

DLSG Week 3: Neural Network



Thuong Nguyen Canh

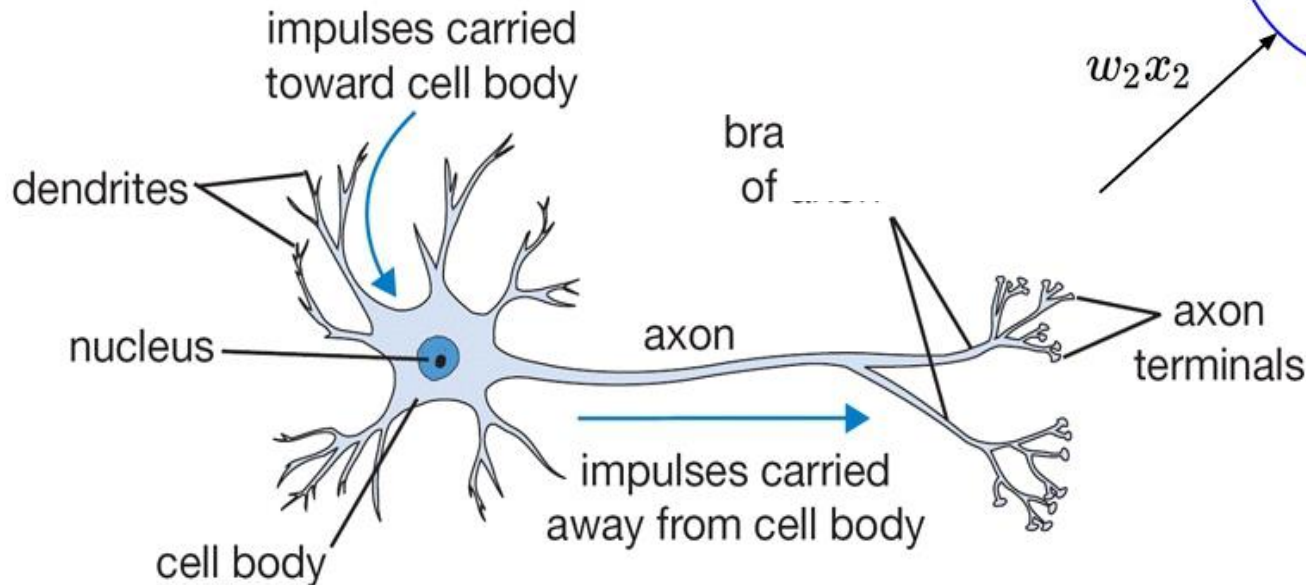
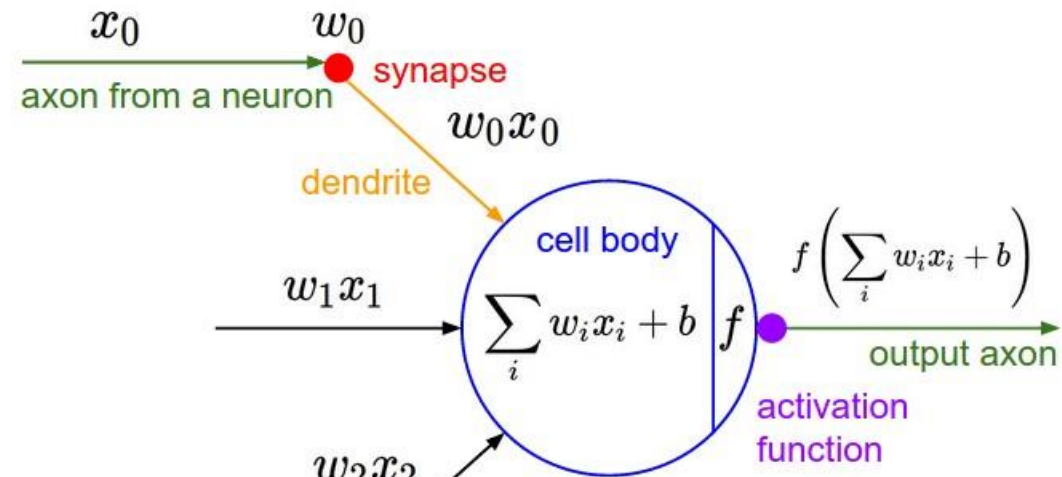
August 2019



