

Support Vector Machines (SVM)

Krittameth Teachasrisaksakul

บทนำ

จนถึงบัดนี้ เราได้เรียนเกี่ยวกับ machine learning models ที่ค่อนข้างง่ายที่จะวิเคราะห์ และหาค่าที่เหมาะสมที่สุดได้ (optimal) เมื่อ **assumptions (สมมติฐาน)** ของมันเป็นจริง

ตัวอย่าง เช่น Gaussian Discriminant Analysis (GDA) มีสมมติฐานว่า conditional distribution (การแจกแจงแบบมีเงื่อนไข) $p(x | y)$ เป็นแบบ multivariate Gaussian (การแจกแจงแบบปกติหลายตัวแปร)

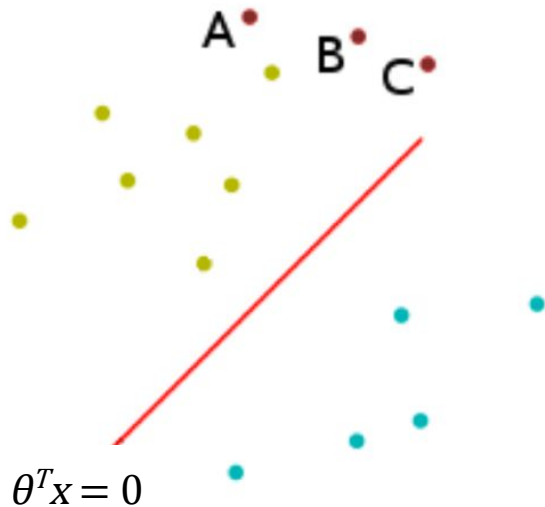
จะเกิดอะไรขึ้นเมื่อ assumptions ถูกละเมิด (ไม่เป็นจริง) ?

เราจะเรียนเกี่ยวกับ **support vector machines (SVMs)** ซึ่ง flexible มากกว่า และสามารถประยุกต์ใช้ได้อย่างกว้างขวางมากกว่าวิธีที่เรียนไปแล้ว

แม้ว่า deep neural networks ได้รับความสนใจมากที่สุดเมื่อเร็ว ๆ นี้ ผู้คนมากมายยังคงเชื่อว่า SVM เป็น supervised classifiers (ตัวแยกประเภทที่เรียนรู้แบบมีผู้สอน / จากตัวอย่างข้อมูล) ที่หาได้ง่าย / มีพร้อมใช้ (off-the-shelf) ที่ดีที่สุด

บทนำ

SVMs เกิดมาจากแนวคิด **maximum margin classification** (การแยกประเภทโดยทำให้ margin สูงที่สุด)

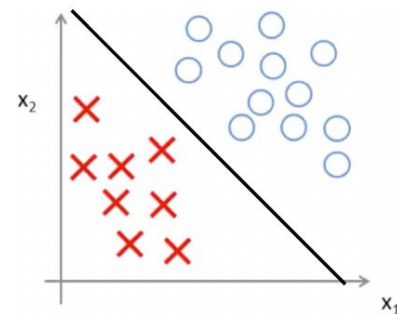
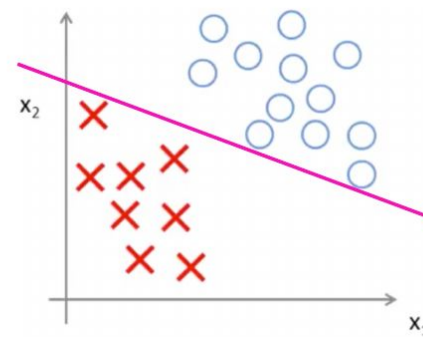
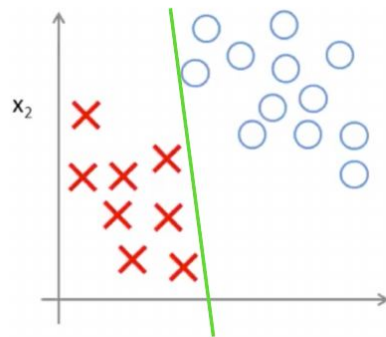
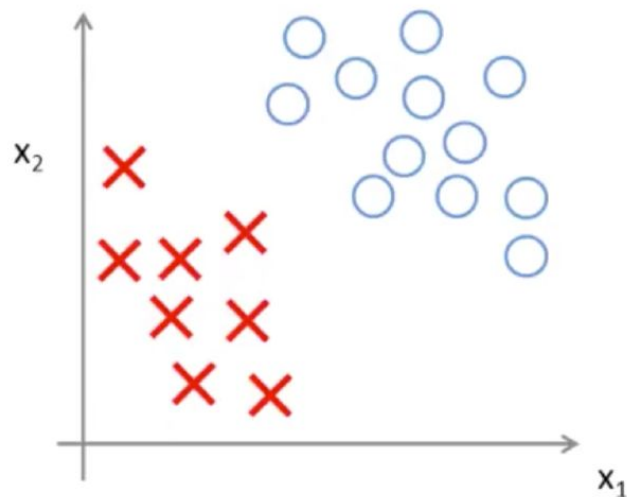


- สมมติ จุดสีเหลืองเป็นข้อมูล training data จาก class 1 และ จุดสีฟ้าเป็นข้อมูล training data จาก class 0
- $\theta^T x = 0$ เป็น ระนาบ hyperplane ที่แบ่งข้อมูลทั้ง 2 class หรือ decision boundary (ขอบเขตตัดสินใจ) ระหว่าง 2 class
- จุด A อยู่ไกลที่สุดจาก decision boundary เราสามารถทำนาย (อย่างมั่นใจ) ว่าจุด A เป็น class 1 แต่จุด C กำกวมมากกว่าว่าจะเป็น class ไດ

ข้อสังเกตนี้ ทำให้เกิดหลักการ maximizing the margin (ทำให้ margin สูงที่สุด) !

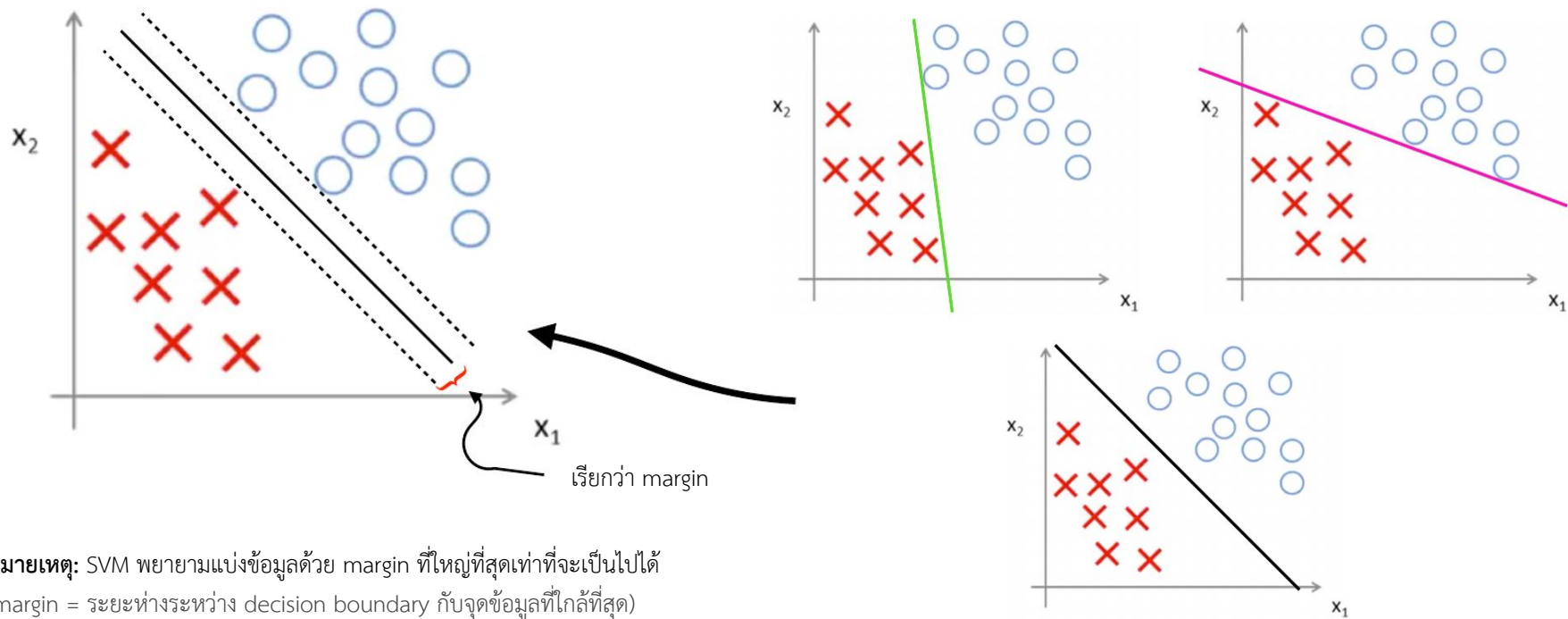
SVM Decision Boundary: Linearly Separable Case

(กรณีที่สามารถแบ่ง class ได้ด้วยขอบเขตตัดสินใจที่เป็นเชิงเส้น)



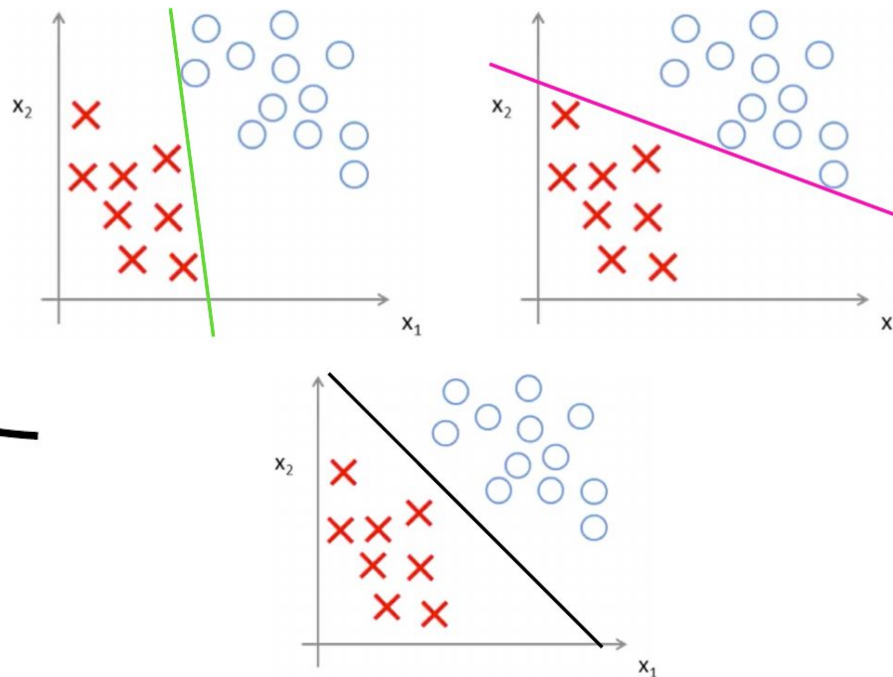
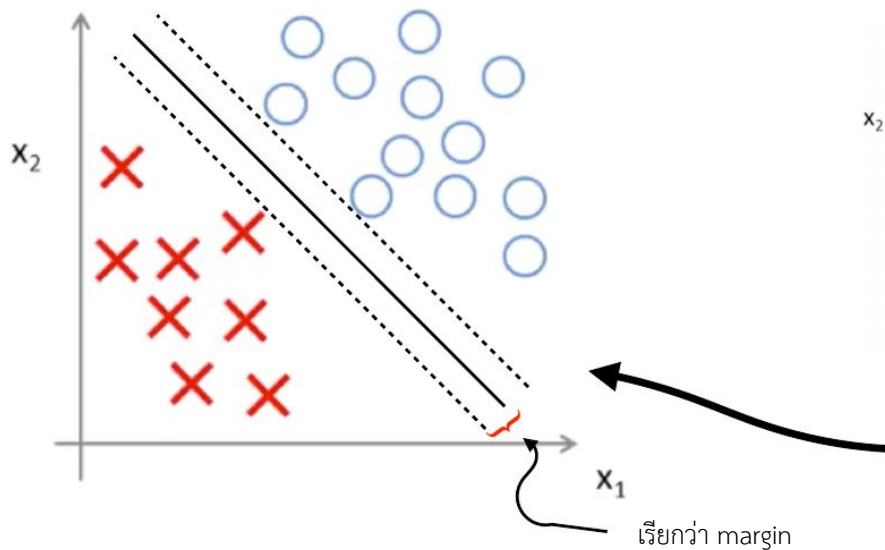
SVM Decision Boundary: Linearly Separable Case

(กรณีที่สามารถแบ่ง class ได้ด้วยขอบเขตตัดสินใจที่เป็นเชิงเส้น)



ใช้ SVM เป็น Large Margin Classifier

(ตัวแยกประเภทที่มี margin ขนาดใหญ่)



หมายเหตุ: SVM พยายามแบ่งข้อมูลด้วย margin ที่ใหญ่ที่สุดเท่าที่จะเป็นไปได้
(margin = ระยะห่างระหว่าง decision boundary กับจุดข้อมูลที่ใกล้ที่สุด)

Support Vector Machines (SVM)

Optimization Objective

Krittameth Teachasrisaksakul

อีกมุมมองหนึ่งของ Logistic Regression

บททวน: Logistic regression เราทำแบบจำลอง (model) ของ $p(y = 1 \mid x; \theta)$ ด้วย

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

กฎการแยกประเภทแบบ logistic (logistic classification rule) คือ:

$$y^{\text{pred}}(x) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5 \\ 0, & \text{otherwise} \end{cases}$$

เป้าหมายของเรา ควรเป็น การหา θ ที่ทำให้

- $\theta^T x^{(i)} \gg 0$ สำหรับ i ทุกตัวที่มี $y^{(i)} = 1$ และ
- $\theta^T x^{(i)} \ll 0$ สำหรับ i ทุกตัวที่มี $y^{(i)} = 0$

