

# Introduction to Neural Networks (NN)

# Why Neural Net

- Superior to computer vision techniques
- Auto feature engineering
- No need for manual inspection

# Why Neural Net 2

- Able to recognize patterns of complex / unexpected inputs
- Image / Text / Audio / Video Search

# Goals

- Application of NN
- Software Ecosystem
- Logistic / Linear Regression as example of NN
- Multi Hidden Layer NN

# MNIST Dataset



# Why Neural Net 2

Classifier	Test Error (%)	Reference
Large/deep conv. net	99.77	Multi-column Deep Neural Networks for Image Classification
Virtual SVM, deg-9 poly, 2-pixel jittered	99.44	Training Invariant Support Vector Machines
product of stumps on Haar f.	99.13	Boosting products of base classifiers

# Why not Neural Net

- Hard to tune
- Computationally intensive
- **Lack of feature introspection**
  - ★ [arxiv.org/abs/1503.02531](https://arxiv.org/abs/1503.02531)

# Ecosystem (May, 2015)

- **Python**

- ★ **Theano (Low level + GPU)**

- ★ [deeplearning.net/software/theano/](http://deeplearning.net/software/theano/)

- ★ **Lasagne (Abstraction of Theano)**

- ★ [lasagne.readthedocs.org/en/latest/index.html](http://lasagne.readthedocs.org/en/latest/index.html)

- ★ **Others:** <https://www.cbinsights.com/blog/python-tools-machine-learning/>

- **Java / Scala**

- ★ **DL4J (Run on distributed systems):** [deeplearning4j.org/](http://deeplearning4j.org/)

- **C++**

- ★ **Caffe (GPU):** [caffe.berkeleyvision.org/](http://caffe.berkeleyvision.org/)



