

Term explanation

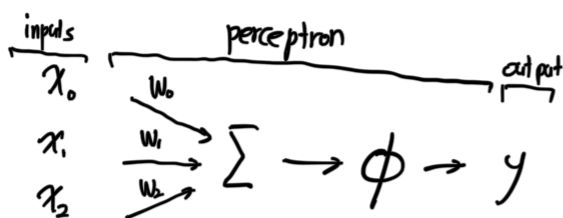
P_n^l : n -th perception of P_n^l -th layer or the output value of it.

W_{nm}^l : the weight about m -th input value of P_n^l .

x_{nm} : m -th input of n -th perception in 0 -th layer.

How does a single perceptron calculate:

There's only one perception so we can omit the perception number and the layer number.



The Σ and the ϕ of the above image are functions.

The Σ calculates $\sum_{i=0}^N x_i w_i$ when N is the number of the inputs.

The ϕ is a function called 'activation function' which makes linear function to non-linear.

Sigmoid function ($y = \frac{1}{1+e^{-x}}$) and ReLU ($y = \begin{cases} x & (x \geq 0) \\ 0 & (x < 0) \end{cases}$) are the examples of the activation function.

So we can write the formula of above perceptron like this:

$$\phi(\Sigma(w, x)) = \phi\left(\sum_{i=0}^{N-1} w_i x_i\right) \quad (w \text{ and } x \text{ are sequences of } \underbrace{\text{weights}}_{(w_0, w_1, \dots, w_{N-1})} \text{ and } \underbrace{\text{inputs}}_{(x_0, x_1, \dots, x_{N-1})})$$

How does a single perceptron learn

Consider following case:

The perceptron have to print \hat{y} as the output when we give 2 specific number x_0 and x_1 as the input.

And the real output of the perceptron is y .

There is difference between \hat{y} (the value we expect) and y (the value of perceptron's calculation).

We can evaluate this difference by cost function (J) which calculate the degree of error.

So $J(y, \hat{y})$ is the degree of error between y and \hat{y} .

Then we can adjust the specific weight w_k like this:

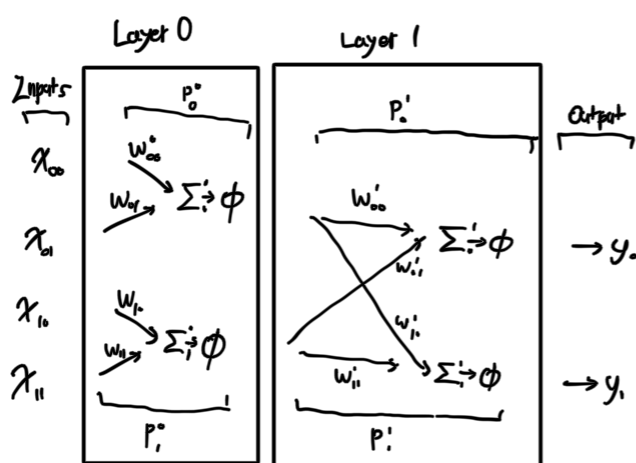
$$\underbrace{w_k}_{\text{adjusted weight}} \leftarrow \underbrace{w_k}_{\text{old weight}} - \frac{J(y, g)}{dw_k}$$

) * This is from 'gradient descent', the algorithm for adjusting the weights.
I omitted the explanation of this because this isn't main subject.

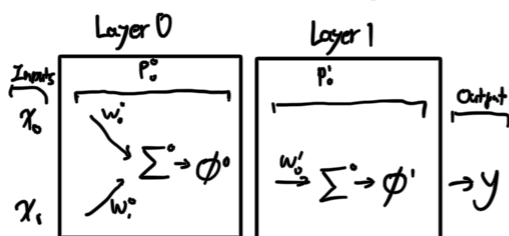
The above process is what we call 'learning of perceptron'.

How does multi layer perceptron learn

Multi layer perceptron (model) looks like this:



Or we can think about following simple one:



P_0^o calculates like this: ; P_0^i calculates like this:

$$P_0^o = \phi'(\sum_{i=0}^n x_i w_{io})$$

$$P_0^i = \phi'(\sum_{i=0}^n P_0^o w'_{io})$$

And the whole perceptron model calculates like this:

$$y = \phi'(\sum_{i=0}^n (\phi'(\sum_{j=0}^n x_j w_{ij}) w'_{io}))$$

And the error level is like this:

$$J(y, g)$$

We have to get the value of $\frac{dJ(y, g)}{dw'_0}$, $\frac{dJ(P_0^o, \hat{P}_0^o)}{dw'_0}$, $\frac{dJ(P_0^i, \hat{P}_0^i)}{dw'_0}$ to train the model.

We can easily get $\frac{dJ(y, g)}{dw'_0}$ but how can we get $\frac{dJ(P_0^o, \hat{P}_0^o)}{dw'_0}$ and $\frac{dJ(P_0^i, \hat{P}_0^i)}{dw'_0}$ when we don't know \hat{P}_0^o and \hat{P}_0^i , which mean the ideal value of P_0^o and P_0^i ?

To do this, $\frac{dJ(y, \hat{y})}{dw_0}$ and $\frac{dJ(y, \hat{y})}{dw_i}$ can be used instead of them.

Let me calculate $\frac{dJ(y, \hat{y})}{dw_i}$ as an example.

$$\begin{aligned} \frac{dJ(y, \hat{y})}{dw_i} &= \frac{d}{dw_i} J\left(\Phi\left(\sum_{\lambda=0}^0 \left(\Phi\left(\sum_{j=0}^1 \pi_j w_j\right) w_{\lambda}\right), \hat{y}\right) = \frac{dJ}{d\phi'} \cdot \frac{d}{dw_i} \Phi\left(\sum_{\lambda=0}^0 \left(\Phi\left(\sum_{j=0}^1 \pi_j w_j\right) w_{\lambda}\right)\right) \\ &= \frac{dJ}{d\phi'} \cdot \frac{d\phi'}{d\Sigma'} \cdot \frac{d}{dw_i} \sum_{\lambda=0}^0 \left(\Phi\left(\sum_{j=0}^1 \pi_j w_j\right) w_{\lambda}\right) = \dots \\ &= \frac{dJ}{d\phi'} \cdot \frac{d\phi'}{d\Sigma'} \cdot \frac{d\Sigma'}{d\phi^0} \cdot \frac{d\phi^0}{d\Sigma^0} \cdot \frac{d\Sigma^0}{dw_i} \end{aligned}$$

The above process is 'back propagation'.

After doing this, differentiate with J respect to w of the perceptrons in the hidden layers can be convert to what is easy to calculate.