

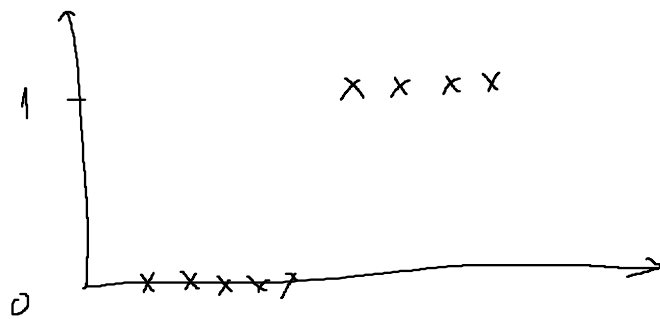
VD: Email: spam / not spam?

Tiền bầu: lành tính / ác tính?

$y \in \{0, 1\}$ 0. negative
1. Positive class

$y \in \{0, 1, 2, 3, 4, \dots\} \rightarrow$ multiple classification

* Đồ thị



Đồ thị sử dụng linear regression nhưng không chính xác

Classification: 0 or 1

$h_{\theta}(x)$: can be > 1 or < 0

Logistic regression: $0 \leq h_{\theta}(x) \leq 1$

Model

mong muốn: $0 \leq h_{\theta}(x) \leq 1$

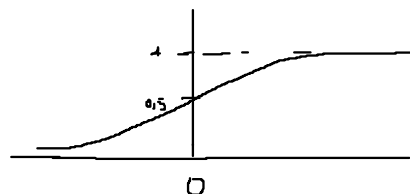
$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid function

Logistic function

$$\Leftrightarrow h_{\theta}(x) = \frac{1}{1 + e^{-(\theta^T x)}}$$



giải ví dụ toán học

$h_{\theta}(x)$ = xác suất để $y=1$ tại mẫu x

VĐ: $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \end{bmatrix}$

$$h_{\theta}(x) = 0,7 = 70\%$$

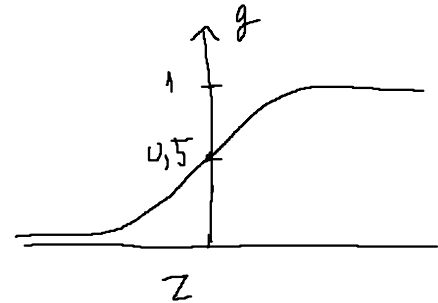
\Rightarrow 70% là xác suất

$h_{\theta}(x) = P(y=1|x, \theta)$ - tỉ lệ để $y=1$ tại x với tham số là θ

$$h_{\theta}(x) = g(\theta^T x) = P(y=1|x;\theta)$$

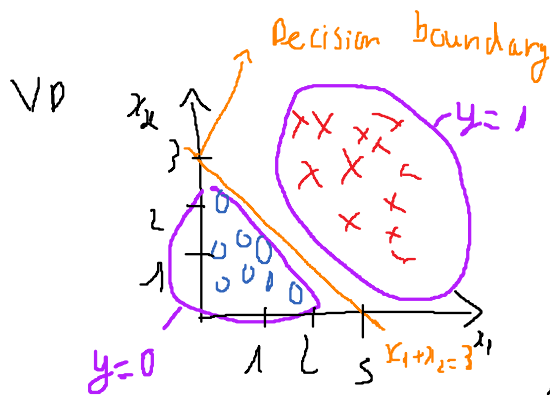
$$g(u) = \frac{1}{1+e^{-u}}$$

do đó
 $y=1$ nếu $h_{\theta}(x) \geq 0,5$
 $y=0$ nếu $h_{\theta}(x) < 0,5$



$$g(z) \geq 0,5 \text{ khi } z \geq 0$$

$$\Leftrightarrow h_{\theta}(x) \geq 0,5 \text{ hay } \theta^T x \geq 0$$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\text{giá trị} \quad \begin{matrix} -3 & 1 & 1 \end{matrix}$$

$$\Leftrightarrow \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

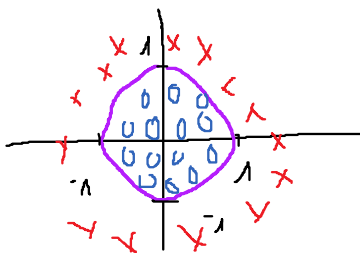
$$y=1 \Leftrightarrow h_{\theta}(x) \geq 0,5 \Leftrightarrow -3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 \geq 3$$

$$x_1 + x_2 = 3 \Leftrightarrow h_{\theta}(x) = 0,5$$

↓
decision boundary

non-linear decision boundaries



$$h_{\theta}(x) = g(-1 + x_1^2 + x_2^2)$$

decision boundary $x_1^2 + x_2^2 = 1$

x

Training set $\{(x^0, y^0), (x^1, y^1), \dots, (x^m, y^m)\}$

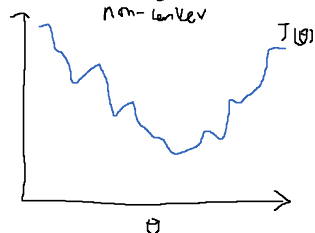
m examples $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$ $y \in \{0, 1\}$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Linear regress. $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^i) - y^i)^2$
cost $(h_{\theta}(x^i), y^i)$

