

Activation Functions

$$f(x) = \textit{sigmoid}(x)$$

$$z_2 = XW_1$$

$$a_2 = f(z_2)$$

$$z_3 = a_2W_2$$

$$\hat{y} = a_3 = f(z_3)$$

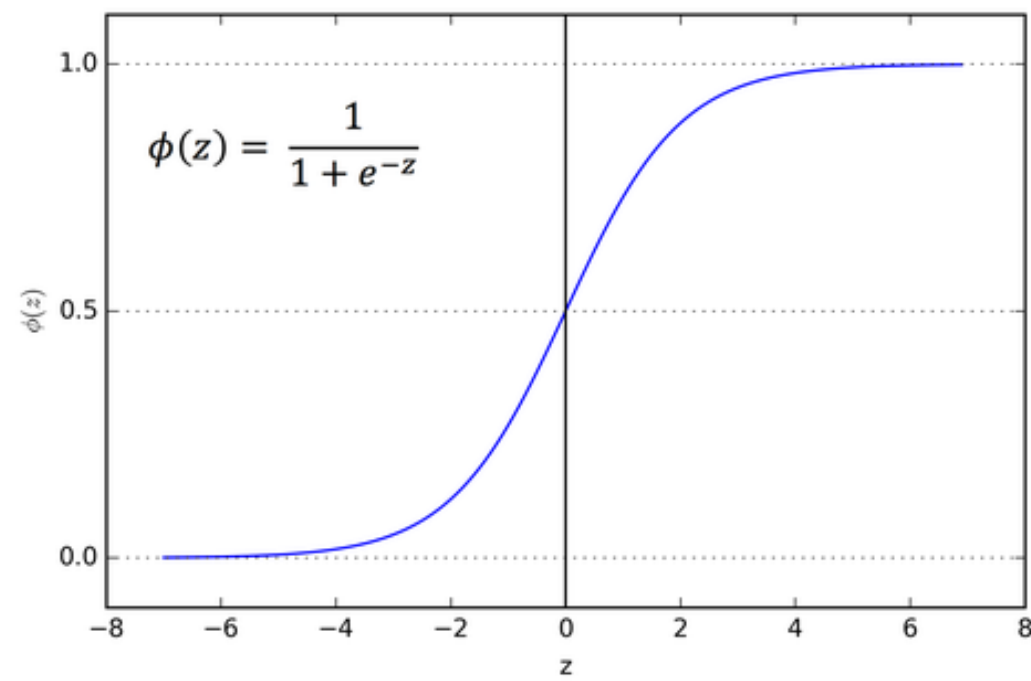


Image Credit:

<https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

IBM Watson IoT ™ © Copyright IBM Corp. 2017 Romeo Kienzler

- add non-linearity to the model
- Universal Function approximator with single hidden layer

