

## บทที่ 1

### บทนำ

#### 1.1 ความเป็นมาและความสำคัญของปัญหา

แบบจำลองจักรกลเรียนรู้ (Machine Learning models) นั้นถูกใช้อย่างกว้างขวางในปัจจุบัน อย่างไรก็ตาม แบบจำลองใดๆ นั้นอาจมีความผิดพลาดต่อการทำการโจมตีประสงค์ร้าย (Adversarial attacks) เพื่อจงใจให้ผลลัพธ์ที่แบบจำลองนั้นคาดเดามีความผิดพลาดจากผลลัพธ์ที่ควรจะเป็น

ในการเรียนรู้เชิงตัวแปรเสริม (parameter-based learning) นั้น ตัวแปรเสริม (parameters) ค่าน้ำหนัก (weights) บนแบบจำลองการเรียนรู้เชิงลึก (deep Learning models) เป็นตัวกำหนดความฉลาดของแบบจำลอง อาจมีตัวแปรเสริมบางชุดที่ทำให้แบบจำลองมีช่องโหว่ต่อการโจมตีประสงค์ร้าย การโจมตีนั้นอาจเกิดจากการเพิ่มสัญญาณรบกวนซึ่งผ่านการคำนวณ (calculated artefacts) เข้าสู่ข้อมูลรับเข้า (inputs) ซึ่งทำให้ความผิดพลาดของแบบจำลองในการพยากรณ์คำตอบนั้นเปลี่ยนไปอย่างชัดเจน

โครงการวิศวกรรมคอมพิวเตอร์นี้มุ่งหวังจะนำตัวแปรเสริมบนแบบจำลองมาสร้างภาพแสดง (visualise) ถึงจุดโหว่ในการพยากรณ์ใดๆ ของแบบจำลอง เพื่อลดความเสียหายอันอาจเกิดขึ้นได้จากการโจมตีแบบจำลองขณะถูกใช้งานจริง

#### 1.2 วัตถุประสงค์ของการศึกษา

โครงการนี้มีวัตถุประสงค์และเป้าหมายดังนี้

1. สร้างแบบจำลองเชิงลึก (Deep Learning models) ซึ่งสามารถถูกโจมตีประสงค์ร้าย (Adversarial attacks) ได้
2. นำแบบจำลองในข้อ (1) มาสร้างเป็นรูปภาพแสดง (visualisation) เพื่อหาจุดโหว่ต่อการโจมตี รวมถึงคาดเดาแนวโน้มการโจมตีที่เป็นไปได้
3. ใช้ความรู้ในข้อ (2) สร้างแบบจำลองที่ทนทาน (prone) ต่อการโจมตีมากขึ้น

#### 1.3 ขอบเขตของการทำโครงการ

โครงการนี้มีขอบเขตการดำเนินงานดังนี้

1. สร้างแบบจำลองเชิงลึก (Deep Learning models) ซึ่งสามารถถูกโจมตีประสงค์ร้าย (Adversarial attacks) ได้
2. นำแบบจำลองในข้อ (1) มาสร้างเป็นรูปภาพแสดง (visualisation) เพื่อหาจุดโหว่ต่อการโจมตี รวมถึงคาดเดาแนวโน้มการโจมตีที่เป็นไปได้
3. ใช้ความรู้ในข้อ (2) สร้างแบบจำลองที่ทนทาน (prone) ต่อการโจมตีมากขึ้น

#### 1.4 ระยะเวลาและแผนดำเนินงาน

ในช่วงแรกของการทำโครงการ แผนการดำเนินงานนั้นจะใช้ในรูปแบบของรวนวนซ้ำ (iteration) ตามกรรมวิธีการดำเนินงานแบบเอจิล (agile) ซึ่งประกอบไปด้วยขั้นตอนการรวนวนดังนี้...

