

3. Neural networks and Deep learning

2017년 1월 6일 금요일 오후 3:59

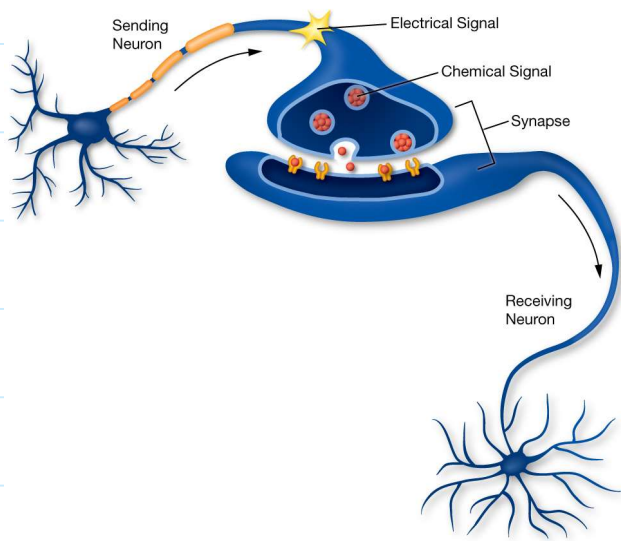
Reference : Michael A. Nielson. (2015), "Neural Networks and Deep Learning".

<http://www.neuralnetworksanddeeplearning.com>.

Mark F. Bear et. al. (2016). Neuroscience: exploring the brain 4th ed. Wolters Kluwer.

3.1. How do neurons work?

For individual neuron,



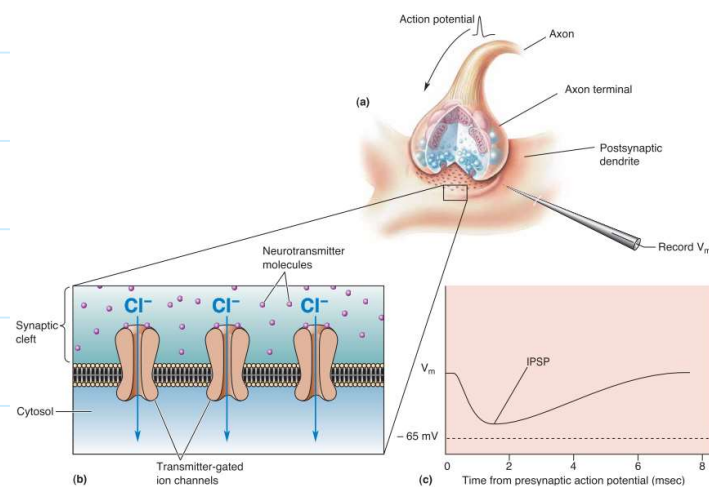
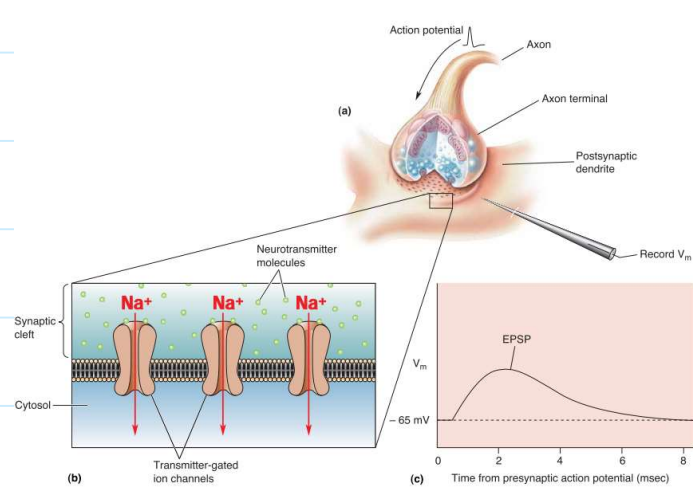
electrical signal over threshold.

→ action potential

→ synaptic transmission.

→ transmit signals to the following neurons.