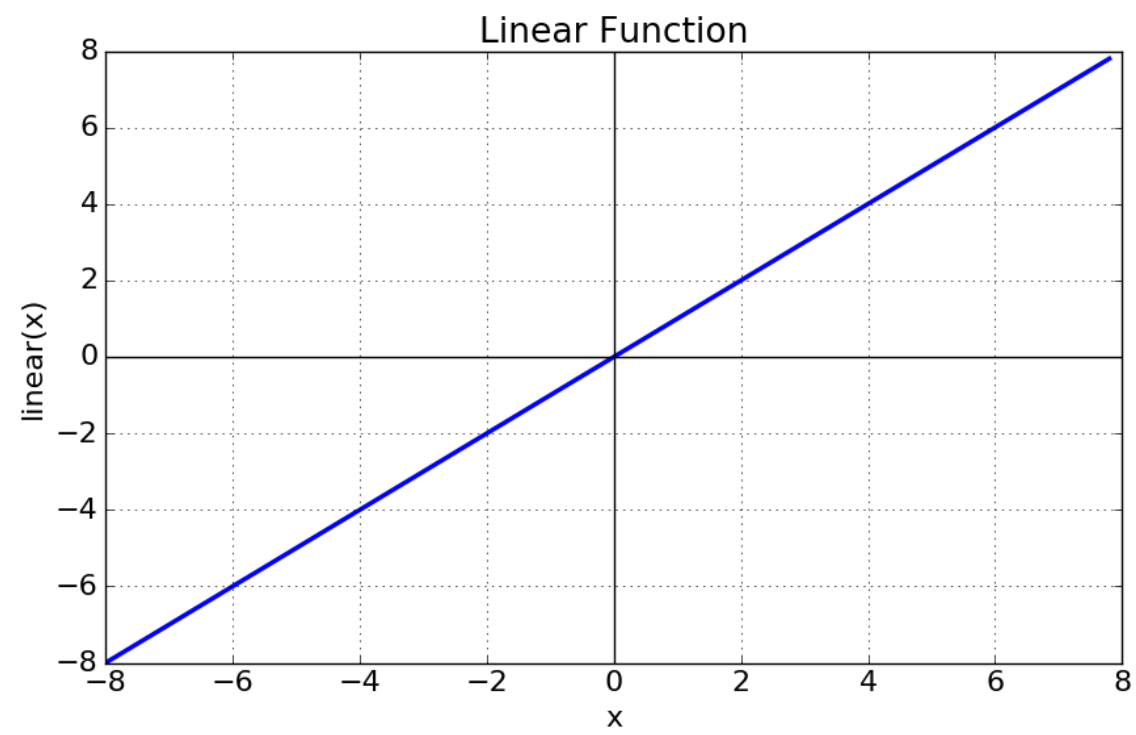


Image Credit:

<https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

IBM Watson IoT ™ © Copyright IBM Corp. 2017 Romeo Kienzler

- but can add as many layers with linear activation, never get UFA property

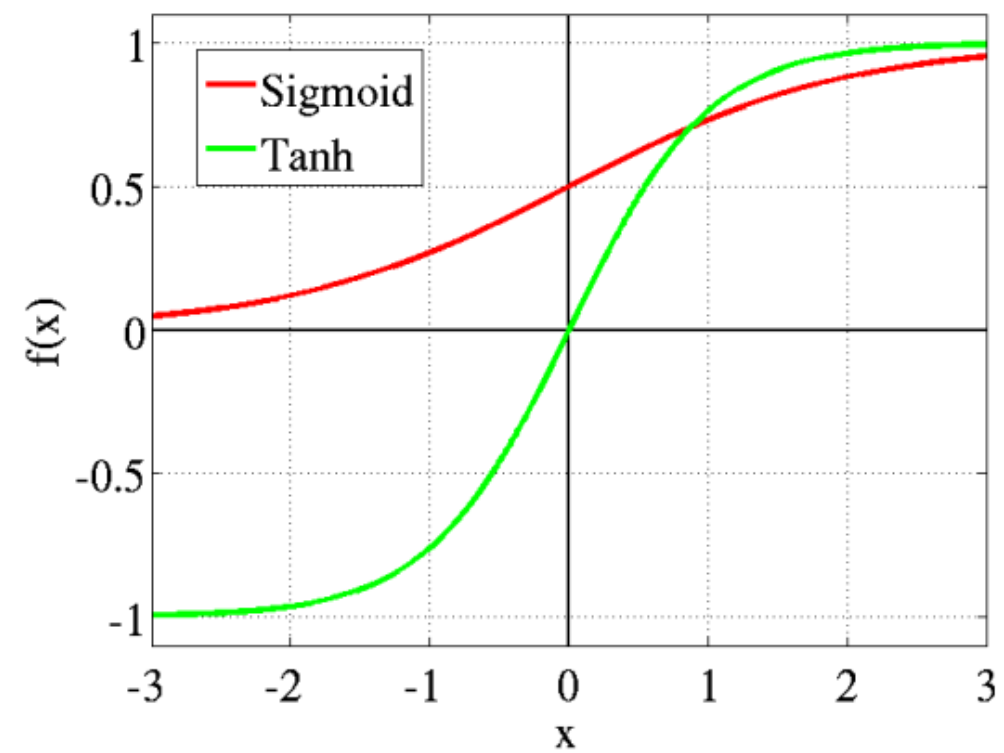


Image Credit:
<https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

IBM Watson IoT ™ © Copyright IBM Corp. 2017 Romeo Kienzler

- better than logistic sigmoid since also covers negative range
- mostly used in output layer in binary classification