#### 1.5 ประโยชน์ที่คาดว่าจะได้รับ

1. เข้าใจถึงพื้นฐาน หลักการทำงาน และระบบจักรกลเรียนรู้แบบต่างๆ
2. เข้าใจถึงจุดอ่อนของระบบจักรกลเรียนรู้ในแต่ละกรณี
3. สามารถโจมตีระบบจักรกลเรียนรู้ เพื่อสร้างระบบจักรกลเรียนรู้ที่ทนทานต่อการโจมตีได้

#### 1.6 คำนิยามศัพท์เฉพาะ

- **จักรกลเรียนรู้ (machine learning)** คือระบบ หรือโค้ด หรือโปรแกรมคอมพิวเตอร์ที่เรียนรู้โครงสร้างของชุดคำถาม และคำตอบโดยมีจำเป็นต้องทำการโปรแกรมลำดับการทำงานอย่างชัดเจน (explicitly)
- **การเรียนรู้เชิงโจมตี (adversarial learning)** หมายถึงการศึกษาถึงการโจมตีแบบจำลอง (model) ของจักรกลเรียนรู้ (machine learning)

## บทที่ 2

### ทฤษฎีและงานวิจัยที่เกี่ยวข้อง

#### 2.1 จักรกลเรียนรู้

ระบบจักรกลเรียนรู้ (machine learning) อาจนิยามได้ว่าเป็นระบบที่ไม่ต้องมีการป้อนข้อมูล หรือวิธีทำงาน เข้าไป ยังโค้ดโปรแกรมอย่างชัดเจน (explicitly) โดยระบบดังกล่าวจะถูกฝึกสอนด้วยชุดของข้อมูลหรือประสบการณ์ (experience) และปรับตัวเองให้ส่งออกคำตอบซึ่งอิงจากประสบการณ์ที่ตนเองเคยได้เรียนรู้มา

หากจะกล่าวให้ละเอียด เราสามารถนิยามโปรแกรมซึ่งสามารถทำการเรียนรู้ได้ดังนี้ [1]

**บทนิยาม 2.1.1.** โปรแกรมใดๆ เรียน (learn) จากประสบการณ์ (experience)  $E$  บนงาน (task)  $T$  และการวัดประสิทธิผล (performance measurement)  $P$  หากประสิทธิผลบน  $T$  ซึ่งถูกวัดโดย  $P$  เพิ่มขึ้นตามประสบการณ์  $E$

#### 2.2 การเรียนรู้เชิงลึก

การเรียนรู้เชิงลึก (Deep Learning) คือความพยายามในการจำลองเซลล์ประสาทของมนุษย์ให้อยู่ในรูปแบบจำลอง คณิตศาสตร์ ด้วยความเชื่อทางหลักประสาทวิทยา (neurosciences) ว่าความฉลาดของสมองมนุษย์เกิดขึ้นได้จากโครงข่าย ประสาทจำนวนมาก ที่เชื่อมเข้าถึงกัน [2]

##### 2.2.1 เพอร์เซปตรอน (Perceptron)

เพอร์เซปตรอน (Perceptron) [3] เป็นแบบจำลองทางคณิตศาสตร์ของเซลล์สมองหนึ่งเซลล์ โดยมีคุณสมบัติดังนี้

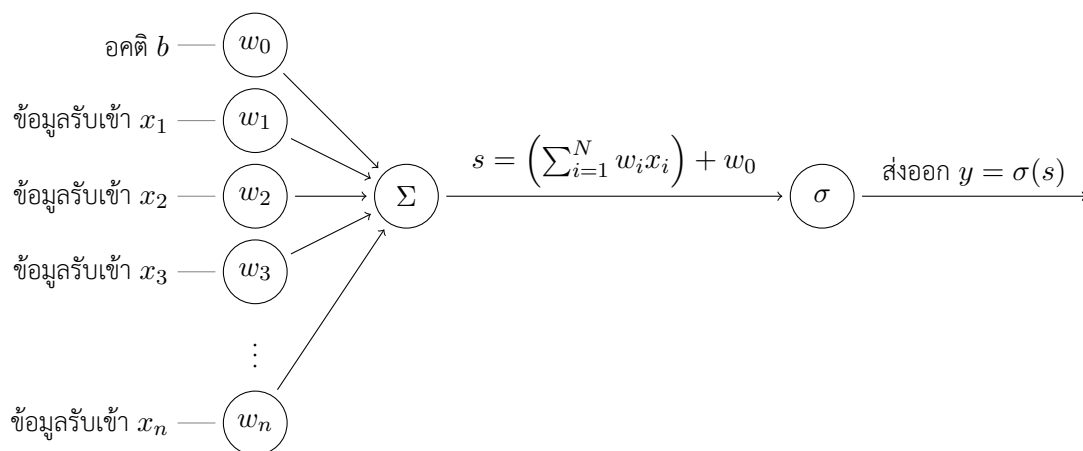
- รับเข้าข้อมูลมาในเซลล์จากหลายแหล่ง และให้น้ำหนักกับข้อมูลนั้นต่างกันไป
- ส่งออกข้อมูลเพียงค่าเดียว