อีกมุมมองหนึ่งของ Logistic Regression

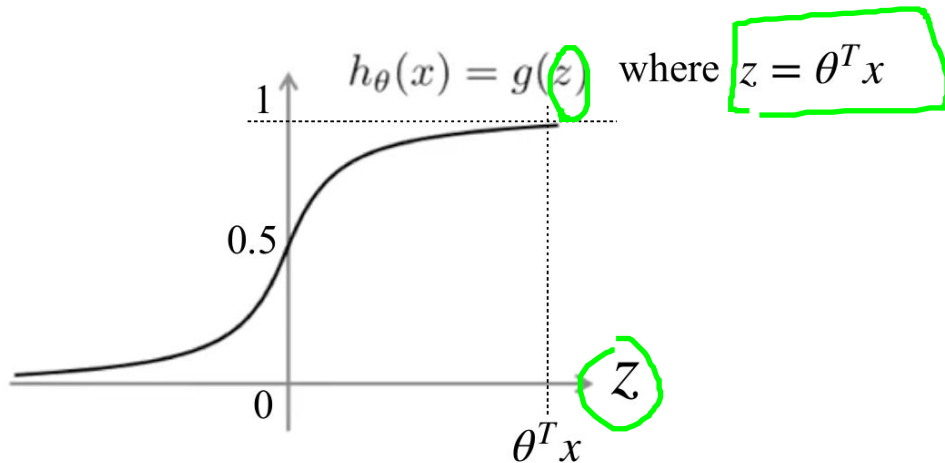
Hypothesis function:

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

ปรับ logistic regression เพื่อให้ได้ SVM

ลองคิดเกี่ยวกับ สิ่งที่เราอยากให้ logistic regression ทำ:

- ถ้า $y = 1$ แล้วเราอยากให้ $h_{\theta}(x) \approx 1$ ก็คือ $\theta^T x \gg 0$ หรือ $z \gg 0$ (มากกว่า 0 มากๆ)
- ถ้า $y = 0$ แล้วเราอยากให้ $h_{\theta}(x) \approx 0$ ก็คือ $\theta^T x \ll 0$ หรือ $z \ll 0$



sigmoid activation function

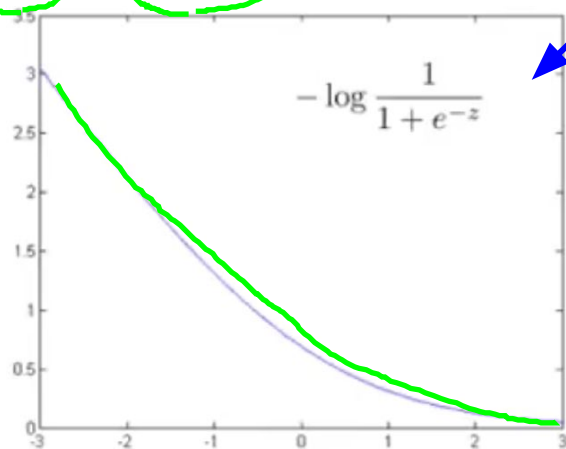
อีกมุมมองหนึ่งของ Logistic Regression

Cost ของ example (ตัวอย่าง)

$$\rightarrow -(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x)))$$

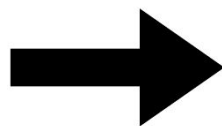
$$\Leftrightarrow -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$

กรณี 1 $(y = 1 \Rightarrow \theta^T x \gg 0)$:



แค่พจน์แรกที่มีผล เพราะ พจน์หลัง = 0

in SVM



Plot: function ของ Z vs. Z

ถ้า Z หรือ $\theta^T x$ มาก

\rightarrow cost function น้อย

เมื่อ logistic regression เจอ positive example ($y=1$)

มันจะทำให้ $z = \theta^T x$ มีค่าใหญ่ ($\gg 0$) เพราะจะทำให้พจน์ใน cost function มีค่าน้อย

อีกมุมมองหนึ่งของ Logistic Regression

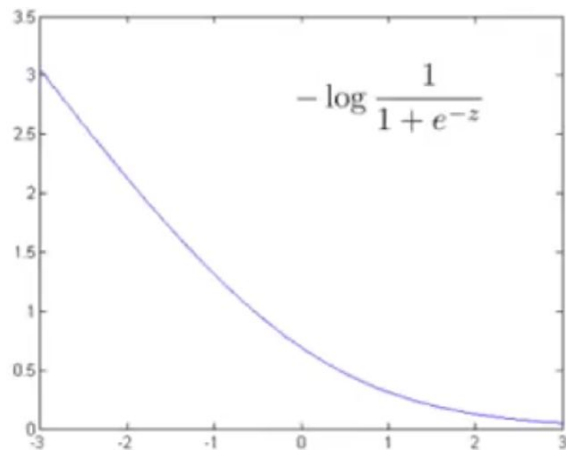
Cost ของ example (ตัวอย่าง)

$$-(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x)))$$

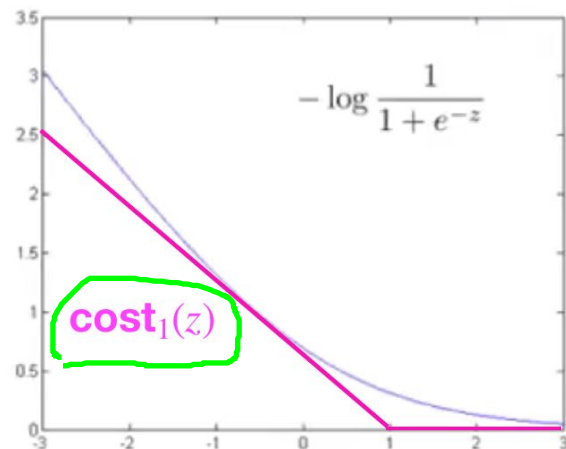
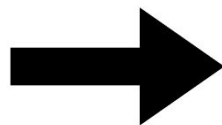
$$\Leftrightarrow -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log(1 - \frac{1}{1 + e^{-\theta^T x}})$$

$$z = \theta^T x$$

กรณี 1 ($y = 1 \Rightarrow \theta^T x \gg 0$):



in SVM

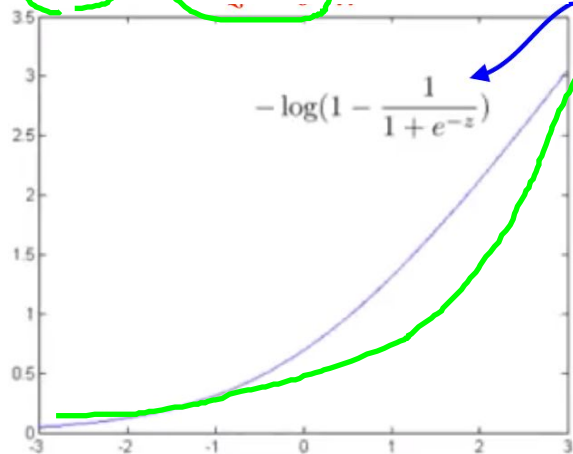


อีกมุมมองหนึ่งของ Logistic Regression

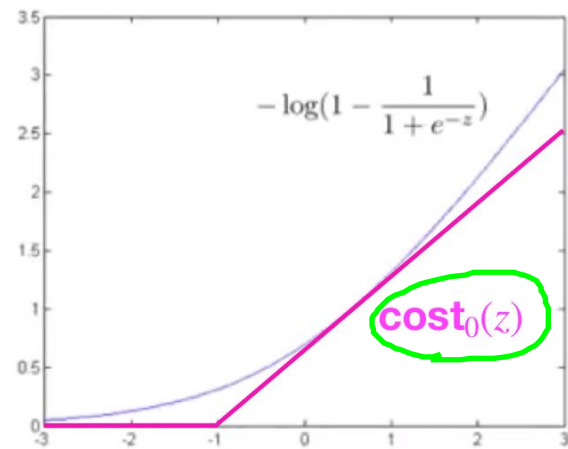
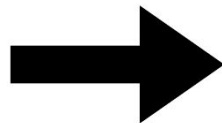
Cost ของ example (ตัวอย่าง)

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

กรณี 2 ($y = 0 \Rightarrow \theta^T x \ll 0$):



in SVM



Support Vector Machine

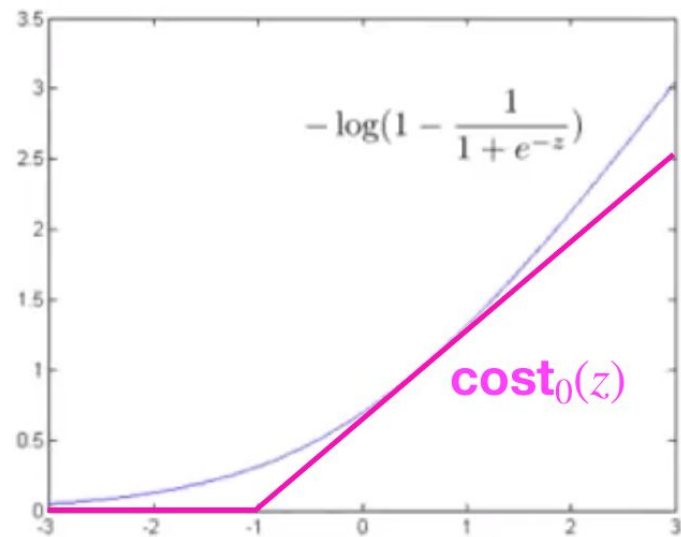
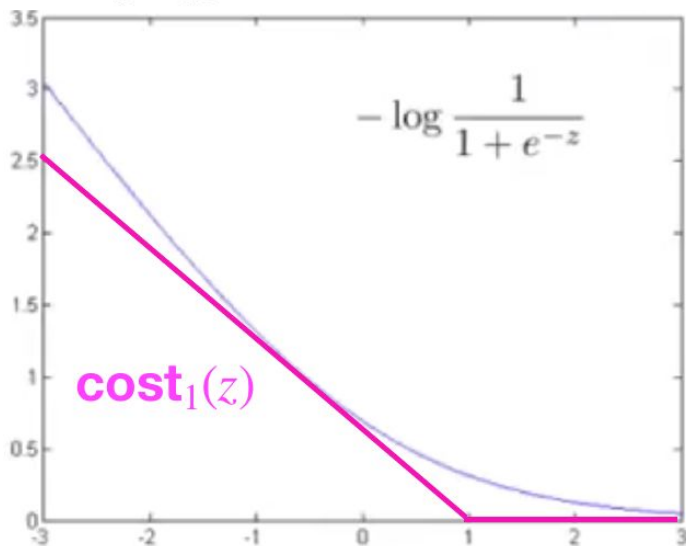
Logistic regression:

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

Support Vector Machine

Logistic regression:

$$J(\theta) = \min_{\theta} \frac{1}{m} \left[\underbrace{\sum_{i=1}^m y^{(i)} \left(-\log h_{\theta}(x^{(i)}) \right)}_{\text{cost}_1(\theta^T x^{(i)})} + \underbrace{\sum_{i=1}^m (1 - y^{(i)}) \left(-\log(1 - h_{\theta}(x^{(i)})) \right)}_{\text{cost}_0(\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$



Support Vector Machine

Logistic regression:

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



Support vector machine:

$$J(\theta) = \min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

เขียนสมการด้านบนใหม่ เพื่อให้มันสอดคล้องกับ convention ของ SVM

Support Vector Machine

Cost function:

1. เอา 'm' ออกจากสมการ

$$J(\theta) = \min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

ทำไมนี่จึงสมเหตุสมผล ?

ลองคิดเกี่ยวกับ:

- $\min_u \left((u + 5)^2 + 1 \right) \Rightarrow 2(u - 5) \cdot 1 = 0 \Leftrightarrow u = 5$
- $\min_u 10 \left((u - 5)^2 + 1 \right) \Rightarrow 2 \cdot 10 \cdot (u - 5) \cdot 1 = 0 \Rightarrow u = 5$

$$\frac{d}{du} \left((u - 5)^2 + 1 \right) = 2(u - 5) = 0$$

Support Vector Machine

Cost function:

2. จัดระเบียบสมการเล็กน้อย

$$J(\theta) = \min_{\theta} \frac{1}{n} \left[\underbrace{\sum_{i=1}^m y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)})}_A \right] + \underbrace{\frac{\lambda}{2n} \sum_{j=1}^n \theta_j^2}_B$$

ก็คือ ใน logistic regression เราเขียน: $A + \lambda B$

ในขณะที่ ใน SVM เราเขียน $CA + B$

สังเกตว่า C ไม่จำเป็นต้องเป็น $1 / \lambda$

Support Vector Machine

Logistic regression:

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



Support vector machine:

$$J(\theta) = \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Question

พิจารณาปัญหา minimization ดังนี้

$$J(\theta) = \min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$$J(\theta) = \min_{\theta} C \left[\sum_{i=1}^m y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

ปัญหา optimization 2 ปัญหานี้ จะให้ค่า θ เดียวกัน (ก็คือ ค่า θ ค่าเดียวกันเป็นคำตอบที่เหมาะสม (optimal solution) ของปัญหาทั้ง 2 ปัญหา

- | | |
|-----------------------|----------------------|
| (i) $C = \lambda$ | (ii) $C = -\lambda$ |
| (iii) $C = 1/\lambda$ | (iv) $C = 2/\lambda$ |

Support Vector Machine

ไม่เหมือนกับ logistic regression : SVM ไม่มีการตีความในรูปของความน่าจะเป็น (probabilistic interpretation)

Cost function:

Hypothesis function..

$$J(\theta) = \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$h_{\theta}(x) = \begin{cases} 1, & \text{if } \theta^T x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

