



# THUẬT TOÁN VÀ LÝ THUYẾT MÁY HỌC

## CÁC BÀI TOÁN HỒI QUY TUYẾN TÍNH

TS. Lương Ngọc Hoàng



# Nội dung

1. Hồi quy tuyến tính đơn giản
2. Hồi quy tuyến tính đa biến
3. Hồi quy đa thức
4. Quá khớp (overfitting) và các kỹ thuật điều chuẩn (regularization)



# HỒI QUY TUYẾN TÍNH ĐƠN GIẢN

## SIMPLE LINEAR REGRESSION



# Hồi quy tuyến tính đơn giản – Mô hình (model)

Một **mô hình (model)** hồi quy tuyến tính đơn giản (simple linear regression) thực hiện dự đoán đầu ra (output) với **một hàm tuyến tính (linear function)** của **một** đặc trưng đầu vào (input feature)  $x$ :

$$f(x; w_0, w_1) = w_0 + w_1 x$$

$w_0, w_1$ : **các tham số (parameters)** của mô hình.

Để tìm ra giá trị phù hợp của  $w_0$  và  $w_1$ , ta sử dụng một **tập dữ liệu (dataset)** gồm các mẫu dữ liệu với cặp thông tin đầu vào (input) – đầu ra (output) tương ứng:

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$$

Tập dữ liệu cũng có thể được biểu diễn như sau:  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$ .

Ta tìm giá trị phù hợp cho  $w_0$  và  $w_1$  dựa vào dữ liệu như thế nào? Ta cần một độ đo để đánh giá tính phù hợp của các giá trị tham số với dữ liệu.



# Hồi quy tuyến tính đơn giản – Hàm mất mát

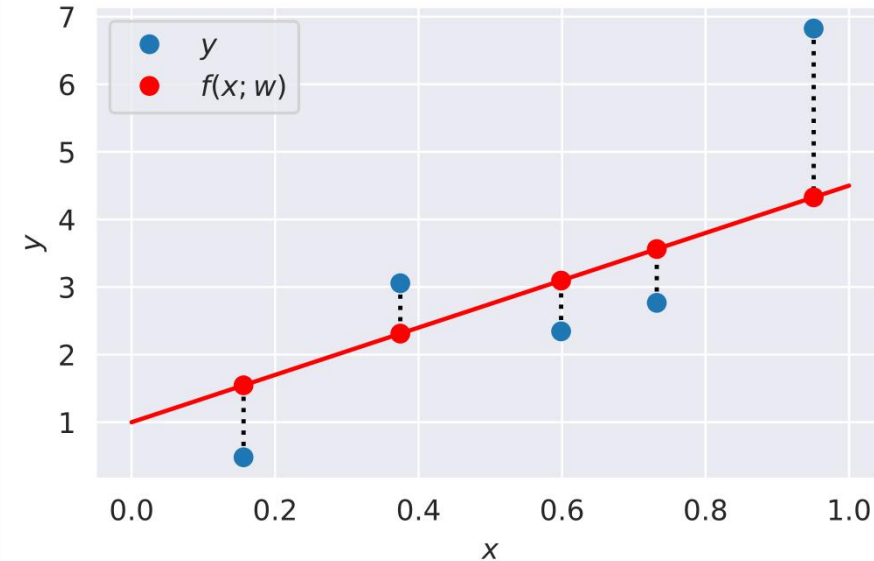
**Hàm mất mát (loss function)**, hay cũng có thể gọi là hàm chi phí (cost function), phản ánh sự phù hợp của các giá trị tham số  $w_0, w_1$  đối với tập dữ liệu:

$$L(w_0, w_1) = \sum_{i=1}^N (y^{(i)} - f(x^{(i)}; w_0, w_1))^2$$

Hàm mất mát trên được gọi là **hàm mất mát bình phương (squared loss)** hoặc tổng bình phương phần dư (Residual Sum of Squares – RSS).

Ta có thể định nghĩa hàm mất mát **sai số bình phương trung bình (Mean Squared Error – MSE)** như sau:

$$L(w_0, w_1) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(x^{(i)}; w_0, w_1))^2$$





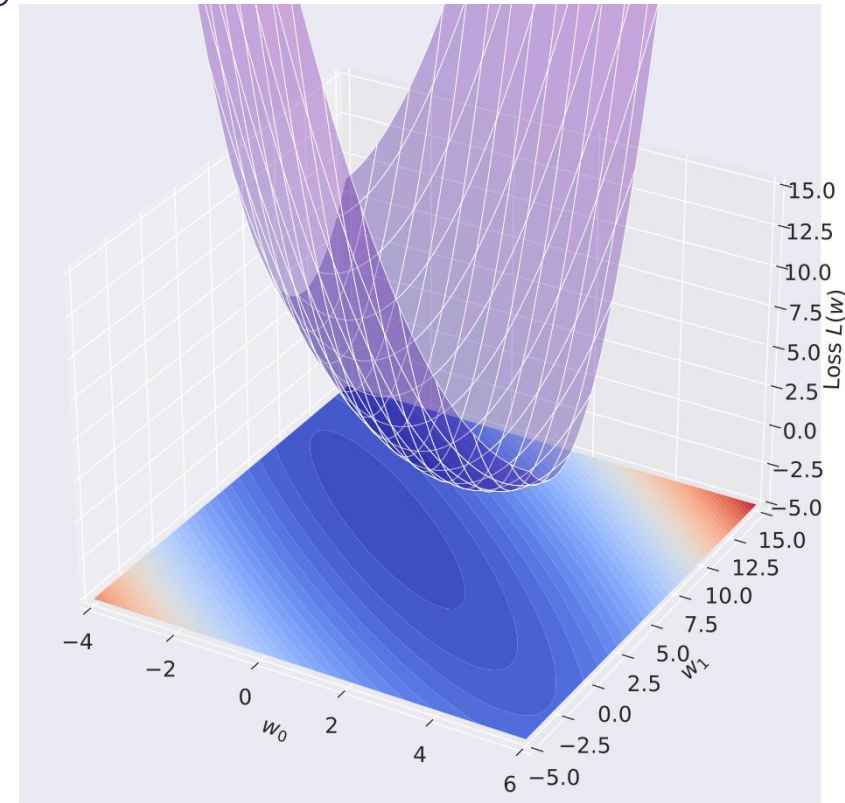
# Hồi quy tuyến tính đơn giản – Tối ưu hóa