ดังนั้น แบบจำลองทางคณิตศาสตร์สามารถเขียนออกมาจากหลักการสองข้อดังกล่าวได้ด้วยสมการ

$$y = f(W^T X + b)$$

เมื่อ  $W$  และ  $X$  เป็นเมทริกซ์ขนาด  $1 \times n$  (โดย  $n$  เป็นจำนวนข้อมูลรับเข้า),  $b$  เป็นค่าสัมประสิทธิ์คงที่ (อคติ: bias) และ  $f$  เป็นฟังก์ชันกระตุ้น (activation function) ซึ่งอาจเขียนรูปร่างของเพอร์เซปตรอนให้มีลักษณะรูปคล้ายเซลล์ สมองได้ในลักษณะรูปที่ 2.1

ยกตัวอย่างการใช้เพอร์เซปตรอนในการแก้ปัญหาอย่างง่ายได้ในที่นี้



รูปที่ 2.1: เพอร์เซปตรอน

### การคาดเดาราคาส่งหาหมักรัพย์

หากสำรวจราคาอสังหาริมทรัพย์แล้วพบว่า

- ราคาอสังหาริมทรัพย์จะเพิ่มขึ้นตามที่ดิน โดยเพิ่มขึ้นทุก 10,000 บาทต่อตารางวา
- ราคาอสังหาริมทรัพย์จะเพิ่มขึ้นตามจำนวนห้องนอน โดยเพิ่มขึ้นทุก 200,000 บาทต่อห้องนอน
- ราคาอสังหาริมทรัพย์จะลดลงตามจำนวนอายุปี โดยลดลงทุก 7,000 บาทต่ออายุของอสังหาริมทรัพย์
- ราคาที่กำหนดจริง (fixed cost) ของอสังหาริมทรัพย์ อยู่ที่ 100,000 บาท

จะสามารถเขียนเพอร์เซปตรอนเพื่อคาดเดาราคาส่งหาหมักรัพย์ได้โดย

$$y = \sigma(W^T X)$$

เมื่อ  $W$  ซึ่งเป็นค่าสัมประสิทธิ์แสดงถึงความสัมพันธ์ข้อมูลรับเข้า ซึ่งเขียนได้จากความสัมพันธ์ดังแสดงด้านล่าง

$$W^T = \begin{bmatrix} 100000 & 10000 & 200000 & -7000 \end{bmatrix}$$

หากต้องการคาดเดาราคาบ้านที่มี 3 ห้องนอน เนื้อที่ 100 ตารางวา และมีอายุ 7 ปี จะสามารถเขียนเมทริกซ์  $X$  ได้เป็น

$$X = \begin{bmatrix} 1 \\ 3 \\ 100 \\ 7 \end{bmatrix}$$

โปรดสังเกตว่า  $x_0 = 1$  เนื่องจากผลคูณของเทอม  $w_0$  และ  $x_0$  เป็นค่าที่เรียกว่าค่าอคติ (bias) ของแบบจำลอง

เนื่องจากเพอร์เซปตรอนตัวนี้ถูกใช้ในการทำนายราคา ซึ่งกล่าวว่ามีความสัมพันธ์กันกับตัวแปรที่กำหนดข้างต้นในเชิงเส้น ดังนั้นจะกล่าวได้ว่าฟังก์ชันกระตุ้น (activation function) ที่เลือกใช้ จะเลือกใช้ฟังก์ชันเส้นตรง (linear function)  $\sigma(x) = x$

ดังนั้น ผลการทำนายราคาบ้านคำนวณได้จาก

$$\begin{aligned}
 y &= \sigma(W^T X + b) \\
 &= \sigma \left( \begin{bmatrix} 100000 & 10000 & 200000 & -7000 \end{bmatrix} \times \begin{bmatrix} 1 \\ 3 \\ 100 \\ 7 \end{bmatrix} \right) \\
 &= \sigma(100000 + 30000 + 20000000 + (-49000)) = f(20981000) \\
 &= 20981000
 \end{aligned}$$