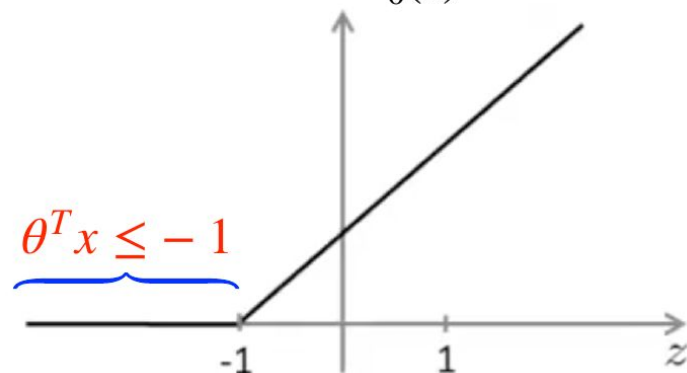
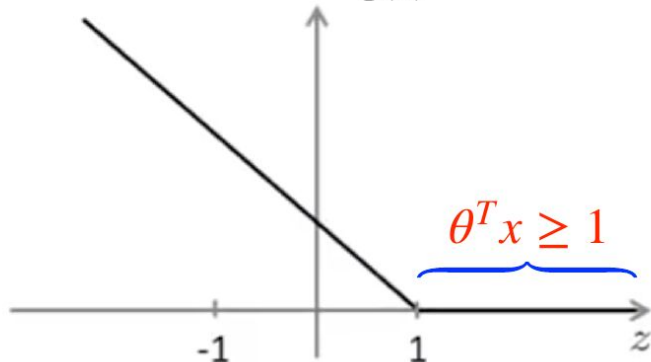
Support Vector Machines (SVM)

Large Margin Classifier (ตัวแยกประเภทที่มี margin ขนาดใหญ่)

Krittameth Teachasrisaksakul

ความเข้าใจพื้นฐาน

Cost function:
$$J(\theta) = \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

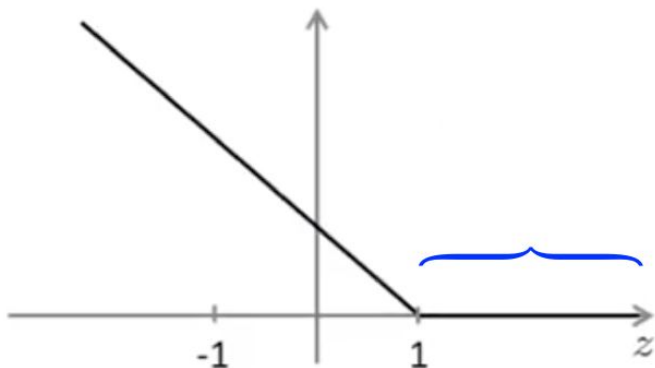


- ถ้า $y = 1$ เราอยากให้ $\theta^T x \geq 1$ (ไม่ใช่แค่ ≥ 0)
- ถ้า $y = 0$ เราอยากให้ $\theta^T x \leq -1$ (ไม่ใช่แค่ < 0)
- ต่อไป ถ้าเราอยากตั้งค่า C เป็นค่าที่สูงมาก เช่น 100,000 ?

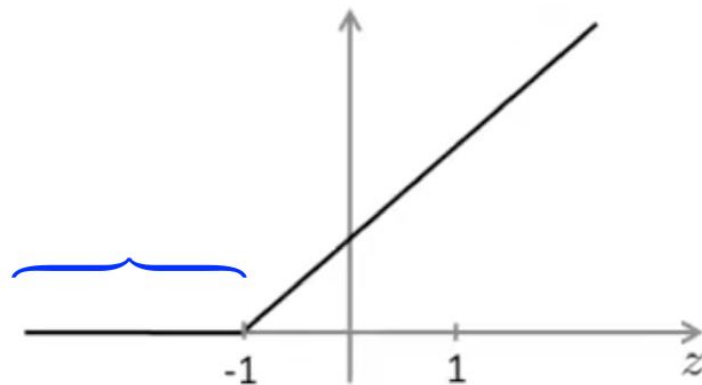
Decision Boundary ของ SVM

$$\min_{\theta} \underbrace{C \sum_{i=1}^m}_{10^5} \underbrace{\left[y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right]}_{\approx 0} + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

เมื่อใดก็ตามที่ $y^{(i)} = 1$ จะได้ $\theta^T x^{(i)} \geq 1$



เมื่อใดก็ตามที่ $y^{(i)} = 0$ จะได้ $\theta^T x^{(i)} \leq -1$



Decision Boundary ของ SVM

$$\min_{\theta} \underbrace{C \sum_{i=1}^m}_{10^5} \underbrace{\left[y^{(i)} \mathbf{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T x^{(i)}) \right]}_{\approx 0} + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

เมื่อใดก็ตามที่ $y^{(i)} = 1$ จะได้ $\theta^T x^{(i)} \geq 1$

เมื่อใดก็ตามที่ $y^{(i)} = 0$ จะได้ $\theta^T x^{(i)} \leq -1$

$$\min_{\theta} C \cdot 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

subject to $\theta^T x^{(i)} \geq 1$ if $y^{(i)} = 1$
- $\theta^T x^{(i)} \leq -1$ if $y^{(i)} = 0$

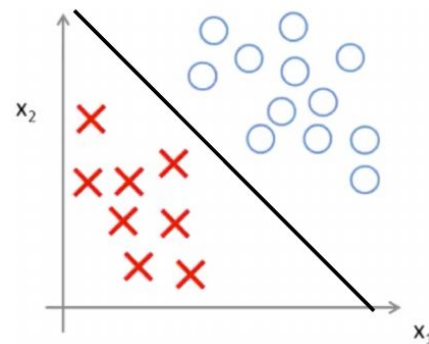
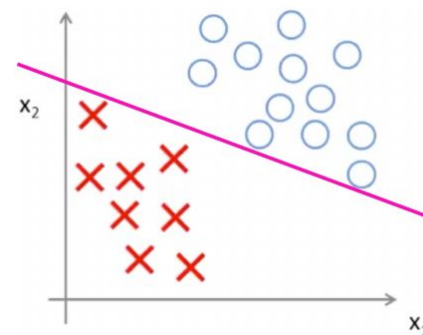
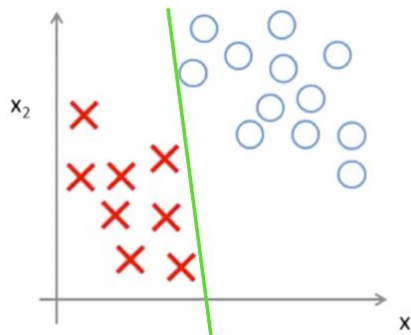
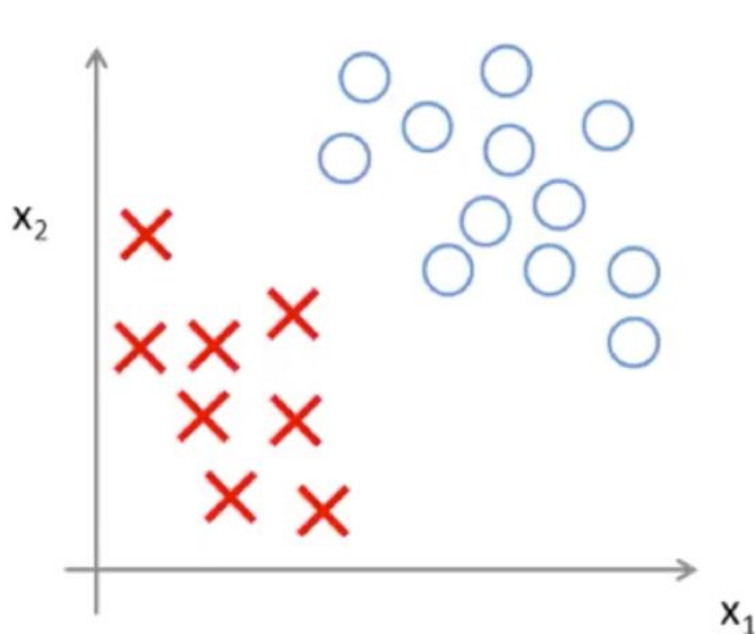
เมื่อแก้ปัญหา optimization นี้

เราจะได้ decision boundary ที่น่าสนใจ !

- ก็คือ เมื่อ C สูงมาก \rightarrow SVM จะอ่อนไหวต่อ outlier (ค่า/ข้อมูลผิดปกติ)
- ลดค่า C จะทำให้ \rightarrow SVM อ่อนไหว น้อยลง

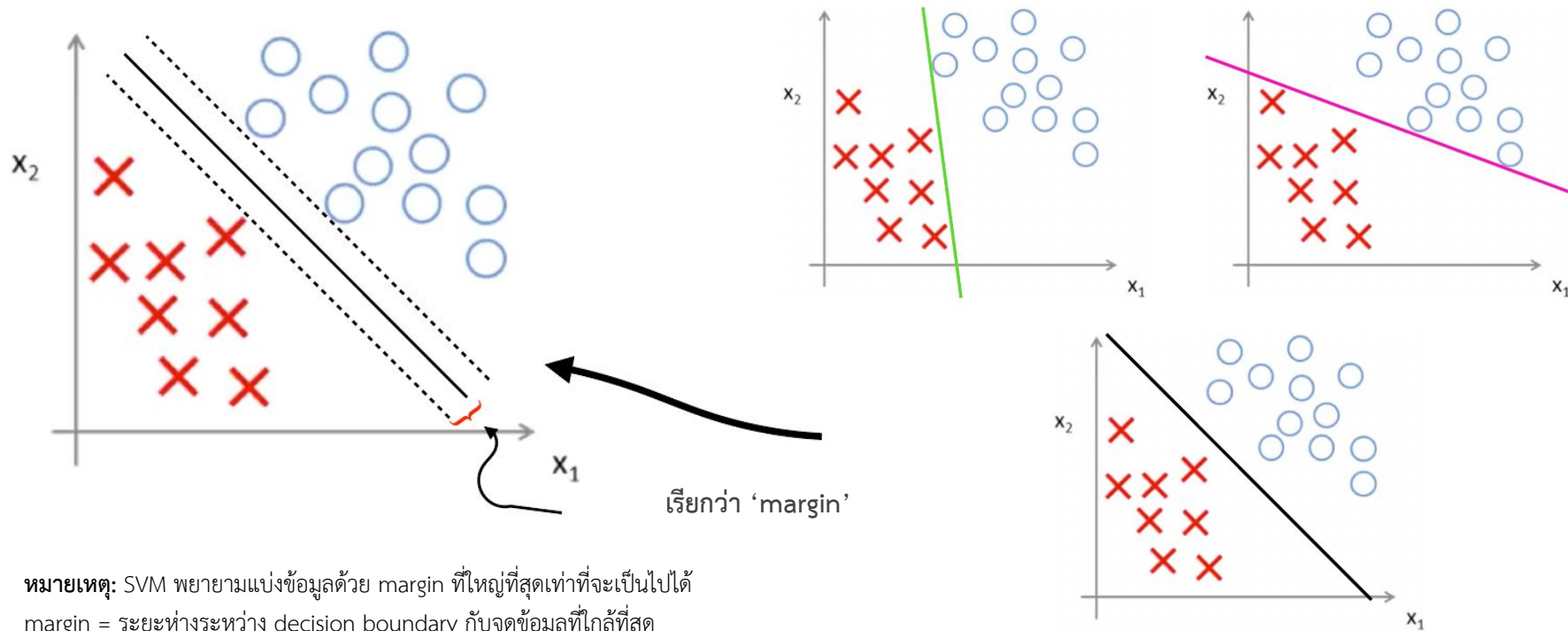
SVM Decision Boundary: Linearly Separable Case

กรณีที่สามารถแบ่ง class ได้ด้วยขอบเขตตัดสินใจที่เป็นเชิงเส้น



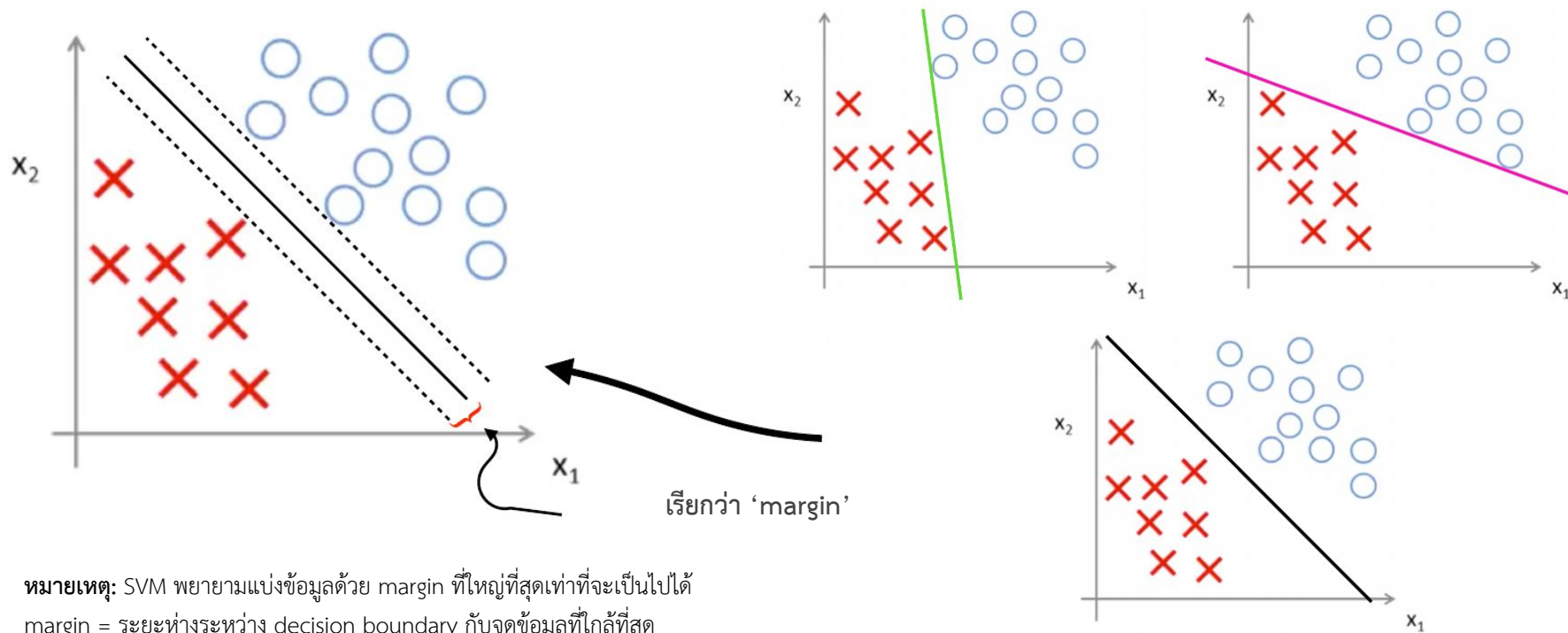
SVM Decision Boundary: Linearly Separable Case

