

บทที่ 4

การเรียนรู้แบบย้อนกลับ

4.1 ประวัติของการเรียนรู้แบบย้อนกลับ

ความคิดพื้นฐานแรกเริ่มของการเรียนรู้แบบย้อนกลับ (Backpropagation หรือ BP ซึ่งนิยมเรียกกันโดยทั่วไป) ถูกนำเสนอในปี ค.ศ. 1974 โดย Paul Werbos จากนั้นก็มีการนำเสนอใหม่ โดย David Parker ในปี ค.ศ. 1982 และ ถูกนำเสนอต่อผู้อ่านอย่างกว้างขวาง โดย Rumelhart And McClelland ในปี ค.ศ. 1986 ในหนังสือที่ชื่อ Parallel Distributed Processing ซึ่งกล่าวถึง ศักยภาพของนิวรัล เนทเวิร์ค และ Back-Error Propagation Network

งานประยุกต์แรกเริ่มของการเรียนรู้แบบย้อนกลับ สร้างโดย Terry Sejnowski และคณะในมหาวิทยาลัย Johns Hopkins โดยเป็นโปรแกรม NETTALK สร้างโดย Terry Sejnowski และ Rosenberg

การเรียนรู้แบบย้อนกลับ เป็นรูปแบบหนึ่งของนิวรัล เนทเวิร์ค ที่ใช้กันอย่างกว้างขวาง และถูกนำมาประยุกต์ใช้กับการศึกษางานประเภทต่างๆ เช่น ทางการทหาร การช่วยในการวินิจฉัยทางการแพทย์ การจดจำคำพูด การจดจำตัวอักษร หุ่นยนต์

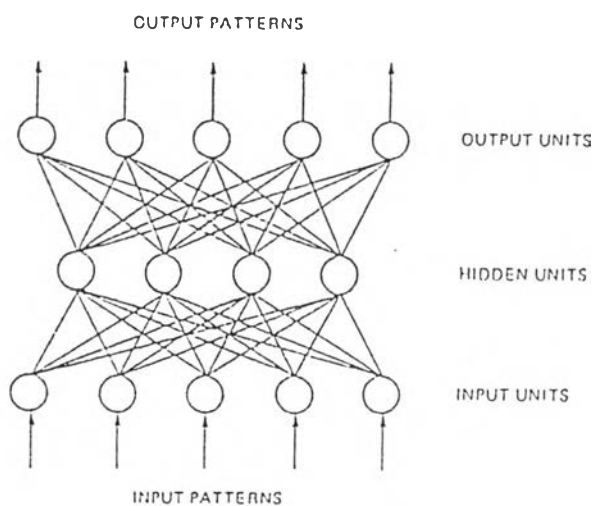
การเรียนรู้แบบย้อนกลับ สามารถแก้ปัญหาที่ต้องการรูปแบบ โดยการป้อนรูปแบบเข้าไป นิวรัล เนทเวิร์ค ก็จะผลิตรูปแบบผลลัพธ์ที่เกี่ยวข้องกันออกมา (Judith E. Dayhoff, 1990)

การเรียนรู้แบบย้อนกลับ เป็นวิธีการหนึ่งของ นิวรัล เนทเวิร์ค ที่ง่ายต่อการเข้าใจ เนื่องจากกระบวนการเรียนรู้และปรับปรุงแก้ไขนั้นเป็นไปด้วยตัวมันเอง ถ้าเน็ตเวิร์คให้คำตอบที่ผิด ดังนั้น ค่าน้ำหนัก จะถูกแก้ไขจนกว่า ค่าความผิดพลาดจะลดน้อยลง หรืออยู่ในเกณฑ์ที่ยอมรับได้ นั่นก็คือค่าได้ในครั้งต่อไปจะมีความถูกต้องมากยิ่งขึ้น

4.2 โครงสร้างของการเรียนรู้แบบย้อนกลับ

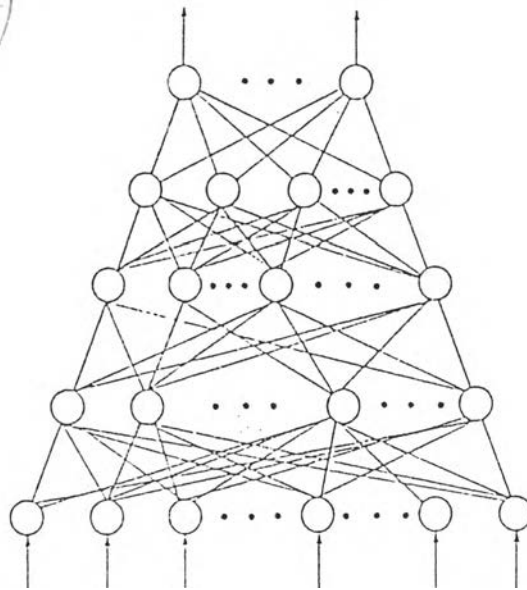
การเรียนรู้แบบย้อนกลับ มีโครงสร้างเป็นชั้นๆ โดยมีโครงสร้างอย่างง่าย 3 ชั้น คือ ชั้นข้อมูลเข้า ชั้นแอบแฝง ชั้นแสดงผลลัพธ์ แต่ละชั้น จะติดต่อกันอย่างสมบูรณ์

ดังรูป 4-1



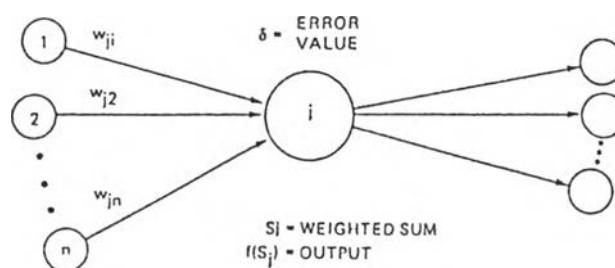
รูป 4-1 แสดงการเรียนรู้แบบย้อนกลับ ชนิด 3 ชั้น

จากรูป 4-1 แสดงรูปแบบ(Topology) ของการเรียนรู้แบบย้อนกลับ ชนิด 3 ชั้น ชั้นล่างคือชั้นข้อมูลเข้า ซึ่งนำเข้าข้อมูลจากภายนอก ชั้นถัดมาคือชั้นแอบแฝง ซึ่งจะติดต่อกับทั้ง ชั้นข้อมูลเข้า และ ชั้นผลลัพธ์ ที่อยู่ชั้นบน