# Terminology

**Neural Network == Artificial Neural Network**

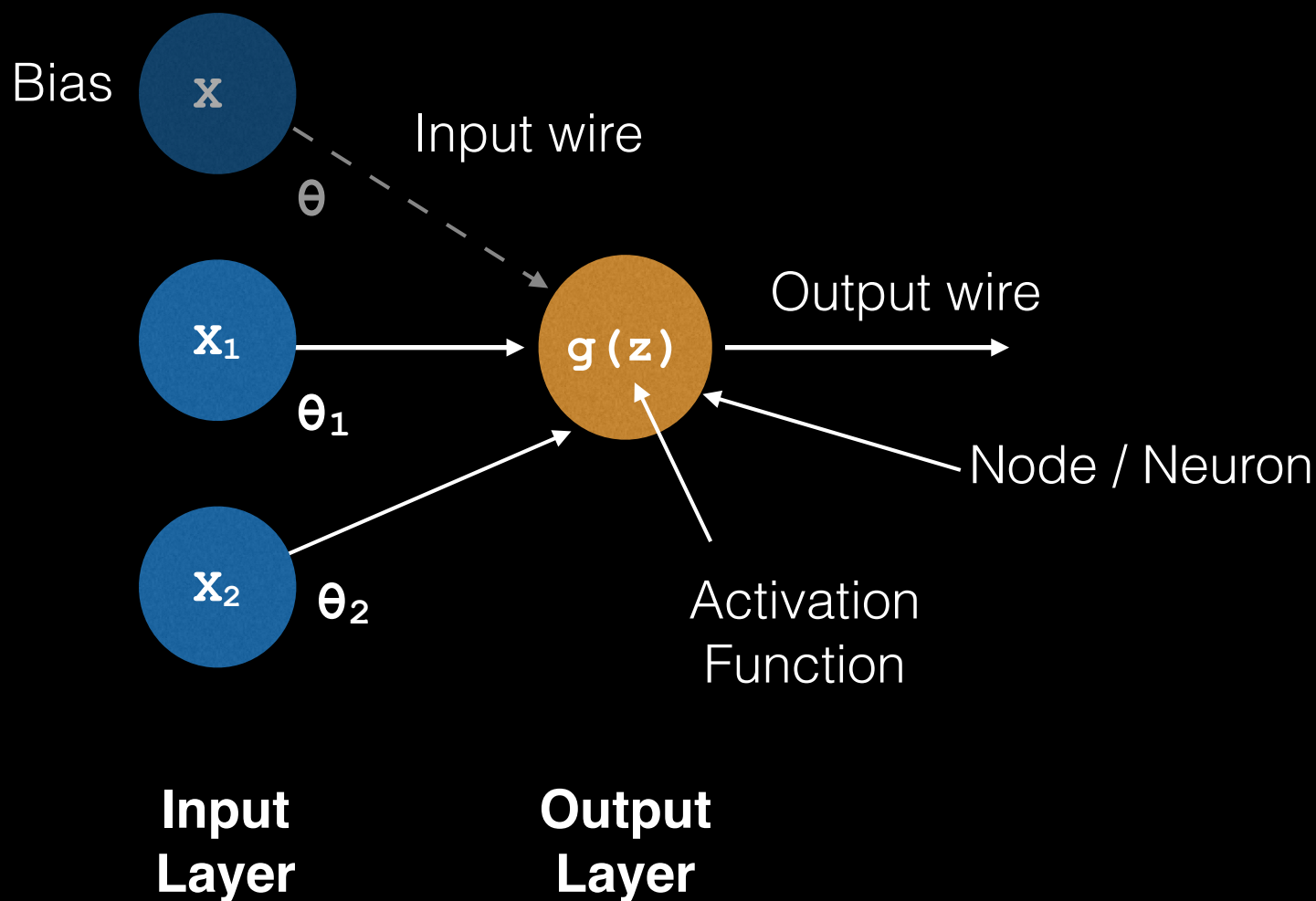
## **Deep Learning**

- 3 or more hidden layers in a NN

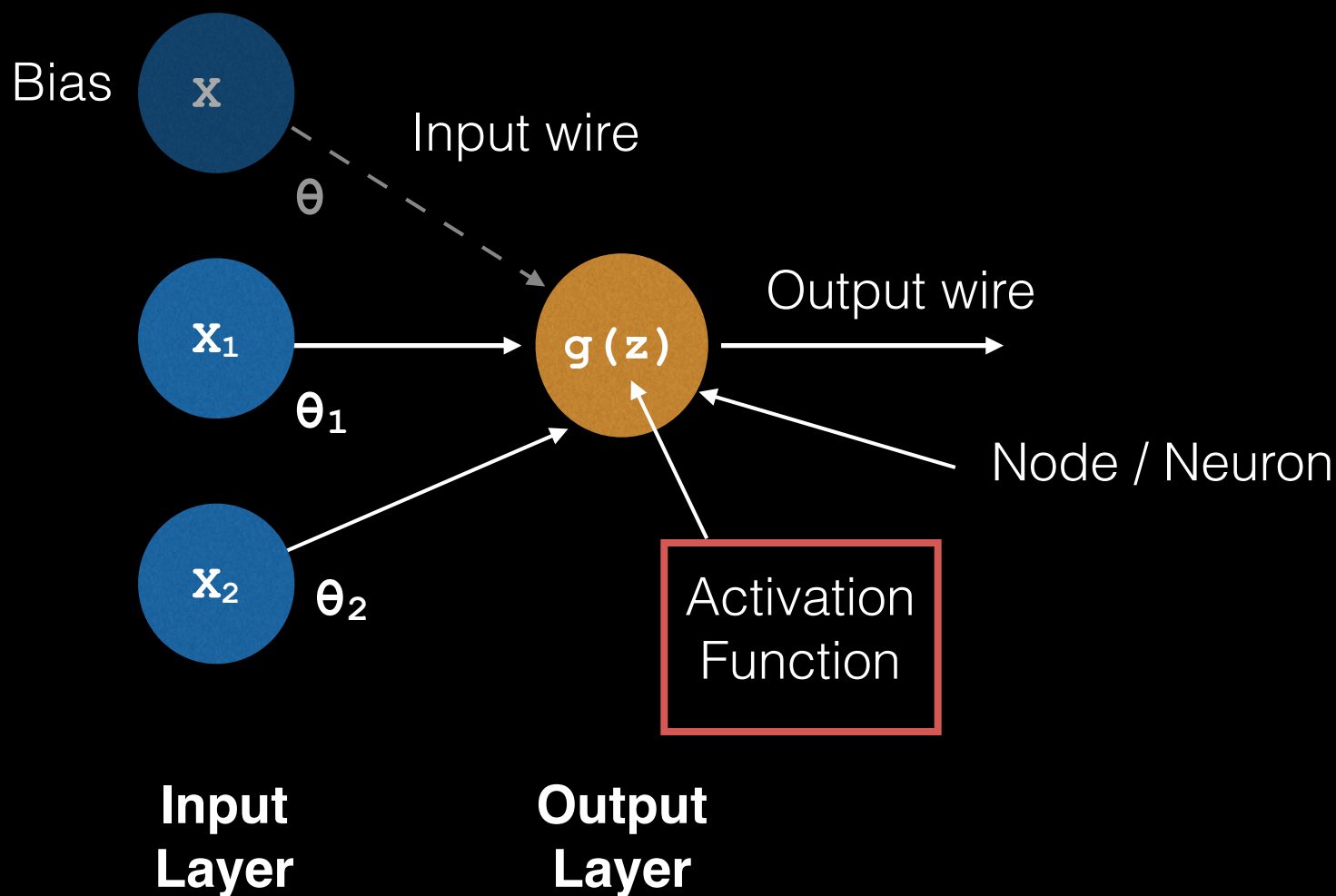
# Feed-Forward NN

- Information moves in **1 direction**