Giá trị tham số  $w_0$  và  $w_1$  phù hợp nhất đối với tập dữ liệu giúp mô hình đạt giá trị hàm mất mát nhỏ nhất trên tập dữ liệu:

$$\hat{w}_0, \hat{w}_1 = \arg \min_{w_0, w_1} L(w_0, w_1)$$

Để giải quyết bài toán **tối ưu hoá (optimization)** trên, ta có thể tính đạo hàm của hàm mất mát theo từng tham số  $\frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}$ , và tìm nghiệm của các phương trình  $\frac{\partial L}{\partial w_0} = 0$  và  $\frac{\partial L}{\partial w_1} = 0$ .

$$\begin{aligned} L(w_0, w_1) &= \sum_{i=1}^N (y^{(i)} - f(x^{(i)}; w_0, w_1))^2 \\ &= \sum_{i=1}^N (y^{(i)} - (w_0 + w_1 x^{(i)}))^2 \end{aligned}$$





# Hồi quy tuyến tính đơn giản – Tối ưu hóa

$$L(w_0, w_1) = \sum_{i=1}^N (y^{(i)} - (w_0 + w_1 x^{(i)}))^2$$

$$\begin{aligned} \frac{\partial L}{\partial w_0} &= \frac{\partial}{\partial w_0} \sum_{i=1}^N (y^{(i)} - (w_0 + w_1 x^{(i)}))^2 \\ &= \sum_{i=1}^N 2 (y^{(i)} - (w_0 + w_1 x^{(i)})) (-1) \end{aligned}$$

$$\text{Giải } \frac{\partial L}{\partial w_0} = 0 = \sum_{i=1}^N 2 (y^{(i)} - (w_0 + w_1 x^{(i)})) (-1)$$

$$\sum_{i=1}^N w_0 = \sum_{i=1}^N y^{(i)} - w_1 \sum_{i=1}^N x^{(i)}$$

$$N w_0 = \sum_{i=1}^N y^{(i)} - w_1 \sum_{i=1}^N x^{(i)}$$

$$w_0 = \frac{1}{N} \sum_{i=1}^N y^{(i)} - w_1 \frac{1}{N} \sum_{i=1}^N x^{(i)} = \bar{y} - w_1 \bar{x}$$

$\bar{y}$ : trung bình giá trị đích trong tập dữ liệu.

$\bar{x}$ : trung bình giá trị đặc trưng trong tập dữ liệu.

$$\frac{\partial L}{\partial w_1} = \frac{\partial}{\partial w_1} \sum_{i=1}^N (y^{(i)} - w_0 - w_1 x^{(i)})^2$$

$$= \frac{\partial}{\partial w_1} \sum_{i=1}^N (y^{(i)} - \bar{y} + w_1 \bar{x} - w_1 x^{(i)})^2$$

$$= \frac{\partial}{\partial w_1} \sum_{i=1}^N (y^{(i)} - \bar{y} - w_1 (x^{(i)} - \bar{x}))^2$$

$$= \sum_{i=1}^N 2 (y^{(i)} - \bar{y} - w_1 (x^{(i)} - \bar{x})) (-1) (x^{(i)} - \bar{x})$$

$$\text{Đặt } \frac{\partial L}{\partial w_1} = 0 = \sum_{i=1}^N (y^{(i)} - \bar{y} - w_1 (x^{(i)} - \bar{x})) (x^{(i)} - \bar{x})$$

$$\sum_{i=1}^N w_1 (x^{(i)} - \bar{x})^2 = \sum_{i=1}^N (y^{(i)} - \bar{y}) (x^{(i)} - \bar{x})$$

$$w_1 = \frac{\sum_{i=1}^N (y^{(i)} - \bar{y}) (x^{(i)} - \bar{x})}{\sum_{i=1}^N (x^{(i)} - \bar{x})^2}$$

Giá trị  $w_0$  và  $w_1$  được gọi là ước lượng bình phương tối thiểu (least square estimates)



# HỒI QUY TUYẾN TÍNH ĐA BIẾN

## MULTIPLE LINEAR REGRESSION





# Hồi quy tuyến tính đa biến – Mô hình (model)

Một **mô hình (model)** hồi quy tuyến tính đa biến (multiple linear regression) thực hiện dự đoán đầu ra (output) với **một hàm tuyến tính (linear function)** của **các** đặc trưng đầu vào (input features)  $x_1, x_2, \dots, x_D$ :

$$f(x_1, x_2, \dots, x_D; w_0, w_1, w_2, \dots, w_D) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D$$

$w_0, w_1, w_2, \dots, w_D$ : **các tham số (parameters)** của mô hình.

$w_0$ : tham số **bias**.

$w_1, w_2, \dots, w_D$ : **trọng số (weights)** của mỗi đặc trưng.

Ta có thể biểu diễn

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \quad \text{và} \quad \mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}. \text{ Do đó, ta có mô hình } f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}.$$



# Hồi quy tuyến tính đa biến – Hàm mất mát

Hàm mất mát bình phương (squared loss) của mô hình hồi quy tuyến tính đa biến:

$$L(\mathbf{w}) = \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2 = \sum_{i=1}^N (y^{(i)} - (w_0 + w_1 x_1 + \dots + w_D x_D))^2 = \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2$$

Ta có thể định nghĩa:

$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}^{(1)T} & \text{---} \\ \text{---} & \mathbf{x}^{(2)T} & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{x}^{(N)T} & \text{---} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}. \text{ Vector sai số là: } \mathbf{e} = \begin{bmatrix} y^{(1)} - \mathbf{x}^{(1)T} \mathbf{w} \\ y^{(2)} - \mathbf{x}^{(2)T} \mathbf{w} \\ \vdots \\ y^{(N)} - \mathbf{x}^{(N)T} \mathbf{w} \end{bmatrix} = \mathbf{y} - \mathbf{X}\mathbf{w}.$$

Hàm mất mát có thể được biểu diễn như sau:

$$\begin{aligned} L(\mathbf{w}) &= \mathbf{e}^T \mathbf{e} = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &= \mathbf{y}^T \mathbf{y} - \underbrace{\mathbf{y}^T \mathbf{X}\mathbf{w}}_{1 \times 1} - \underbrace{\mathbf{w}^T \mathbf{X}^T \mathbf{y}}_{1 \times 1} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \end{aligned}$$



# Hồi quy tuyến tính đa biến – Tối ưu hóa

Hàm mất mát:

$$L(\mathbf{w}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}$$

Ta tính gradient:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} [\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}] \\ &= -2\mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + (\mathbf{X}^T \mathbf{X})^T) \mathbf{w} \\ &= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} \end{aligned}$$

Ta sử dụng các công thức sau:

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

Ta đặt  $\frac{\partial L}{\partial \mathbf{w}} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{0}$  và giải:

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Công thức trên có tên gọi là **Normal Equations**, trả về ước lượng bình phương tối thiểu.