รูป 4-2 แสดงการเรียนรู้แบบย้อนกลับ ชนิด 5 ชั้น

เป็นตัวอย่างของการเรียนรู้แบบย้อนกลับ ชนิด 5 ชั้น ที่ติดต่อกันอย่างสมบูรณ์ โดยมี
ชั้นแอบแฝง อยู่ 3 ชั้น



รูป 4-3 แสดงแต่ละโหนดภายในชั้นแอบแฝง

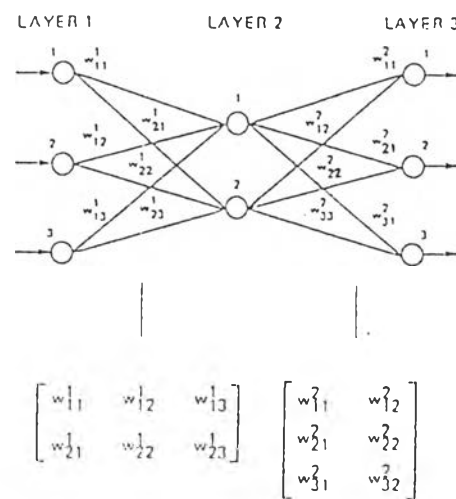
จะเห็นว่าข้อมูลนำเข้าอยู่ทางด้านซ้าย ส่วนทางด้านขวาคือหน่วยที่รับผลลัพธ์
จากโหนด j

หน่วยประมวลผลกลาง(Processing Unit) ประกอบด้วย

- ผลรวมถ่วงน้ำหนักของข้อมูลนำเข้า (S_j)
- ผลลัพธ์ (X_j)
- ค่าความผิดพลาด (δ_j) ซึ่งถูกใช้ระหว่างการการปรับค่าน้ำหนัก

ค่าน้ำหนักที่เกี่ยวข้องเนื่องกับการติดต่อแต่ละชั้น จะถูกปรับแต่งตลอดเวลาระหว่างการ
เรียนรู้ เพื่อให้ลดความแตกต่างระหว่างผลลัพธ์ที่ได้กับผลลัพธ์ที่ต้องการ

กำหนดให้ w_{ji}^i หมายถึง ค่าน้ำหนักจากหน่วยที่ i ไปยัง หน่วยที่ j



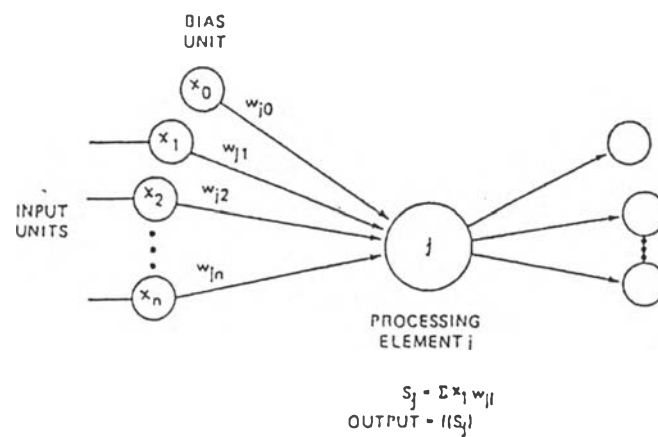
รูป 4-4 แสดงการแทนค่าน้ำหนักด้วยเมตริกซ์

จากรูป เป็นการแทนค่าน้ำหนักด้วย เมตริกซ์ โดยมีค่ากำกับบน (Superscript) บ่งบอกถึงชั้นที่ต่าง ๆ กัน เช่น w_{21}^1 หมายถึงค่าน้ำหนักจากหน่วยประมวลผลที่ 2 ในระดับชั้นที่ 1 ไปยัง หน่วยประมวลผลที่ 1 ในระดับที่ 2

4.3 ขั้นตอนการเรียนรู้ของการเรียนรู้แบบย้อนกลับ

ประกอบด้วย 2 ขั้นตอน

4.3.1 การแพร่เดิหน้า (Forward-Propagation)



รูป 4-5 แสดงขั้นตอนของการแพร่เดิหน้า

จากรูปเป็นการคำนวณผลรวมของผลลัพธ์ที่เข้ามายังหน่วยที่ j ดังสมการ

$$S_j = \sum_i X_i w_{ji} \quad \text{----- (4.1)}$$

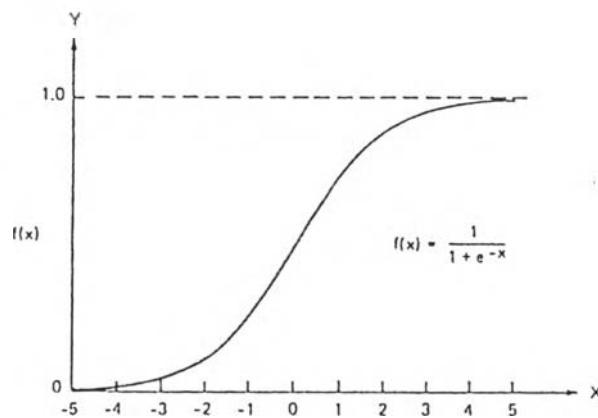
โดยที่ X_i = ข้อมูลจากหน่วยที่ i

w_{ji} = ค่าน้ำหนักจากหน่วยที่ i ไปยังหน่วยที่ j

เมื่อกำหนดจนได้ค่า S_j ก็ทำการคำนวณหาค่า $f(S_j)$ อีกครั้ง

โดย $f(x) = \frac{1}{1 + e^{-x}}$ ซึ่งเป็นสมการของฟังก์ชันซิกมอยด์

$$\text{ดังนั้น } f(S_j) = \frac{1}{1 + e^{-S_j}} \quad \text{----- (4.2)}$$



รูป 4-6 กราฟของฟังก์ชันซิกมอยด์

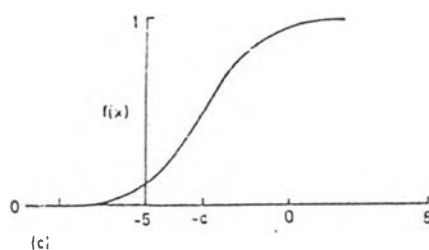
เมื่อได้ค่า $f(S_j)$ แล้ว ค่า $f(S_j)$ ก็จะกลายเป็นผลลัพธ์ของหน่วยที่ j ซึ่งก็คือค่า X_j ดังรูป 4-5 โดยจะส่งออกไปทางหน่วยอื่นๆ ด้วยค่า X_j ที่เท่ากัน