### การสร้างประตูลัษณตรรกะด้วยเปอร์เซปตรอน

เราสามารถสร้างประตูลัษณตรรกะ (logic gates) บางชนิดได้ด้วยเปอร์เซปตรอน เช่นการสร้าง AND และ OR gate

ยกตัวอย่างโครงสร้างของประตูลัษณตรรกะและซึ่งสามารถสร้างได้ด้วยการกำหนดให้

- $X$  เป็นเมทริกซ์ขนาด  $1 \times 2$  กล่าวคือเมื่อรับค่า  $x_1, x_2$  เป็นค่า 0 หรือ 1 แทนสัญญาณจริงหรือเท็จแล้ว

$$X = \begin{bmatrix} 1 \\ a_1 \\ a_2 \end{bmatrix}$$

- กำหนดค่าของเมทริกซ์  $W$  เป็น

$$W^T = \begin{bmatrix} -2 & 1 & 1 \end{bmatrix}$$

- กำหนดฟังก์ชัน  $\sigma(x)$  เป็น step function กล่าวคือ

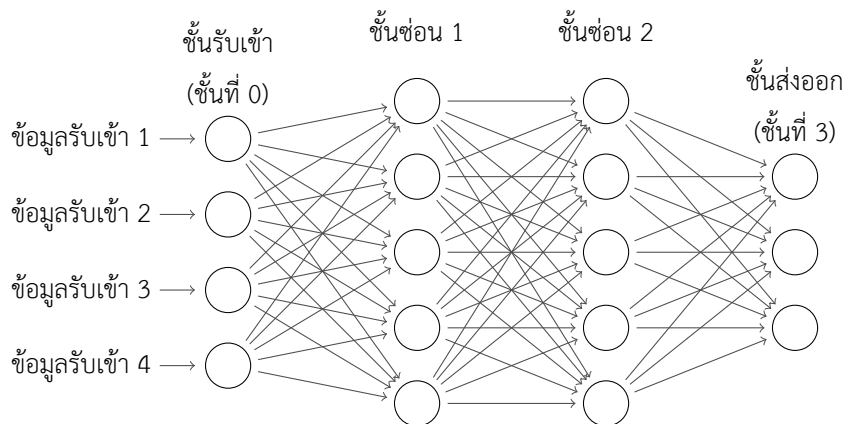
$$\sigma(x) = \begin{cases} 1; & x \geq 0 \\ 0; & \text{ในกรณีอื่น} \end{cases}$$

และการสร้างประตูลัษณตรรกะหรือสามารถทำได้ในลักษณะเดียวกันโดยเปลี่ยนชุดน้ำหนัก เป็น

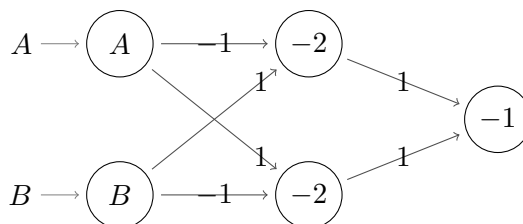
$$W^T = \begin{bmatrix} -1 & 1 & 1 \end{bmatrix}$$

### 2.2.2 เปอร์เซปตรอนแบบหลายชั้น (Multi Layer Perceptron)

เราอาจสังเกตว่าเปอร์เซปตรอนหนึ่งตัวนั้นทำหน้าที่ได้เพียงแยก (classify) หรือถดถอย (regress) ปัญหาที่เป็นปัญหาเชิงเส้น (linear problems) ได้เท่านั้น อย่างไรก็ตามหากเรากำหนดให้ฟังก์ชัน  $f$  เป็นฟังก์ชันที่ไม่ใช่ฟังก์ชันเส้นตรงแล้ว เราอาจสร้างเปอร์เซปตรอนแบบหลายชั้น (Multi Layer Perceptron) ขึ้นมาได้โดยมีลักษณะดังรูปที่ 2.2



