

Credit Assignment Problem: If any distributed system or any distributed network, then it is a common practice to give some credits or blame to each and every internal decisions so that it can be reflected on the overall outcome of the network.

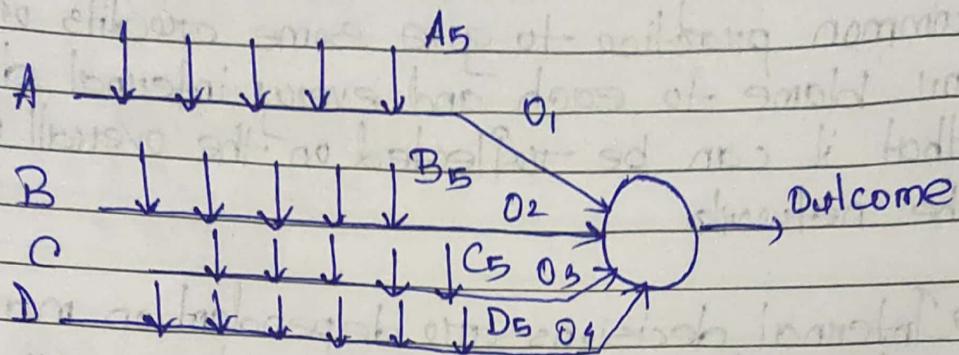
→ Internal decisions are dependent on many actions. Among such, some actions are very prominent to get the outcome of internal decision. In this cases, the credit assignment is divided into two parts:

i) Temporal: If we have given credit or blame for outcomes of the internal decision, then this is called temporal credit assignment problem.

ii) Structural: If we have given credit or blame for outcomes of the actions of the internal decision, then this is called structural credit assignment.

→ Temporal = to get the selected outcomes of the internal decisions.

→ Structural = to get the selected outcome actions of the internal decision.

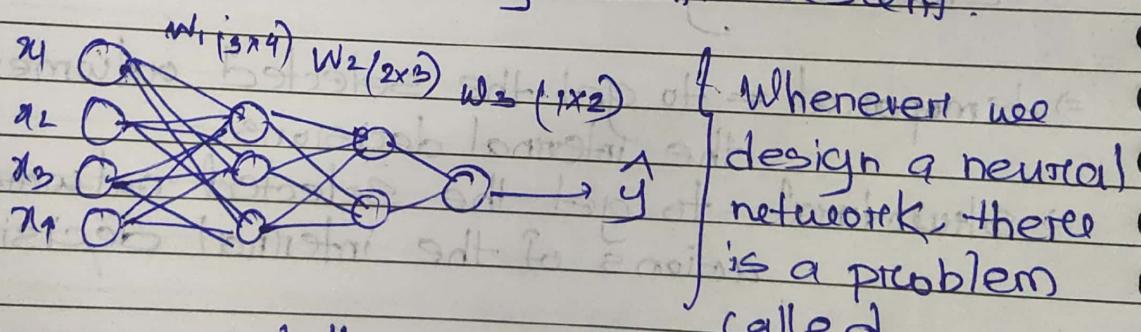


Here, internal decision means

Credit assignment problem, we need to give some credit or blame to our outcomes of internal decisions only. that means $O_1 - O_4$, we need to provide whether it is credit or blame
 ⇒ if credit is there, it will be included in the final outcome
 ⇒ if blame is there, it is not included in final outcome.

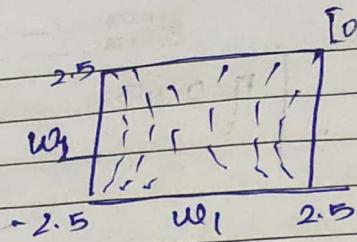
Delta Rule : (Perception learning)

→ derived from gradient descent



When we are initially assigning all the weights equally, the model is updated equally and it is learning nothing. So, we assign weights randomly.
 Good Write

"Symmetry Breaking".



It is forming a saddle shape across the entire pattern. So, the errors in one particular layer is mitigated with the help of adjusting the weights which are put randomly.

Our major goal is of Delta learning is to learn the weights so that weight updation whenever we get the final output that is changed with the help of some rule.

Now,

$$\text{error function, } J = \frac{1}{2} \sum_{d=1}^D (t_d - o_d)^2$$

dependent on the weights and biases

Initially,

$$O = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + w_{n+1} x_{n+1}$$

Here,

$$\frac{\partial J}{\partial w_i} = \left[\frac{\partial J}{\partial w_0}, \frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \dots, \frac{\partial J}{\partial w_n}, \frac{\partial J}{\partial w_{n+1}} \right]$$

$$\text{or, } \frac{\partial J}{\partial w_i} = \frac{1}{2} \sum_{d=1}^D 2(t_d - o_d) \cdot \frac{\partial J}{\partial w_i} (t_d - o_d)$$

$$\text{or, } \Delta \vec{w} = (t_d - o_d) \sum_{d \in D} \frac{\partial J}{\partial w_i} (t_d - o_d)$$

$$= \sum_{d \in D} (t_d - o_d) \left[t_d - \left(\vec{w} \cdot \vec{x}_d + \vec{w} \cdot \vec{x}_d^2 \right) \right]$$

$$= \sum_{d \in D} (t_d - o_d) \left(-\vec{x}_{i,d} - \vec{x}_{i,d}^2 \right)$$

$$= \sum_{d \in D} (t_d - o_d) \cdot \left(\vec{x}_{i,d} - \vec{x}_{i,d}^2 \right)$$