จากรูป 4-5 ค่าเอนเอียง (Bias Unit) เป็นค่าคงที่ที่เราใส่เข้าไป เพื่อให้การเรียนรู้ ของเน็ตเวิร์คเร็วขึ้น หรือที่เรียกว่า Convergence Time เร็วขึ้น กล่าวคือ ค่าเอนเอียง ยังมีผลต่อระดับการกระตุ้น (Threshold) ด้วย เช่น ให้ค่า $C = 5$ ทำให้กราฟขยับไปทางซ้าย 5 หน่วย ดังรูป 4-7

$$C = W_{j0}$$

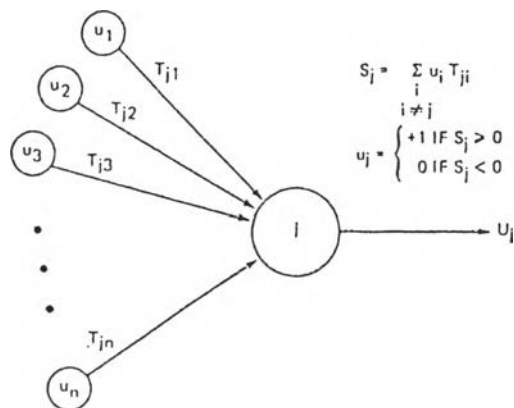
ให้ $S = \sum x_i W_{ji}$

ดังนั้น ผลรวม $= S + C$ ----- (4.3)



รูป 4-7 แสดงผลของค่าเอนเอียงที่มีต่อการเรียนรู้

สาเหตุที่ต้องใช้ฟังก์ชันซิกมอยด์ ก็เนื่องจาก ต้องการให้เป็นระดับการกระตุ้นแบบอ่อน (Soft Threshold) มากกว่าการกระตุ้นแบบรุนแรง (Hard Threshold) (ดังเช่น ฟังก์ชันขั้นบันได ดังรูป 4-8) นั่นคือ ฟังก์ชันซิกมอยด์ให้ค่าที่ต่อเนื่องกัน



รูป 4-8 แสดงนิวรัลเน็ตเวิร์คแบบที่ใช้ฟังก์ชันขั้นบันได

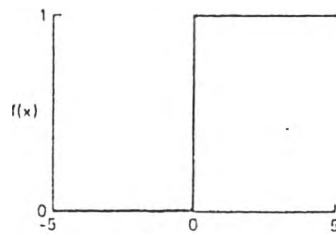
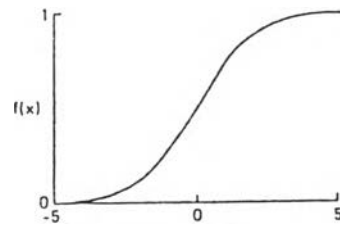
จากรูป 4-8 ฟังก์ชันขั้นบันได คือ

$$S_j \geq 0 \quad \text{ถ้า } U_j = 1$$

$$S_j < 0 \quad \text{ถ้า } U_j = 0$$

$$S_j = \sum_{i \neq j} U_i T_{ji}$$

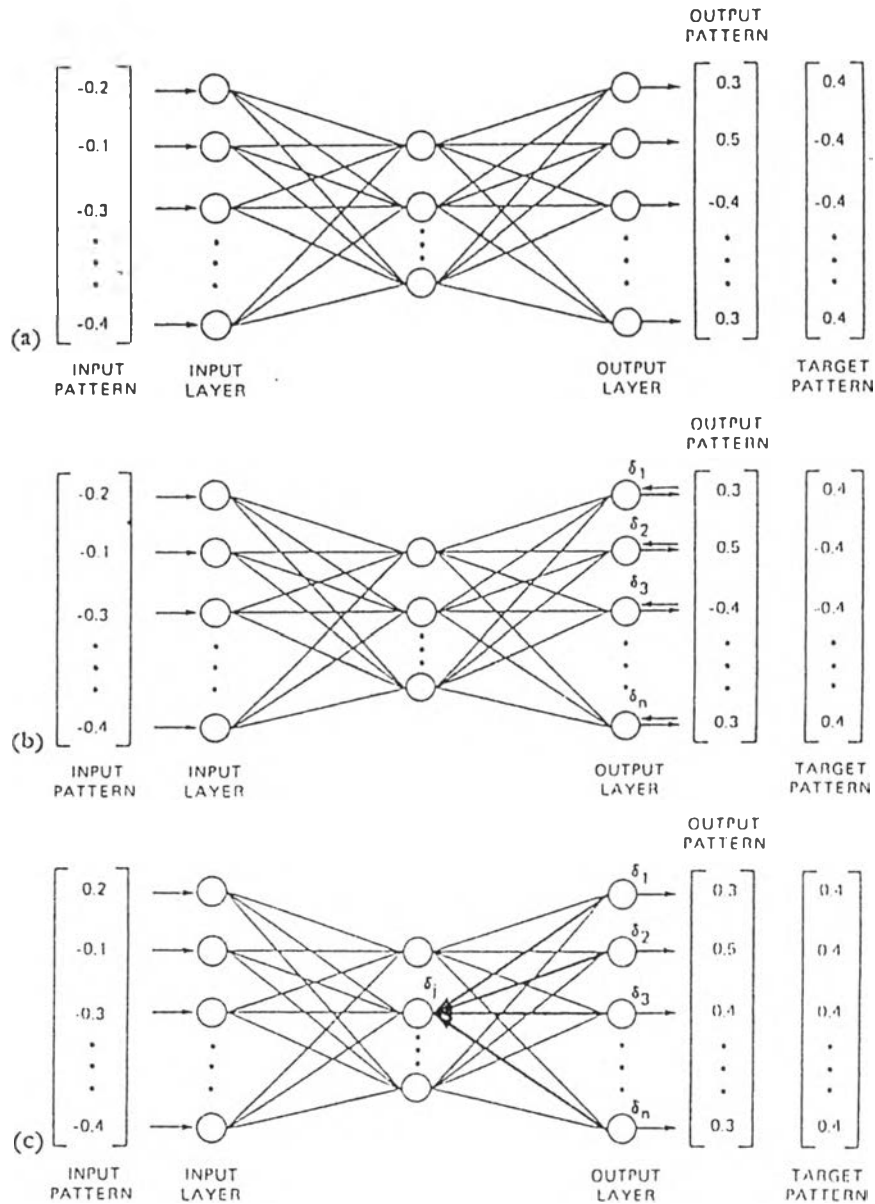
ซึ่งให้ผลลัพธ์ U_j 2 ค่า เท่านั้น คือ 0 , 1