- Input → Output

# Components of NN (2-layer)



# Components of NN (2-layer)



# Activation Function

- **Transforms** certain input to an other domain

Input

Parameter

Output

$x$

$\theta^T$

Linear  
Activation

$\hat{y}$

Input

Parameter

Output

$x$

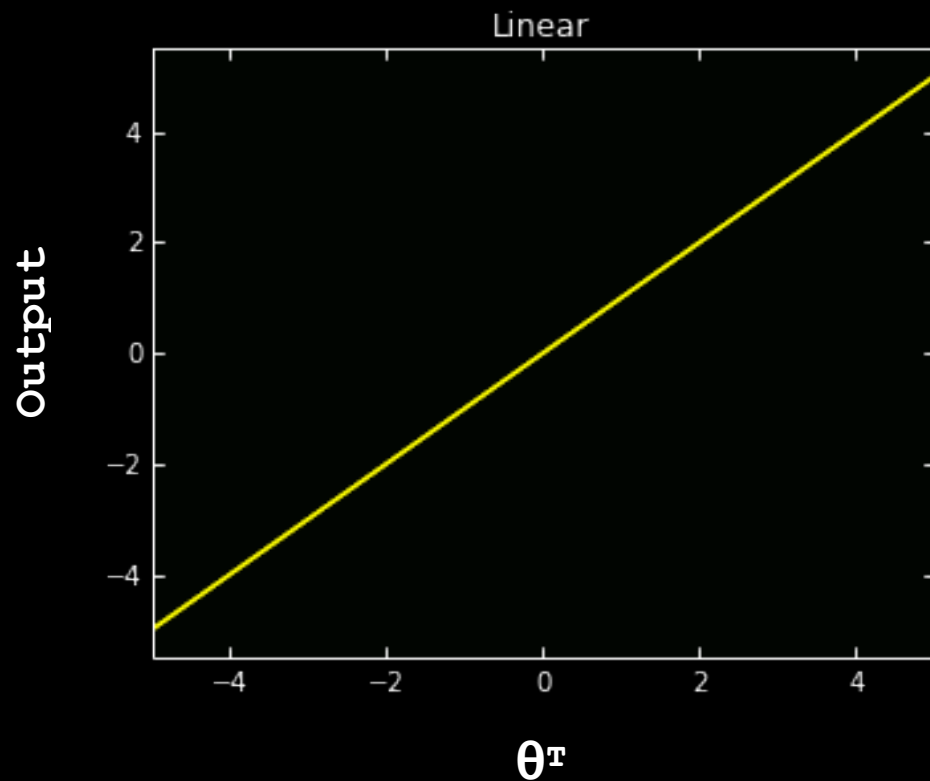
$\theta^T$

Logistic  
Activation

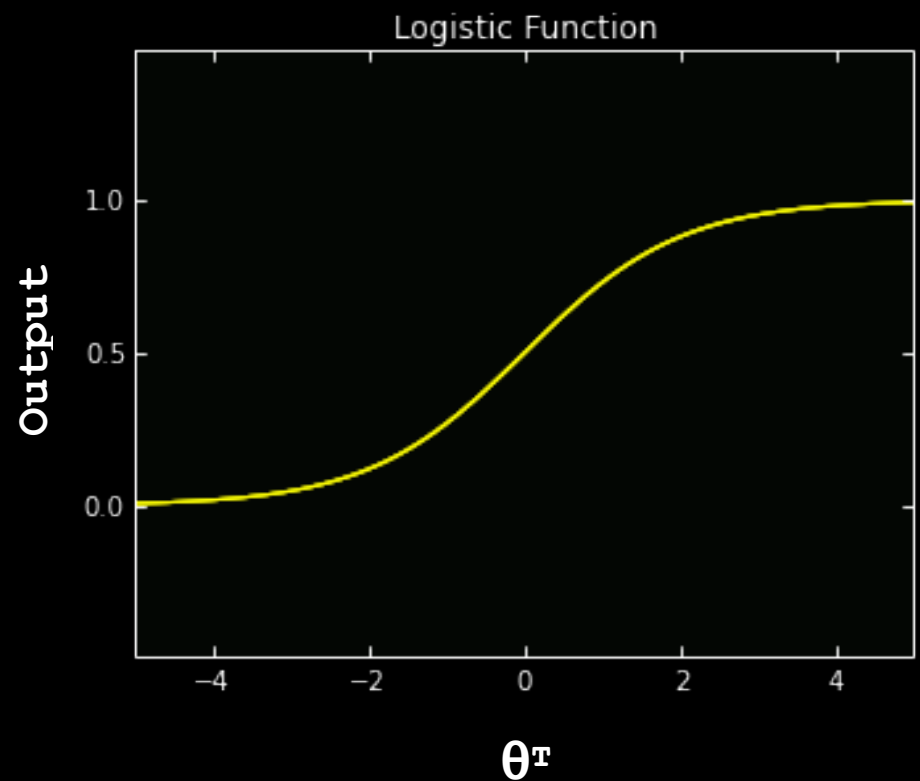
$P(Y=1 \mid x; \theta)$

# Output Layer Activation Function

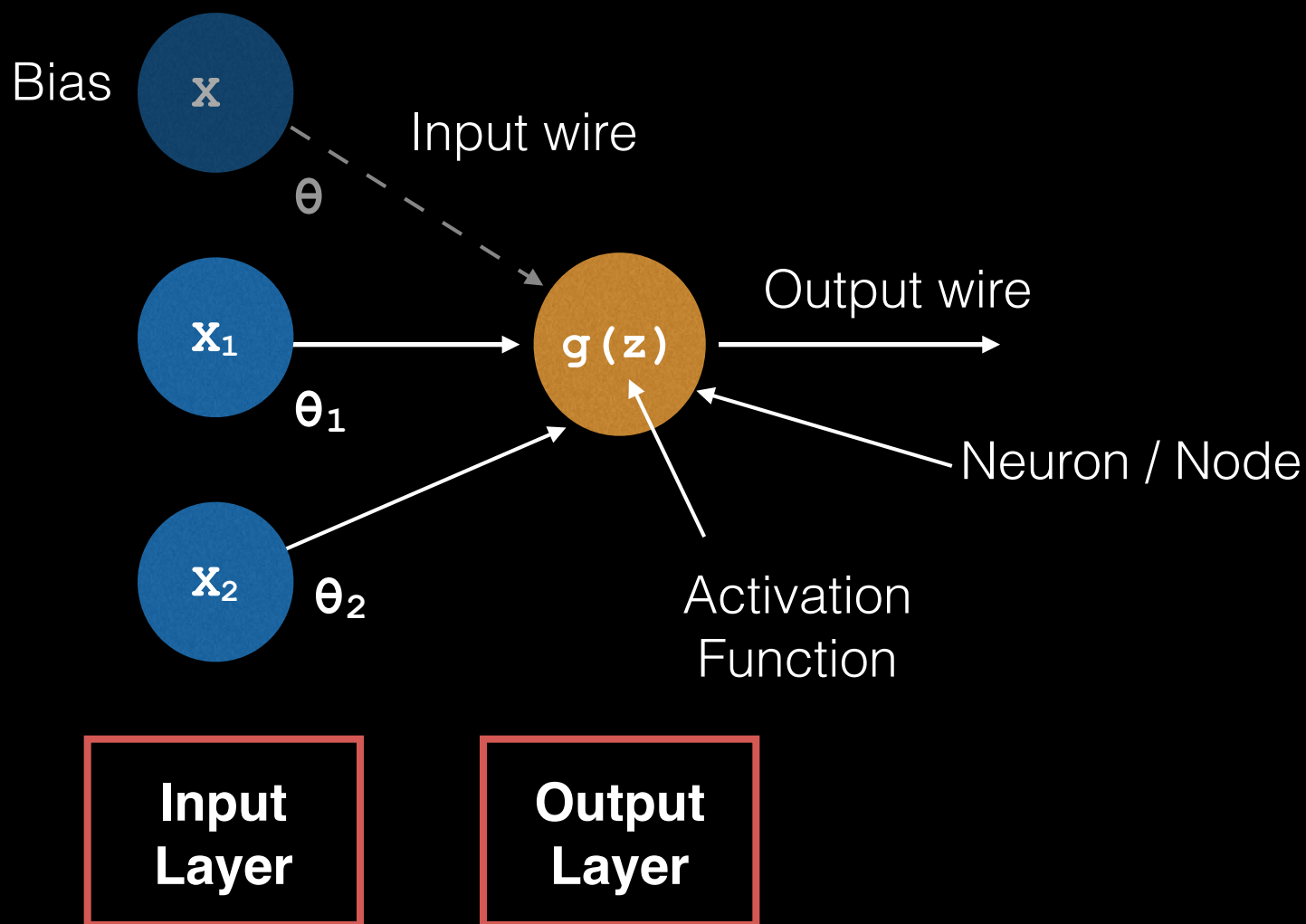
(Regression)



