

## 6. Regularization

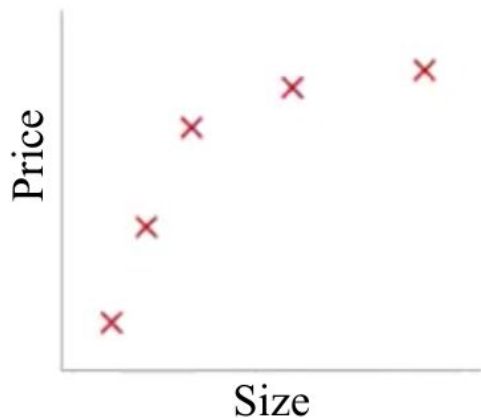
### 6.1 ปัญหา Overfitting

Krittameth Teachasrisaksakul

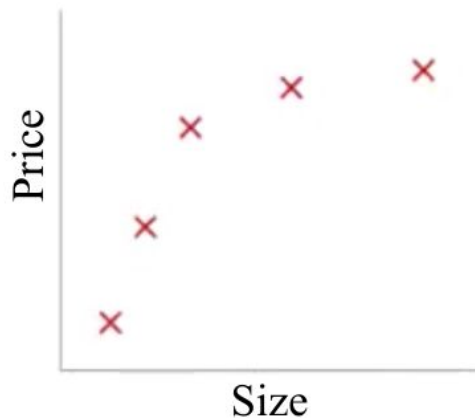
## Summary about supervised learning

- If you have continuous  $\mathcal{X}$  and continuous  $\mathcal{Y}$ , your first go-to model should be **linear regression**. Also, consider non-linear transformation of the inputs.
- If you have continuous  $\mathcal{X}$  and discrete  $\mathcal{Y}$  but don't know much about  $p(\mathbf{x} | y)$ , your first go-to model should be **logistic or softmax regression**, or may come up with a new **GLM** from scratch.
- If you have continuous  $\mathcal{X}$  and discrete  $\mathcal{Y}$  and know something about  $p(\mathbf{x} | y)$ , you should model the distribution accurately, as a **Gaussian (GDA)** or build a new **generative** model from scratch.
- If you have discrete  $\mathcal{X}$  and  $\mathcal{Y}$ , you should probably start with **naive Bayes** and build up from there.

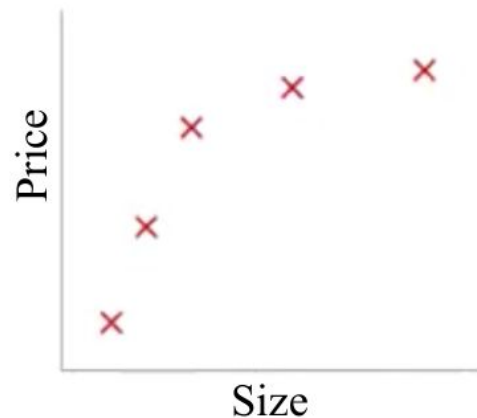
ตัวอย่าง: Linear regression (ราคาบ้าน / housing prices)



$$\theta_0 + \theta_1 x$$



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

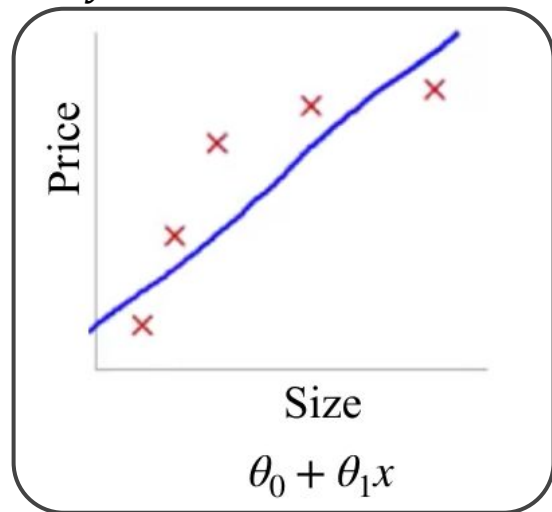


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

## ตัวอย่าง: Linear regression (ราคาบ้าน / housing prices)

(1) ใช้โมเดลเส้นตรง  $\rightarrow$  จุดข้อมูลไม่อยู่บนเส้นตรง  $\rightarrow$  เส้นตรง ไม่เหมาะกับ ข้อมูลชุดนี้

ทำนายค่า  $y$  จาก  $x \in \mathbb{R}$  :

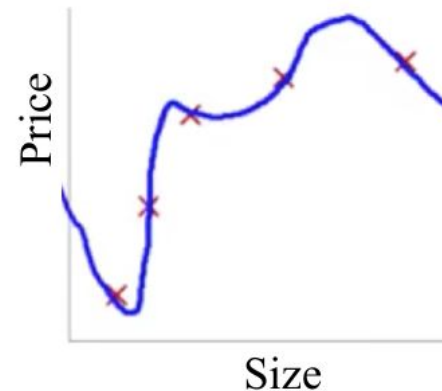
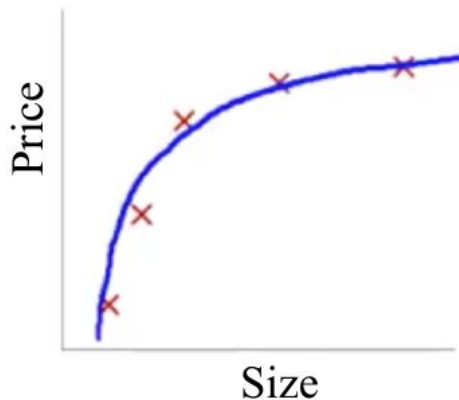


‘Underfit’