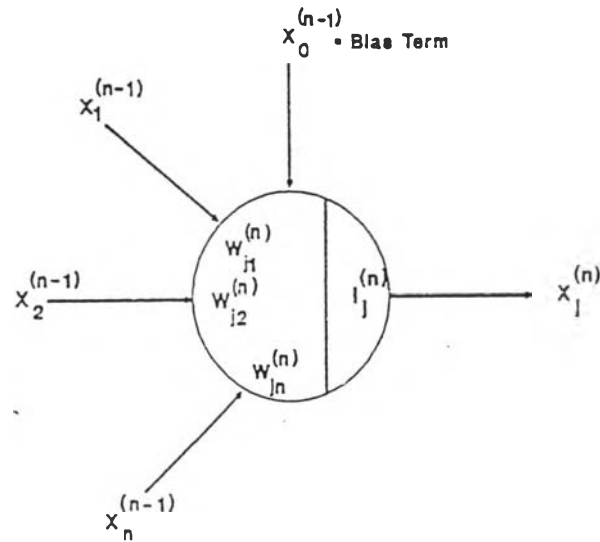
รูป 4-9 กราฟเปรียบเทียบ ฟังก์ชันซิกมอยด์ และ ฟังก์ชันขั้นบันได



4.3.2 การแพร่ย้อนกลับ (Backward-Propagation)



รูป 4-10 แสดง การแพร่ย้อนกลับ



รูป 4-11 รายละเอียดในแต่ละโนด

ในขั้นตอนของการแพร่ขยายเงินหน้า สามารถพิจารณาแต่ละ โหนดหรือหน่วย
ประมวลผลโดยละเอียดได้ดังรูป และอธิบายได้ดังนี้

ให้ X_t = ชุดข้อมูลเข้าลำดับที่ t

$$X_t = \{ x_{t,1}, x_{t,2}, \dots, x_{t,N} \} \quad \text{----- (4.4)}$$

Y_t = ผลลัพธ์ที่เกิดขึ้นจริงจากการใส่ชุดข้อมูลเข้า X_t

$$Y_t = \{ y_{t,1}, y_{t,2}, \dots, y_{t,N} \} \quad \text{----- (4.5)}$$

D_t = ผลลัพธ์ที่ต้องการจากการใส่ชุดข้อมูลเข้า X_t

$$D_t = \{ d_{t,1}, d_{t,2}, \dots, d_{t,N} \} \quad \text{----- (4.6)}$$

$X_{t,j}^{(n)}$ คือ ผลลัพธ์จากโหนดที่ j จากการใส่ชุดข้อมูลเข้าลำดับที่ t ในชั้นที่ n

$W_{j,i}^{(n)}$ คือ ค่าน้ำหนักจากโหนดที่ i ในชั้นที่ $(n-1)$ ไปยังโหนดที่ j ในชั้นที่ n

$I_{t,j}^{(n)}$ หรือ $S_{t,j}^{(n)}$ คือ ผลรวมของโหนดที่ j ในชั้นที่ n จากการใส่ชุดข้อมูลเข้าลำดับที่ t

หมายเหตุ

การใช้สัญลักษณ์ $I_{t,j}^{(n)}$ แทน $S_{t,j}^{(n)}$ เนื่องจากค่าผลรวมที่ได้นั้นจะ

กลายเป็นข้อมูลเข้า (Input) ของโหนดในชั้นถัดไป และต่อไปนี้จะใช้ $I_{t,j}^{(n)}$ แทน $S_{t,j}^{(n)}$ ในการพิสูจน์ในตอนต่อไป

ดังนั้น ค่าผลลัพธ์ของโหนดที่ j ในชั้นที่ n จากการใส่ชุดข้อมูลเข้าลำดับที่ t สามารถเขียนได้ดังนี้

$$X_{t,j}^{(n)} = f \left[\sum_{i=1}^{N_{n-1}} W_{j,i}^{(n)} X_{t,i}^{(n-1)} \right] = f(S_{t,j}^{(n)}) = f(I_{t,j}^{(n)}) \text{ ----- (4.7)}$$

ซึ่ง N_{n-1} คือ จำนวนโหนดในชั้นที่ $(n-1)$

4.3.2.1 ค่าความผิดพลาดท้องถิ่นของการแพร่ย้อนกลับ

สมมติว่า เน็ตเวิร์ค มีค่าความผิดพลาดโดยรวม (Global Error) คือ E ซึ่งสามารถหาอนุพันธ์ได้ทุกค่าในเน็ตเวิร์ค และให้ $\sigma_{t,j}^{(n)}$ คือ ค่าความผิดพลาดท้องถิ่น (Local Error) ที่โหนด j ในชั้นที่ n ดังนั้นค่า $\sigma_{t,j}^{(n)}$ สามารถเขียนได้ดังนี้

$$\text{ให้ } \sigma_{t,j}^{(n)} = -\frac{\partial E}{\partial I_{t,j}^{(n)}} \quad \text{----- (4.8)}$$

$$\text{โดยกฎลูกโซ่จะได้ } \sigma_{t,j}^{(n)} = -\frac{\partial E}{\partial X_{t,j}^{(n)}} \frac{\partial X_{t,j}^{(n)}}{\partial I_{t,j}^{(n)}} \quad \text{----- (4.9)}$$

ดังนั้น ค่าความผิดพลาดท้องถิ่นของโหนดที่ j คือผลลัพธ์ของ 2 เทอมคือ

$$\text{เทอมที่ 1 } \frac{\partial E}{\partial X_{t,j}^{(n)}} \quad \text{คือ ค่าการเปลี่ยนแปลงค่า } E \text{ เปรียบเทียบกับค่า } X_{t,j}^{(n)}$$

$$\text{เทอมที่ 2 } \frac{\partial X_{t,j}^{(n)}}{\partial I_{t,j}^{(n)}} \quad \text{คือ ค่าการเปลี่ยนแปลงค่า } X_{t,j}^{(n)} \text{ เปรียบเทียบกับค่า } I_{t,j}^{(n)}$$

$$\text{ดังนั้นจากสมการ (4.7) จะได้ } \frac{\partial X_{t,j}^{(n)}}{\partial I_{t,j}^{(n)}} = f_j' (I_{t,j}^{(n)}) \quad \text{----- (4.10)}$$