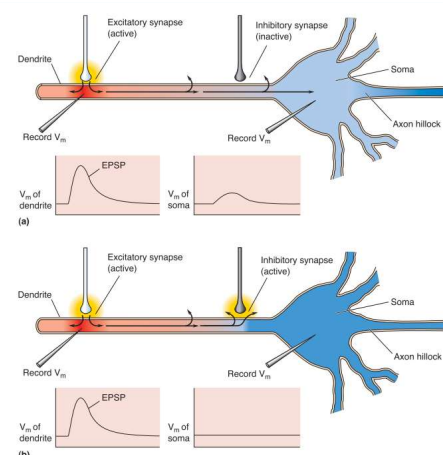
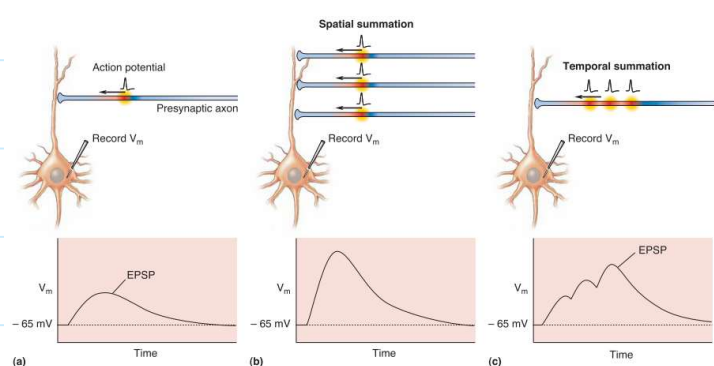
In more detailed and macro view, there are two types of synaptic connections



< excitatory post-synaptic potential >

< inhibitory post-synaptic potential >

⇒ Synaptic integration,

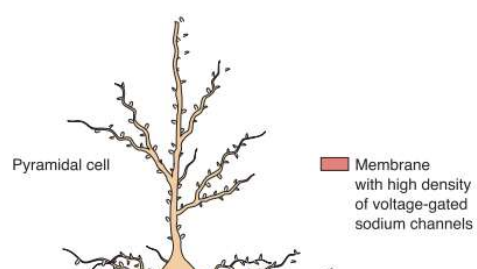


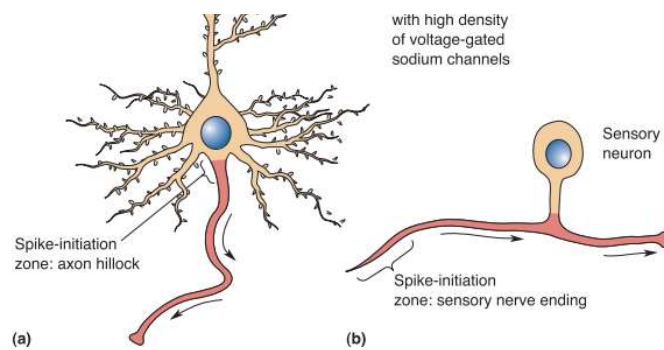
< EPSP summation >

< Shunting inhibition >

⇒ integrated in spike-initiation zone.

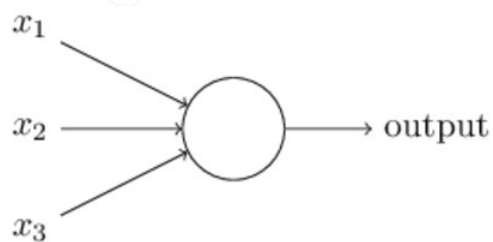
(usually in axon hillock).





Conclusion : action potential is an activation from a weighted sum of preceding neural inputs.

3.2. Perceptron : Artificial neural network.



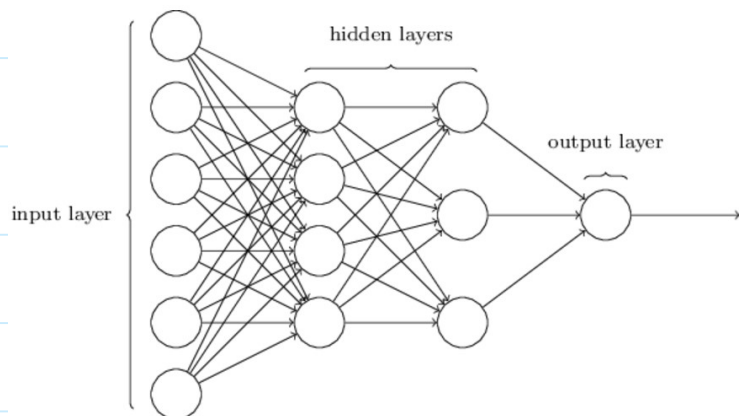
inputs : x_1, x_2, \dots, x_n $\rightarrow \sum_j w_j x_j$: weighted sum.

weights : w_1, w_2, \dots, w_n

$\Rightarrow \text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq \text{threshold.} \\ 1 & \text{if } \sum_j w_j x_j > \text{threshold.} \end{cases} = \underbrace{\sigma(w^T x)}_{\text{activation function.}}$

3.2.1. Multilayer Perceptron. (OR feedforward. network?)

- multilayer perceptron



simple decision making \rightarrow complex decision making.