รูปที่ 2.2: เพอร์เซปตรอนแบบหลายชั้น



รูปที่ 2.3: เพอร์เซปตรอนแบบหลายชั้นซึ่งทำหน้าที่เป็นประตูสัญญาณเฉพาะหรือ

เราอาจเขียนแทนน้ำหนักของโครงข่ายจากเพอร์เซปตรอนชั้นที่  $i$  ไปยังชั้นที่  $j$  ( $j = i + 1$ ) ได้เป็น

$$\mathbf{W}_{ij} = \begin{bmatrix} w_{10} & w_{20} & \dots & w_{n_i 0} \\ w_{11} & w_{21} & \dots & w_{n_i 1} \\ w_{12} & w_{22} & \dots & w_{n_i 2} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1n_j} & w_{2n_j} & \dots & w_{n_j n_i} \end{bmatrix}$$

เมื่อจำนวนเพอร์เซปตรอนในชั้นที่  $k$  เขียนแทนด้วย  $n_k$

ยกตัวอย่างเช่น เราจะสามารถสร้างประตูสัญญาณเฉพาะหรือ (XOR gate) ได้จากเพอร์เซปตรอนแบบหลายชั้น ดังแสดงในรูปที่ 2.3 โดยเลขในแต่ละเพอร์เซปตรอนแทนค่าอคติ ( $b$ ) และเลขบนเส้นเชื่อมแทนค่าน้ำหนัก ( $w$ ) และกำหนดให้ฟังก์ชันกระตุ้น  $\sigma$  เป็นฟังก์ชันขั้นบันได (step function) กล่าวคือ

$$\sigma(x) = \begin{cases} 1; & x \geq 0 \\ 0; & \text{ในกรณีอื่น} \end{cases}$$

เพอร์เซปตรอนดังกล่าว เมื่อรับค่า  $A$  และ  $B$  เป็น 0 หรือ 1 จะส่งออกมา  $A \oplus B$

## 2.3 ฟังก์ชันกระตุ้นและความฉลาดของโครงข่ายประสาทเทียม

### 2.3.1 ทฤษฎีบทตัวประมาณฟังก์ชันครอบจักรวาล

เหตุผลที่โครงข่ายประสาทเทียมสามารถทำงานได้ดี เนื่องจากมีการพิสูจน์ว่าโครงข่ายประสาทเทียมนั้นสามารถทำหน้าที่เป็นตัวประมาณฟังก์ชันครอบจักรวาล [4] (universal function approximator) กล่าวคือโครงข่ายประสาทเทียม  $N : \mathbb{R}^k \rightarrow \mathbb{R}^n$  ที่มีความซับซ้อนมากเพียงพอ (ซึ่งจะกล่าวถึงความซับซ้อนนี้ในภายหลัง) สามารถที่จะจำลองฟังก์ชัน  $f : \mathbb{R}^k \rightarrow \mathbb{R}^n$  (กล่าวคือฟังก์ชันที่มีโดเมน และเรนจ์ เป็นจำนวนจริงใดๆ ในมิติที่เหมือนกับมิติข้อมูลรับเข้าและข้อมูลส่งออกของโครงข่ายประสาทเทียม  $N$ ) ได้ [5] [6] [7]

บทพิสูจน์ของทฤษฎีบทนี้ทั้งในรูปแบบของกรณีไม่ตีกรอบความกว้าง (unbounded width case) และกรณีตีกรอบความกว้าง (bounded width case) สามารถศึกษาได้จากแหล่งอ้างอิง รวมถึงแหล่งอ้างอิงเพิ่มเติมที่ใช้การแสดงทัศนภาพ (visualisation) เพื่อการพิสูจน์ทฤษฎีบทดังกล่าว [8]

### 2.3.2 ข้อสังเกตต่อฟังก์ชันกระตุ้นและความฉลาด

บทพิสูจน์ที่ได้กล่าวถึงไปก่อนหน้านี้สำหรับกรณีไม่ตีกรอบความกว้าง และตีกรอบความกว้าง เป็นบทพิสูจน์ที่ใช้ฟังก์ชันกระตุ้นเป็นฟังก์ชันซิกมอยด์ (sigmoid) และฟังก์ชันรีลู (ReLU) ตามลำดับ