กรณีที่สามารถแบ่ง class ได้ด้วยขอบเขตตัดสินใจที่เป็นเชิงเส้น



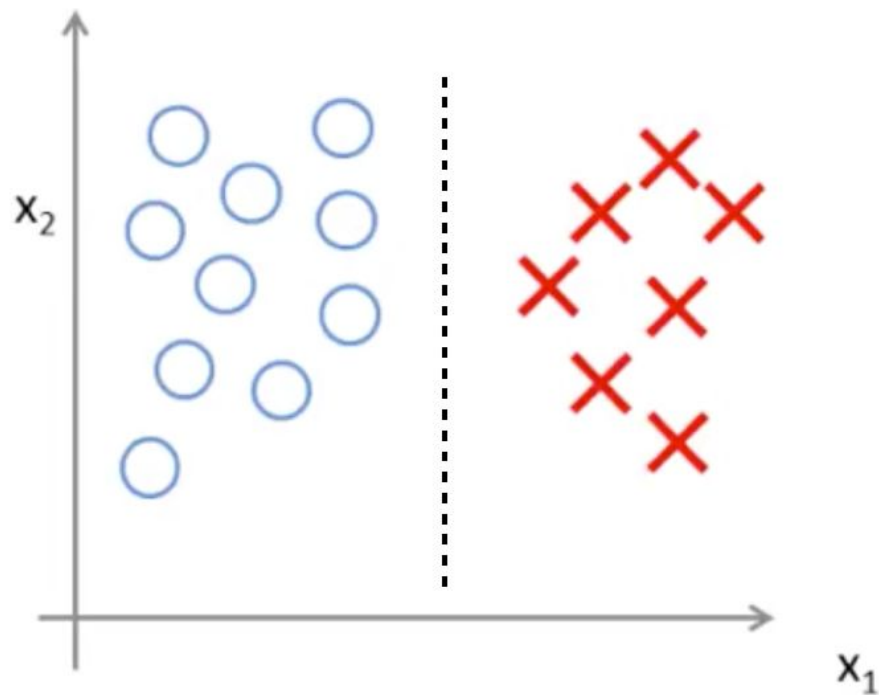
หมายเหตุ: SVM พยายามแบ่งข้อมูลด้วย margin ที่ใหญ่ที่สุดเท่าที่จะเป็นไปได้
margin = ระยะห่างระหว่าง decision boundary กับจุดข้อมูลที่ใกล้ที่สุด

SVM เป็น Large Margin Classifier (ตัวแยกประเภทที่มี margin ขนาดใหญ่)



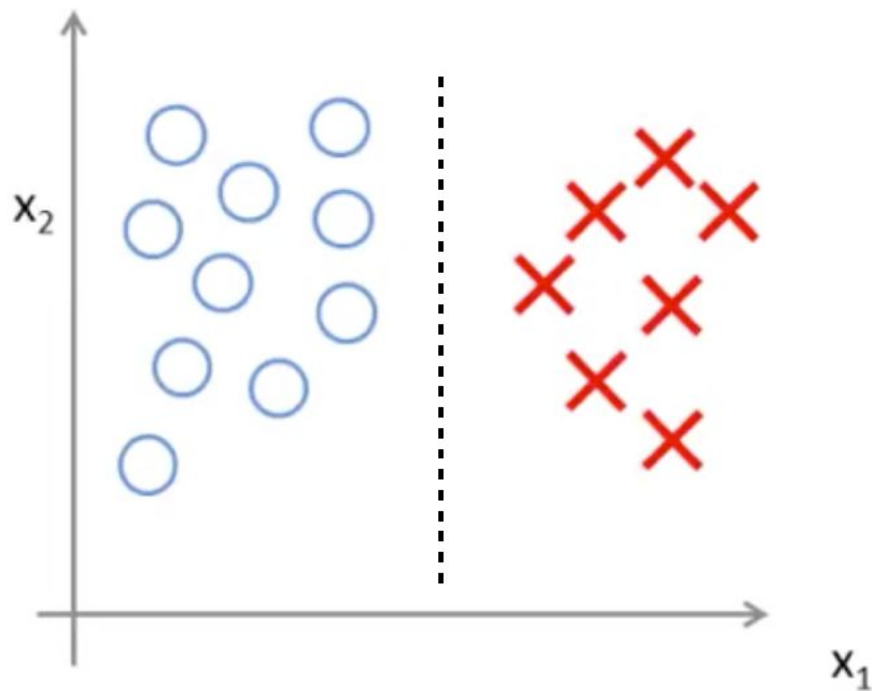
หมายเหตุ: SVM พยายามแบ่งข้อมูลด้วย margin ที่ใหญ่ที่สุดเท่าที่จะเป็นไปได้
margin = ระยะห่างระหว่าง decision boundary กับจุดข้อมูลที่ใกล้ที่สุด

Large Margin Classifier เมื่อมี Outliers (ค่า/ข้อมูลผิดปกติ)



SVM พยายามหา decision boundary ที่มีระยะห่าง (*margin*) จาก
ตัวอย่าง / sample สูงสุด

Large Margin Classifier เมื่อมี Outliers (ค่า/ข้อมูลผิดปกติ)



ลองเพิ่ม outlier



เมื่อ C สูงมาก : SVM จะอ่อนไหวต่อ outlier

เมื่อลดค่า C จะทำให้ SVM อ่อนไหวน้อยลง

ถ้าข้อมูลไม่ 'linearly sample' (ไม่สามารถแยกชนิด / class ด้วย decision boundary ที่เป็นเส้นตรง)

แล้ว SVM ยังคงทำงานได้ถูกต้อง

Question

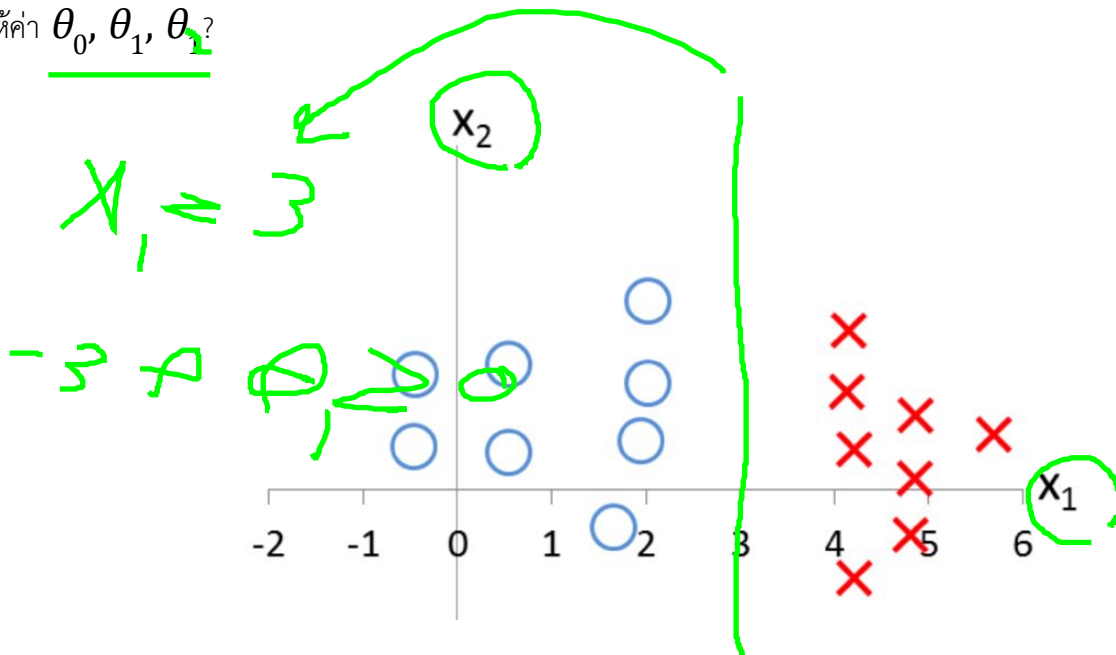
พิจารณาชุดข้อมูล training set เมื่อ 'x' แทน ตัวอย่างที่ ($y = 1$) และ 'o' แทน ตัวอย่าง negative ที่ ($y = 0$) สมมติฝึก/สร้าง (train) SVM (ที่ทำนาย 1 เมื่อ $\theta_0 + \theta_1 x_1 + \theta_2 x_2 \geq 0$) SVM อาจจะให้ค่า $\theta_0, \theta_1, \theta_2$?

(i) $\theta_0 = 3, \theta_1 = 1, \theta_2 = 0$

(ii) $\theta_0 = -3, \theta_1 = 1, \theta_2 = 0$

(iii) $\theta_0 = 3, \theta_1 = 0, \theta_2 = 1$

(iv) $\theta_0 = -3, \theta_1 = 0, \theta_2 = 1$



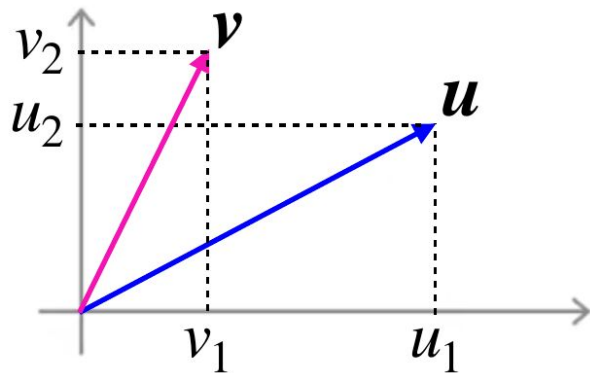
Support Vector Machines (SVM)

คณิตศาสตร์เบื้องหลัง

Large Margin Classification

Krittameth Teachasrisaksakul

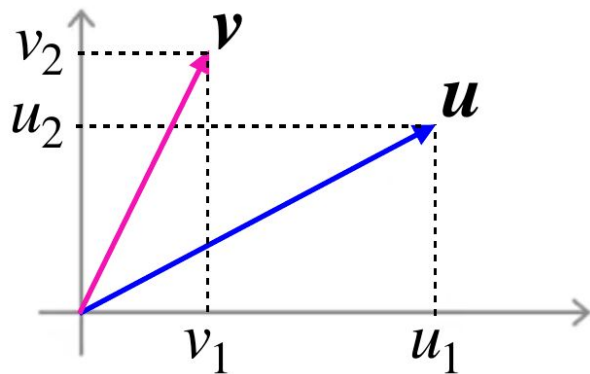
Preliminary: Vector Inner Product



ถ้ามี $u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$ กับ $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$ ให้คำนวณ $u^T v = ?$

($u^T v$ เรียกว่า 'vector inner product')

Preliminary: Vector Inner Product



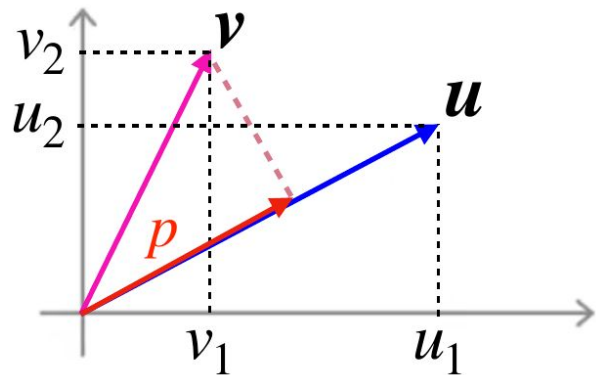
ถ้ามี $u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$ กับ $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$ ให้คำนวณ $u^T v = ?$
($u^T v$ เรียกว่า 'vector inner product')

สามารถ quantify (วัดปริมาณ) vector โดยหาค่า euclidean length หรือ norm ของมัน

$$\|u\| = \sqrt{u_1^2 + u_2^2} \quad (\text{by Pythagoras theorem})$$

\in
 \mathbb{R}

Preliminary: Vector Inner Product



ถ้ามี $u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$ กับ $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$ ให้คำนวณ $u^T v = ?$
($u^T v$ เรียกว่า 'vector inner product')

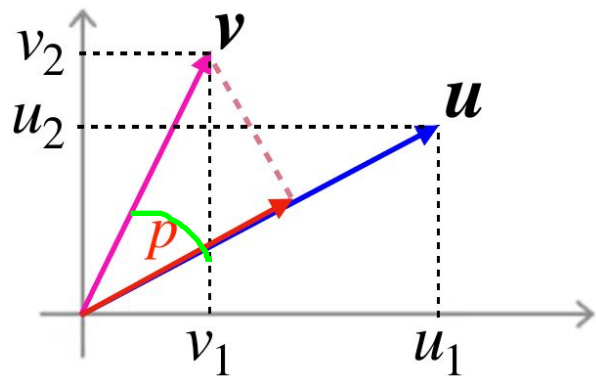
จาก figure :

p = ความยาวของ projection (การฉายภาพ) ของ v บน u

เป็นไปได้ที่จะแสดงให้เห็นว่า

$$\begin{aligned} u^T v &= p \cdot \|u\| = u_1 v_1 + u_2 v_2 \\ &= v^T u \end{aligned}$$

Preliminary: Vector Inner Product



ถ้ามี $\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$ กับ $\mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$ ให้คำนวณ $\mathbf{u}^T \mathbf{v} = ?$
($\mathbf{u}^T \mathbf{v}$ เรียกว่า ‘vector inner product’)