# Hồi quy tuyến tính – Hình học

Chiếu trực giao vector  $b$  lên đường thẳng đi qua vector  $a$ .  
Bởi vì  $e$  trực giao với  $a$ , nên ta có:

$$a^T e = 0$$

$$a^T (b - p) = 0$$

$$a^T (b - xa) = 0$$

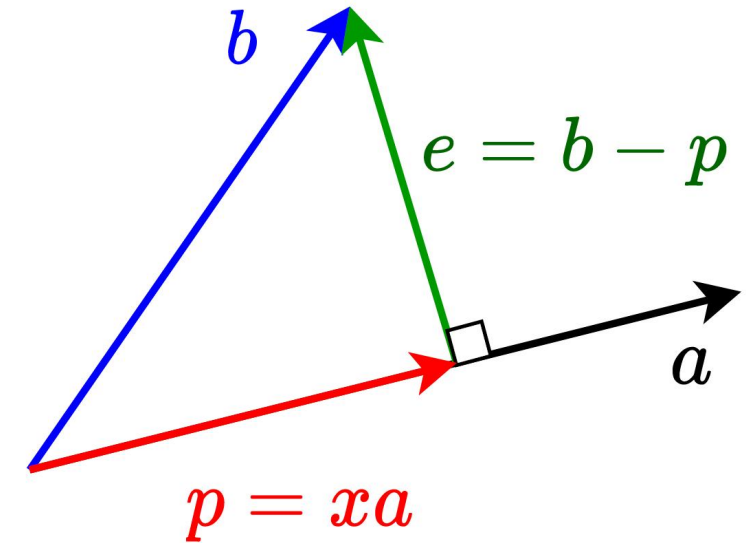
$$xa^T a = a^T b$$

$$x = \frac{a^T b}{a^T a}$$

Do đó, hình chiếu  $p = xa = ax = a \frac{a^T b}{a^T a}$ .

Ta có ma trận chiếu (projection matrix)  $P$  sao cho  $p = Pb$ .

$$P = \frac{aa^T}{a^T a}$$



- Hình chiếu  $p$  của  $b$  nằm trên đường thẳng đi qua  $a$ .  
Do đó:  $p = xa$ .
- Vector  $e$  thể hiện sự khác biệt giữa  $b$  và hình chiếu  $p$ .  
Do đó:  $e = b - p$ .



# Hồi quy tuyến tính – Hình học

Chiếu trực giao vector  $\mathbf{b}$  lên mặt phẳng chứa 2 vector  $\mathbf{a}_1$  và  $\mathbf{a}_2$ .

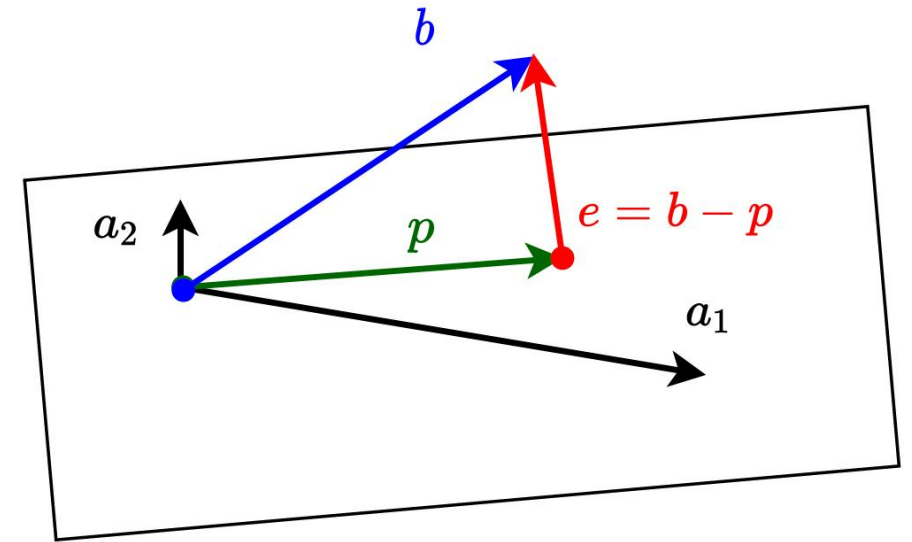
- Mặt phẳng chứa  $\mathbf{a}_1$  và  $\mathbf{a}_2$  là không gian cột (column space) của  $A$ .
- Hình chiếu  $\mathbf{p}$  của  $\mathbf{b}$  do đó là một tổ hợp tuyến tính của  $\mathbf{a}_1$  và  $\mathbf{a}_2$ .

$$\mathbf{p} = x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2$$

$$\mathbf{p} = A\mathbf{x} = \begin{bmatrix} | & | \\ \mathbf{a}_1 & \mathbf{a}_2 \\ | & | \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

- Ta cần tìm  $\mathbf{x}$ .
- Vector  $\mathbf{e}$  thể hiện sự khác biệt giữa  $\mathbf{b}$  và  $\mathbf{p}$ . Do đó:

$$\mathbf{e} = \mathbf{b} - \mathbf{p}$$



- Chọn 2 vector  $\mathbf{a}_1$  và  $\mathbf{a}_2$  độc lập với nhau trong mặt phẳng để tạo thành một cơ sở.

$$A = \begin{bmatrix} | & | \\ \mathbf{a}_1 & \mathbf{a}_2 \\ | & | \end{bmatrix}$$





# Hồi quy tuyến tính – Hình học

- $e = b - p$  trực giao với mặt phẳng chứa  $a_1$  và  $a_2$  nên ta có:

$$\begin{bmatrix} - & a_1^T & - \\ - & a_2^T & - \end{bmatrix} \begin{bmatrix} | \\ e \\ | \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$A^T e = 0$$