Neural Network Definition

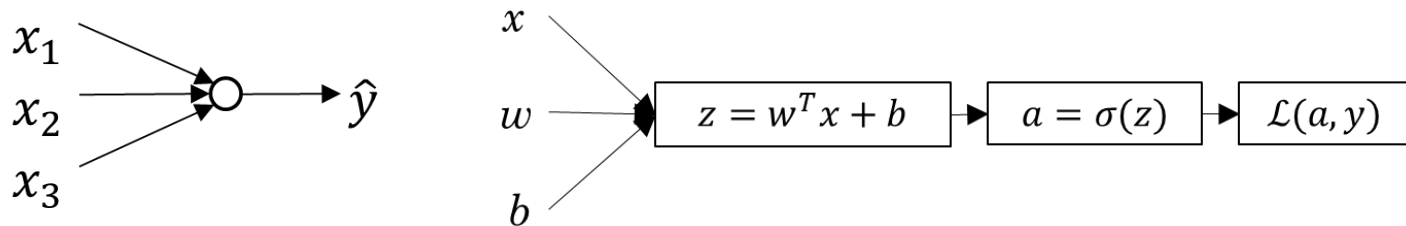
- Biological motivations

*This is a simplified connection.
Neural network does not represent
how human brain works*

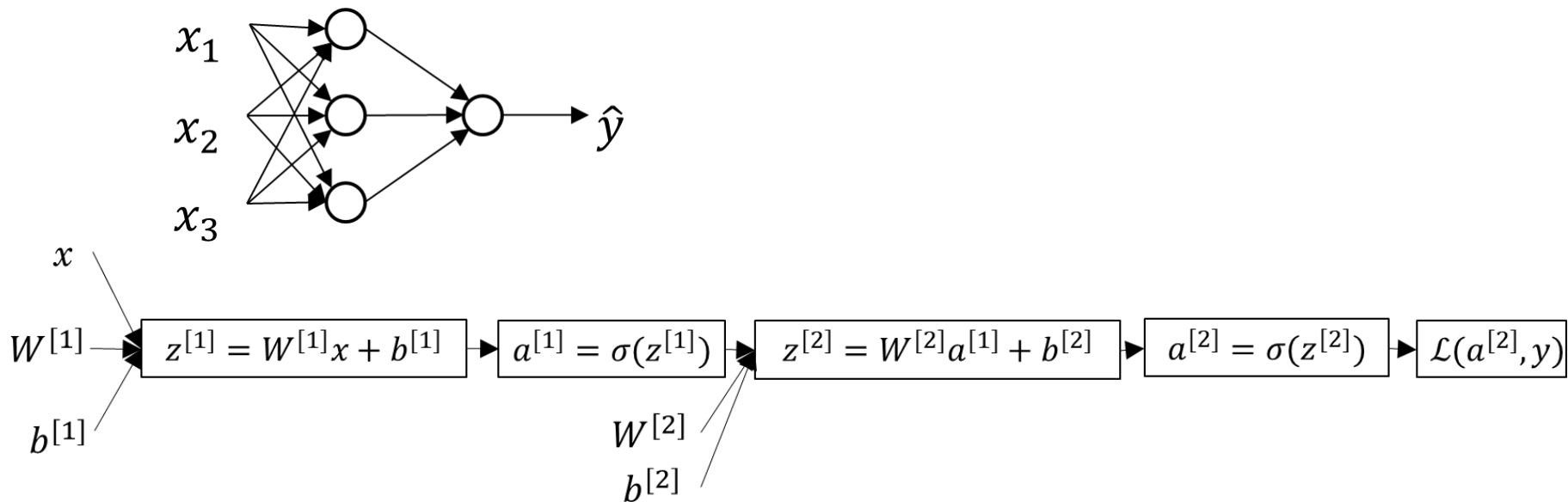


Neural Network Definition (3)