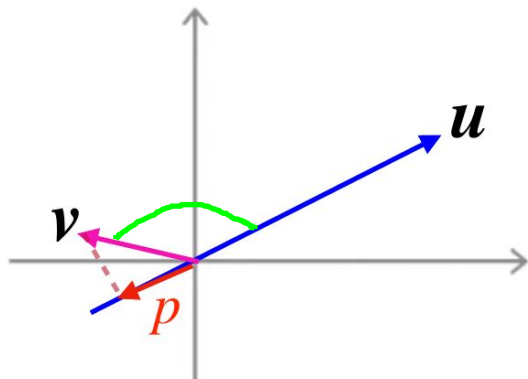
จาก figure :

p = ความยาวของ projection (การฉายภาพ) ของ \mathbf{v} บน \mathbf{u}

เป็นไปได้ที่จะแสดงให้เห็นว่า

$$\mathbf{u}^T \mathbf{v} = p \cdot \|\mathbf{u}\|$$

หมายเหตุ: p เป็นลบ ถ้าขนาดมุมระหว่าง \mathbf{u} กับ $\mathbf{v} > 90$



Decision Boundary ของ SVM

ทำให้สมการง่ายขึ้น เพื่อที่จะวิเคราะห์มันได้ง่ายขึ้น

$$\blacktriangleright J(\theta) = \frac{1}{2} \sum_{j=1}^n \theta_j^2 \text{ s.t. } \begin{aligned} & - \theta^T x^{(i)} \geq 1 \text{ if } y^{(i)} = 1 \\ & - \theta^T x^{(i)} \leq -1 \text{ if } y^{(i)} = 0 \end{aligned}$$

$$\blacktriangleright \theta_0 = 0$$

$$\blacktriangleright n = 2$$

Decision Boundary ของ SVM

ทำให้สมการง่ายขึ้น เพื่อที่จะวิเคราะห์มันได้ง่ายขึ้น

$$\begin{aligned} \blacktriangleright J(\theta) &= \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\underbrace{\sqrt{\theta_1^2 + \theta_2^2}}_{=\|\theta\|})^2 \\ \text{s.t. } & - \theta^T x^{(i)} \geq 1 \text{ if } y^{(i)} = 1 \\ & - \theta^T x^{(i)} \leq -1 \text{ if } y^{(i)} = 0 \end{aligned} \quad \text{where } \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$\blacktriangleright \theta_0 = 0$$

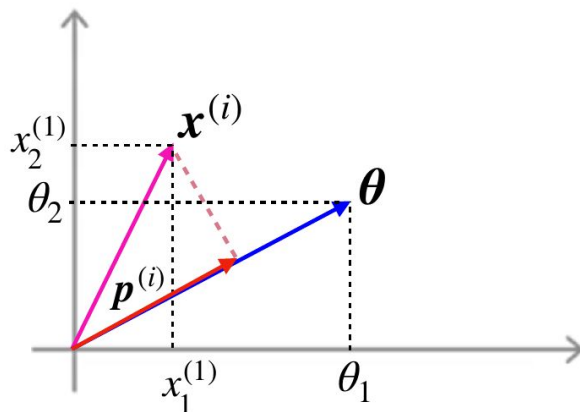
$$\blacktriangleright n = 2$$

Decision Boundary ของ SVM

$$\begin{aligned} \blacktriangleright J(\theta) &= \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2 \\ \text{s.t. } & \theta^T \mathbf{x}^{(i)} \geq 1 \text{ if } y^{(i)} = 1 \\ & -\theta^T \mathbf{x}^{(i)} \leq -1 \text{ if } y^{(i)} = 0 \end{aligned}$$

$$\blacktriangleright \theta_0 = 0$$

$$\blacktriangleright n = 2$$



ลองคำนวณ $\theta^T \mathbf{x}^{(i)} = ?$

$$\begin{aligned} \theta^T \mathbf{x}^{(i)} &= p^{(i)} \cdot \|\theta\| \\ &= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \end{aligned}$$

Decision Boundary ของ SVM

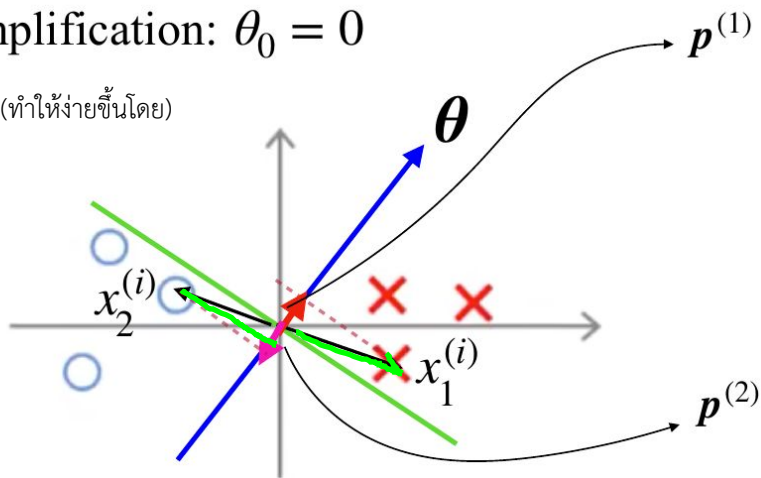
Objective function

ของ SVM จะกลายเป็น

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \iff \min_{\theta} \sum_{j=1}^n \frac{1}{2} \|\theta\|^2 \quad \text{s.t.} \quad \begin{aligned} & - p^{(i)} \cdot \|\theta\| \geq 1 \quad \text{if } y^{(i)} = 1 \\ & - p^{(i)} \cdot \|\theta\| \leq -1 \quad \text{if } y^{(i)} = 0 \end{aligned}$$

Simplification: $\theta_0 = 0$

(ทำให้ง่ายขึ้นโดย)



เมื่อ $p^{(i)}$ เป็น projection ของ $x^{(i)}$ บน vector θ

เพราะ $p^{(1)}$ น้อยมาก

ดังนั้น $\|\theta\|$ ต้องเยอะมาก

(ทำไม นี่จึงสมเหตุสมผล ?)

เพราะ $p^{(2)}$ น้อยมาก

ดังนั้น $\|\theta\|$ ต้องเยอะมาก

(ทำไม นี่จึงสมเหตุสมผล ?)

Decision Boundary ของ SVM

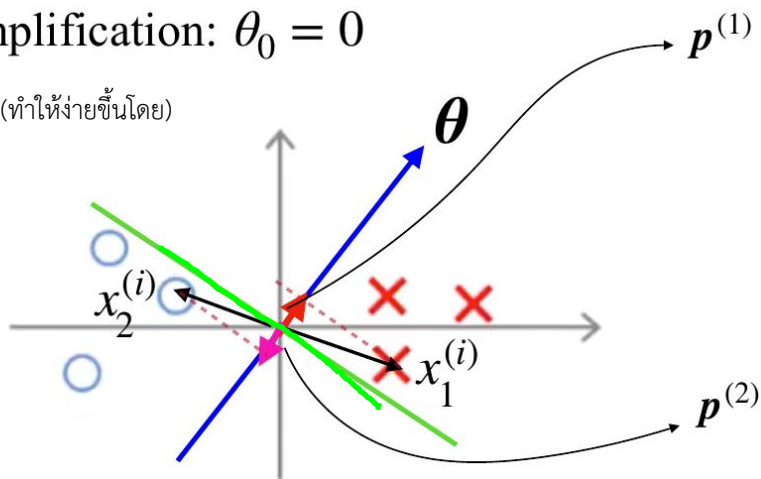
Objective function

ของ SVM จะกลายเป็น

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \iff \min_{\theta} \sum_{j=1}^n \frac{1}{2} \|\theta\|^2 \quad \text{s.t.} \quad \begin{aligned} & - \mathbf{p}^{(i)} \cdot \|\theta\| \geq 1 \quad \text{if } y^{(i)} = 1 \\ & - \mathbf{p}^{(i)} \cdot \|\theta\| \leq -1 \quad \text{if } y^{(i)} = 0 \end{aligned}$$

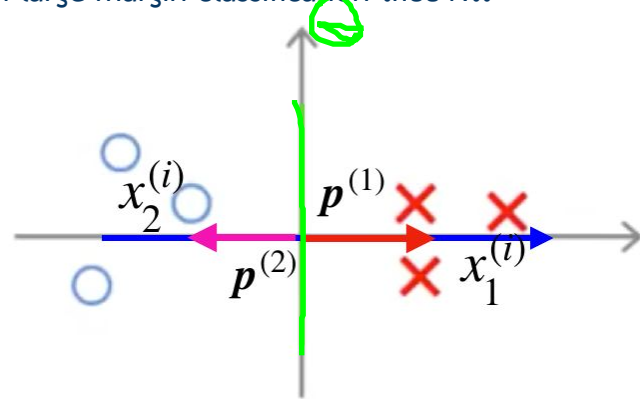
Simplification: $\theta_0 = 0$

(ทำให้ง่ายขึ้นโดย)



เมื่อ $p^{(i)}$ เป็น projection ของ $x^{(i)}$ บน vector θ

นี่อธิบายว่า SVM ทำให้เกิด large margin classification ได้อย่างไร !



Decision Boundary ของ SVM

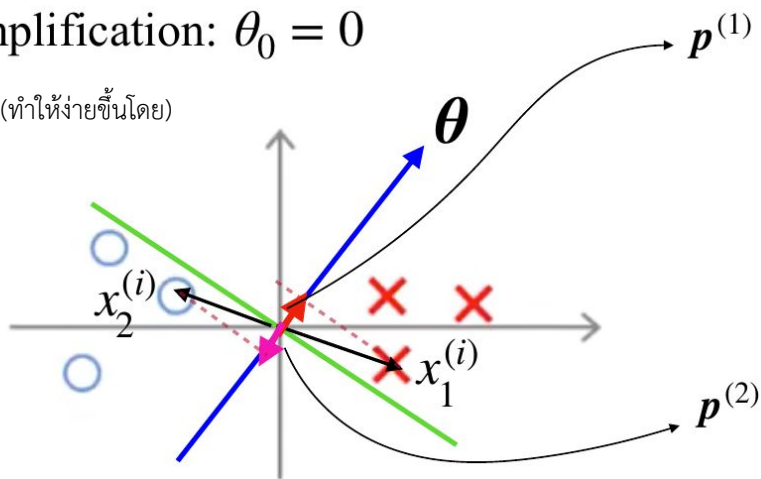
Objective function

ของ SVM จะกลายเป็น

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \iff \min_{\theta} \sum_{j=1}^n \frac{1}{2} \|\theta\|^2 \quad \text{s.t.} \quad \begin{aligned} & - \mathbf{p}^{(i)} \cdot \|\theta\| \geq 1 \quad \text{if } y^{(i)} = 1 \\ & - \mathbf{p}^{(i)} \cdot \|\theta\| \leq -1 \quad \text{if } y^{(i)} = 0 \end{aligned}$$

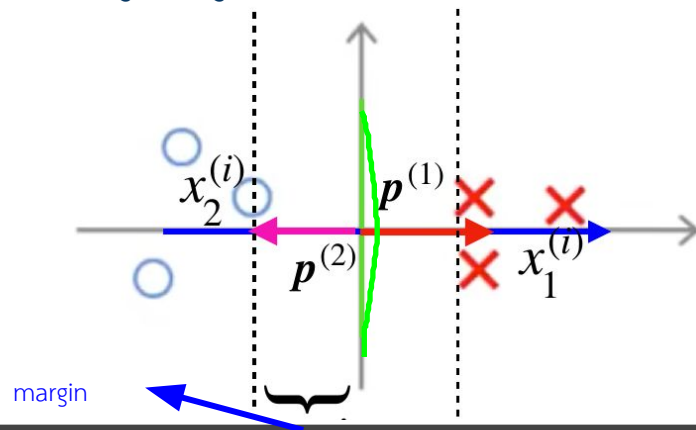
Simplification: $\theta_0 = 0$

(ทำให้ง่ายขึ้นโดย)



เมื่อ $\mathbf{p}^{(i)}$ เป็น projection ของ $\mathbf{x}^{(i)}$ บน vector θ

นี่อธิบายว่า SVM ทำให้เกิด large margin classification ได้อย่างไร !



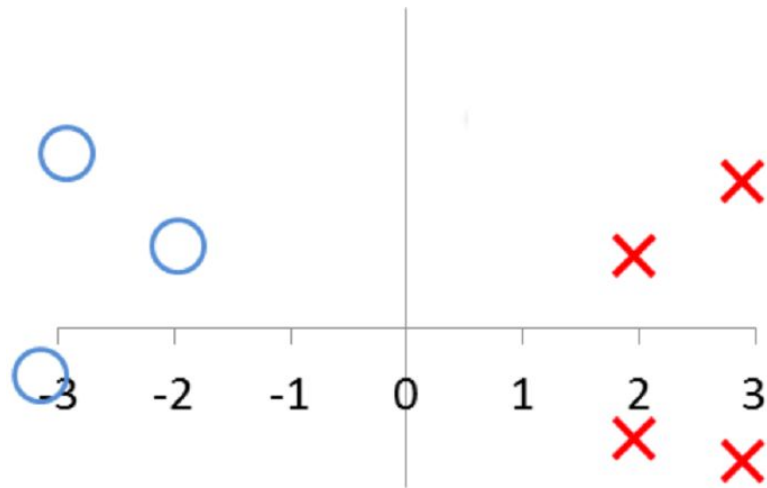
Question

ปัญหา SVM optimization ที่เราใช้ คือ:

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \text{ s.t. } \begin{aligned} & - \mathbf{p}^{(i)} \cdot \|\boldsymbol{\theta}\| \geq 1 \quad \text{if } y^{(i)} = 1 \\ & - \mathbf{p}^{(i)} \cdot \|\boldsymbol{\theta}\| \leq -1 \quad \text{if } y^{(i)} = 0 \end{aligned}$$

เมื่อ $\mathbf{p}^{(i)}$ เป็น projection (signed / ที่มีเครื่องหมาย) ของ $\mathbf{x}^{(i)}$ ลงบน $\boldsymbol{\theta}$ พิจารณาชุดข้อมูล training set ด้านขวา ที่ค่าที่เหมาะสม (optimal value) ของ $\boldsymbol{\theta}$ ค่าของ $\|\boldsymbol{\theta}\|$ จะเป็นเท่าไร?