ซึ่งคือค่าอนุพันธ์ของ f_j สำหรับโหนดที่ j

เทอมที่ 1 มี 2 กรณีคือ

กรณีที่ 1 สมมติว่าโหนดที่ j คือ โหนดผลลัพธ์ของผลลัพธ์ของเน็ตเวิร์ค และให้ค่าความผิดพลาดโดยรวม E คือ ค่า Mean Square Error (MSE) ของผลลัพธ์ที่ต้องการและผลลัพธ์ที่เกิดขึ้นจริง สามารถแสดงได้ดังนี้

$$E = 1/2 \sum_{j=1}^{N_y} (d_{t,j} - y_{t,j})^2 \quad \text{----- (4.11)}$$

N_y คือ จำนวนโหนดในชั้นผลลัพธ์ ดังนั้นค่าอนุพันธ์ที่ได้เมื่อเทียบกับ $y_{t,j}$ คือ

$$\frac{\partial E}{\partial y_{t,j}} = - (d_{t,j} - y_{t,j}) \quad \text{----- (4.12)}$$

จากสมการ (4.9) เมื่อรวมสมการ (4.10) และ (4.12) จะได้ ค่า δ สำหรับโหนดที่ j ใดๆ คือ

$$\delta_{t,j} = (d_{t,j} - y_{t,j}) f'_j(I_{t,j}) \quad \text{----- (4.13)}$$

กรณีที่ 2 ถ้าโหนดที่ j ไม่ใช่โหนดผลลัพธ์ของเน็ตเวิร์ค ค่าความผิดพลาดท้องถิ่นสามารถแทนด้วย

$$\begin{aligned} \sum_k \frac{\partial E}{\partial I_{t,k}^{(n+1)}} \frac{\partial I_{t,k}^{(n+1)}}{\partial X_{t,j}^{(n)}} &= \sum_k \frac{\partial E}{\partial I_{t,k}^{(n+1)}} \frac{\partial (\sum_i W_{k,i}^{(n+1)} X_{t,i}^{(n)})}{\partial X_{t,j}^{(n)}} \\ &= \sum_k \frac{\partial E}{\partial I_{t,k}^{(n+1)}} W_{k,j}^{(n+1)} = -\sum_k \delta_{t,k} W_{k,j}^{(n+1)} \quad \text{----- (4.14)} \end{aligned}$$

ดังนั้น จากสมการ (4.9) (4.10) จะได้

$$\delta_{t,j}^{(n)} = f_j' (I_{t,j}^{(n)}) \sum_k \delta_{t,k}^{(n+1)} w_{k,j}^{(n+1)} \quad \text{----- (4.15)}$$

ซึ่งเรียกว่า สมการ Generalized δ Rule (Rumelhart & McClelland, 1986)

เนื่องจาก f_j คือ ฟังก์ชันซิกมอยด์

จาก ฟังก์ชันซิกมอยด์ $f(z) = 1 / (1 + e^{-z})$ ----- (4.16)

$$f'(z) = e^{-z} / (1 + e^{-z})^2$$

$$= [1 / (1 + e^{-z})] - [1 / (1 + e^{-z})^2]$$

$$= 1 / (1 + e^{-z}) [1 - 1 / (1 + e^{-z})]$$

ดังนั้น $f'(z) = f(z) [1 - f(z)]$ ----- (4.17)

จากสมการ (4.15) สามารถเขียนได้ดังนี้

$$\delta_{t,j}^{(n)} = X_{t,j}^{(n)} (1 - X_{t,j}^{(n)}) \sum_k (\delta_k^{(n+1)} w_{k,j}^{(n+1)}) \quad \text{----- (4.18)}$$

ดังนั้น จุดประสงค์สำคัญของการเรียนรู้แบบย้อนกลับ คือ การแพร่คืนหน้าของ X_t จากชั้นข้อมูลเข้าไปยังชั้นผลลัพธ์ จากนั้นก็ทำการคำนวณค่าความผิดพลาด โดยเปรียบเทียบค่าผลลัพธ์ที่เกิดขึ้นจริงกับผลลัพธ์ที่ต้องการ และแพร่ค่าความผิดพลาดย้อนกลับจากชั้นผลลัพธ์ย้อนกลับมายังชั้นข้อมูลเข้า

4.3.2.2 การหาค่าความผิดพลาดโดยรวมให้น้อยที่สุด

จากแนวความคิดของ ค่าความผิดพลาดท้องถิ่นที่แต่ละโหนด หรือแต่ละหน่วยประมวลผล ค่าความผิดพลาดโดยรวม E สามารถทำให้ลดน้อยที่สุด โดยการปรับปรุงชุดของน้ำหนัก โดยใช้กฎของ Gradient Descent ดังนี้

$$\Delta W_{j,i}^{(n)} = -\eta \left(\frac{\partial E}{\partial W_{j,i}^{(n)}} \right) \quad \text{----- (4.19)}$$

η อ่านว่า อีตา (Eta) คือ ค่าสัมประสิทธิ์การเรีบนรู้

$$\text{จากกฎลูกโซ่จะได้} \quad \frac{\partial E}{\partial W_{j,i}^{(n)}} = \frac{\partial E}{\partial I_{t,j}^{(n)}} \frac{\partial I_{t,j}^{(n)}}{\partial W_{j,i}^{(n)}} \quad \text{----- (4.20)}$$

พิจารณาเทอมที่ 2

$$\begin{aligned} \frac{\partial I_{t,j}^{(n)}}{\partial W_{j,i}^{(n)}} &= \frac{\partial}{\partial W_{j,i}^{(n)}} \left[\sum_k W_{j,k} X_{t,k} \right] \\ &= X_{t,i} \end{aligned} \quad \text{----- (4.21)}$$