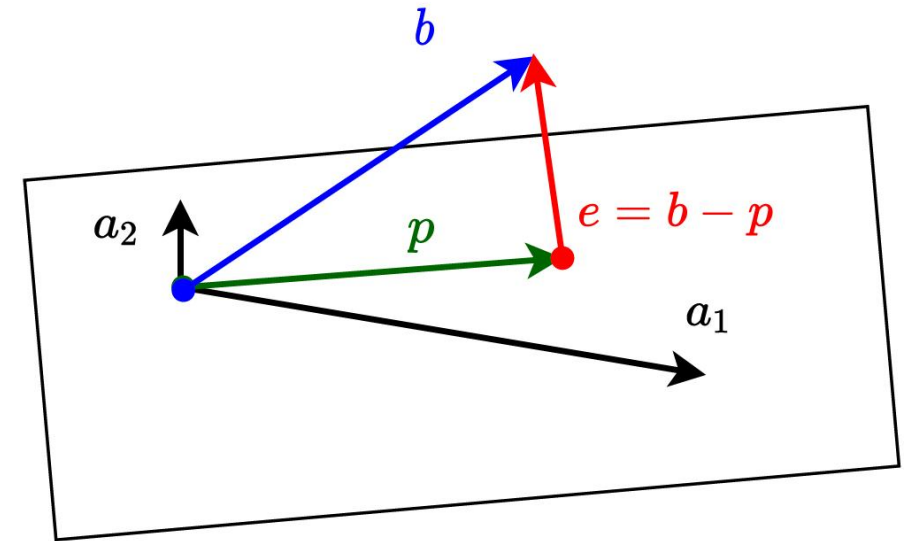
$$A^T (b - Ax) = 0$$

$$A^T Ax = A^T b$$

$$x = (A^T A)^{-1} A^T b$$

- Vector hình chiếu  $p = Ax = A(A^T A)^{-1} A^T b$ .
- Ta có ma trận chiếu (projection matrix)  $P$  là:

$$P = A(A^T A)^{-1} A^T$$



- Chọn 2 vector  $a_1$  và  $a_2$  độc lập với nhau trong mặt phẳng để tạo thành một cơ sở.

$$A = \begin{bmatrix} | & | \\ a_1 & a_2 \\ | & | \end{bmatrix}$$

- $p$  là một tổ hợp tuyến tính của  $a_1$  và  $a_2$ .



# Hồi quy tuyến tính – Hình học

$$X = \begin{bmatrix} \text{---} & \mathbf{x}^{(1)T} & \text{---} \\ \text{---} & \mathbf{x}^{(2)T} & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{x}^{(N)T} & \text{---} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix}$$

- **Hàm mất mát bình phương (squared loss)** của mô hình hồi quy tuyến tính:

$$L(\mathbf{w}) = \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2 = \sum_{i=1}^N (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2 = (\mathbf{y} - X\mathbf{w})^T (\mathbf{y} - X\mathbf{w}) = \|\mathbf{y} - X\mathbf{w}\|^2$$

- **Bộ giá trị trọng số tối ưu  $\mathbf{w}^*$**  cực tiểu hóa (minimize) hàm mất mát trên:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} L(\mathbf{w}) = \arg \min_{\mathbf{w}} \|\mathbf{y} - X\mathbf{w}\|^2$$

- Dự đoán của mô hình  $X\mathbf{w}$  tương ứng với một tổ hợp tuyến tính các cột của ma trận  $X$ , là một vector trong không gian cột của  $X$ .
- Ta muốn tìm bộ giá trị trọng số tối ưu của  $\mathbf{w}$  sao cho sự sai lệch giữa dự đoán của mô hình  $X\mathbf{w}$  và vector giá trị đích  $\mathbf{y}$  là nhỏ nhất có thể.
- $X\mathbf{w}^*$  **tối ưu** do đó sẽ là **hình chiếu trực giao** của  $\mathbf{y}$  vào không gian cột của  $X$ .

$$\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$$



# HỒI QUY ĐA THỨC

## POLYNOMIAL REGRESSION



# Hồi quy đa thức

Hồi quy tuyến tính đa biến:

$$f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} = w_0 + w_1 x_1 + \dots + w_D x_D$$

Huấn luyện trên tập dữ liệu  $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$  sử dụng hàm mất mát:

$$\begin{aligned} L(\mathbf{w}) &= \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2 \\ &= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \end{aligned}$$

với

$$\mathbf{X} = \begin{bmatrix} \text{---} & \mathbf{x}^{(1)T} & \text{---} \\ \text{---} & \mathbf{x}^{(2)T} & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{x}^{(N)T} & \text{---} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}$$

Lời giải:  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Hồi quy đa thức:

Nếu ta muốn sử dụng mô hình

$$f(x; \mathbf{w}) = w_0 + w_1 x + w_2 x^2$$

ta có thể định nghĩa:

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix}$$

Ta có thể biểu diễn  $f(x; \mathbf{w}) = \mathbf{w}^T \phi(x)$  và giải tương tự như hồi quy tuyến tính đa biến với  $\phi(x)$  đóng vai trò như vector đặc trưng  $\mathbf{x}$ .

Mã trận dữ liệu trở thành:

$$\Phi = \begin{bmatrix} \text{---} & \phi(x^{(1)})^T & \text{---} \\ \text{---} & \phi(x^{(2)})^T & \text{---} \\ & \vdots & \\ \text{---} & \phi(x^{(N)})^T & \text{---} \end{bmatrix} = \begin{bmatrix} 1 & x^{(1)} & x^{(1)2} \\ 1 & x^{(2)} & x^{(2)2} \\ & \vdots & \\ 1 & x^{(N)} & x^{(N)2} \end{bmatrix}$$

Lời giải:  $\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$



# Hồi quy đa thức

Một **mô hình (model)** hồi quy đa thức (polynomial regression) thực hiện dự đoán đầu ra (output) với một **hàm đa thức (polynomial function)** của các đặc trưng đầu vào (input features). Với một đặc trưng  $x$ , ta có thể có mô hình sau:

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Dx^D$$

trong đó,  $\phi(x) = [1, x, x^2, \dots, x^D]^T$ .

Hoặc, ví dụ với các đặc trưng  $\mathbf{x} = (x_1, x_2)$ , ta có thể có mô hình sau:

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2 + w_4x_1^2 + w_5x_2^2$$

trong đó,  $\phi(\mathbf{x}) = [1, x_1, x_2, x_1x_2, x_1^2, x_2^2]^T$ . Đây là một ví dụ **chế tác đặc trưng (feature engineering)**.