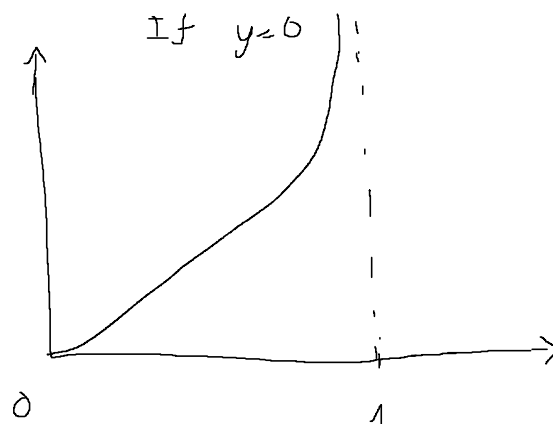
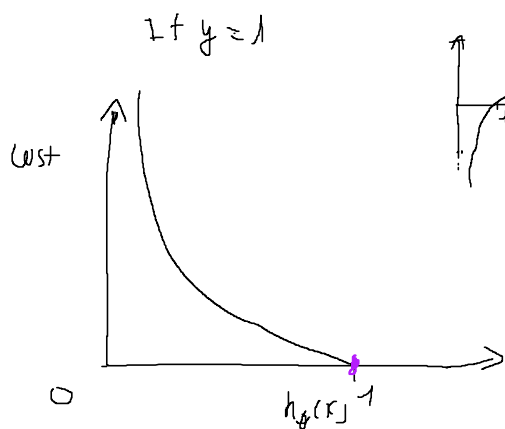
$$\text{cost}(h_{\theta}(x^i), y^i) = \frac{1}{2} (h_{\theta}(x^i) - y^i)^2$$

nếu vẽ đúng $J(\theta)$ của linear cho logistic thì sẽ là một hàm non-convex



Cost Function for logistic

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



cost = 0 khi $y=1$ và $h_{\theta}(x)=1$

nhưng $h_{\theta}(x) \rightarrow 0 \Rightarrow \text{large cost}$
 $\text{cost} \rightarrow \infty$

note. nếu $h_{\theta}(x)=0$, (Prob $P(y=1|x;\theta)$)
 $y=1$ thì sẽ trở nên lớn

note. nếu $h_\theta(x) = 0$, (hoặc $P(y=1|x;\theta)$)
 $y=1$ thì sẽ trở giá lớn

$$\text{cost}(h_\theta(x), y) = 0 \text{ if } h_\theta(x) = y$$

$$\text{cost}(h_\theta(x), y) \rightarrow \infty \text{ if } \begin{cases} y=0, & h_\theta(x) \rightarrow 1 \\ y=1, & h_\theta(x) \rightarrow 0 \end{cases}$$

Simplified cost function

$$\text{cost}(h_\theta(x), y) = -y \cdot \log(h_\theta(x)) - (1-y) \cdot (\log(1-h_\theta(x)))$$

$$\Rightarrow J(\theta) = -\frac{1}{m} \sum_{i=1}^m y \log(h_\theta(x)) + (1-y) (\log(1-h_\theta(x))) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x), y)$$

↓
 tìm bằng maximum likelihood estimation (MLE)
 * còn nhiều cost function khác

$\min_{\theta} J(\theta)$ tìm θ để dự đoán với x mới

output: $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}} \quad P(y=1|x;\theta)$

$$J_\theta(x) = -\frac{1}{m} \sum_{i=1}^m y^i \log h_{\theta^i} + (1-y^i) \cdot \log(1-h_{\theta^i})$$

đạo hàm vectơ
 $h = g(x\theta)$

$$J(\theta) = \frac{1}{m} \cdot (-y^T \log h - (1-y)^T \log(1-h))$$

Lặp lại đến khi hội tụ

$$\theta_j = \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\Leftrightarrow \theta_j = \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta^i} - y^i) x^i$$

vectơ

$$\theta = \theta - \frac{\alpha}{m} X^T (g(X\theta) - \bar{y})$$

$$\text{cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m (y \log(h_{\theta}(x_i)) + (1-y) \log(1-h_{\theta}(x_i)))$$