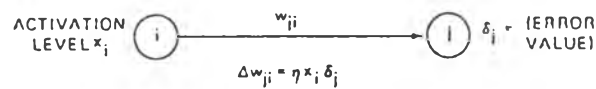
แทนค่า สมการ (4.21) ลงใน (4.20) จะได้

$$\frac{\partial E}{\partial W_{j,i}^{(n)}} = -\delta_{t,j} X_{t,i} \quad \text{----- (4.22)}$$

แทนค่า สมการ (4.22) ลงใน (4.19) จะได้

$$\Delta_t w_{j,i}^{(n)} = \eta \delta_{t,j}^{(n)} x_{t,j}^{(n-1)} \quad \text{----- (4.23)}$$

เป็นสมการการปรับปรุงค่าน้ำหนัก



รูป 4-12 แสดงการปรับปรุงค่าน้ำหนัก

4.3.2.3 ฟังก์ชันความผิดพลาดโดยรวม (The Global Error Function)

ให้ X_t คือ เวกเตอร์ข้อมูลเข้า (Input Vector)

D_t คือ เวกเตอร์ข้อมูลผลลัพธ์ที่ต้องการ (Desired Output Vector)

Y_t คือ เวกเตอร์ข้อมูลผลลัพธ์ที่เกิดขึ้นจริง (Actual Output Vector)

ค่าความผิดพลาดโดยรวมสามารถแสดงได้ดังเช่นสมการ (4.11) คือ

$$E = 1/2 \sum_{j=1}^{N_y} (d_j - y_j)^2 \quad \text{----- (4.24)}$$

$$\begin{aligned}
 \text{จากสมการ (4.8) } \delta_j &= - \frac{\partial E_j}{\partial I_j} \\
 &= - \frac{\partial E_j}{\partial y_j} \frac{\partial y_j}{\partial I_j} \\
 &= (d_j - y_j) f'(I_j)
 \end{aligned}$$

ซึ่งมีค่าเท่ากับสมการ (4.13) คือมีค่าเท่ากับความผิดพลาดท้องถิ่น

4.3.2.4 การใช้ค่าโนดเอนเอียง

จุดประสงค์ก็เพื่อให้การเรียนรู้ของเน็ตเวิร์คเร็วขึ้น และมีความยืดหยุ่นมากขึ้น

$$y_j = f \left[\sum_{i=1}^{N_x} w_{ji} X_i + \theta_j \right] \quad \text{----- (4.25)}$$

ถ้าเป็นสมการเชิงเส้นจะได้

$$y_j = \sum_{i=1}^{N_x} w_{ji} X_i + \theta_j \quad \text{----- (4.26)}$$

4.3.2.5 การใช้ค่าโมเมนตัม

นอกจากค่า η แล้ว ยังมีค่าโมเมนตัม ใช้สัญลักษณ์ α อ่านว่า อัลฟา (Alpha) ที่ช่วยให้การเรียนรู้เร็วขึ้น

ค่าโมเมนตัมนั้น คิดค้นโดย Rumelhart Hinton และ William (1986) โดยได้สร้างเทคนิคช่วยให้เวลาในการสอนเน็ตเวิร์คของการเรียนรู้แบบย้อนกลับเร็วขึ้น โดยการใส่ค่าโมเมนตัม เข้าไป ทำให้เน็ตเวิร์คมีความคงตัวมากขึ้น และได้สมการดังนี้

$$\Delta w_{ji}^{(n)} = \eta \delta_j x_i + \alpha \Delta w_{ji}^{(n-1)} \quad \text{----- (4.30)}$$

ดังนั้นจะได้

$$\Delta w_{ji}^{(n)} = \sum_{t=1}^T \Delta_t w_{ij}^{(n)} \quad \text{----- (4.31)}$$

ค่า k คือ ดัชนีเวลา (Time Index) หรือ จำนวนรอบการปรับปรุงค่าน้ำหนัก

T คือ จำนวนเวกเตอร์ข้อมูลเข้า หรือ จำนวนเวกเตอร์ผลลัพธ์

4.4 สรุปขั้นตอนของการเรียนรู้แบบย้อนกลับมาตรฐาน

สามารถสรุปเป็นขั้นตอนได้ดังนี้

ขั้นตอนที่ 1 กำหนดค่าเริ่มต้นให้กับค่าน้ำหนักและค่าความเอนเอียง

โดยการกำหนดค่าน้ำหนักทั้งหมดและค่าความเอนเอียงให้เป็นค่าเล็กๆ

โดยปกติอยู่ในช่วง $[-0.1, 0.1]$

ขั้นตอนที่ 2 กำหนดชุดข้อมูลเข้าและผลลัพธ์ที่ต้องการ

โดยการกำหนดชุดข้อมูลเข้าและระบุผลลัพธ์ที่ต้องการ ชุดข้อมูลเข้าสามารถกำหนดขึ้นใหม่ในการสอนแต่ละรอบ หรือ เลือกจากชุดการสอน โดยจะมีการกำหนดเป็นรอบๆ จนกระทั่งได้ค่าน้ำหนักที่คงที่

ขั้นตอนที่ 3 คำนวณผลลัพธ์ที่เกิดขึ้นจริง

โดยใช้ ฟังก์ชัน ซิกมอยด์ $f(z) = 1 / (1 + e^{-z})$

ในการคำนวณผลลัพธ์

ชั้นแอบแฝงที่ 1

$$X_j^{(1)} = f \left[\sum_{i=1}^N w_{ij}^{(1)} X_i^{(0)} - \theta_j^{(1)} \right] \quad ; \quad 1 \leq j \leq N_1$$

ชั้นแอบแฝงที่ 2

$$X_k^{(2)} = f \left[\sum_{j=1}^{N_2} W_{kj}^{(2)} X_j^{(1)} - \theta_k^{(2)} \right] \quad ; \quad 1 \leq k \leq N_2$$

ชั้นแอบแฝงที่ n

$$X_k^{(n)} = f \left[\sum_{j=1}^{N_n} W_{kj}^{(n)} X_j^{(n-1)} - \theta_k^{(n)} \right] \quad ; \quad 1 \leq k \leq N_n$$

ชั้นผลลัพธ์

$$Y_j = f \left[\sum_{k=1}^{N_y} W_{jk}^{(n)} X_k^{(n)} - \theta_j^{(n)} \right] \quad ; \quad 1 \leq j \leq N_y$$

โดย N_x = จำนวนโหนดของข้อมูลเข้า

N_y = จำนวนโหนดผลลัพธ์

N_1, N_2, \dots, N_n = จำนวนโหนดของชั้นแอบแฝงที่ $1, 2, \dots, n$

$\theta_j^{(n)}$ = ค่าความเอนเอียง

ขั้นตอนที่ 4 ปรับปรุงค่าน้ำหนัก

ก. ใช้ขั้นตอนย้อนกลับ (Recursive Algorithm) โดยเริ่มต้นจากชั้นผลลัพธ์ย้อนกลับมาชั้นแอบแฝงที่อยู่ถัดมา และปรับปรุงค่าน้ำหนัก

$$\text{โดย } W_{ji}^{(k+1)} = W_{ji}^{(k)} + \eta \delta_j X_i$$

ซึ่ง $w_{ji}(k)$ คือ ค่าน้ำหนักของโหนดที่ i ในชั้นที่ $(n-1)$ ไปยังโหนดที่ j ในชั้นที่ n ณ เวลา k

[.]