- (i) 1 / 4 (ii) 1 / 2
(iii) 1 (iv) 2



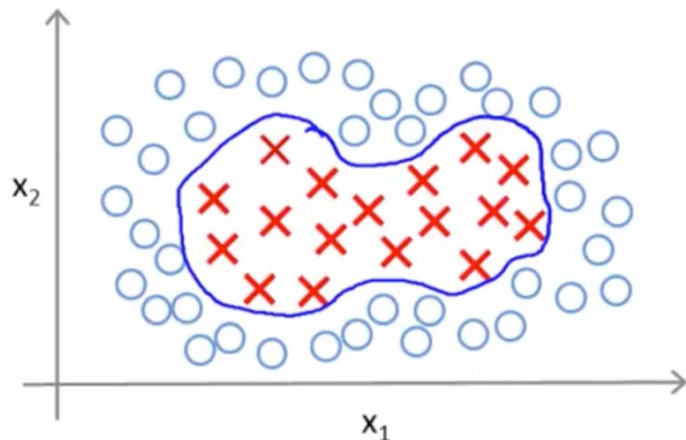
Support Vector Machines (SVM)

Kernels (Part 1)

Krittameth Teachasrisaksakul

ความเข้าใจพื้นฐาน

ขอบเขตตัดสินใจที่ไม่เป็นเชิงเส้น (Non-linear decision boundary)



ทำนาย $y = 1$ ถ้า

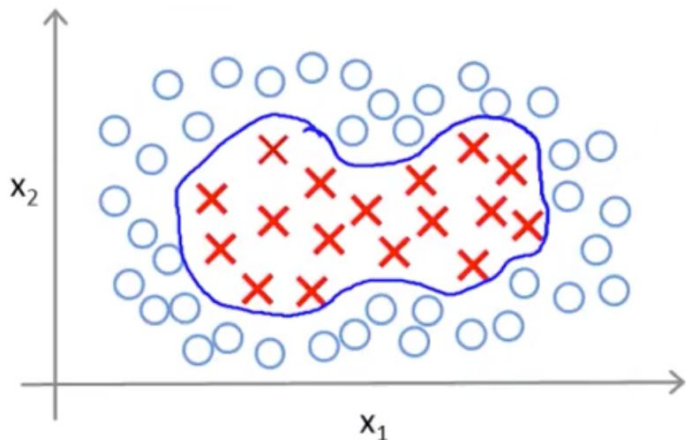
$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$$

ก็คือ

$$h_{\theta}(x) = \begin{cases} 1, & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

ความเข้าใจพื้นฐาน

ขอบเขตตัดสินใจที่ไม่เป็นเชิงเส้น (Non-linear decision boundary)



ทำนาย $y = 1$ ถ้า

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$$

ก็คือ

$$h_\theta(x) = \begin{cases} 1, & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

สัญลักษณ์ใหม่ (New notation):

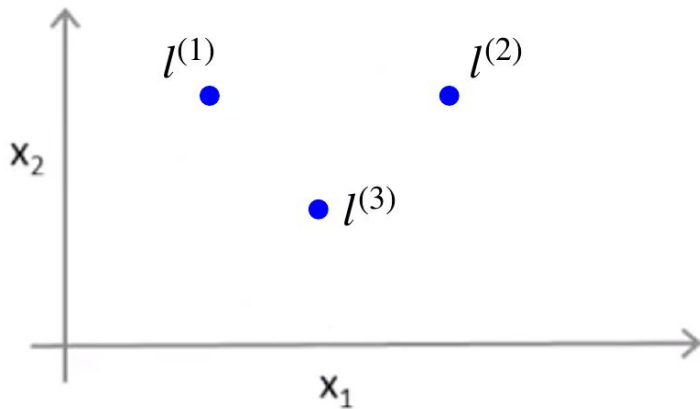
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \dots$$

เช่น $f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, f_5 = x_2^2, \dots$

คำถาม: มีตัวเลือก features f_1, f_2, f_3 ที่ต่างไป หรือ ดีกว่าหรือไม่?

Kernels : ความเข้าใจพื้นฐาน

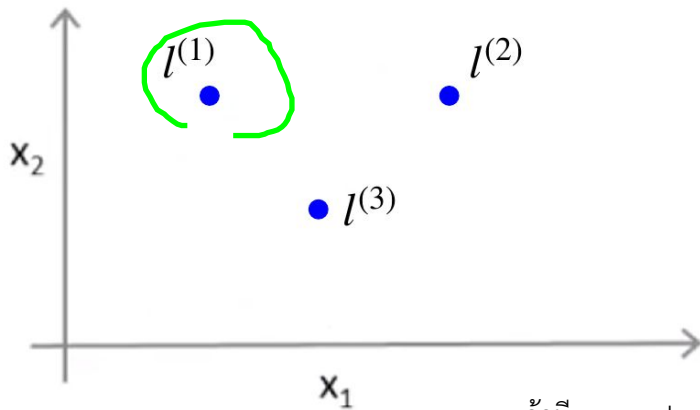
นี่คือ แนวคิดใหม่ในการนิยาม feature ใหม่ f_1, f_2, f_3



ถ้ามี x จำนวน feature ใหม่ ที่ขึ้นอยู่กับ proximity (ความใกล้ชิด) ของ landmarks $l^{(1)}, l^{(2)}, l^{(3)}$

Kernels : ความเข้าใจพื้นฐาน

นี่คือ แนวคิดใหม่ในการนิยาม feature ใหม่ f_1, f_2, f_3



ถ้ามี x จำนวน feature ใหม่ ที่ขึ้นอยู่กับ proximity (ความใกล้ชิด) ของ landmarks $l^{(1)}, l^{(2)}, l^{(3)}$

ถ้ามี example (ตัวอย่าง) $x \rightarrow$

$$f_1 := \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

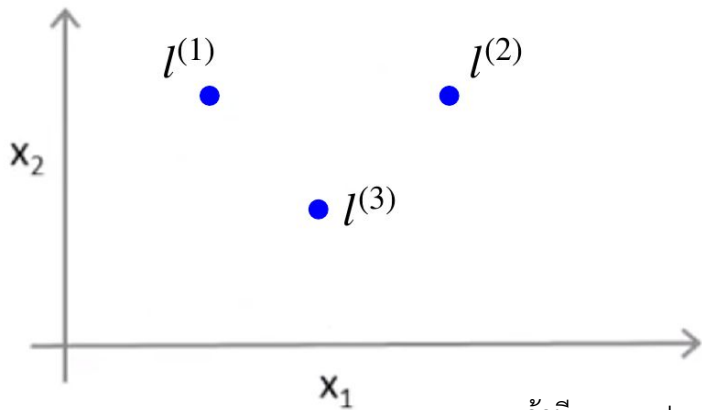
$$f_2 := \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

\vdots

kernels
(Gaussian kernels)

Kernels : ความเข้าใจพื้นฐาน

นี่คือ แนวคิดใหม่ในการนิยาม feature ใหม่ f_1, f_2, f_3



ถ้ามี example (ตัวอย่าง) $x \rightarrow$

ถ้ามี x จำนวน feature ใหม่ ที่ขึ้นอยู่กับ proximity (ความใกล้ชิด) ของ landmarks $l^{(1)}, l^{(2)}, l^{(3)}$

$$f_1 := \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$f_2 := \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

⋮

kernels
(Gaussian kernels)

Kernels และ Similarity

$$f_1 := \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right),$$

$$\text{where } \|x - l\|^2 = \sum_{j=1}^n (x_j - l_j)^2$$

ถ้า $x \approx l^{(1)}$ แล้ว

ถ้า x อยู่ไกลจาก $l^{(1)}$ แล้ว

Kernels และ Similarity

$$f_1 := \mathbf{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right),$$

$$\text{where } \|x - l\|^2 = \sum_{j=1}^n (x_j - l_j)^2$$

ถ้า $x \approx l^{(1)}$ แล้ว

$$f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

ถ้า x อยู่ไกลจาก $l^{(1)}$ แล้ว

$$f_1 \approx \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

สังเกตว่า landmark แต่ละจุด ทำให้เกิด (นิยาม) feature ใหม่

Kernels และ Similarity

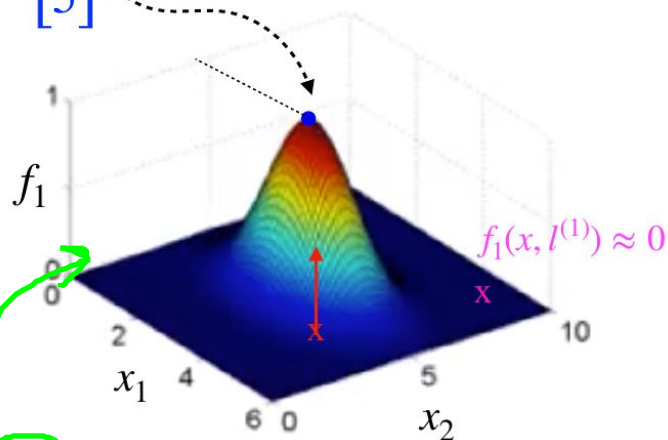
ตัวอย่าง: ถ้ามี

$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

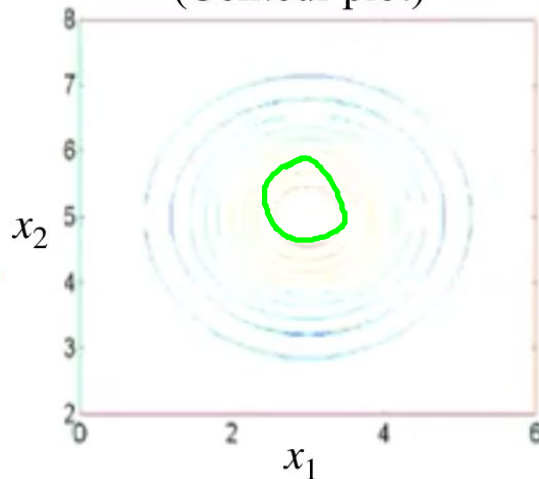
$$f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$\sigma^2 = 1$$



(Contour plot)



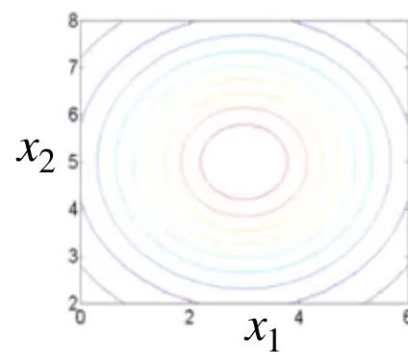
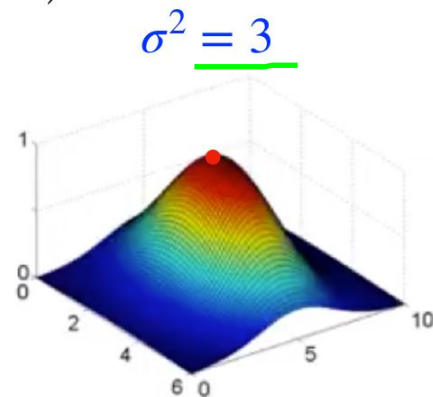
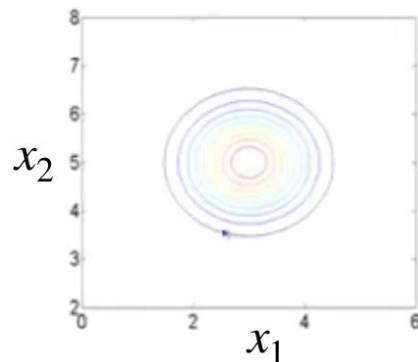
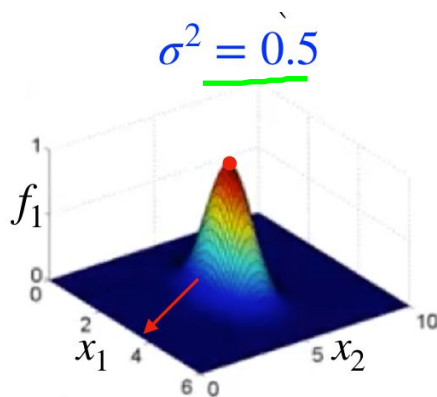
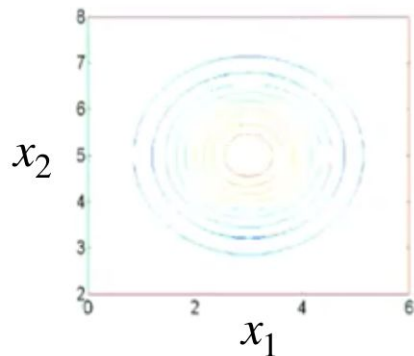
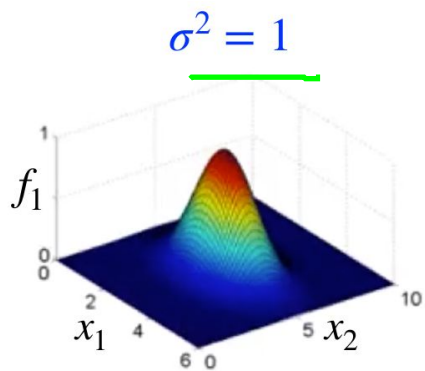
ต่อไป ลองดูผลของ σ^2 (ซึ่งเป็น parameter ของ Gaussian kernel)