Gradient descent

minimize $J(\theta)$.
 repeat {
 $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) = \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_j^{(i)}$
 }

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$

Prove

$$\begin{aligned} \frac{\partial}{\partial \theta} J(\theta) &= \frac{\partial}{\partial \theta} \left(-y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x)) \right) \quad \text{with } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad \text{and } h'_{\theta}(x) = x(1-h_{\theta}(x))h_{\theta}(x) \\ &= -y \cdot x(1-h_{\theta}) + (1-y) \cdot h_{\theta}(x) \cdot x \\ &= (-y + y h_{\theta}(x) - y h_{\theta}(x) + h_{\theta}(x)) \cdot x \\ &= (h_{\theta}(x) - y) \cdot x \end{aligned}$$

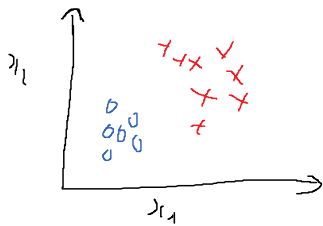
optimization

- conjugate gradient
- BFGS
- left BFGS

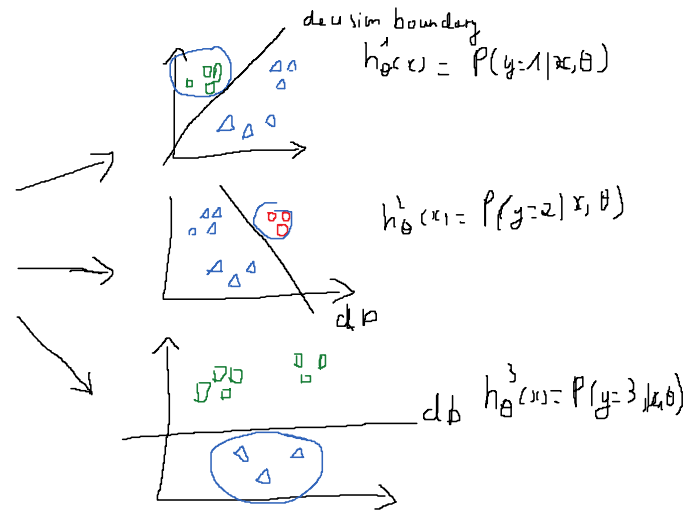
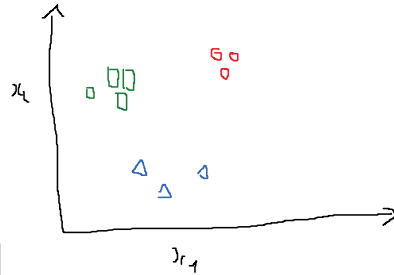
multi class classification

eg: email, work, friend, family, hobby
 $y=1$ 2 3 4

binary classification



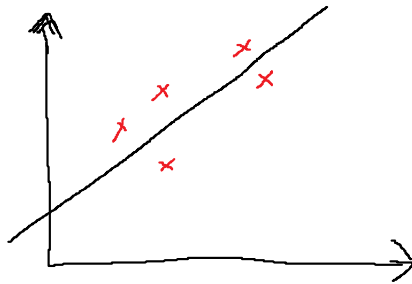
multi class



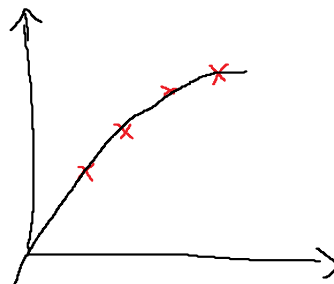
+ train model by $h_{\theta}^{(i)}$ cho tất cả y để dự đoán $y=i$

$$\text{Max}_j h_{\theta}^{(j)}(x)$$

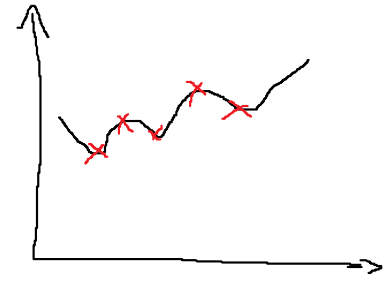
overfitting



$\theta_0 + \theta_1 x$
under fit high bias



$\theta_0 + \theta_1 x + \theta_2 x^2$



$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
over fitting high variance

- Under fit: mô hình dự đoán sai ở cả train set và test set
- Over fit: mô hình dự đoán tốt ở training set nhưng tệ ở test set
 - Khi có quá nhiều feature và model fit với data set nhưng dự đoán sai với data mới
- cách xử lý: 1 - giữ lại những feature quan trọng
 - thuật toán lựa chọn mô hình