\Rightarrow engage sophisticated decision making.

- feedforward vs. recurrent neural network.

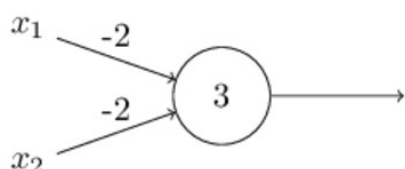
feed-forward network : output from one layer is used as an output of the next layer.
(no loops in the network).

recurrent neural network \leftarrow cascade of neural firing.

\hookrightarrow loops doesn't cause problems.

3.2.2. Computational universality of MLP.

MLP learns NAND gate. \Rightarrow can immitate all hard-coded programs (logical functions)



Furthermore than just mimicking logical functions,

we can devise learning algorithms by tuning the parameter automatically.

↳ Obtained by optimization.

3.3. Training feedforward neural networks.

3.3.1. Gradient descent.

cost function : $C(\phi) = \mathbb{E}_x(C_x(\phi; y))$ while C_x : individual cost on an example x .

x : input, y : desired output, ϕ : parameters.

Target: "minimize the cost function by gradient descent"

e.g. error rate, -log-likelihood.

For parameters, $\phi = (\phi_1, \phi_2, \dots, \phi_n)$

$$\Delta C = \frac{\partial C}{\partial \phi_1} \Delta \phi_1 + \frac{\partial C}{\partial \phi_2} \Delta \phi_2 + \dots + \frac{\partial C}{\partial \phi_n} \Delta \phi_n$$

$$\text{while } \nabla C := \left(\frac{\partial C}{\partial \phi_1}, \dots, \frac{\partial C}{\partial \phi_n} \right)$$

$$\Delta C = \nabla C \cdot \Delta \phi$$

↳ $\Delta \phi = -\eta \nabla C \rightarrow$ reduce ΔC the most.

$\phi \rightarrow \phi' = \phi - \eta \nabla C$: gradient descent algorithm.

For MLP, parameters are weights or biases.

$$w_k \rightarrow w_k' = w_k - \eta \cdot \frac{\partial C}{\partial w_k}$$

$$b_e \rightarrow b_e' = b_e - \eta \cdot \frac{\partial C}{\partial b_e}$$

3.3.2. Backpropagation.

(Notation). (Gradient). For a tensor T with dimension n , and cost function C

$$\nabla_T C = \left(\frac{\partial C}{\partial T_{i_1 \dots i_n}} \right) \text{ i.e. } (\nabla_T C)_{i_1 \dots i_n} = \frac{\partial C}{\partial T_{i_1 \dots i_n}}$$

(Derivative). For a vector $x \in \mathbb{R}^n$ and $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$

$\frac{df}{dx}$ is defined as $f(x+dx) \approx f(x) + \frac{df}{dx} dx$.

As a result, if $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $\frac{df}{dx} = (\nabla_x f)^T$

activation of l^{th} layer, a^l is related to the activations in the $(l-1)^{\text{th}}$ layer, a^{l-1}

$$a^l = \sigma(w^l a^{l-1} + b^l)$$

When $\frac{dC}{da^l} (:= \nabla_{a^l} C^T)$ is given.

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{dC}{da^l} \cdot \frac{da^l}{dz^l} \cdot \frac{\partial z^l}{\partial w_{ij}^l} \text{ while } z^l = w^{lT} a^{l-1} + b^l$$

$$= \frac{dC}{da^l} \cdot \sigma'(z^l) \cdot a^{l-1} e_j \text{ (Note that } \sigma'(z^l) \text{ is a diagonal matrix of } \sigma'(z_i^l) \text{'s)}$$

$$= \frac{dC}{da_j^l} \cdot \sigma'(z_j^l) \cdot a_j^{l-1}$$

$$\Rightarrow \nabla_w C = a^{l-1} \cdot (\sigma'(z^l) \cdot \nabla_{a^l} C)^T$$

$$\frac{dC}{dz^L} = \frac{dC}{da^L} \cdot \frac{da^L}{dz^L} = \frac{dC}{da^L} \frac{da^L}{dz^L} \frac{dz^L}{db^L} = \frac{dC}{da^L} \sigma'(z^L)$$

$$\Rightarrow \nabla_{\theta^L} C = \sigma'(z^L) \cdot \nabla_{a^L} C$$

$$\frac{dC}{da^{L-1}} = \frac{dC}{da^L} \frac{da^L}{dz^L} \frac{dz^L}{da^{L-1}} = \frac{dC}{da^L} \sigma'(z^L) \cdot w^L$$

$$\Rightarrow \nabla_{a^{L-1}} C = w^{L^T} \cdot \sigma'(z^L) \cdot \nabla_{a^L} C$$

Recursively, we can calculate all $\nabla_{a^L} C$ and $\nabla_{\theta^L} C$.