$\rightarrow$  hypothesis function  $h$  เข้ากับข้อมูลได้ไม่ดี

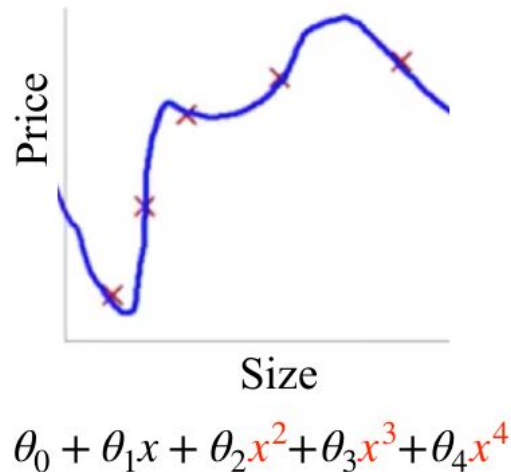
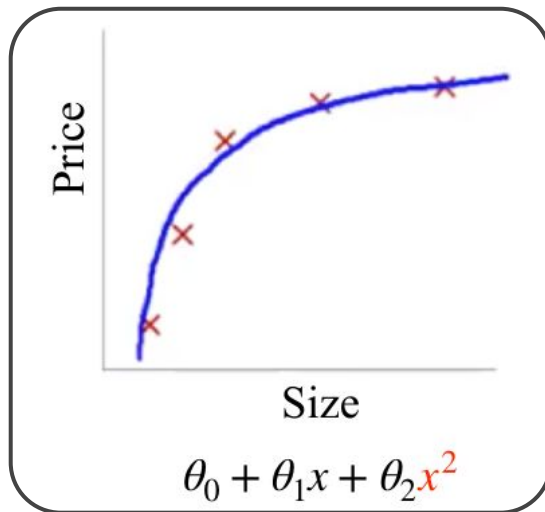
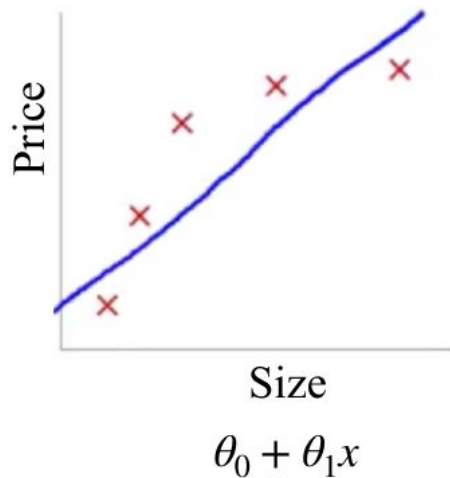
‘High bias’

$\rightarrow$  สาเหตุ: function เรียบง่าย (simple) เกินไป หรือ ใช้ feature จำนวนน้อยเกินไป



## ตัวอย่าง: Linear regression (ราคาบ้าน / housing prices)

ทำนายค่า  $y$  จาก  $x \in \mathbb{R}$  :

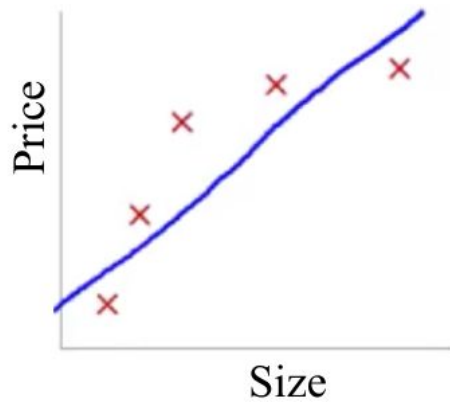


(2) เพิ่ม feature 1 ตัว ( $x^2$ )  $\rightarrow$  ได้โมเดล ที่เข้ากับข้อมูล มาก  
ขึ้นเล็กน้อย  $\rightarrow$  เหมือนว่า ยังมี feature มาก ยังได้ผลดีขึ้น

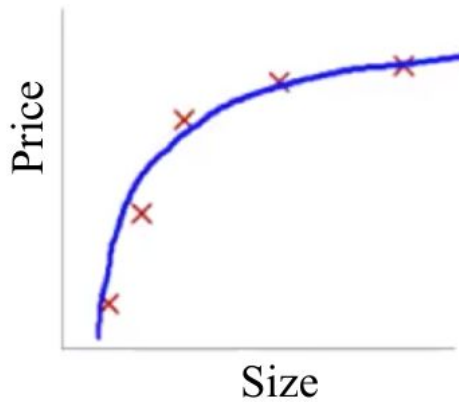
## ตัวอย่าง: Linear regression (ราคาบ้าน / housing prices)

ทำนายค่า  $y$  จาก  $x \in \mathbb{R}$  :

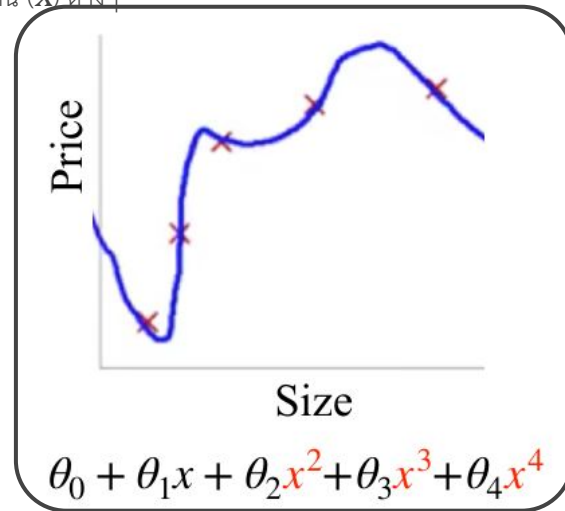
(3) ใช้ฟังก์ชันพหุนาม (5<sup>th</sup> order polynomial)  $\rightarrow$  โมเดล / เส้นโค้ง จะลากผ่านจุดข้อมูลทุกจุดอย่าง perfect  $\rightarrow$  แต่  
โมเดลนี้อาจไม่ได้เป็น ตัวทำนายที่ดีของ ราคา ( $y$ ) ของบ้าน ( $x$ ) ต่างๆ



$$\theta_0 + \theta_1 x$$



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



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

สาเหตุ: ใช้ feature จำนวนมากเกินไป หรือ function ที่ซับซ้อน จะมีส่วนโค้งและมุม ที่ไม่สอดคล้อง (unrelated) กับข้อมูล