X_i = ค่าผลลัพธ์ของโหนด i หรือเป็นข้อมูลของโหนด j

η = ค่าสัมประสิทธิ์การเรียนรู้

(n)

δ_j = ค่าความผิดพลาดของโหนด j

โดยถ้า กรณี 1 โหนด j เป็นโหนดผลลัพธ์จะได้

$$\delta_j = y_j (1 - y_j) (d_j - y_j)$$

d_j คือ ผลลัพธ์ที่ต้องการของโหนด j

y_j คือ ผลลัพธ์ที่เกิดขึ้นจริงของโหนด j

กรณี 2 โหนด i เป็นโหนดในชั้นแอบแฝง จะได้

$$\delta_j = X_j (1 - X_j) \sum_k \delta_k w_{kj}$$

k คือ จำนวนโหนดทั้งหมดในชั้นถัดจาก j

ข. หยุดเมื่อค่าน้ำหนักทั้งหมดมีค่าคงที่ ซึ่งหมายความว่า กระบวนการ
เรียนรู้ได้สิ้นสุดแล้ว

ขั้นตอนที่ 5 วนซ้ำย้อนกลับไปขั้นตอนที่ 2

ค่า η โดยปกติ มีค่าอยู่ระหว่าง 0.25-0.75 โดยถูกกำหนดโดยผู้ใช้
ถ้ากำหนดค่า η มากเกินไปจะทำให้เกิดความไม่แน่นอนขึ้นในเน็ตเวิร์ค และทำให้การเรียนรู้
ของเน็ตเวิร์คไม่ดีเท่าที่ควร หรือถ้ากำหนดค่าน้อยเกินไปจะทำให้การเรียนรู้ช้ามาก นอก
จากค่า η แล้ว ยังมีค่าโมเมนตัม α ที่ช่วยให้การเรียนรู้เร็วขึ้น และทำให้เน็ตเวิร์ค
มีความคงตัวมากขึ้น นอกจากนี้ ยังช่วยไม่ให้เน็ตเวิร์คตกอยู่ภายใต้จุดค่าสุดท้องถิ่น (Local
Minima)

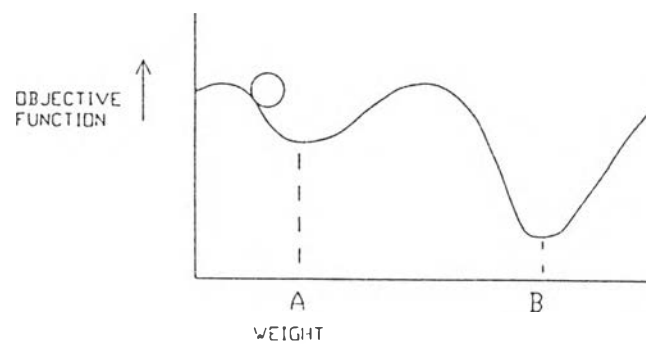
$$\Delta w_{ji}(k+1) = \eta \delta_j x_i + \alpha [\Delta w_{ji}(k)] \quad \text{----- (4.32)}$$

$$\text{ดังนั้น} \quad w_{ji}(k+1) = w_{ji}(k) + \Delta w_{ji}(k+1) \quad \text{----- (4.33)}$$

ค่า α ที่เหมาะสมอยู่ในช่วง 0.00 - 1.0 (Phillip D.Wasserman,1989)
และมีค่าเท่ากับ 0.9 (John Herz,Anders Krogh,Richard G.
Palmer,1991)

ค่า η ที่เหมาะสมอยู่ในช่วง 0.01 - 1.0 (Phillip D.Wasserman,1989)
และ 0.25 - 0.75 (Judith E.Dayhoff,1990)

ซึ่งการจะกำหนดค่าของ η และ α นั้นขึ้นกับแต่ละงานประยุกต์ Rumelhart และ McClelland กำหนดให้ตอนเริ่มการสอนเน็ตเวิร์ค มีค่า $\eta = 0.5$ และ $\alpha = 0.9$ จะเห็นว่า การกำหนดค่า η และ α นั้นควรเริ่มจากค่ามากๆ แล้วค่อยลดลงเรื่อยๆ จะทำให้เน็ตเวิร์คสามารถเข้าสู่จุดต่ำสุดโดยรวม (Global Minimum) เร็วขึ้น (แทนที่จะไปตกอยู่ที่จุดต่ำสุดท้องถิ่น) และ มีความคงตัว ดังรูป 4-13



รูป 4-13 ปัญหาของจุดต่ำสุดท้องถิ่น

4.4 การสอนเน็ตเวิร์ค

การสอนเน็ตเวิร์คชนิดการเรียนรู้แบบย้อนกลับ เรียกว่า การสอนแบบมีครู โดยมี การกำหนดรูปแบบข้อมูลเข้าควบคุม กับรูปแบบเป้าหมาย โดยมีการกำหนดชุดการสอนหลายๆ รูปแบบ เพื่อให้เน็ตเวิร์คสามารถ เรียนรู้ เน็ตเวิร์คสามารถแยกแยะได้ว่า ลักษณะ สัญญาณความถี่นั้นผิดปกติกจากความถี่ธรรมชาติ (Natural Frequency) หรือไม่