Kernels และ Similarity

ตัวอย่าง: ถ้ามี

$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$



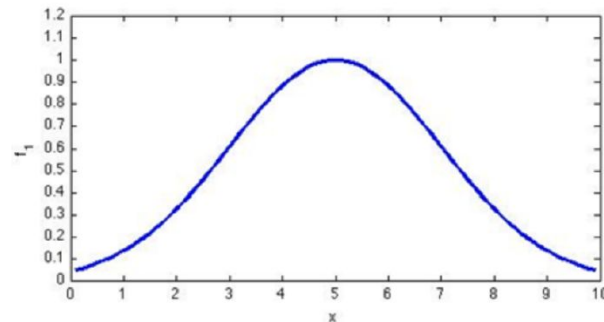
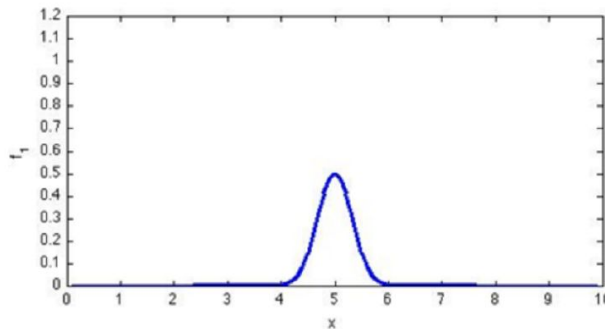
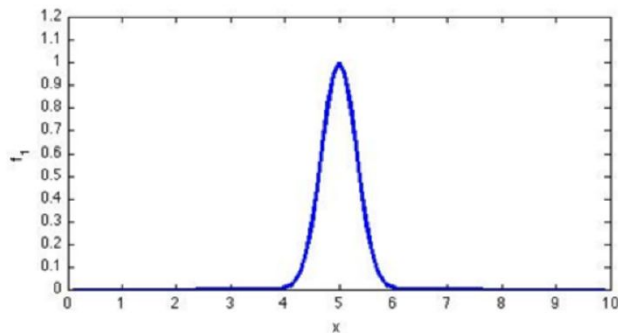
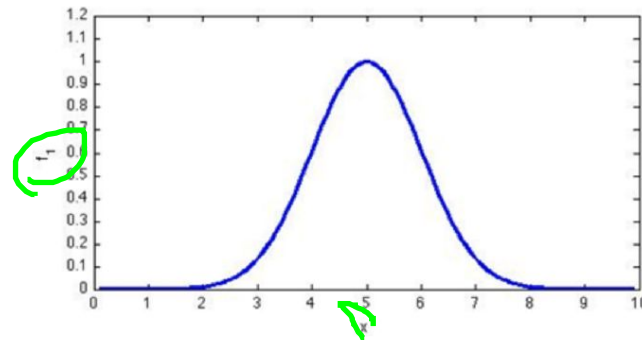
Question

พิจารณาตัวอย่าง 1 มิติ (1-D example) ที่มี 1 feature X_1 และสมมติ $l^{(1)} = 5$

ภาพด้านบน คือ plot ของ

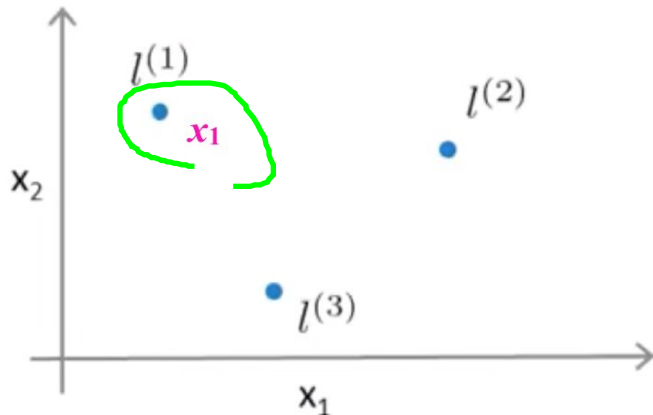
$$f_1 = \exp\left(-\frac{\|x_1 - l^{(1)}\|^2}{2\sigma^2}\right)^2 = 1$$

สมมติ เราเปลี่ยนให้ $\sigma^2 = 4$ ภาพใดที่เป็น plot ของ f_1 ที่มีค่า σ^2 ค่าใหม่?



Kernels และ Similarity

ถ้ามี features (คุณลักษณะ) เหล่านี้ ลองดูว่าจะเรียนรู้ hypothesis function อะไรได้



$$\because f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$$

ก็คือ $\theta_0 + \theta_1 \cdot 1 + \theta_2 \cdot 0 + \theta_3 \cdot 0 = -0.5 + 1 = 0.5 \geq 0$

\therefore ทำนาย $y = 1$ สำหรับ x_1

Hypothesis function:

ทำนาย '1' เมื่อ

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

ถ้ามี example x_1

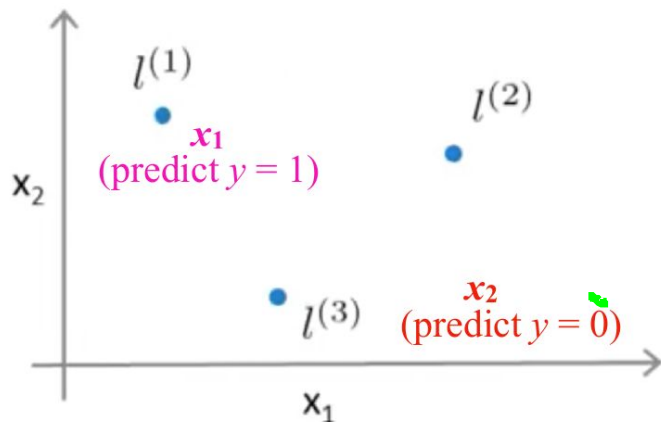
เราจะคำนวณ: f_1, f_2, f_3

สมมติ เราทำนาย

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

Kernels และ Similarity

ถ้ามี features (คุณลักษณะ) เหล่านี้ ลองดูว่าจะเรียนรู้ hypothesis function อะไรได้



$$\because f_1 \approx 0, f_2 \approx 0, f_3 \approx 0$$

ก็คือ $\theta_0 + \theta_1 \cdot 1 + \theta_2 \cdot 0 + \theta_3 \cdot 0 = \underline{\underline{-0.5 < 0}}$

\therefore ทำนาย $y = 0$ สำหรับ x_2

Hypothesis function:

ทำนาย '1' เมื่อ

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

ถ้ามี example x_2

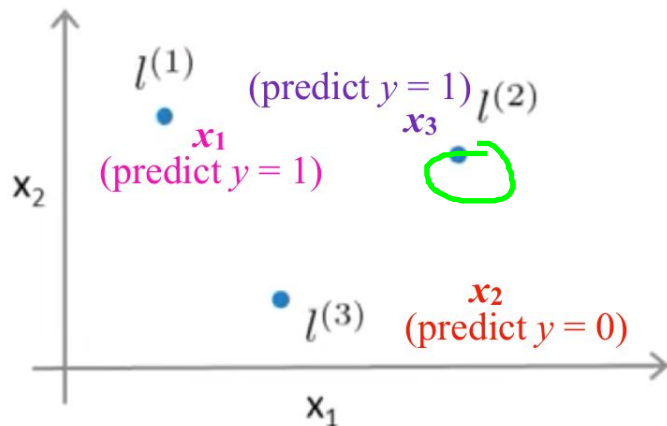
เราจะคำนวณ: f_1, f_2, f_3

สมมติ เราทำนาย

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

Kernels และ Similarity

ถ้ามี features (คุณลักษณะ) เหล่านี้ ลองดูว่าจะเรียนรู้ hypothesis function อะไรได้



$$\because f_1 \approx 0, f_2 \approx 0, f_3 \approx 0$$

$$\text{ก็คือ } \theta_0 + \theta_1 \cdot 1 + \theta_2 \cdot 0 + \theta_3 \cdot 0 = -0.5 < 0$$

\therefore ทำนาย $y=0$ สำหรับ x_2

Hypothesis function:

ทำนาย '1' เมื่อ

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

ถ้ามี example x_2

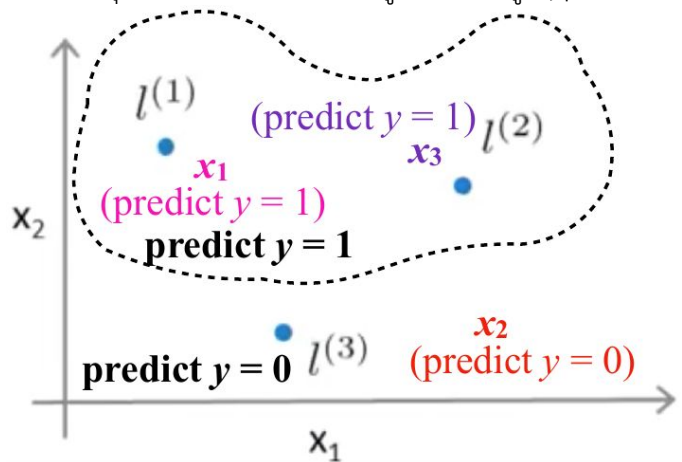
เราจะคำนวณ: f_1, f_2, f_3

สมมติ เราทำนาย

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

Kernels และ Similarity

ถ้ามี features (คุณลักษณะ) เหล่านี้ ลองดูว่าจะเรียนรู้ hypothesis function อะไรได้



Hypothesis function:

ทำนาย '1' เมื่อ

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

ถ้ามี example \mathbf{x}_2

เราจะคำนวณ: f_1, f_2, f_3

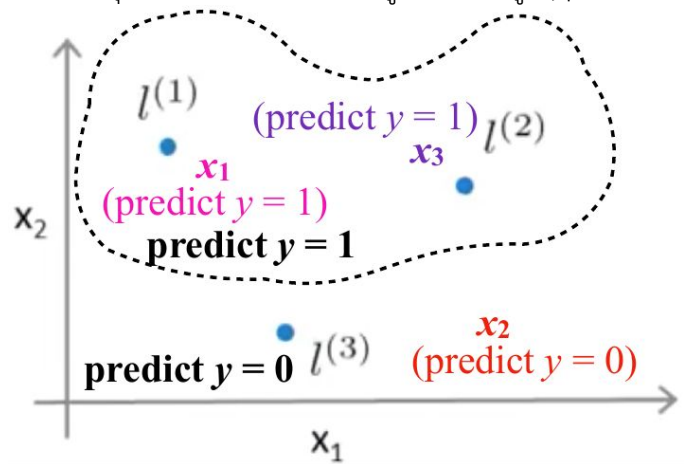
สมมติ เราทำนาย

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

นี้ให้ความเข้าใจพื้นฐานเกี่ยวกับว่า นิยามของ landmarks และ kernel function ทำให้เรียนรู้ non-linear decision boundary ที่ค่อนข้างซับซ้อนได้อย่างไร !

Kernels และ Similarity

ถ้ามี features (คุณลักษณะ) เหล่านี้ ลองดูว่าจะเรียนรู้ hypothesis function อะไรได้



Hypothesis function:

ทำนาย '1' เมื่อ

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

ถ้ามี example \mathbf{x}_2

เราจะคำนวณ: f_1, f_2, f_3

สมมติ เราทำนาย

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$

คำถาม: หา landmarks เหล่านี้ได้อย่างไร ? เลือก landmarks เหล่านี้ได้อย่างไร ? Similarity function อื่นๆ มีอะไรบ้าง?
เป็นต้น

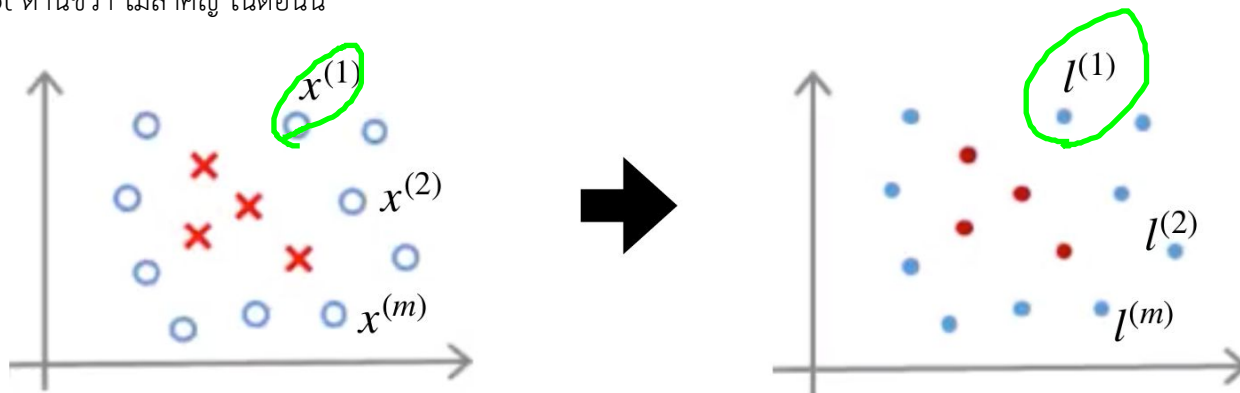
Support Vector Machines (SVM)

Kernels (Part 2)

Krittameth Teachasrisaksakul

การเลือก landmarks

สำหรับ X ทุกตัว ถ้า X อยู่ในชุดข้อมูล เรากำหนดให้ X เป็น landmarks
สี่แต่ละสีของจุดใน plot ด้านขวา ไม่สำคัญ ในตอนนี้



นิยาม: ถ้ามีตัวอย่าง example m ตัว
เลือก / กำหนดให้ $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$
 $l^{(1)} := x^{(1)}, l^{(2)} := x^{(2)}, \dots, l^{(m)} := x^{(m)}$

SVM ที่ใช้ Kernels

ถ้ามีตัวอย่าง \mathbf{X} (จากชุดข้อมูล training / cross validation / testing):


$$\left. \begin{array}{l} f_1 = \text{similarity}(x, l^{(1)}) \\ f_2 = \text{similarity}(x, l^{(2)}) \\ \vdots \end{array} \right\} \mathbf{f} = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad \text{เมื่อ } \underline{f_0 = 1 \text{ (interceptor)}}$$

SVM ที่ใช้ Kernels

ถ้ามีตัวอย่าง \mathbf{X} (จากชุดข้อมูล training / cross validation / testing):

$$\left. \begin{array}{l} f_1 = \mathbf{similarity}(x, l^{(1)}) \\ f_2 = \mathbf{similarity}(x, l^{(2)}) \\ \vdots \end{array} \right\} \mathbf{f} = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad \text{เมื่อ } f_0 = 1 \text{ (interceptor)}$$

ตัวอย่าง: สำหรับตัวอย่างจากชุดข้อมูล training $(x^{(i)}, y^{(i)})$ เราสามารถสร้าง vector

$$x^{(i)} \rightarrow \begin{array}{l} f_1^{(i)} = \mathbf{similarity}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \mathbf{similarity}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} = \mathbf{similarity}(x^{(i)}, l^{(m)}) \end{array} \quad \leftarrow \quad \begin{array}{l} f_i^{(i)} = \mathbf{similarity}(x^{(i)}, l^{(i)}) \\ \hline = \exp\left(-\frac{0}{2\sigma^2}\right) = 1 \end{array}$$


SVM ที่ใช้ Kernels

ถ้ามีตัวอย่าง \mathbf{X} (จากชุดข้อมูล training / cross validation / testing):

$$\left. \begin{array}{l} f_1 = \text{similarity}(x, l^{(1)}) \\ f_2 = \text{similarity}(x, l^{(2)}) \\ \vdots \end{array} \right\} \mathbf{f} = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad \text{เมื่อ } f_0 = 1 \text{ (interceptor)}$$

ตัวอย่าง: สำหรับตัวอย่างจากชุดข้อมูล training ($x^{(i)}, y^{(i)}$) เราสามารถสร้าง vector

$$x^{(i)} \rightarrow \begin{array}{l} f_1^{(i)} = \text{similarity}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{similarity}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} = \text{similarity}(x^{(i)}, l^{(m)}) \end{array}$$

แทนที่จะใช้ $\mathbf{x}^{(i)} \in \mathbb{R}^{n+1}$

เราจะเขียนแทน $\mathbf{x}^{(i)}$ แต่ละตัวด้วย feature vector

$$\mathbf{f}^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$$

SVM ที่ใช้ Kernels

สมมติ เรามี parameter ที่เรียนรู้แล้ว $\theta (\theta \in \mathbb{R}^{m+1})$

Hypothesis: ถ้ามี X เราคำนวณ features $f \in \mathbb{R}^{m+1}$

แล้ว เราทำนาย 'y=1' ถ้า

$$\theta^T f \geq 0 \iff \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m \geq 0$$

จะหาค่า parameter θ ได้อย่างไร?