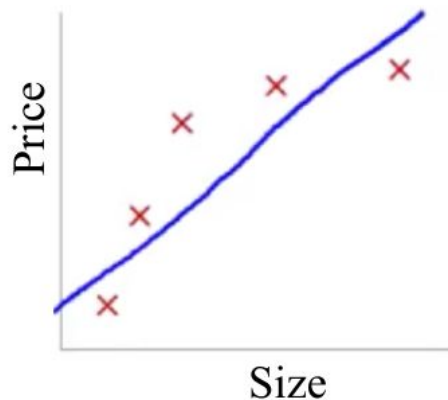
‘Overfit’

‘High variance’

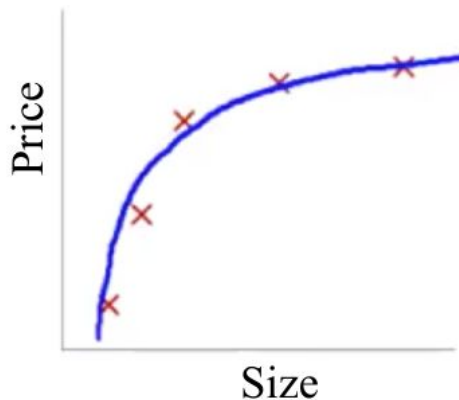
# Overfitting คืออะไร?

$$(J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0)$$

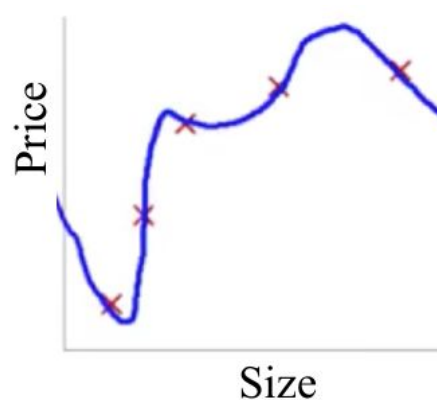
ตัวอย่าง: Linear regression (ราคาบ้าน / housing prices)



$$\theta_0 + \theta_1 x$$



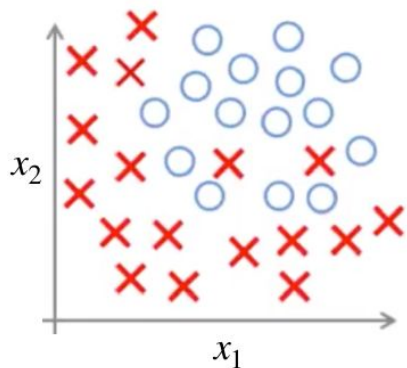
$$\theta_0 + \theta_1 x + \theta_2 x^2$$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

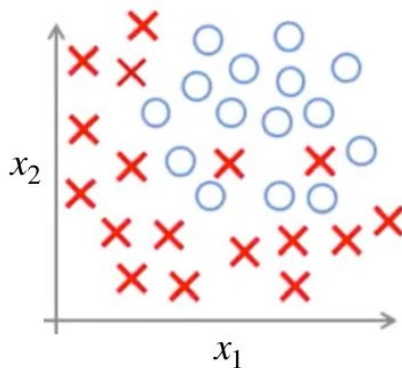
**Overfitting:** ผล  $\rightarrow$  hypothesis ที่เรียนรู้ อาจประมาณค่า ชุดข้อมูล training set ได้ดีมาก แต่ไม่สามารถประมาณค่า (generalize) ข้อมูลใหม่ ที่ไม่เคยเจอมาก่อน (ก็คือ ทำนายราคาของตัวอย่างใหม่)

## Overfitting ใน Logistic Regression

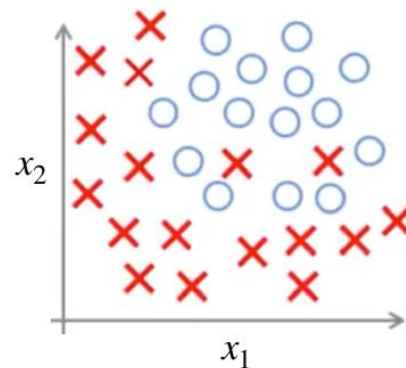


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

( $g$  is a sigmoid function)



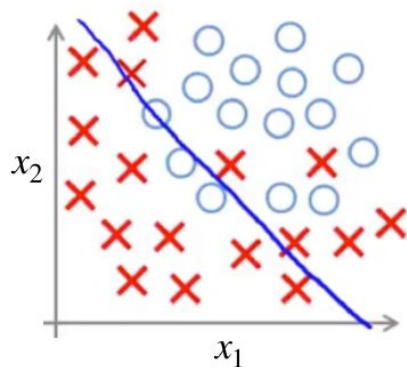
$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$