ในกระบวนการสอน เพื่อให้เน็ตเวิร์คสามารถเรียนรู้ได้ดีนั้น จะต้องใช้ การวนซ้ำ การสอนหลายๆรอบ ซึ่งอาจเป็น 1,000 - 30,000 รอบ ก็ได้ เพื่อให้ค่า RMS. ลดลงต่ำกว่า 0.1 ซึ่งก็หมายความว่า เน็ตเวิร์คได้เกิดการเรียนรู้แล้ว ดังนั้นจำเป็นต้อง ใช้ ฮาร์ดแวร์ที่มีประสิทธิภาพสูงเข้าช่วย เช่นต้องมี ตัวช่วยประมวลคณิตศาสตร์ หรือ อาจ ต้องใช้เครื่องเวิร์คสเตชันประมวลผล เพื่อต้องการผลลัพธ์ที่เร็วขึ้น

แต่ถ้าเน็ตเวิร์ค ตกอยู่ภายใต้จุดต่ำสุดท้องถิ่น นั่นคือ เน็ตเวิร์คได้หยุดการ เรียนรู้แล้ว หรือ ไม่สามารถระลึก ชุดการทดสอบได้เลย แม้ว่าจะระลึก ชุดการสอนได้



รูป 4-14 แสดงจุดต่ำสุดโดยรวม และ จุดต่ำสุดท้องถิ่น

ถ้าเน็ตเวิร์ค ตกอยู่ที่จุดต่ำสุดโดยรวม นั่นคือ ค่า RMS. ของเน็ตเวิร์คได้ถึงจุดต่ำสุดแล้ว ซึ่งมันสามารถระลึก ได้ทั้งชุดการสอน และ ชุดการทดสอบ ได้ในระดับที่พอใจ มีวิธีการแก้ปัญหาไม่ให้เน็ตเวิร์ค ตกอยู่ภายใต้จุดต่ำสุดท้องถิ่น คือ

4.4.1 การปรับค่า η และ α โดยเริ่มต้นจากค่ามากๆ ซึ่งน้อยกว่า 1 แล้วค่อยๆ ลดลงเรื่อยๆ ตามจำนวนรอบที่ใช้สอนเน็ตเวิร์ค

4.4.2 การเพิ่มจำนวนโนดแอบแฝง

4.4.3 การเปลี่ยนฟังก์ชันการแปลงค่า รูปแบบหรือโครงสร้าง (Model)

4.4.4 การเปลี่ยนชุดการสอนเสียใหม่ เนื่องจากชุดการสอนเดิมไม่ได้เป็นตัวแทนที่ดีของประเภทอาการนั้นๆ

4.5 การประเมินผลการสอนนิวรัล เน็ตเวิร์ค

การสอนนิวรัล เน็ตเวิร์ค จะสำเร็จหรือไม่ เราวัดจากค่า RMS. ดังสมการ

$$\text{RMS} = \sqrt{\frac{\sum_p \sum_j (t_{jp} - x_{jp})^2}{n_p n_o}} \quad \text{-----} \quad (4.34)$$

โดย n_p = จำนวนของรูปแบบในชุดการสอน

n_o = จำนวนของหน่วยในชั้นผลลัพธ์

t_{jp} = ค่าเป้าหมายของหน่วยที่ j หลังจากที่ได้เสนอรูปแบบ p

x_{jp} = ค่าผลลัพธ์ที่ออกมาโดยหน่วยที่ j หลังจากทีเสนอรูปแบบ p

โดยปกติ ถ้าค่า RMS. ที่ได้ < 0.1 แสดงว่านิเวศ เนทเวิร์ค ได้เกิดการ
เรียนรู้แล้ว

4.6 ข้อดีและข้อจำกัดของการเรียนรู้แบบย้อนกลับ (Judith E. Dayhoff, 1990)

4.6.1 ข้อดี

คือ มีความสามารถในการจดจำรูปแบบของปัญหา(Pattern-Mapping) ซึ่งการเรียนรู้แบบย้อนกลับ สามารถที่จะเรียนรู้ความสัมพันธ์ของรูปแบบได้มากมาย โดยการเรียนรู้แบบย้อนกลับ ต้องการตัวอย่างรูปแบบที่จะเรียนรู้ ความยืดหยุ่นของการเรียนรู้แบบย้อนกลับ อยู่ที่การมีความหลากหลายในการออกแบบทางเลือกต่างๆ เช่น จำนวนชั้นเส้นเชื่อมโยง จำนวนของโนด ค่า ส.ป.ส. การเรียนรู้ ค่าโมเมนตัม ที่เรากำหนดขึ้น และการแทนรูปแบบของข้อมูล

ความยืดหยุ่นดังกล่าว ทำให้การเรียนรู้แบบย้อนกลับ สามารถแก้ปัญหา งานประยุกต์ได้อย่างมากมาย

4.6.2 ข้อจำกัด

คือ การใช้เวลาอย่างมากเพื่อให้การสอนนิเวศ เนทเวิร์ค สามารถเรียนรู้ได้ หรือที่เรียกว่า Convergence Time ซึ่งบางครั้งอาจต้องการการกรวนซ้ำถึง 100-1,000 รอบ สำหรับการเรียนรู้การ แก้ปัญหานั้นๆ ทำให้บางครั้งต้องใช้เวลา 1 วัน หรือมากกว่านั้นในการสอนนิเวศ เนทเวิร์ค

อุปสรรคอีกข้อหนึ่ง คือ การต้องใช้ชุดการสอนอย่างมากมาย ซึ่งบางครั้งอาจต้องใช้ถึง 1,000 ตัวอย่าง ด้วยเหตุผลดังกล่าว เพื่อที่จะลดเวลาการเรียนรู้ ทำให้ต้องใช้ฮาร์ดแวร์ ที่มีกำลังและประสิทธิภาพสูง กล่าวคือ จำเป็นต้องใช้ บอร์ดควมเร่งชนิดพิเศษ เครื่องประมวลแบบคู่ขนาน

นอกจากความสามารถทางด้านฮาร์ดแวร์แล้ว อาจมีการเปลี่ยนชุดการสอนใหม่ การปรับเปลี่ยนค่า ส.ป.ส. การเรียนรู้ ก็มีส่วนช่วยให้การเรียนรู้เร็วขึ้น โดยปกติค่า η จะเริ่มจากค่ามากและลดค่าลงเรื่อยๆ เพื่อให้การเรียนรู้เป็นไปอย่างมีประสิทธิภาพ ค่า η ที่เหมาะสมจะมีค่าระหว่าง 0.25-0.75