**Lưu ý:**

- Mô hình hồi quy đa thức là **phi tuyến tính** theo các đặc trưng đầu vào  $\mathbf{x}$ .
- Mô hình hồi quy đa thức là **tuyến tính** theo các trọng số  $\mathbf{w}$  của mô hình. Do đó, ta giải quyết bài toán hồi quy đa thức tương tự như hồi quy tuyến tính đa biến (multiple linear regression).

$$\Phi = \begin{bmatrix} \text{--} \phi(x^{(1)})^T \text{--} \\ \text{--} \phi(x^{(2)})^T \text{--} \\ \vdots \\ \text{--} \phi(x^{(N)})^T \text{--} \end{bmatrix}. \text{ Ta có: } \hat{\mathbf{w}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}.$$



# Hồi quy đa thức – Ví dụ

## Tập dữ liệu (Dataset)

Phát sinh một tập dữ liệu (dataset) gồm  $N$  điểm dữ liệu

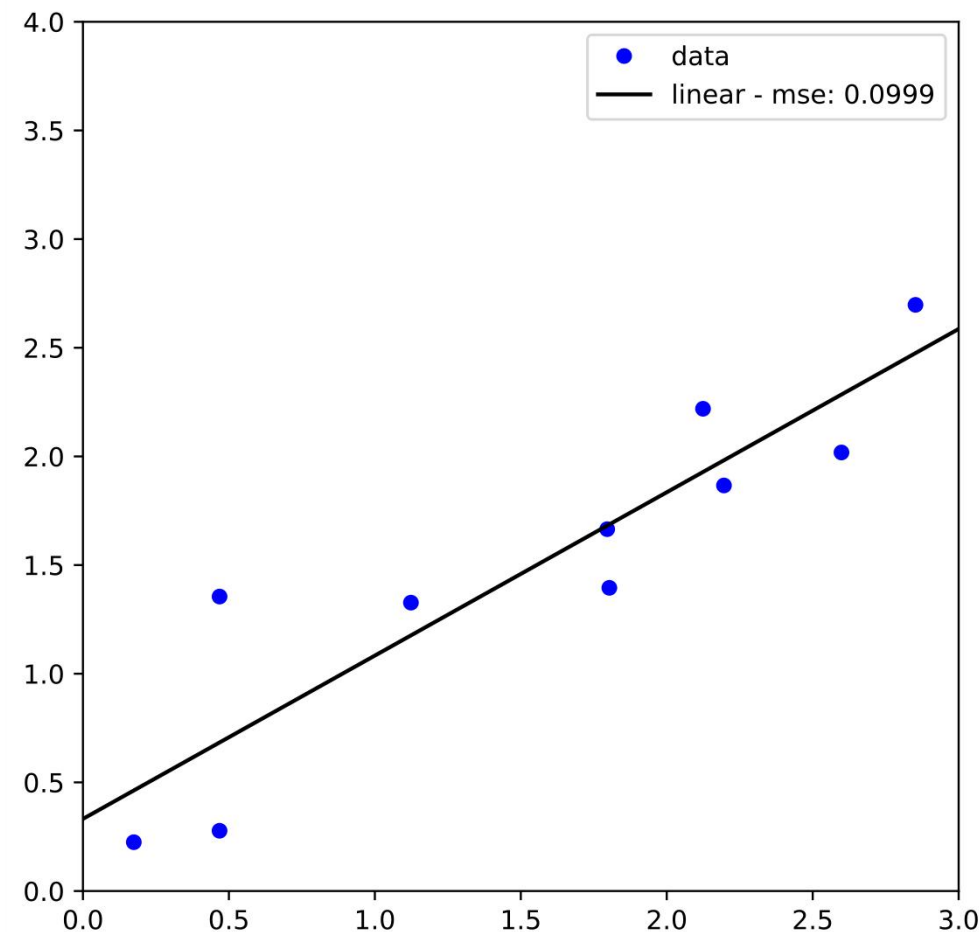
(data points)  $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$  với:

- $x^{(i)} \in (0,1)$  là **đặc trưng đầu vào** (input feature) của điểm dữ liệu thứ  $i$ .
- $y^{(i)} \in \mathbb{R}$  là **giá trị đích** (target value) của điểm dữ liệu thứ  $i$ . Ta mô phỏng giá trị  $y^{(i)}$  thu thập được có một lượng nhiễu nhỏ cho mỗi điểm dữ liệu.

## Hồi quy tuyến tính đơn giản

Sử dụng một mô hình hồi quy tuyến tính đơn giản, ta có:

$$f(x; \mathbf{w}) = w_0 + w_1 x$$



```
from sklearn.linear_model import LinearRegression
np.random.seed(42)
N = 10
X_train = 3.0*np.random.rand(N, 1)
y_train = 1.0 + 0.5*X_train + np.random.randn(N,1)/2.0
lin_reg = LinearRegression()
lin_reg.fit(X_train,y_train)
```

# Hồi quy đa thức – Ví dụ

## Hồi quy tuyến tính đơn giản

- Sử dụng một mô hình hồi quy tuyến tính đơn giản:

$$f(x; \mathbf{w}) = w_0 + w_1x$$

## Hồi quy đa thức

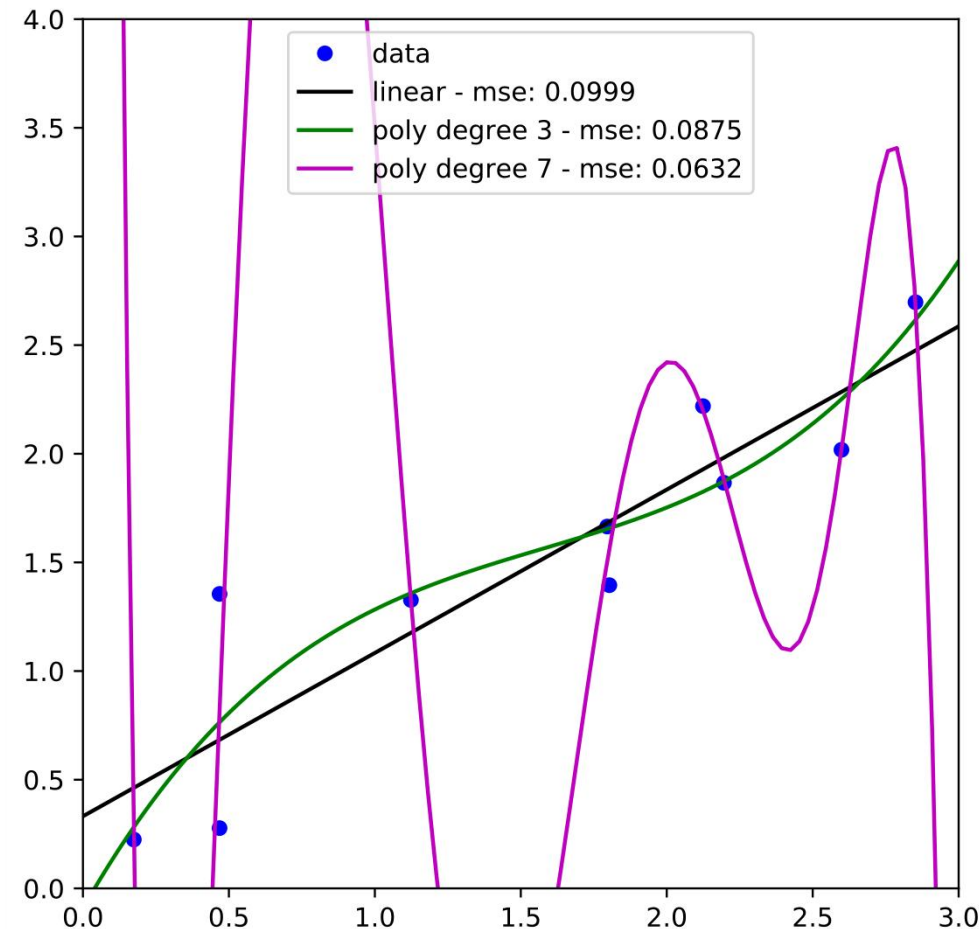
- Sử dụng một mô hình hồi quy đa thức bậc 3:

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3$$

- Sử dụng một mô hình hồi quy đa thức bậc 7:

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 + \dots + w_7x^7$$

- Mô hình nào là tốt nhất đối với tập dữ liệu đang xét?



```
from sklearn.preprocessing import PolynomialFeatures
deg = 7
rng = np.random.RandomState(42)
poly = PolynomialFeatures(degree=deg,
                           include_bias=False)
X_train = 3.0 * np.random.rand(N, 1)
y_train = 1.0 + 0.5 * X_train + np.random.randn(N, 1) / 2.0
poly_reg = LinearRegression()
poly_reg.fit(X_train, y_train)
```

# Hồi quy đa thức – Ví dụ

- Hồi quy tuyến tính đơn giản:

$$f(x; \mathbf{w}) = w_0 + w_1x$$

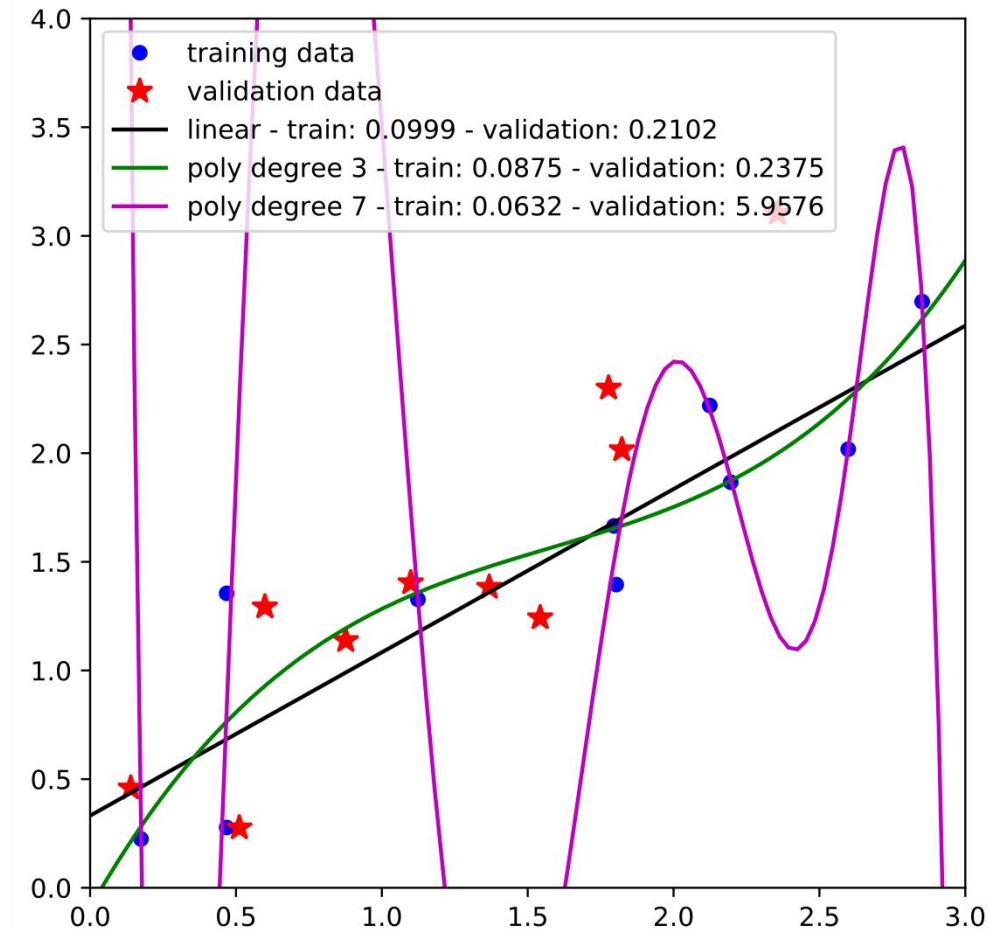
- Hồi quy đa thức bậc 3:

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3$$

- Hồi quy đa thức bậc 7:

$$f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3 + \dots + w_7x^7$$

- Mô hình nào là tốt nhất đối với tập dữ liệu đang xét?
- Ta cần đánh giá các mô hình sau khi huấn luyện trên **một tập dữ liệu khác** với tập dữ liệu đã dùng để huấn luyện các mô hình.
- Mô hình phức tạp có thể có **độ lỗi nhỏ trên tập huấn luyện** nhưng có **độ lỗi lớn trên tập dữ liệu đánh giá**.



```
from sklearn.preprocessing import PolynomialFeatures
N_train = 3.0*np.random.rand(N, 1)
X_train = PolynomialFeatures(degree=7).fit_transform(N_train/2.0,
include_bias=False)
N_test = 3.0*np.random.rand(N_test, 1)
X_test = PolynomialFeatures(degree=7).fit_transform(N_test/2.0,
include_bias=False)
```