# Overfitting ใน Logistic Regression

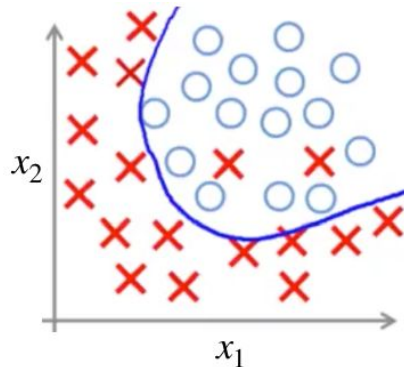


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

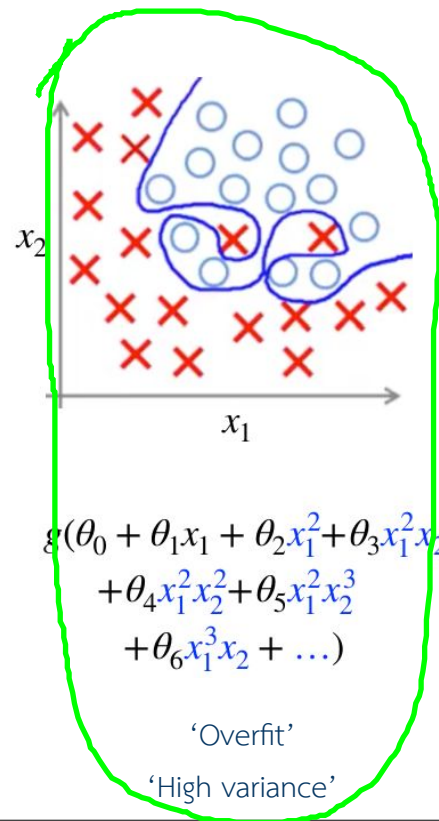
( $g$  is a sigmoid function)

‘Underfit’

‘High bias’



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

‘Overfit’

‘High variance’

# คำถาม

พิจารณาปัญหาการวินิจฉัยโรค (medical diagnosis) ที่แบ่งประเภทเนื้องอกเป็น ร้าย (malignant) หรือ ไม่ร้าย (benign) ถ้า hypothesis  $h_\theta(\mathbf{x})$  overfit ชุดข้อมูล training set หมายความว่าอะไร ?

- (i) มันทำนายตัวอย่างใน training set ได้ แม่นยำ และ generalize ได้ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้แม่นยำด้วย
- (ii) มันทำนายตัวอย่างใน training set ได้ **ไม่**แม่นยำ และ generalize ได้ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้แม่นยำด้วย
- (iii) มันทำนายตัวอย่างใน training set ได้ แม่นยำ และ generalize ได้ **ไม่**ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้ **ไม่**แม่นยำ
- (iv) มันทำนายตัวอย่างใน training set ได้ **ไม่**แม่นยำ และ generalize ได้ **ไม่**ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้ **ไม่**แม่นยำ

# คำถาม

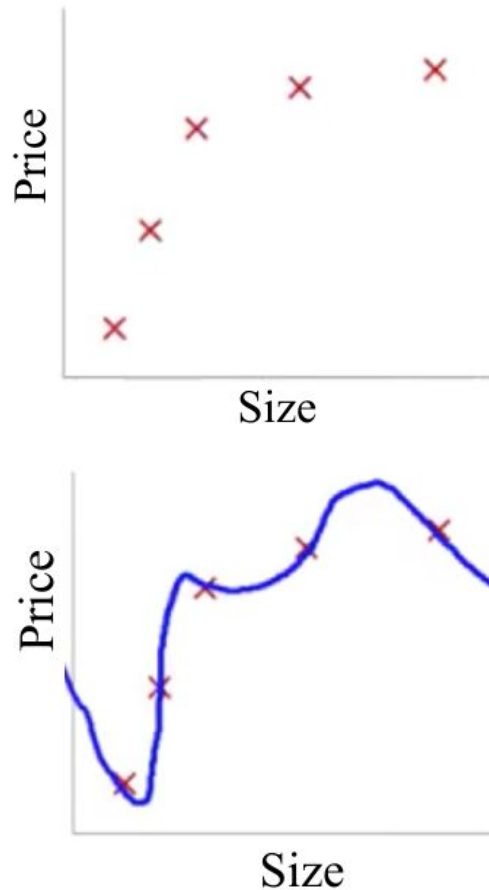
พิจารณาปัญหาการวินิจฉัยโรค (medical diagnosis) ที่แบ่งประเภทเนื้องอกเป็น ร้าย (malignant) หรือ ไม่ร้าย (benign) ถ้า hypothesis  $h_\theta(x)$  overfit ชุดข้อมูล training set หมายความว่าอะไร ?

- (i) มันทำนายตัวอย่างใน training set ได้ แม่นยำ และ generalize ได้ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้แม่นยำด้วย
- (ii) มันทำนายตัวอย่างใน training set ได้ **ไม่**แม่นยำ และ generalize ได้ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้แม่นยำด้วย
- (iii) มันทำนายตัวอย่างใน training set ได้ แม่นยำ และ generalize ได้ **ไม่**ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้ **ไม่**แม่นยำ
- (iv) มันทำนายตัวอย่างใน training set ได้ **ไม่**แม่นยำ และ generalize ได้ **ไม่**ดี ทำให้ทำนายตัวอย่างใหม่ ที่ไม่เคยเจอ ได้ **ไม่**แม่นยำ

# จัดการปัญหา Overfitting

ใช้ features จำนวนมากเกินไป อาจทำให้เกิด overfitting

- $X_1$  = ขนาดพื้นที่บ้าน
- $X_2$  = จำนวน ห้องนอน
- $X_3$  = จำนวน ชั้น
- $X_4$  = อายุบ้าน
- $X_5$  = รายได้เฉลี่ยของบริเวณใกล้เคียง
- $X_6$  = ขนาดพื้นที่ที่ห้องครัว
- $\vdots$
- $X_{100}$



# จัดการปัญหา Overfitting

## 1. ลดจำนวน features

- เลือก features ที่จะเก็บไว้ ด้วยมือ
- ใช้ algorithm ที่ทำการเลือก model (Model selection algorithm)