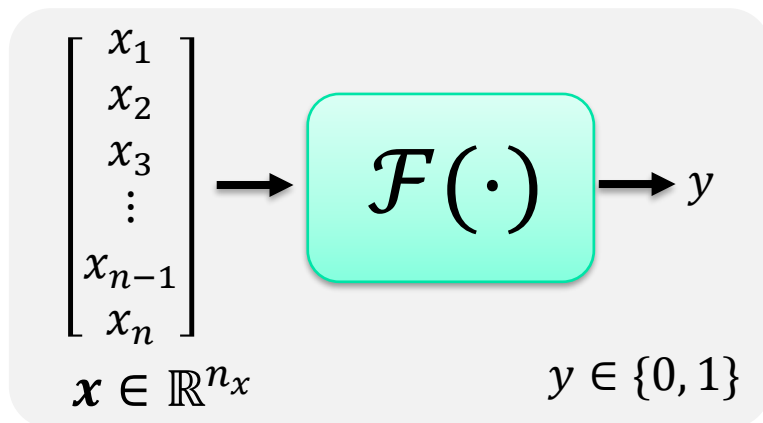
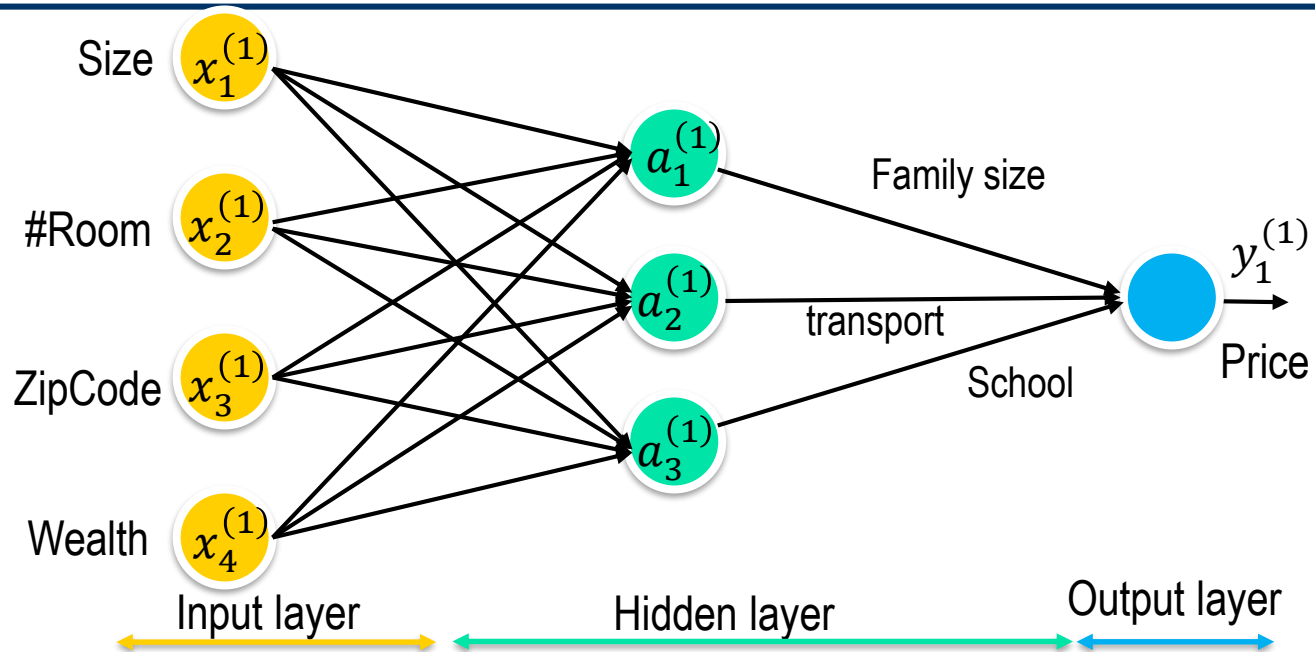
- Single neuron as a linear classifier