$$\sigma : \mathbb{R}^K \rightarrow [0, 1]^K$$
$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

softmax

- mostly used in output layer in multiclass-classification

$$\sigma : \mathbb{R}^K \rightarrow [0, 1]^K$$

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

$$\sigma : \mathbb{R}^K \rightarrow [0, 1]^K$$
$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

$$\sigma : \mathbb{R}^K \rightarrow [0, 1]^K$$
$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

so it has the property that the sum of all classes is one, and the most probable class approximates one whereas values for the other classes approximate zero. in case training was successful ... with a one hot encoded vector .. a one hot encoded vector is a way to express class membership by just setting one element to one, the rest to zero

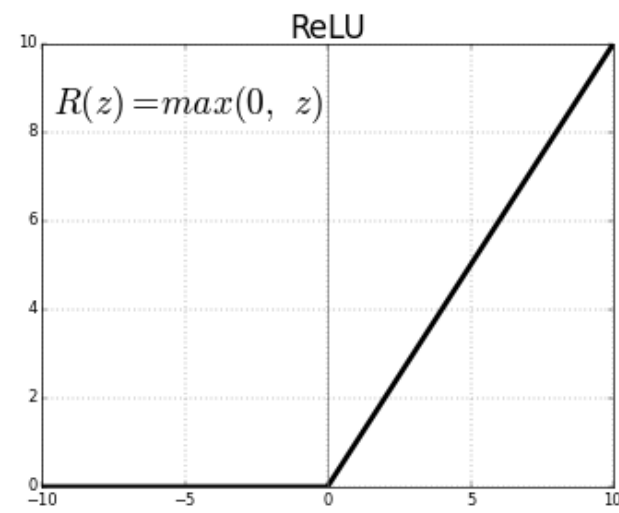
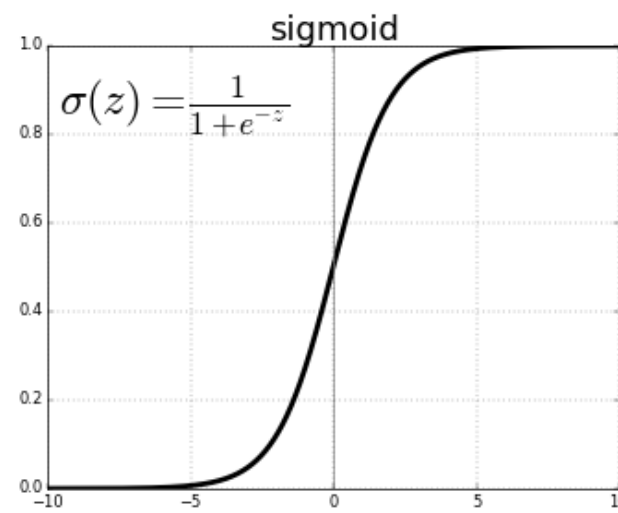


Image Credit:

<https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

IBM Watson IoT ™ © Copyright IBM Corp. 2017 Romeo Kienzler

- so doesn't this function look linear?
- not, not a straight line
- with multiple relu functions summed up we can build arbitrary functions
- mostly used activation function today ...