## 2. Regularization

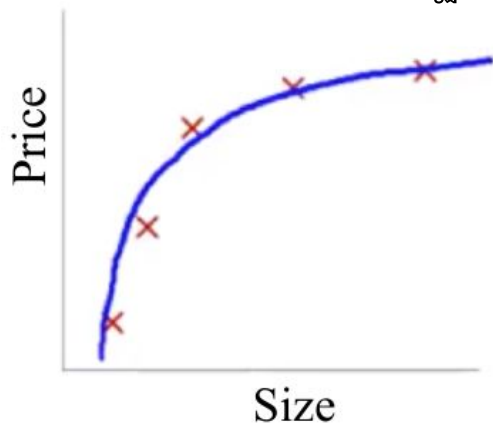
- เก็บ features ทั้งหมดไว้ แต่ลดขนาด (magnitude) หรือ ค่าของ parameters  $\theta_j$
- ทำงานได้ดีเมื่อเรามีจำนวน features มากๆ โดยที่ feature แต่ละตัวส่งผลให้ทำนาย  $y$  ได้

## 6. Regularization

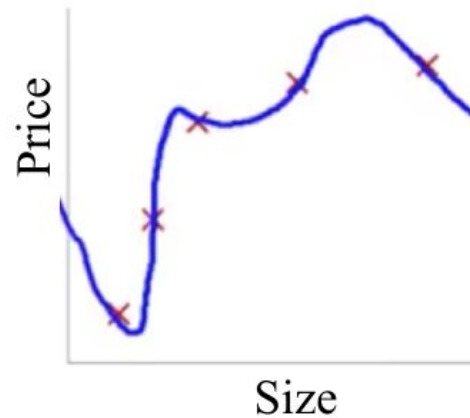
### 6.2 Cost Function

Krittameth Teachasrisaksakul

## Regularization : ความเข้าใจพื้นฐาน

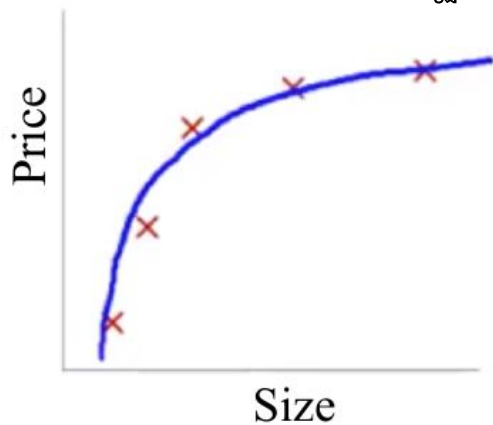


$$\theta_0 + \theta_1 x + \theta_2 x^2$$

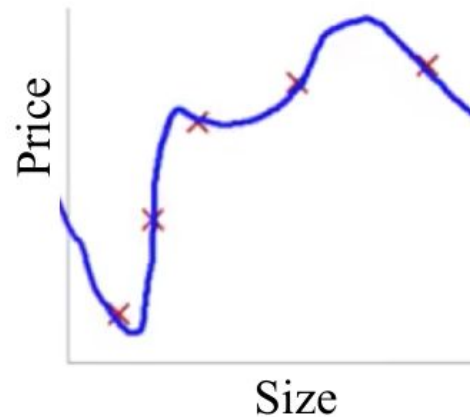


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

## Regularization : ความเข้าใจพื้นฐาน



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



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

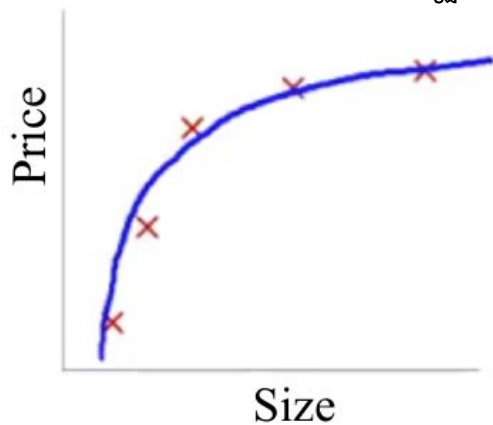
### Optimization objective:

(เป้าหมายของ optimization หรือ การปรับค่า parameter เพื่อหาค่าที่เหมาะสม)

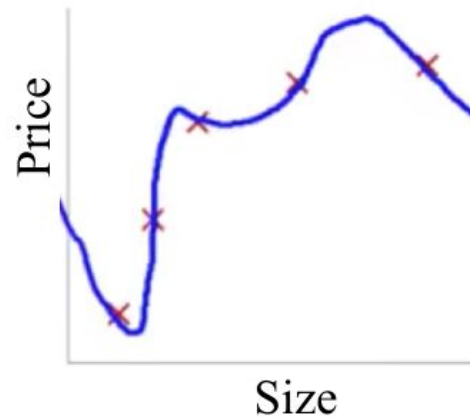
$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



## Regularization : ความเข้าใจพื้นฐาน



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



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

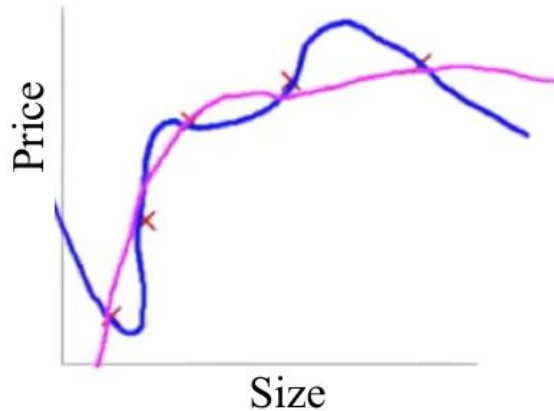
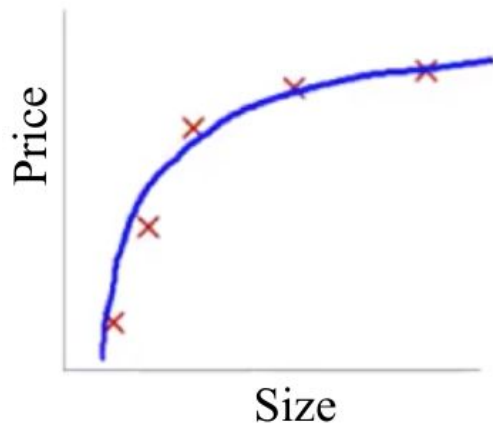
สมมติเรา penalize (ทำโทษ) และทำให้  $\theta_3, \theta_4$  น้อยมากๆ (ก็คือ ไม่สนับสนุนให้ใช้  $\theta_3, \theta_4$ )