- Multiple neurons as a classifier



Neural Network Representation



Training pairs $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$

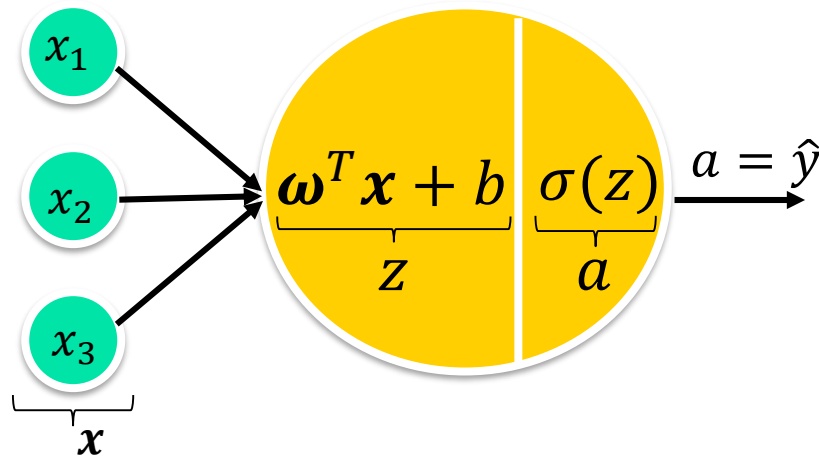
In matrix form

$$\mathbf{X} = \begin{bmatrix} | & | & & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(m)} \\ | & | & & | \end{bmatrix} \quad \mathbf{y} = [y^{(1)}, \dots, y^{(m)}]$$

Inputs

Labels

Neural Network Representation



$$z = \omega^T x + b$$

$$a = \sigma(z)$$

$$W^{(1)} = \begin{bmatrix} \text{---} \omega_1^{(1)T} \text{---} \\ \dots \\ \text{---} \omega_4^{(1)T} \text{---} \end{bmatrix}$$

$$z_1^{(1)} = \omega_1^{(1)T} x + b_1^{(1)}, \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

$$z_2^{(1)} = \omega_2^{(1)T} x + b_2^{(1)}, \quad a_2^{(1)} = \sigma(z_2^{(1)})$$

$$z_3^{(1)} = \omega_3^{(1)T} x + b_3^{(1)}, \quad a_3^{(1)} = \sigma(z_3^{(1)})$$

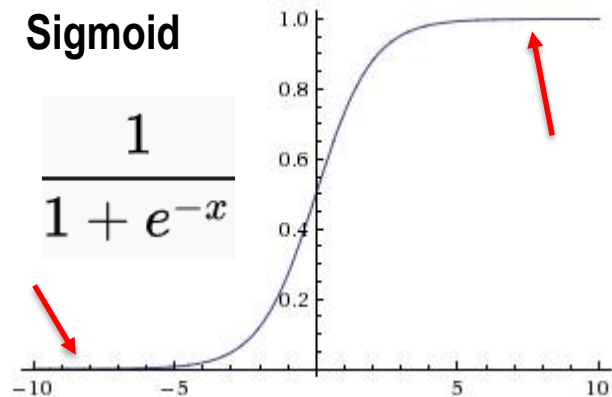
$$\underbrace{z_4^{(1)}}_{z^{(1)}} = \underbrace{\omega_4^{(1)T} x}_{W^T x} + \underbrace{b_4^{(1)}}_{b}, \quad \underbrace{a_4^{(1)}}_{a^{(1)}} = \sigma(\underbrace{z_4^{(1)}}_{z^{(1)}})$$

$$z^{(1)} = W^T x + b, \quad a^{(1)} = \sigma(z^{(1)})$$

Activation Functions (1)

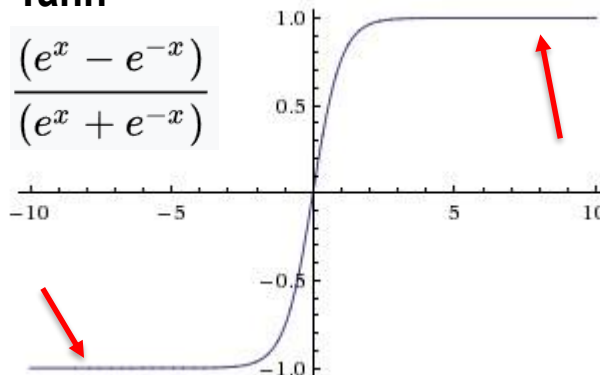
- Activation function defines the output of given input