Good Write

$$= - \sum_{d \in D} (t_d - o_d) (\vec{x}_{i,d} + \vec{w}_{i,d}). \quad | \quad \eta = 0.05$$

$$= - \eta \nabla J[\vec{w}]$$

$$\therefore \Delta \vec{w}_i = - \eta \frac{\partial J}{\partial w_i}$$

Derivation of Back Propagation Algorithm :

⇒ Back Propagation : (training example, η , n_{in} , n_{out} , n_{hidden})

i) training example is a form (x, t) where,
 x is the vector of network input values,
and t is the vector of target network values.

- ii) η = learning rate (0.05)
- iii) n_{in} = no. of network inputs
- iv) n_{out} = no. of output units
- v) n_{hidden} = no. of units in the hidden layer

= The input from unit i into unit j is denoted x_{ij} , the weight from unit i to unit j is denoted w_{ij} .

Algorithm:

Step 01: Create a feed forwarded network with n_i inputs, n_h hidden hidden units, and n_o output units.

Step 02: Initialize all network weights to small random numbers.

Step 03: Until the termination condition is met,
Do, For each training set (x, t) in
Do,

→ Propagate the input forward through the network.

→ Propagate the errors backward through the network.

→ For each network k , calculate error term, $\delta_k \leftarrow s_k o_k (1 - o_k) (t_k - o_k)$

→ For each network unit h , calculate it's error term.

$$\delta_h = o_h (1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k$$

Step 04: Update each network's weights.

$$w_{ji} \leftarrow w_{ji} + \eta \Delta w_{ji}$$

$$\text{where, } \Delta w_{ji} = \eta \delta_j x_j$$

Derivation:

- ⇒ To derive the equation for updating weights in back propagation algorithm, we use Stochastic gradient descent rule.
- ⇒ Stochastic gradient descent involves iterating through the training examples one at a time, for each training example (d) descending the gradient of the error $\rightarrow (E_d)$ with respect to this single example.
- ⇒ In other words, for each training example d every weight w_{ij} is updated by adding to it Δw_{ij} .

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

where,

$$\Delta w_{ij} = -\eta \frac{\partial E_d}{\partial w_{ij}}$$

where, E_d is the error on training example d that is half the square difference between the target output and the actual gain output over all output units in the network.

$$E_d(\vec{w}) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

Here, outputs is the set of output units in the network, t_k is the target value of unit k for training example d , and o_k = output of unit k .

t_k = target value

o_k = output of unit k .

Good Write

Notation used:

x_{ji} = the i^{th} input to unit j .

w_{ji} = the weight associated with the i^{th} input to unit j

$\text{net}_j = \sum_i w_{ji} x_{ji}$ (the weighted sum of inputs for unit j)

o_j = the output computed by unit j .

t_j = the target output for unit j .

σ = the sigmoid function.

outputs = the set of units in the final layer of the network

$$w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$$

where,

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

i) Notice that weight w_{ji} can influence the rest of the network only through net_j .

\therefore We can use the chain rule to write,

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial f_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}}$$

$$\text{net}_j = \sum_i w_{ji} x_{ji}$$

$$\frac{\partial \text{net}_j}{\partial w_{ji}} = \sum_i x_{ji}$$

$$\text{or}, \frac{\partial E_d}{\partial w_{ji}} = \frac{\partial f_d}{\partial \text{net}_j} \cdot x_{ji}$$

$$\therefore \Delta w_{ji} = -\eta \frac{\partial f_d}{\partial \text{net}_j} \cdot x_{ji}$$

To derive a convenient expression for $\frac{\partial E_d}{\partial \text{net}_j}$

We consider two cases in turn:

case 1: where unit j is an output unit for the network.

case 2: where unit j is an internal unit of the network.

Case 1: Training Rule for Output Unit Weights.

Just as w_{ij} can influence the rest of the network only through net_j , net_j can influence the network only through o_j .