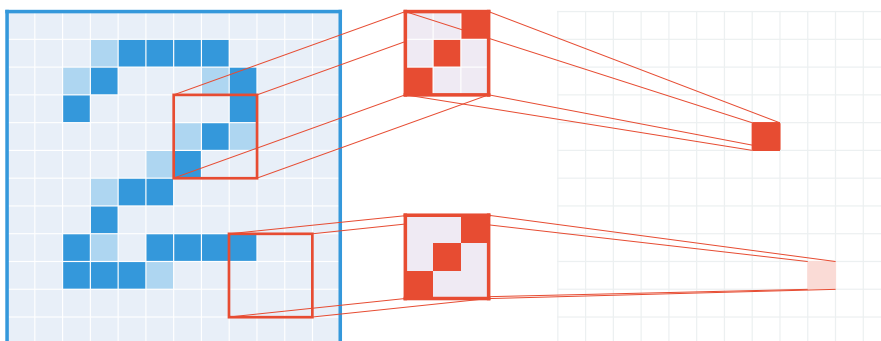
อย่างไรก็ดี หากพิจารณาโครงข่ายประสาทเทียมใดๆ ที่ใช้ฟังก์ชันกระตุ้นเป็นฟังก์ชันเชิงเส้น  $f(x) = x$  เราจะพบว่าโครงข่ายประสาทเทียมใดๆ จะสามารถยุบให้อยู่ในรูปของเปอร์เซปตรอนเพียงตัวเดียว และทำให้ไม่สามารถตัดสินใจปัญหาได้มากกว่าปัญหาที่แบ่งแยกเชิงเส้นได้ (linearly separable problems)

ดังนั้น อาจกล่าวด้วยการพิจารณา (intuition) ในลักษณะดังกล่าวได้ว่า ส่วนหนึ่งของความเป็นไปได้ของการที่โครงข่ายประสาทเทียมใดๆ สามารถทำหน้าที่เป็นตัวประมาณฟังก์ชันครอบจักรวาลได้ ส่วนหนึ่งมาจากการที่ฟังก์ชันกระตุ้นทำหน้าที่เป็นตัวบีบ (squeezer) ช่วงของข้อมูลรับเข้าบนโดเมนจำนวนจริงใดๆ ( $\mathbb{R}$ ) ให้กลายเป็นช่วงจำกัดช่วงอื่น (เช่น ช่วง  $(0, 1)$  ของฟังก์ชันซิกมอยด์ หรือช่วง  $[0, \infty)$  ของฟังก์ชันรีลู)

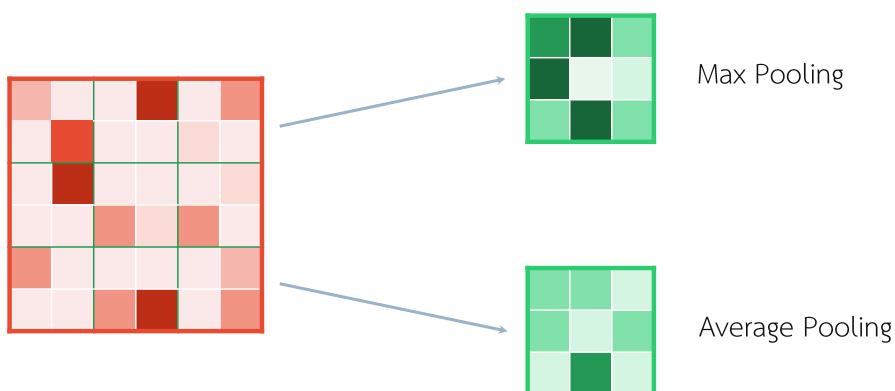
## 2.4 โครงข่ายประสาทเทียมแบบสังวัฒนาการ

โครงข่ายประสาทเทียมแบบสังวัฒนาการ (Convolutional Neural Networks: CNN) [9] เป็นโครงข่ายประสาทเทียมซึ่งมักถูกใช้กับข้อมูลภาพ [10] โดยคร่าวแล้วโครงข่ายประสาทเทียมในลักษณะดังกล่าวมักประกอบด้วยชั้นประสาทเทียมในลักษณะดังนี้