Sigmoid



- Enforce output range to $[0, 1]$
- Saturation kills the gradients
- Outputs are not zero means

Tanh



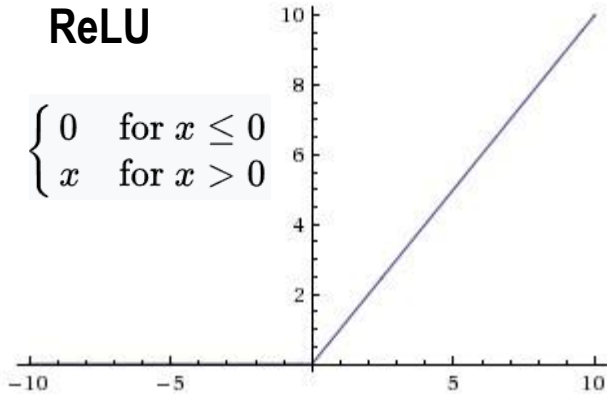
- Enforce output range to $[-1, 1]$
- Saturation kills the gradients

Activation Functions (2)

- Activation function defines the output of given input

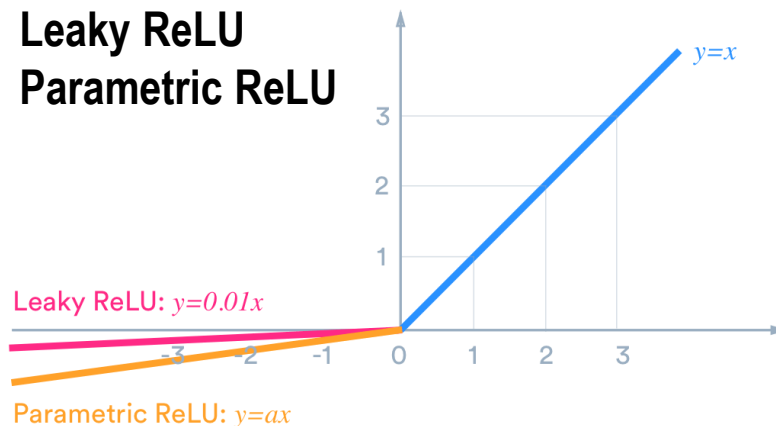
ReLU

$$\begin{cases} 0 & \text{for } x \leq 0 \\ x & \text{for } x > 0 \end{cases}$$



- Not saturate
- Computation efficient
- Faster convergence
- Dead ReLU

Leaky ReLU Parametric ReLU

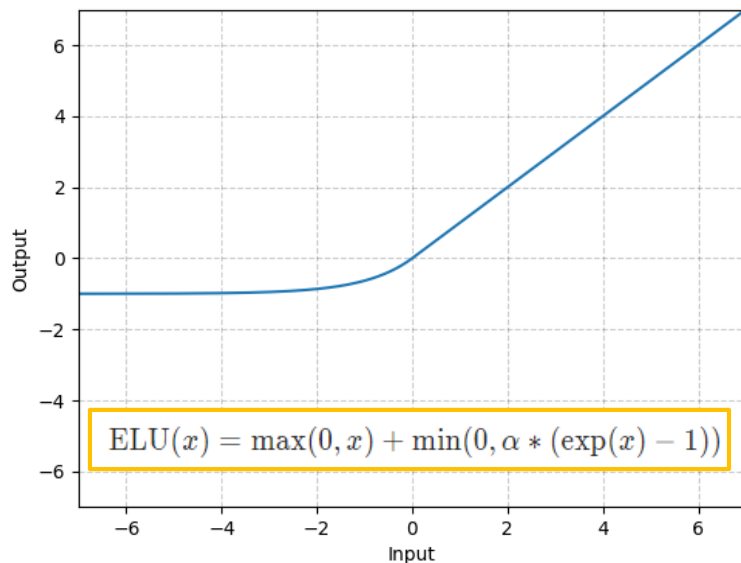


- Not saturated
- Computation efficient
- Faster convergence
- Will not "die"