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \cdot \frac{\partial o_j}{\partial net_j}$$

$$\frac{\partial E_d}{\partial o_j} = \frac{2}{2} \times \frac{1}{2} \sum_{k \text{ outputs}} (t_k - o_k)^2$$

$$\text{or, } \frac{\partial E_d}{\partial o_j} = \frac{2}{2} \times \frac{1}{2} (t_j - o_j)^2$$

$$= \frac{1}{2} \times 2 (t_j - o_j) \times \frac{\partial (t_j - o_j)}{\partial o_j}$$

$$= -(t_j - o_j)$$

$$\frac{\partial o_j}{\partial (net_j)} = \frac{\partial \sigma (net_j)}{\partial (net_j)}$$

$$= \sigma (net_j) (1 - \sigma (net_j))$$

$$= o_j (1 - o_j)$$

$$\frac{\partial \sigma}{\partial (x)} = \sigma(x) (1 - \sigma(x))$$

$$\frac{\partial E_d}{\partial net_j} = -(t_j - o_j) o_j (1 - o_j)$$

$$\text{Now, } \Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial net_j} x_{ji}$$

$$\frac{\partial E_d}{\partial net_j} = -(t_j - o_j) o_j (1 - o_j)$$

Good Write

$$\therefore \Delta w_{ij} = \eta (t_j - o_j) o_j (1 - o_j) x_{ji}$$

$$S_j = (t_j - o_j) o_j (1 - o_j)$$

$$\Delta w_{ji} = \eta S_j x_{ji}$$

Case 2: Training Rule For Hidden Unit Weights

$$\text{on } i, \frac{\partial E_d}{\partial net_j} = \sum_{k \in \text{Downstream}} \frac{\partial E_d}{\partial net_k} \cdot \frac{\partial net_k}{\partial net_j} \quad \frac{\partial E_d}{\partial net_j} = -(t_j - o_j) o_j (1 - o_j)$$

$$\text{on } i, \frac{\partial E_d}{\partial net_j} = \sum_{k \in \text{Downstream}} -S_k \frac{\partial net_k}{\partial o_j} \cdot \frac{\partial o_j}{\partial net_j} \quad \frac{\partial net_k}{\partial o_j} = \frac{\partial o_j w_{kj}}{\partial o_j}$$

$$\text{on } i, \frac{\partial E_d}{\partial net_j} = \sum_{k \in \text{Downstream}} -S_k w_{kj} \times \frac{\partial o_j}{\partial net_j} \quad = \frac{\partial o_j w_{kj}}{\partial o_j}$$

$$\text{on } i, \frac{\partial E_d}{\partial net_j} = \sum_{k \in \text{Downstream}} -S_k w_{kj} o_j (1 - o_j) \quad \frac{\partial o_j}{\partial net_j} = \frac{\partial o_j}{\partial (net_j)} \cdot \frac{\partial o_j (net_j)}{\partial (net_j)}$$

$$= \sigma'(net_j) (1 - \sigma(net_j))$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial net_j} x_{ji}$$

$$= \cancel{\eta} \\ = o_j (1 - o_j)$$

$$\Delta w_{ji} = \eta o_j (1 - o_j) \sum_{k \in D_i} S_k w_{kj} x_{ji}$$

$$\therefore \Delta w_{ji} = \eta S_j x_{ji}$$

Kolmogorov theorem:

→ Chapman-Kolmogorov equation was derived independently by both the British mathematician Sydney Chapman and Russian mathematician Andrey Kolmogorov.

Equation:

Suppose, $x_n, n = 0, 1, 2, \dots$ is a homogeneous Markov chain.

$$P_{ij}^{(m+n)} = \sum_k P_{ik}^{(m)} P_{kj}^{(n)}$$

P_{ij} = conditional probability that the Markov chain goes from state i to j in $m+n$ steps.

P_{ik} = conditional probabilities of reaching an intermediate state k in m steps.

P_{kj} = conditional probabilities from k reaching state j in n steps.

We know that, n -step transition probability:

$$P_{ij}^{(n)} = P(X_{n+1} = j | X_0 = i)$$

That is, $P(X_{m+n} = j | X_0 = i) = \sum P(X_m = k | X_0 = i)$

$$P(X_{m+n} = j | X_m = k)$$

Good Write

$$P_{ij}^{(m+n)} = P(X_{m+n}=j, X_m=k \mid X_0=i)$$

$$\text{or, } P_{ij}^{(m+n)} = \sum_K P(X_{m+n}=j, X_m=k \mid X_0=i)$$

$$= \sum_K \frac{P(X_{m+n}=j, X_m=k, X_0=i)}{P(X_0=i)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$= \sum_K \frac{P(X_{m+n}=j, X_m=k, X_0=i)}{P(X_0=i)} \cdot \frac{P(X_m=k \mid X_0=i)}{P(X_0=i)}$$

$$= \sum_K P(X_{m+n}=j \mid X_m=k, X_0=i) \cdot P(X_m=k \mid X_0=i)$$

$$= \sum_K P(X_{m+n}=j \mid X_m=k) \cdot P(X_m=k \mid X_0=i).$$

$$= \sum_K P_{kj}^{(m)} \cdot P_{ik}^{(m)}$$

In matrix notation,

$$\boxed{P^{(m+n)} = P^{(m)} \cdot P^{(n)}} \quad \text{or}$$

$$\boxed{P^{(m+n)} = P(m) \hat{P}(n)}$$