# QUÁ KHỚP VÀ KỸ THUẬT ĐIỀU CHUẨN

---

## OVERFITTING & REGULARIZATION

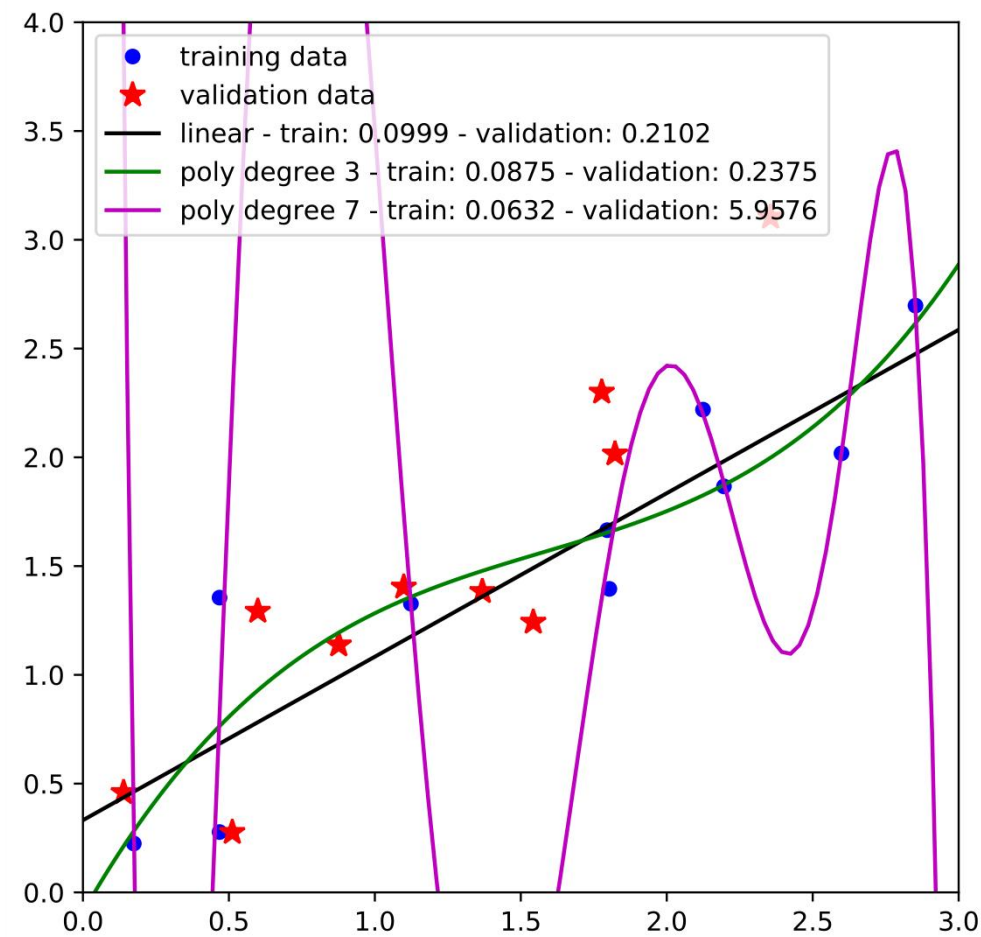


# Hồi quy đa thức – Quá khớp

➤ Mô hình nào là tốt nhất đối với tập dữ liệu đang xét?

```
---Simple Linear Regression---  
[0.33181066] [[0.75155622]]  
Training Loss: 0.09991531905633781  
Validation Loss: 0.21020128681013045  
  
---Polynomial Regression degree = 3---  
[-0.0916451] [[ 2.35120734 -1.23811156 0.26177074]]  
Training Loss: 0.08750190066671261  
Validation Loss: 0.23747048298743456  
  
---Polynomial Regression degree = 7---  
[ 3 8 . 0 6 9 8 7 1 0 1 ] [ [ - 3 9 9 . 6 0 8 7 3 8 9 3 1 3 6 4 . 6 9 5 7 9 4 8 8 -  
2093.21939536 1665.82047801  
-716.24860925 158.03695712 -14.0390417 ] ]  
Training Loss: 0.06322418953671403  
Validation Loss: 5.9575569104420065
```

➤ **Quá khớp (overfitting):** hiện tượng mô hình dự đoán rất chính xác trên dữ liệu huấn luyện, nhưng lại không thể dự đoán tốt trên dữ liệu mới (không tổng quát hóa tốt).







# Hiện tượng quá khớp (Overfitting)

Giả sử ta sử dụng các kỹ thuật **chế tác đặc trưng (feature engineering)** với các phép biến đổi  $\Phi$  cho mô hình hồi quy tuyến tính. Giả sử số điểm dữ liệu huấn luyện  $N = 10$ .

Với tập dữ liệu trên, ta huấn luyện mô hình để tìm bộ giá trị trọng số  $w$  sao cho mô hình dự đoán chính xác nhất có thể:  $y \approx \Phi w$ .

- Nếu ta sử dụng  $\Phi$  với 2 đặc trưng thì ta có:  $\overset{10 \times 1}{\tilde{y}} \approx \overset{10 \times 2}{\tilde{\Phi}} \overset{2 \times 1}{\tilde{w}}$ .
- Nếu ta sử dụng  $\Phi$  với 10 đặc trưng thì ta có:  $\overset{10 \times 1}{\tilde{y}} \approx \overset{10 \times 10}{\tilde{\Phi}} \overset{10 \times 1}{\tilde{w}}$ .
- Ta có thể giải chính xác giá trị tối ưu của  $w$  nếu  $\Phi$  khả nghịch:  $w = \Phi^{-1}y$ .
- Giá trị của hàm mất mát trên tập dữ liệu huấn luyện là bao nhiêu?
- Mô hình này liệu có hoạt động tốt khi dự đoán trên các mẫu dữ liệu mới không nằm trong số  $N$  điểm dữ liệu huấn luyện?



# Điều chuẩn (Regularization)

Ta muốn sử dụng các mô hình bậc cao và có tính linh hoạt cao để có kết quả huấn luyện tốt trên tập dữ liệu huấn luyện, nhưng ta cũng muốn kiểm soát “độ phức tạp” của mô hình để giảm thiểu khả năng quá trình huấn luyện dẫn đến hiện tượng quá khớp.

Ta sử dụng hàm mất mát được **điều chuẩn (regularized)** để huấn luyện mô hình:

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \{L(\mathbf{w}) + \text{penalty}(\mathbf{w})\}$$

- $L(\mathbf{w})$  : hàm mất mát thông thường, ví dụ như hàm mất mát bình phương.
- **penalty**( $\mathbf{w}$ ): hàm phạt được sử dụng để giới hạn độ lớn của các tham số  $\mathbf{w}$ . Hai phương pháp phạt phổ biến là:
  - Ridge –  $L_2$  regularization
  - Lasso –  $L_1$  regularization



# Điều chuẩn (Regularization)

## Ridge – $L_2$ regularization

$$L_{\lambda}(\mathbf{w}) = \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2 + \lambda \sum_{k=1}^D (w_k)^2$$

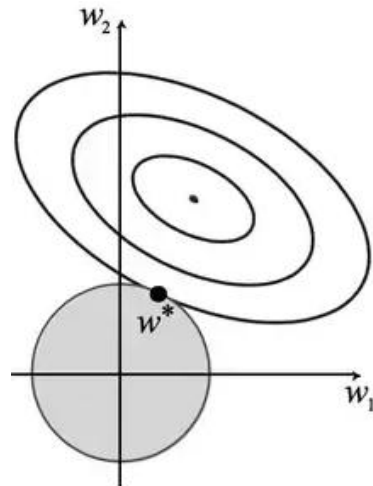
Ta thường không cần phải điều chuẩn  $w_0$ . Tại sao?