SVM ที่ใช้ Kernels

สมมติ เรามี parameter ที่เรียนรู้แล้ว $\theta (\theta \in \mathbb{R}^{m+1})$

Hypothesis: ถ้ามี X เราคำนวณ features $f \in \mathbb{R}^{m+1}$

แล้ว เราทำนาย ' $y=1$ ' ถ้า

$$\theta^T f \geq 0 \iff \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m \geq 0$$

Objective function:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \mathbf{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

สังเกตว่า สำหรับปัญหา optimization นี้ เรามี $n = m$

SVM ที่ใช้ Kernels

สมมติ เรามี parameter ที่เรียนรู้แล้ว $\theta (\theta \in \mathbb{R}^{m+1})$

Hypothesis: ถ้ามี X เราคำนวณ features $f \in \mathbb{R}^{m+1}$

แล้ว เราทำนาย ' $y=1$ ' ถ้า

$$\theta^T f \geq 0 \iff \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m \geq 0$$

Objective function:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \mathbf{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \mathbf{cost}_0(\theta^T f^{(i)}) + \boxed{\frac{1}{2} \sum_{j=1}^n \theta_j^2}$$

$$\because \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \boxed{\frac{1}{2} \|\theta\|^2}$$

$$\because \underline{\sum_j \theta_j^2 = \theta^T \theta}$$

เพิ่มมาเพื่อ rescaling !

ในบาง implementation : คำนวณ $\theta^T M \theta$ เพื่อ computational efficiency (ประสิทธิภาพในการคำนวณ)

การเลือก Parameter ของ SVM

C (เช่น $\frac{1}{\lambda}$)

C มาก : bias ต่ำลง, variance สูง ($\approx \lambda$ น้อย)

มีแนวโน้ม overfitting

C น้อย : bias สูงขึ้น, variance ต่ำ ($\approx \lambda$ มาก)

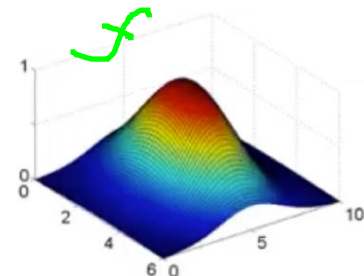
มีแนวโน้ม underfitting

σ^2 σ^2 มาก : features f_i เปลี่ยนอย่าง smooth
bias สูงขึ้น, variance ต่ำลง

มีแนวโน้ม underfitting

$$\sigma^2 = 3$$

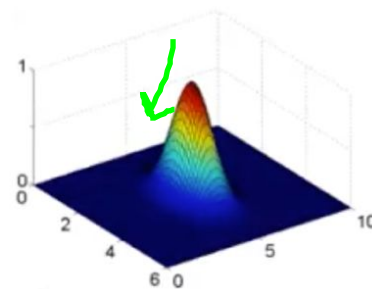
ทบทวน:
$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$



σ^2 น้อย : features f_i เปลี่ยนอย่าง smooth น้อยลง
bias ต่ำลง, variance สูงขึ้น

มีแนวโน้ม overfitting

$$\sigma^2 = 0.5$$



Question

สมมติ เรา train SVM และพบว่ามัน overfit ชุดข้อมูล training

ข้อใดต่อไปนี้เป็นขั้นตอนต่อไปที่เหมาะสม ? วงทุกข้อที่ถูกต้อง

(i) เพิ่ม C

(ii) ลด C

(iii) เพิ่ม σ^2

(iv) ลด σ^2

Support Vector Machines (SVM)

การใช้ SVM

Krittameth Teachasrisaksakul

การประยุกต์ใช้

1. ใช้ SVM software package (เช่น scikit-learn, libsvm, ...) เพื่อแก้หาค่า parameter
2. จำเป็น ต้องกำหนดค่า:
 - i. การเลือกค่า parameter C
 - ii. การเลือก kernel ก็คือ similarity function เช่น
 - ไม่มี kernel (หรือ 'linear kernel')

ทำนาย $y = 1$ ถ้า $\theta^T x \geq 0$ ($\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$)

Linear kernel บอกเป็นนัยว่าเราใช้ 'standard linear classifier' (ตัวแยกประเภทที่ใช้ฟังก์ชันเชิงเส้นแบบมาตรฐาน)

คำถาม: เมื่อใดที่เราจำเป็นต้องใช้มัน?

เมื่อ n มาก และ m น้อย ก็คือ $x_i \in \mathbb{R}^{n+1}$

ในกรณีนี้ เรามีข้อมูลไม่เพียงพอ และอยากหลีกเลี่ยง overfitting !

การประยุกต์ใช้

1. ใช้ SVM software package (เช่น scikit-learn, libsvm, ...) เพื่อแก้หาค่า parameter
2. จำเป็น ต้องกำหนดค่า:

- i. การเลือกค่า parameter C
- ii. การเลือก kernel ก็คือ similarity function เช่น
 - ไม่มี kernel (หรือ 'linear kernel')

ทำนาย $y = 1$ ถ้า $\theta^T x \geq 0$ ($\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$)

- Gaussian kernel (ต้องเลือกค่า σ^2):

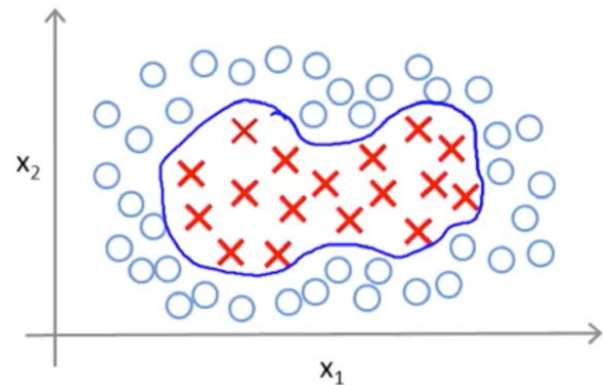
$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$

เมื่อ

$$l^{(i)} = x^{(i)}$$

คำถาม: เมื่อใดที่เราจะใช้ Gaussian kernel ?

เมื่อ $x \in \mathbb{R}^n$, n น้อย, และ m มาก



การประยุกต์ใช้

เตือนความจำเกี่ยวกับ implementation

ต้องทำ feature scaling ก่อนใช้ Gaussian kernel

เพราะอะไร ?

$$\because \|x - l\|^2 = (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2 \quad (x \in \mathbb{R}^{n+1})$$

ตัวอย่าง (Housing domain):

$$x_1 \in [0, 1000] \text{ feet}^2$$

$$x_2 \in \{1, 2, 3, 4, 5\}$$

ถ้าในกรณีนี้ : ขนาดพื้นที่บ้าน (size) จะมีอิทธิพลมากกว่า ระยะห่าง (distance) เหล่านี้

Kernel แบบอื่นๆ

- ไม่ใช่ similarity function ทุกอัน **similarity**(x, l) ที่จะเป็น kernel ที่ถูกต้อง / ใช้ได้ (valid)
- เพื่อยอมรับว่าเป็น kernel ที่ถูกต้อง ต้องสอดคล้องกับ Mercer's theorem

- เพื่อให้แน่ใจว่า optimization ของ SVM package ทำงานอย่างถูกต้อง และ
- ไม่ diverge !

- ตัวเลือกอื่นของ off-the-shelf kernels (ที่ไม่ต้องพัฒนาเอง)

- Polynomial kernel: **similarity**(x, l) := $(x^T l + \epsilon)^d$ เช่น

$$(x^T l)^2 \quad (x^T l)^3 \quad (x^T l + 1)^3 \quad (x^T l + 5)^4$$

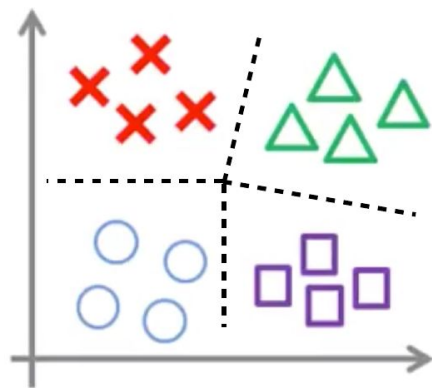
- Kernel ที่เข้าใจยากขึ้น: string kernel, chi-square kernel, histogram, intersection kernel, ...

Question

สมมติ เราพยายามตัดสินใจเลือกระหว่าง kernel ไม่กี่ตัว และเลือกค่า parameter เช่น C, σ^2 เป็นต้น เราควรเลือกอย่างไร?

- (i) เลือกอะไรก็ตามที่ ทำงานมีประสิทธิภาพสูงสุดกับข้อมูล training
- (ii) เลือกอะไรก็ตามที่ ทำงานมีประสิทธิภาพสูงสุดกับข้อมูล cross-validation
- (iii) เลือกอะไรก็ตามที่ ทำงานมีประสิทธิภาพสูงสุดกับข้อมูล test
- (iv) เลือกอะไรก็ตามที่ ทำให้มี SVM margin มากที่สุด

Multi-class Classification (การแยกประเภท มากกว่า 2 ประเภท)



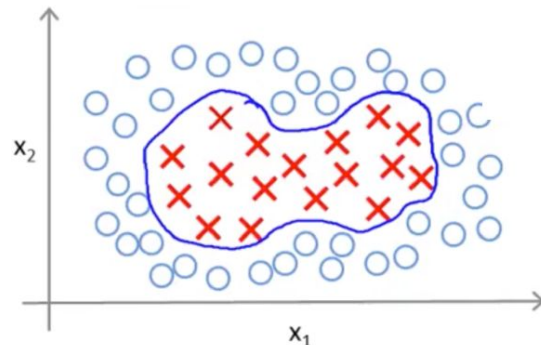
$$y \in \{1, 2, 3, \dots, K\}$$

- SVM package หลายๆอัน มี multi-class classification functionality (ที่ built-in / มีพร้อมใช้ได้)
- ไม่อย่างนั้น ใช้วิธี one-vs-all (แยก 1 class ออกจาก class ที่เหลือทั้งหมด) ก็คือ
 - Train SVM K ตัว แต่ละตัวใช้เพื่อแยก $y = i$ ออกจาก class ที่เหลือทั้งหมด เมื่อ $i = 1, 2, \dots, K$
 - หาค่า $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$
 - เลือก class i ที่มีค่า $(\theta^{(i)})^T \mathbf{x}$ มากที่สุด

Logistic Regression vs. SVM

เมื่อใดควรใช้ algorithm แต่ละอัน (เทียบกับอีกอัน) ?

สมมติ n = จำนวน features ($\mathbf{x} \in \mathbb{R}^{n+1}$), m = จำนวน training examples



1. ถ้า n มาก (เทียบกับ m) (เช่น $n \geq m, n = 10,000, 10 \leq m \leq 10,000$)

ใช้ logistic regression หรือ SVM ที่ไม่มี kernel (ก็คือ 'linear kernel')

2. ถ้า n น้อย, m มีค่าปานกลาง (intermediate) (เช่น $1 \leq n \leq 10,000, 10 \leq m \leq 50,000$)

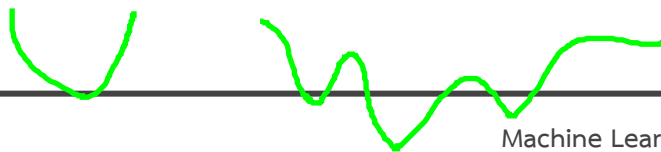
ใช้ SVM ที่ใช้ Gaussian kernel

3. ถ้า n น้อย, m มาก (เช่น $1 \leq n \leq 1,000, m > 50,000$)

สร้าง / เพิ่ม features แล้วใช้ logistic regression หรือ SVM ที่ไม่มี kernel

4. Neural network (NN) มีแนวโน้มที่จะทำงานได้ดี ใน setting ส่วนมากที่พูดถึง แต่อาจ train ได้ช้ากว่า

5. SVM เป็น convex optimization problem ในทางปฏิบัติ ค่า local optima ไม่ใช่ปัญหาใหญ่ เมื่อใช้ neural network แต่เราไม่ต้องกังวลเรื่องนี้ เมื่อใช้ SVM



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>