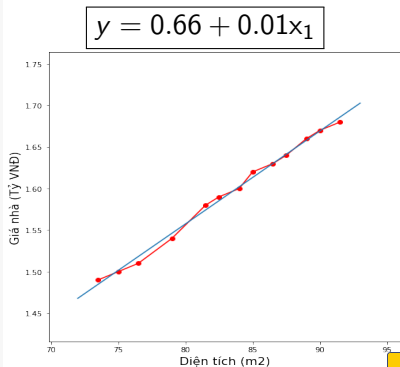


■ Xây dựng hàm tính toán tham số

- Tạo ma trận hệ số điều chỉnh sử dụng **np.ones**
- Ghép ma trận hệ số điều chỉnh và ma trận dữ liệu theo cột sử dụng **np.concatenate**
- Để dễ dàng tính toán, đặt $\mathbf{A} = \bar{\mathbf{X}}^T \bar{\mathbf{X}}$, $\mathbf{b} = \bar{\mathbf{X}}^T \mathbf{y}$

✓
0s

```
[5] """  
    xác định tham số mô hình w1 và w0  
    """  
def model(X, y):  
    # Tạo ma trận hệ số điều chỉnh  
    one = np.ones((X.shape[0],1))  
    # Ghép theo cột với ma trận X (ma trận dữ liệu)  
    Xbar = np.concatenate((one, X), axis = 1)  
    A = np.dot(Xbar.T, Xbar)  
    b = np.dot(Xbar.T, y)  
  
    # Tính nghiệm  
    w = np.dot(np.linalg.pinv(A), b)  
    w_0 = w[0][0]  
    w_1 = w[1][0]  
    #print(w)  
    print("Mô hình:" , w_0 , "+" , w_1 , "*x_1")  
    return w_0, w_1
```



Đánh giá mô hình

- Trong mô hình hồi quy tuyến tính, người ta thường sử dụng chỉ số R-squared để đánh giá chất lượng của mô hình.
- R-squared cho ta biết mức độ các biến đầu vào (biến đầu vào) sẽ giải thích được bao nhiêu phần trăm các biến mục tiêu.
- R-squared càng lớn thì mô hình càng tốt.

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (7)$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (8)$$

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \quad (9)$$



Đánh giá mô hình

Trong đó TSS là tổng bình phương sai số toàn bộ mô hình (Total Sum Squared), RSS là tổng bình phương sai số ngẫu nhiên (Residual Sum Squared), ESS là tổng bình phương sai số được giải thích bởi mô hình (Explained Sum Squared).

$$R^2 = 1 - \frac{RSS}{TSS} \quad (10)$$



PHƯƠNG PHÁP OLS

Ưu điểm	Nhược điểm
Có khả năng tổng quát hóa mô hình	Rất nhạy cảm với các dữ liệu nhiễu, không nhất quán
Thời gian xây dựng và tính toán nhanh	Mô hình không làm việc khi ma trận $\bar{\mathbf{X}}^T \bar{\mathbf{X}}$ là không khả nghịch
Có khả năng dự đoán tốt với các dữ liệu có quy luật	Không hiệu quả với các dữ liệu phân phối không bình thường



PHƯƠNG PHÁP OLS

Ưu điểm	Nhược điểm
Có khả năng tổng quát hóa mô hình	Rất nhạy cảm với các dữ liệu nhiễu, không nhất quán
Thời gian xây dựng và tính toán nhanh	Mô hình không làm việc khi ma trận $\bar{\mathbf{X}}^T \bar{\mathbf{X}}$ là không khả nghịch
Có khả năng dự đoán tốt với các dữ liệu có quy luật	Không hiệu quả với các dữ liệu phân phối không bình thường

Nhận xét: Có thể sử dụng các mô hình hồi quy biến thể khác để hạn chế một số nhược điểm trên:

- Sử dụng mô hình hồi quy Ridge, Lasso
- Sử dụng phương pháp **đạo hàm ngược**
- Sử dụng mô hình hồi quy dạng phi tuyến



Nội dung

- ① Ôn tập kiến thức buổi 1
- ② Khái niệm cơ bản
- ③ Xây dựng mô hình hồi quy
- ④ Triển khai mô hình
- ⑤ Tổng kết



Kiến thức quan trọng

■ Nắm rõ các khái niệm cơ bản

- Mô hình hồi quy, mô hình hồi quy tuyến tính.
- Hiểu được ý nghĩa và mục đích của các mô hình hồi quy và mô hình hồi quy tuyến tính.

■ Hiểu rõ quá trình xây dựng mô hình

- Các kiến thức liên quan tới quá trình học, hàm mất mát.
- Nắm rõ phương pháp OLS và cách biến đổi nghiệm tối ưu của mô hình.

■ Biết cách triển khai mô hình trên các dữ liệu thực tế



Yêu cầu:

- Xem lại mã nguồn và Slide bài giảng trên lớp sau đó tiến hành lựa chọn một bộ dữ liệu mẫu để thực hành.
- Đọc trước tài liệu về mô hình hồi quy Lasso, mô hình hồi quy Ridge và cách đánh giá hiệu quả của mô hình hồi quy tuyến tính.
- Truy cập và đường dẫn dưới đây để tham khảo và đọc trước các tài liệu yêu cầu.

Tài liệu tham khảo: <https://bit.ly/tltkhm>

Mã Nguồn: <https://bit.ly/sourceb1>



THANK YOU