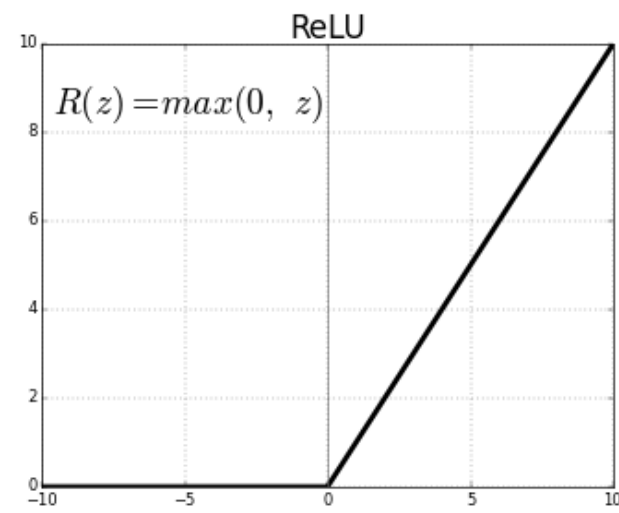
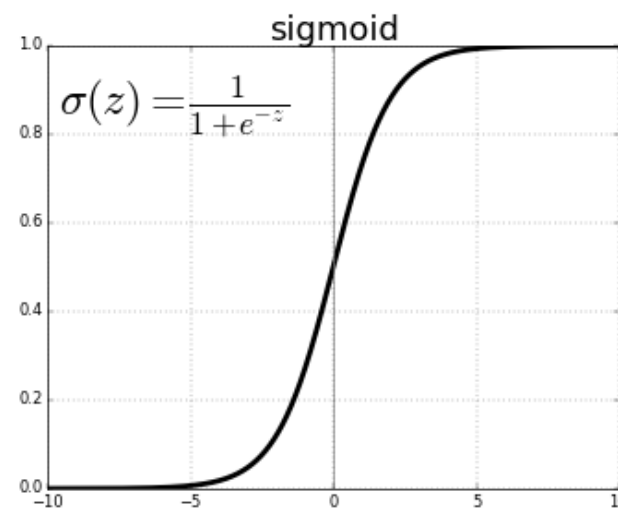


Image Credit:
<https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

IBM Watson IoT ™ © Copyright IBM Corp. 2017 Romeo Kienzler

- so doesn't this function look linear?
- not, not a straight line
- with multiple relu functions summed up we can build arbitrary functions
- mostly used activation function today ...
- NN with relu converge faster
- compu complexity is less than logistic sigmoid or tanh

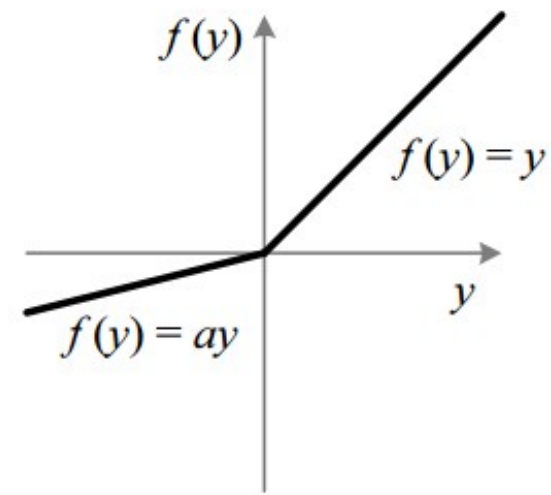
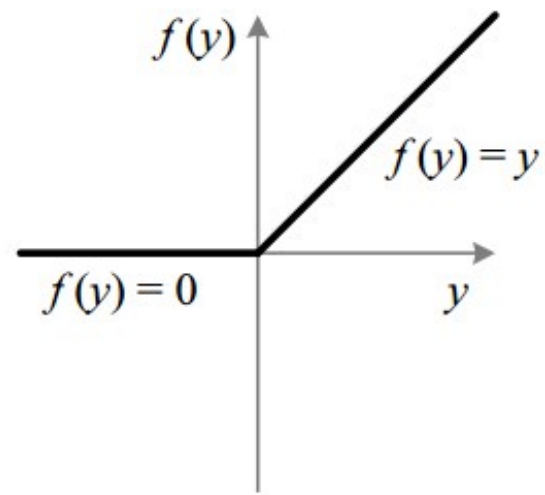


Image Credit:

<https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6>

IBM Watson IoT ™ © Copyright IBM Corp. 2017 Romeo Kienzler

- increases value range of function
- Relu can cause dead neurons
- if NN training not performing, just try it out

Summary

start with RELU whenever possible for input and hidden layers

output layer activation depends on problem and the cost function , but rule of thumb

- regression => linear
- classification => sigmoid / softmax

nobody stops you to try your own combinations

Bias-Variance Tradeoff