ปรับ Optimization objective เป็น:

$$\therefore \theta_3 \approx 0, \theta_4 \approx 0$$

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2$$

# Regularization : ความเข้าใจพื้นฐาน



จากการปรับ optimization objective ทำให้ได้กราฟเส้นใหม่ (สีฟ้าอ่อน)

ถ้ามี hypothesis function ที่ overfitting เมื่อใช้กับข้อมูล

- (1) เพิ่ม 2 พจน์ท้าย
- เพิ่ม cost ของ  $\theta_3, \theta_4$
- ลดค่าน้ำหนัก (weight)  $\theta_3, \theta_4$  ของบางพจน์ใน function
- (2) ถ้าอยากให้ cost function เข้าใกล้ 0  $\rightarrow$  ต้องลดค่า  $\theta_3, \theta_4$  ให้ใกล้ 0

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

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

สมมติเรา penalize (ทำโทษ) และ ทำให้  $\theta_3, \theta_4$  น้อยมากๆ (ก็คือ ไม่สนับสนุนให้ใช้  $\theta_3, \theta_4$ )

ปรับ Optimization objective เป็น:

(2)  $\therefore \theta_3 \approx 0, \theta_4 \approx 0$

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2 \quad (1)$$

# Regularization : อธิบายแบบทางการ

ค่าที่น้อย ของ parameter  $\theta_1, \theta_2, \dots, \theta_n$  จะทำให้เกิด

- Hypothesis ที่ง่ายขึ้น (smooth มากขึ้น)
- มีแนวโน้ม overfitting น้อยลง

ตัวอย่าง การทำนายราคาบ้าน:

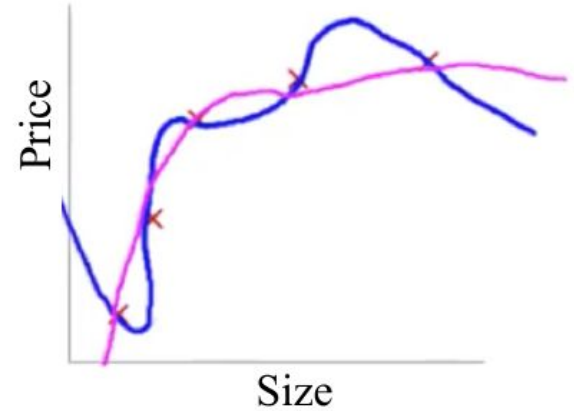
- Features:  $x_1, x_2, \dots, x_{100}$
- Parameters:  $\theta_1, \theta_2, \dots, \theta_{100}$

Cost function (ของ linear regression) :

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

เพิ่มพจน์นี้  $\rightarrow$  ทำให้ค่า output ของ hypothesis function smooth  $\rightarrow$  เพื่อลด overfitting

ถ้า  $\lambda$  มีค่ามากเกินไป  $\rightarrow$  อาจ smooth out function มากเกินไป  $\rightarrow$  ทำให้เกิด underfitting



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

# Regularization : อธิบายแบบทางการ

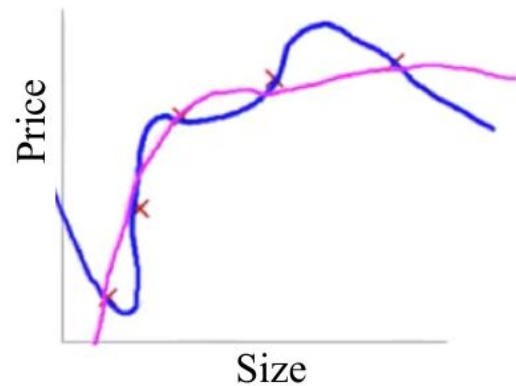
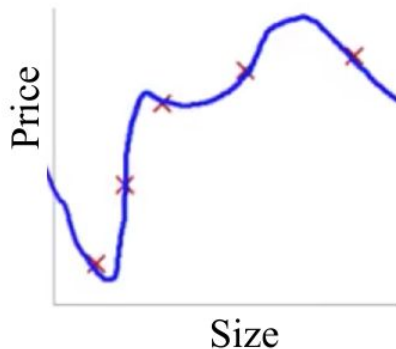
กำหนดว่า cost ของ parameter  $\theta$  ถูกเพิ่มขึ้นเท่าไร

**Regularized cost function:**

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

regularization parameter

**Goal:**  $\min_{\theta} J(\theta)$



## คำถาม

ใน regularized linear regression เราเลือกค่า  $\theta$  เพื่อให้  $J(\theta)$  น้อยที่สุด

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

ถ้า  $\lambda$  ถูกตั้งค่าเป็นค่าที่เยอะมากๆ (อาจมากเกินไปสำหรับปัญหาของเรา สมมติ  $\lambda = 10^{10}$ )

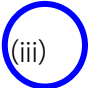
- (i) Algorithm ทำงานได้ดี; ตั้งค่า  $\lambda$  เป็นค่าเยอะมาก ไม่มีผลอะไร
- (ii) Algorithm ไม่สามารถแก้ปัญหา overfitting ได้
- (iii) Algorithm ทำให้เกิด underfitting (ไม่สามารถหา parameter ของ training set ได้)
- (iv) Gradient descent จะไม่ converge

## คำถาม

ใน regularized linear regression เราเลือกค่า  $\theta$  เพื่อให้  $J(\theta)$  น้อยที่สุด