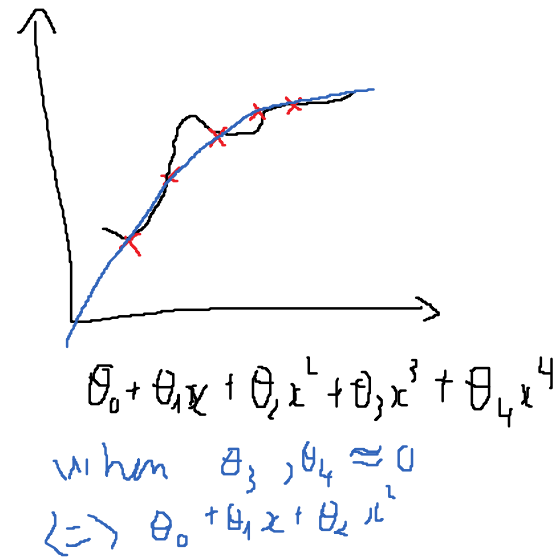
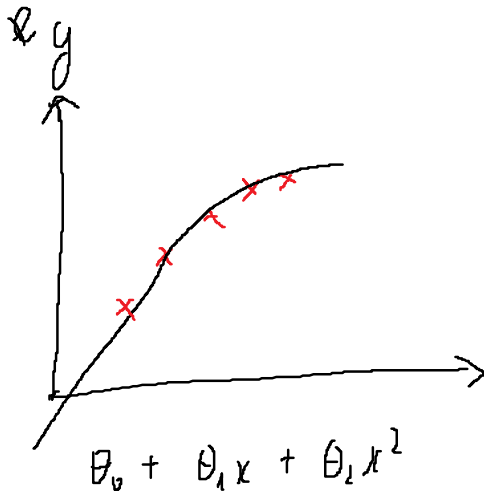
2 regularization

- giữ tất cả feature nhưng giảm độ lớn/giá trị của tham số
- tất cả các feature đóng góp và việc dự đoán

Regularization - cost function

Small value for parameter $\theta_0, \theta_1, \dots, \theta_n$ (Penalize)

- "simpler" hypothesis
- Less prone to overfitting



housing.

- Features x_1, x_2, \dots, x_{100}
- parameter: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

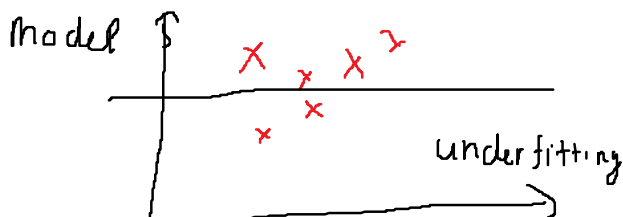
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

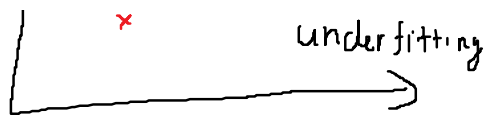
regularization
parameter

$$\sum_{j=1}^m \theta_j^2 = \|\theta\|_2^2$$

$$\min_{\theta} J(\theta)$$

new λ quá lớn $\Leftrightarrow \theta_1, \theta_2, \theta_3, \theta_4, \dots \approx 0$ (quá nhỏ)



A hand-drawn diagram consisting of a vertical line on the left and a horizontal line extending to the right, ending in an arrowhead. A red 'x' is marked in the upper-left quadrant. The word "underfitting" is written in a handwritten style above the horizontal line.

Regularizaion for Linear

Regularization

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \sum_{j=1}^n \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$

Gradient descent

Repeat {

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j = \theta_j - \alpha \left[\frac{1}{m} \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

\Downarrow

$$\frac{\partial}{\partial \theta_j} J(\theta) \quad j = (1, 2, \dots, n)$$

$$\theta_j = \underbrace{\theta_j \left(1 - \alpha \cdot \frac{\lambda}{m} \right)}_{\text{shrinking} < 1} - \alpha \cdot \frac{1}{m} \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Normal equation

$$X = \begin{bmatrix} (x^{(0)})^T \\ (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

$m \times (n+1)$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} \in \mathbb{R}^m$$

$$\theta = \left(X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \right)^{-1} X^T y$$

$$\theta = (X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix})^{-1} X^T y$$

Regularization for logistic

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^i \log(h_{\theta}(x^i)) + (1-y^i) \log(1-h_{\theta}(x^i)) + \frac{\lambda}{m} \sum_{j=1}^m \theta_j^2 \right]$$

Gradient descent

$$\begin{cases} \theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x_0^i \\ \theta_j = \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x_j^i + \frac{\lambda}{m} \theta_j \right] \end{cases}$$