3.4. Modification on gradient descent algorithm.

3.4.1. Difficulties on training gradient algorithms.

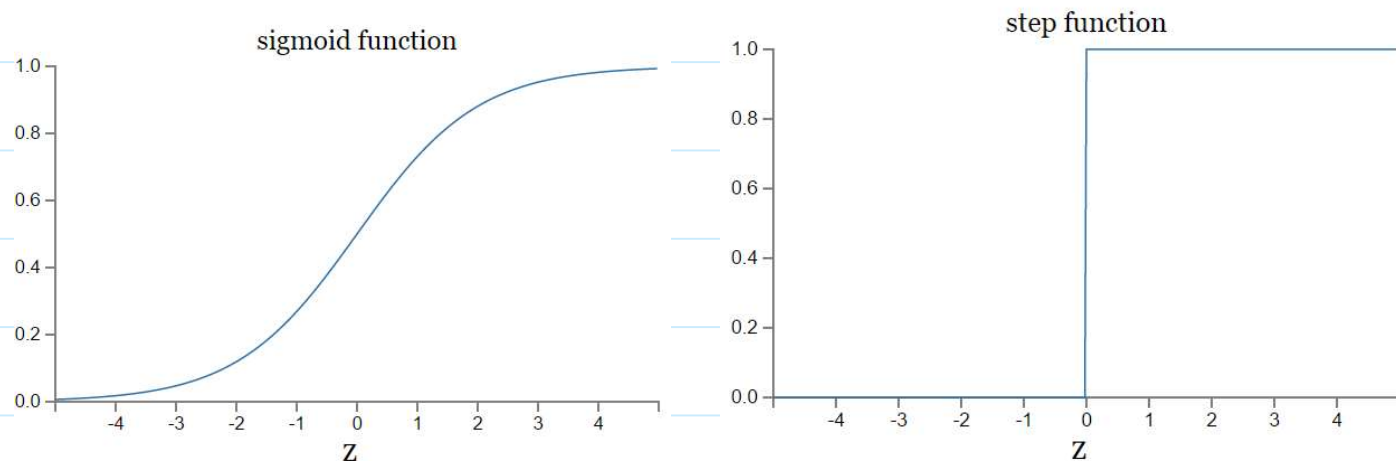
- ① undifferentiability of activation function.
- ② undifferentiability and under-representation of individual loss.
- ③ unavailable loss function

\Rightarrow Solutions : sigmoid neuron, proxy loss function, stochastic gradient descent.

3.4.2. Sigmoid neuron.

sigmoid function : $\sigma(z) = \frac{1}{1+e^{-z}}$

\hookrightarrow sigmoid neuron output = $\sigma(\sum_j w_j x_j + b)$



$w x + b \rightarrow$ linear function.

$\sigma \rightarrow$ activation function (nonlinearity).

\Rightarrow sigmoid neuron is much easier to get trained.

3.4.3. Proxy loss function

$$C_\alpha(\phi; y) = \|y - \hat{y}(x; \phi)\|^2$$

where \hat{y} is the output of the neural network. with parameter ϕ .

\Rightarrow ① differentiable with respect to parameters.

② motivates the learning further than direct loss function as accuracy.

3.4.4. Stochastic gradient descent.

: estimate the gradient ∇C from a small sample.

$$\frac{1}{m} \sum_{j=1}^m \nabla C x_j \approx E[\nabla C_\alpha] = \nabla C$$

$$\theta_k \rightarrow \theta_k - \frac{\eta}{m} \sum_i \frac{\partial C x_i}{\partial \theta_k}$$

'''$w_k \leftarrow w_k - \frac{\eta}{m} \sum_j \frac{\partial C_j}{\partial w_k}$'''

$$w_k \rightarrow w_k - \frac{\eta}{m} \sum_j \frac{\partial C_j}{\partial w_k}$$

$$b_e \rightarrow b_e - \frac{\eta}{m} \sum_j \frac{\partial C_j}{\partial b_e}$$