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

ถ้า  $\lambda$  ถูกตั้งค่าเป็นค่าที่เยอะมากๆ (อาจมากเกินไปสำหรับปัญหาของเรา สมมติ  $\lambda = 10^{10}$ )

- (i) Algorithm ทำงานได้ดี; ตั้งค่า  $\lambda$  เป็นค่าเยอะมาก ไม่มีผลอะไร
- (ii) Algorithm ไม่สามารถแก้ปัญหา overfitting ได้
-  (iii) Algorithm ทำให้เกิด underfitting (ไม่สามารถหา parameter ของ training set ได้)
- (iv) Gradient descent จะไม่ converge

## 6. Regularization

### 6.3 Regularized Linear Regression

Krittameth Teachasrisaksakul

- สามารถใช้ regularization ได้กับ linear regression และ logistic regression
- พิจารณา linear regression ก่อน

## Cost Function (Recap / ทบทวน)

Cost function:

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

Goal:

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



## Gradient Descent (เดิม)

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (j = 0, 1, 2, \dots, n)$$

}

## Gradient Descent สำหรับ Regularized Linear Regression

- ปรับเปลี่ยน gradient descent function เพื่อแยก  $\theta_0$  ออกจาก parameter ตัวอื่นๆ
- เพราะเราไม่ต้องการ penalize  $\theta_0$  (ลางโทษ / ขัดขวางการใช้  $\theta_0$ )

Repeat {

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

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

}

( $j = 1, 2, \dots, n$ )

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

ทำ regularization

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

## Gradient Descent สำหรับ Regularized Linear Regression

Repeat {

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

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

}

( $j = 1, 2, \dots, n$ )

หลังการปรับปรุง กฎในการปรับค่า (update rule) ของ parameter  $\theta$  จะเป็น:

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

## Gradient Descent สำหรับ Regularized Linear Regression

$$\text{Repeat } \left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \boxed{\theta_j} - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \boxed{\theta_j} \right] \end{array} \right. \quad \left. \begin{array}{l} \xrightarrow{\frac{\partial}{\partial \theta_0} J(\theta)} \\ \xrightarrow{(j = \cancel{0}, 1, 2, \dots, n) \quad \frac{\partial}{\partial \theta_j} J(\theta)} \end{array} \right.$$

หลังการปรับปรุง กฎในการปรับค่า (update rule) ของ parameter  $\theta$  จะเป็น:

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

พจน์แรกมาจากการดึงตัว  
ร่วม  $\theta_j$

# Gradient Descent สำหรับ Regularized Linear Regression

Repeat {

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

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

}

$\frac{\partial}{\partial \theta_0} J(\theta)$   
 $\frac{\partial}{\partial \theta_j} J(\theta)$

(j = ~~0~~, 1, 2, ..., n)

หลังการปรับปรุง กฎในการปรับค่า (update rule) ของ parameter  $\theta$  จะเป็น:

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

(1) พจน์แรก  
น้อยกว่า 1 เสมอ

พจน์แรก คือ ลดค่า  $\theta_j$  ลง ทุกครั้งที่  
ปรับค่ามัน

(2) พจน์ที่ 2  
เหมือนกับสมการบน

## คำถาม

สมมติเรากำลังทำ gradient descent กับชุดข้อมูล training set ที่มีตัวอย่างจำนวน  $m > 0$  ตัวอย่าง โดยใช้ learning rate ที่ค่อนข้างน้อย  $\alpha > 0$  และ regularization parameter  $\lambda > 0$  พิจารณา update rule (กฎการปรับค่า parameter)

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

ข้อใดเป็นจริงเกี่ยวกับพจน์

$\left(1 - \alpha \frac{\lambda}{m}\right)$

$$1 - \alpha \frac{\lambda}{m} > 1$$

$$1 - \alpha \frac{\lambda}{m} = 1$$

$$1 - \alpha \frac{\lambda}{m} < 1$$

None of these

## คำถาม

สมมติเรากำลังทำ gradient descent กับชุดข้อมูล training set ที่มีตัวอย่างจำนวน  $m > 0$  ตัวอย่าง โดยใช้ learning rate ที่ค่อนข้างน้อย  $\alpha > 0$  และ regularization parameter  $\lambda > 0$  พิจารณา update rule (กฎการปรับค่า parameter)

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

ข้อใดเป็นจริงเกี่ยวกับพจน์

$$\left(1 - \alpha \frac{\lambda}{m}\right)$$

$$1 - \alpha \frac{\lambda}{m} > 1$$

$$1 - \alpha \frac{\lambda}{m} = 1$$

$$1 - \alpha \frac{\lambda}{m} < 1$$

None of these

## Normal Equation (Recap)

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

$m \times (n + 1)$

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

**Goal:**  $\min_{\theta} J(\theta)$



## Normal Equation (Recap)

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

$m \times (n + 1)$

**Goal:**  $\min_{\theta} J(\theta)$

**Solution:**  $\theta = (X^T X)^{-1} X^T y$

## Normal Equation (Recap)

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}_{m \times (n+1)}$$

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

ทำ regularization ด้วยวิธี Normal Equation: ใช้สมการเดิม + เพิ่ม 1 พจน์ในวงเล็บ (เพื่อทำ regularization)

**Goal:**  $\min_{\theta} J(\theta)$

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

$(n+1) \times (n+1)$

$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0$

## ปัญหา non-invertibility (การหา inverse matrix ไม่ได้)

- $m$  = จำนวน examples / ตัวอย่าง,  $n$  = จำนวน features
- ถ้า  $m < n$  → แล้ว  $(X^T X)$  เป็น non-invertible / singular (หา inverse matrix ไม่ได้)
- ถ้า  $m = n$  → แล้ว  $(X^T X)$  อาจเป็น be non-invertible