- **ชั้นสังวัฒนาการ** (convolution layer) เป็นชั้นที่กระทำตัวดำเนินการสังวัฒนาการ (convolve) ตัวกรอง (filter)  $F$  บนข้อมูลนำเข้า  $I$  ด้วยระยะเคลื่อน (stride)  $S$  ผลลัพธ์จากการสังวัฒนาการนี้จะเรียกว่าแผนที่ลักษณะ (feature map) ยกตัวอย่างการสังวัฒนาการเพื่อหาเส้นเฉียงในรูปที่ 2.4 สังเกตว่าการสังวัฒนาการด้วยตัวกรองเส้นเฉียงบนเส้นเฉียงบริเวณข้อมูลนำเข้า จะให้ค่าส่งออกที่มากกว่าการสังวัฒนาการตัวกรองเส้นเฉียงบนจุดพื้นที่อื่นของข้อมูลนำเข้า (ในที่นี้เขียนแทนด้วยสีแดงเข้ม และสีแดงอ่อน)
- **ชั้นบ่อรวม** (pooling layer) เป็นชั้นที่ทำการสุ่มตัวอย่างแบบลดขนาด (downsampling) เพื่อลดขนาดของข้อมูล ในขณะที่ยังคงไว้ซึ่งชุดคุณสมบัติที่ข้อมูลรับเข้ามี ชั้นบ่อรวมอาจแบ่งเป็นสองประเภทหลัก
  - **ชั้นบ่อรวมแบบมากที่สุด** (maximum pooling layer) เป็นชั้นบ่อรวมที่พบได้บ่อยที่สุด



รูปที่ 2.4: ชั้นสังวัตนาการ ซึ่งแสดงข้อมูลนำเข้าด้วยสีฟ้า และตัวกรองด้วยสีแดง



รูปที่ 2.5: ชั้นบ่อรวม ทั้งแบบบ่อรวมมากที่สุดและแบบบ่อรวมเฉลี่ย โดยพิจารณาบ่อตามขอบเขตสี่เหลี่ยม

- ชั้นบ่อรวมแบบเฉลี่ย (average pooling layer) เป็นชั้นบ่อรวมที่พบในโครงข่ายประสาทเทียมแบบสังวัตนาการบางรูปแบบ เช่น LeNet

การสังวัตนาการของชั้นสังวัตนาการในโครงข่ายประสาทเทียม ทำหน้าที่เป็นตัวตรวจจับคุณสมบัติ (feature detector) เช่นการตรวจจับขอบ (edge detection) และชั้นบ่อรวมทำให้ขนาดของผลลัพธ์จากชั้นสังวัตนาการมีขนาดเล็กลง เพื่อให้จำนวนค่าน้ำหนักของโครงข่ายประสาทเทียมที่ต้องคำนวณนั้นน้อยลง

## 2.5 การเรียนรู้ด้วยวิธีก้าวเคลื่อนถอยหลัง

การเรียนรู้ด้วยวิธีก้าวเคลื่อนถอยหลัง (backpropagation learning)...



## บรรณานุกรม

- [1] T. M. Mitchell, *Machine Learning*. New York: McGraw-Hill, 1997.
- [2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [3] F. Rosenblatt, "The perceptron: A probabilistic model for information storage and organization in the brain.," *Psychological Review*, vol. 65, no. 6, p. 386–408, 1958.
- [4] K. Hornik, "Approximation capabilities of multilayer feedforward networks," *Neural Networks*, vol. 4, no. 2, p. 251–257, 1991.
- [5] G. Cybenko, "Approximation by superpositions of a sigmoidal function," *Mathematics of Control, Signals, and Systems*, vol. 2, pp. 303–314, Dec. 1989.
- [6] M. Leshno, V. Y. Lin, A. Pinkus, and S. Schocken, "Multilayer feedforward networks with a nonpolynomial activation function can approximate any function," *Neural Networks*, vol. 6, pp. 861–867, Jan. 1993.
- [7] A. Kratsios, "Universal approximation theorems," 2019.
- [8] M. A. Nielsen, "Neural networks and deep learning," Jan 1970.
- [9] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," in *Proceedings of the IEEE*, vol. 86, pp. 2278–2324, 1998.
- [10] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1, NIPS'12, (USA)*, pp. 1097–1105, Curran Associates Inc., 2012.