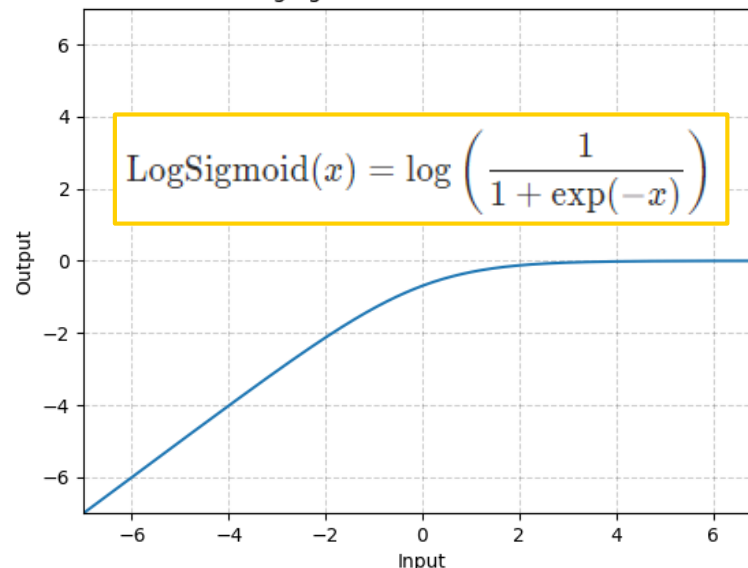


Activation Function (3)