อย่างไรก็ตาม regularization สามารถจัดการกับปัญหา non-invertibility ได้

ถ้า  $\lambda > 0$  แล้ว :

$$\theta = \underbrace{\left( X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \right)}_{\text{invertible (หา inverse matrix ได้)}}^{-1} X^T y$$

$(n+1) \times (n+1)$

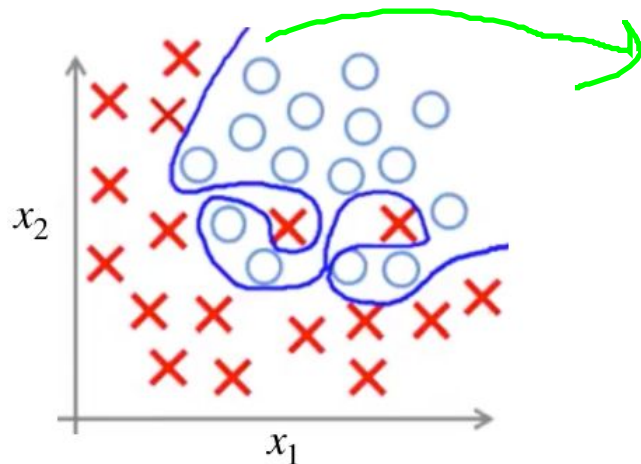
## 6. Regularization

### 6.4 Regularized Logistic Regression

Krittameth Teachasrisaksakul

เพื่อหลีกเลี่ยงปัญหา overfitting → ทำ regularization กับ logistic regression ด้วย วิธีคล้ายๆกับบทก่อนหน้านี้ (ทำ regularization กับ linear regression)

## Logistic Regression (Recap)

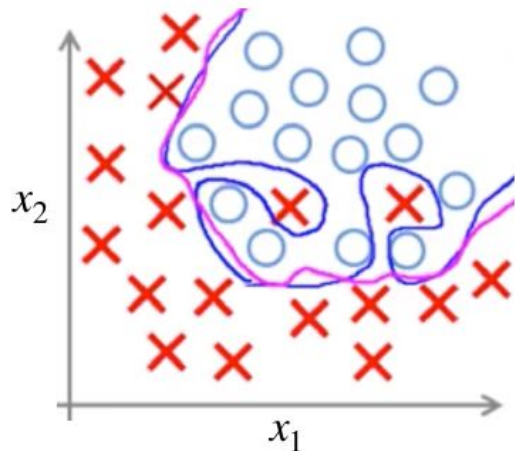


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 \\ + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 \\ + \theta_6 x_1^3 x_2 + \dots)$$

**Cost function:**

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

## Regularized Logistic Regression



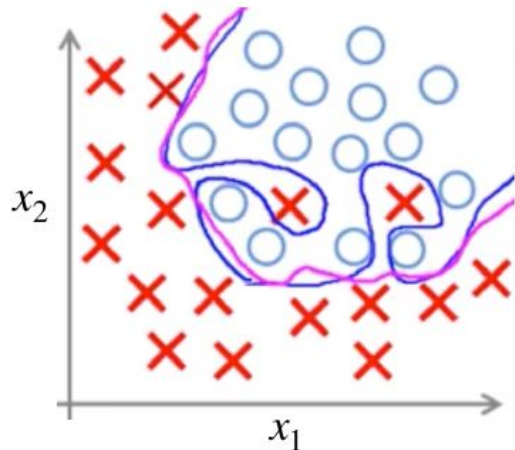
**Cost function:**

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

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

regularized function (เส้นสีชมพู) มีแนวโน้มที่จะ overfit น้อยกว่า non-regularized function (เส้นสีน้ำเงิน):

## Regularized Logistic Regression



**Cost function:**

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 \\ + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 \\ + \theta_6 x_1^3 x_2 + \dots)$$

Regularize โดยเพิ่มพจน์ท้ายสุด

- vector  $\theta$  มี index จาก 0 ถึง  $n$  (มีทั้งหมด  $n+1$  ตัว :  $\theta_0$  ถึง  $\theta_n$ )
- ผลรวมข้าม  $\theta_0$  โดยให้  $j$  เป็น 1 ถึง  $n$  (ข้าม 0)
- ก็คือ แยก พจน์ bias term  $\theta_0$  ออก

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

## Regularized Logistic Regression

Repeat {

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

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

( $j = \cancel{0}, 1, 2, \dots, n$ )

นี่ไม่ใช่ algorithm เดียวกับ gradient descent สำหรับ regularized linear regression เพราะ ...

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



## คำถาม

เมื่อใช้ regularized logistic regression วิธีใดเป็นวิธีที่ดีที่สุดที่จะสังเกตการณ์ว่า gradient descent ทำงานอย่างถูกต้อง ?

- (i) Plot  
เป็น function  $-\left[\frac{1}{m}\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)})\log(1 - h_{\theta}(x^{(i)}))\right]$
- (ii) Plot  
เป็น function  $-\left[\frac{1}{m}\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)})\log(1 - h_{\theta}(x^{(i)}))\right] + \frac{\lambda}{2m}\sum_{j=1}^n \theta_j^2$
- (iii) Plot  
เป็น function ของจำนวน iterations และทำให้แน่ใจว่ามันลดลง  
 $\sum_{j=1}^n \theta_j^2$

## คำถาม

เมื่อใช้ regularized logistic regression วิธีใดเป็นวิธีที่ดีที่สุดที่จะสังเกตการณ์ว่า gradient descent ทำงานอย่างถูกต้อง ?

(i) Plot  
เป็น function  $-\left[\frac{1}{m}\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right]$

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

(iii) Plot  
เป็น function ของจำนวน iterations และทำให้แน่ใจว่ามันลดลง  
 $\sum_{j=1}^n \theta_j^2$

# References

1. Andrew Ng, Machine Learning, Coursera.
2. Teeradaj Racharak, AI Practical Development Bootcamp.
3. What is Machine Learning?, <https://www.digitalskill.org/contents/5>