Ta có thể chuẩn hóa zero-mean dữ liệu trước khi huấn luyện mô hình: các cột của  $\mathbf{X}$  (hoặc  $\Phi$ ) được chuẩn hóa để có giá trị trung bình là 0.

Lời giải dạng đóng  
(closed-form solution) của

Hồi quy Ridge là:

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

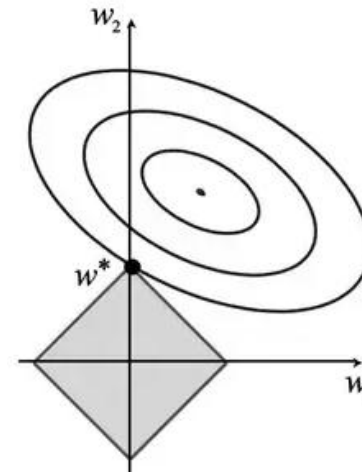


## Lasso – $L_1$ regularization

$$L_{\lambda}(\mathbf{w}) = \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}))^2 + \lambda \sum_{k=1}^D |w_k|$$

$L_1$  regularization không có lời giải dạng đóng vì có thành phần tính giá trị tuyệt đối  $|w_k|$ .

Lasso có thể đưa trọng số  $w_k$  của một số đặc trưng  $x_k$  về giá trị 0. Ta có mô hình đơn giản hơn, loại bỏ đi những đặc trưng không cần thiết.



# Hồi quy đa thức – Quá khớp

---Polynomial Regression degree = 7---

```
[38.06987101] [[ -399.60873893  1364.69579488 -2093.21939536  1665.82047801  
-716.24860925  158.03695712 -14.0390417 ]]
```

Training Loss: 0.06322418953671403

Validation Loss: 5.9575569104420065

--- lambda = 0 ---

```
[38.06839355] [ -399.5931998  1364.64294788 -2093.13814505  1665.75551427  
-716.22051288  158.03071754 -14.03848365]
```

--- lambda = 1e-05 ---

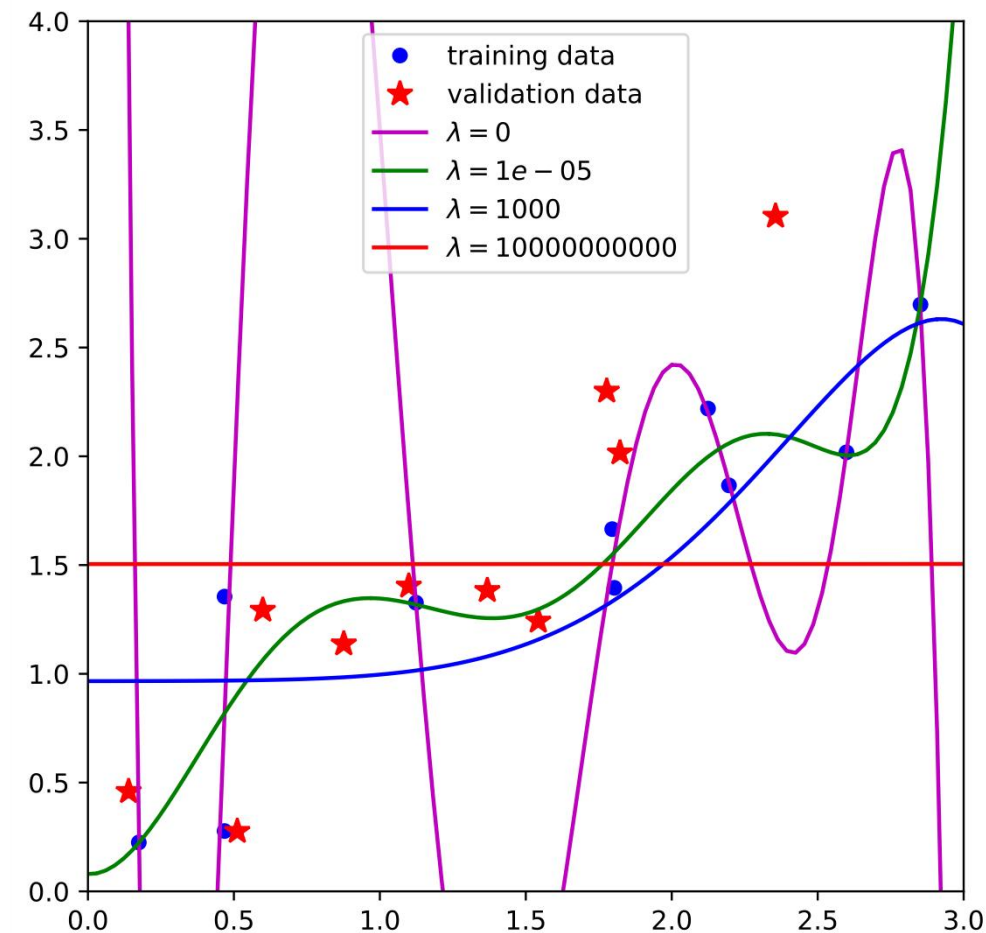
```
[0.07987999] [-0.09370931  6.14945012 -4.77478402 -3.78021829  5.75140457 -2.28394658  
0.29814782]
```

--- lambda = 1000 ---

```
[0.96643294] [ 0.00141423  0.00266913  0.00442294  0.00683627  0.00931761  0.00896911  
-0.00359498]
```

--- lambda = 10000000000 ---

```
[1.50446301] [6.07278017e-10 1.77351966e-09 4.72471362e-09 1.25448975e-08  
3.35844395e-08 9.07463757e-08 2.47260904e-07]
```



```
from sklearn.preprocessing import PolynomialFeatures  
deg = 7  
poly = PolynomialFeatures(degree=deg,  
include_bias=False)  
X_poly = poly.fit_transform(X_train)  
poly_reg = LinearRegression()  
poly_reg.fit(X_poly, y_train)  
ridge_reg = Ridge(alpha=10000000000)  
ridge_reg.fit(X_poly, y_train)
```

➤ Giá trị phù hợp cho siêu tham số (hyperparameter)  $\lambda =$

?

Thực hiện bởi Trường Đại học Công nghệ Thông tin, ĐHQG-HCM