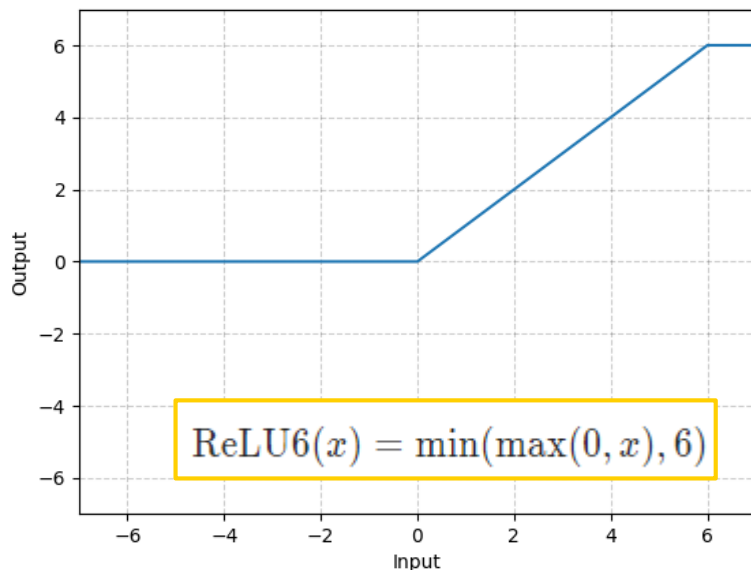
ELU activation function



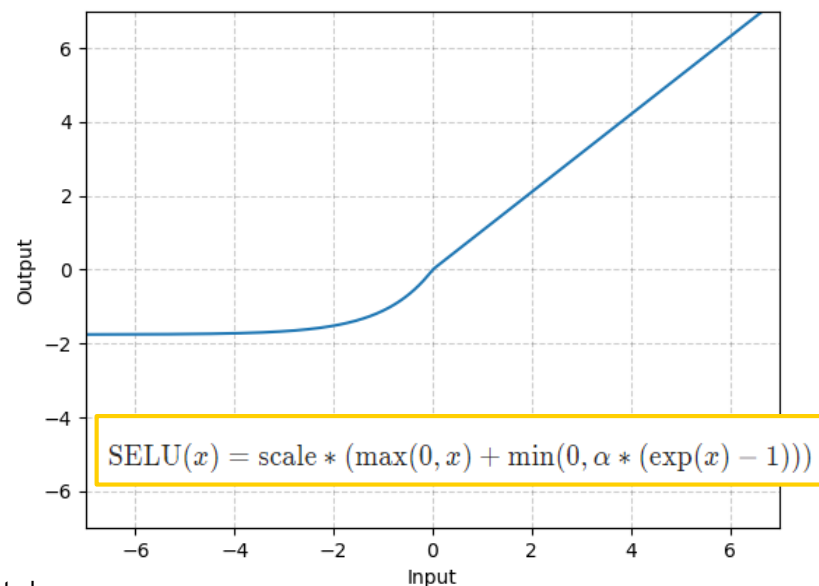
LogSigmoid activation function



ReLU6 activation function



SELU activation function



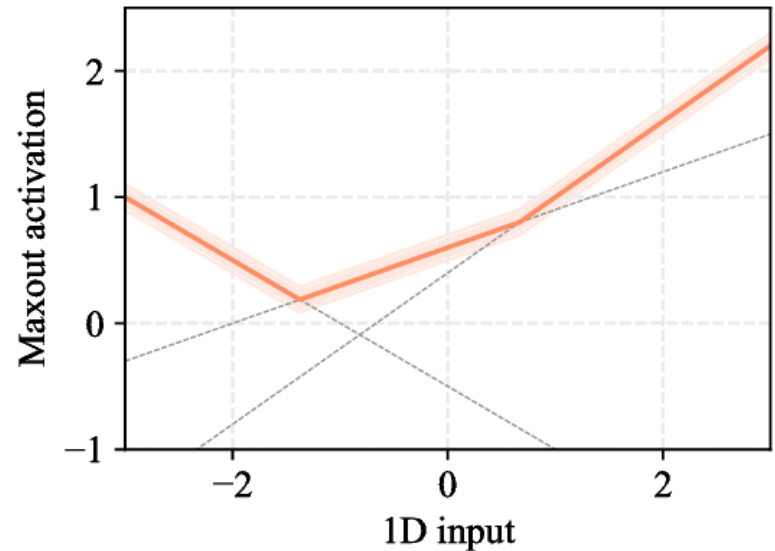
Activation (4)

- Maxout “activation”

- Generalized Relu & Leaky ReLU

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

- Not saturate
- Faster convergence
- Not die
- Double parameters



- Rule of thumb

- Use ReLU. Be careful with your learning rates
- Try out Leaky ReLU / Maxout
- Try out tanh but don't expect much
- Never use sigmoid except for the desired output range.



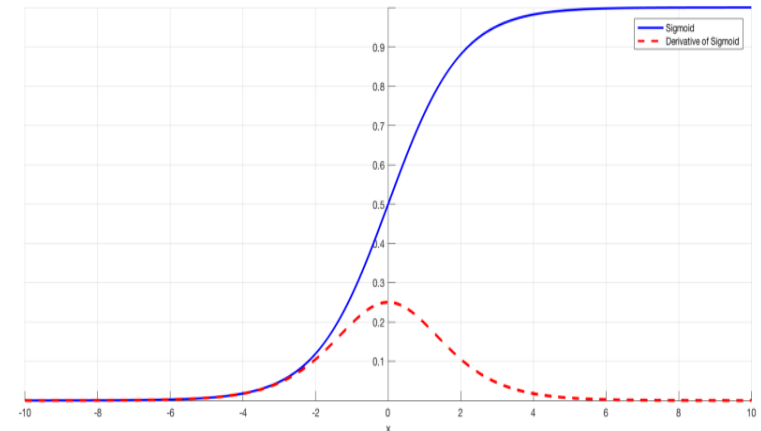
Back Propagation

- See the Lecture note from CS231n note



Random Initialization

- Gradient vanishing (deep network)
 - Certain activation change very slow
 - Gradient get smaller via back propagation
 - Too small gradients,
 - Network changing slowly
 - Slow to converge → cannot train more
- Gradient Exploding
 - Very large gradient can jump to far away
 - Objective function can be much bigger
 - Undo all previous training effort



Random Initialization

- Zero initialization or very small random initialization
 - Network is more likely to stuck at local minimal
 - All neural would follow the same directions due to same weights
- Large random initialization
 - Gradient will converse slowly if using sigmoid
 - Potential of gradient exploding
- Random initialization with a scale (depend on size of parameters)