(Classification)



# Components of NN (2-layer)





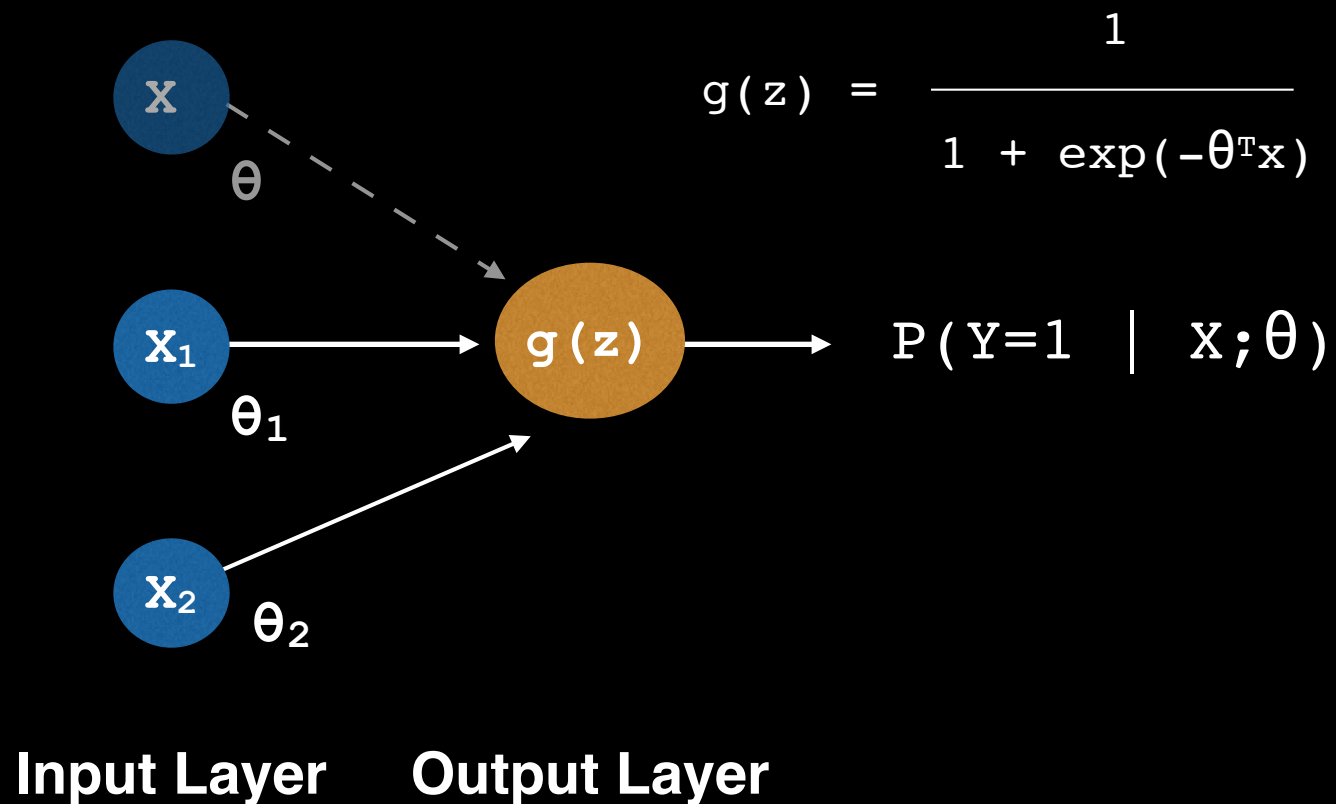
# How many input nodes

- Number of features ( $x_i$ ) + 1 bias term

# How many output nodes

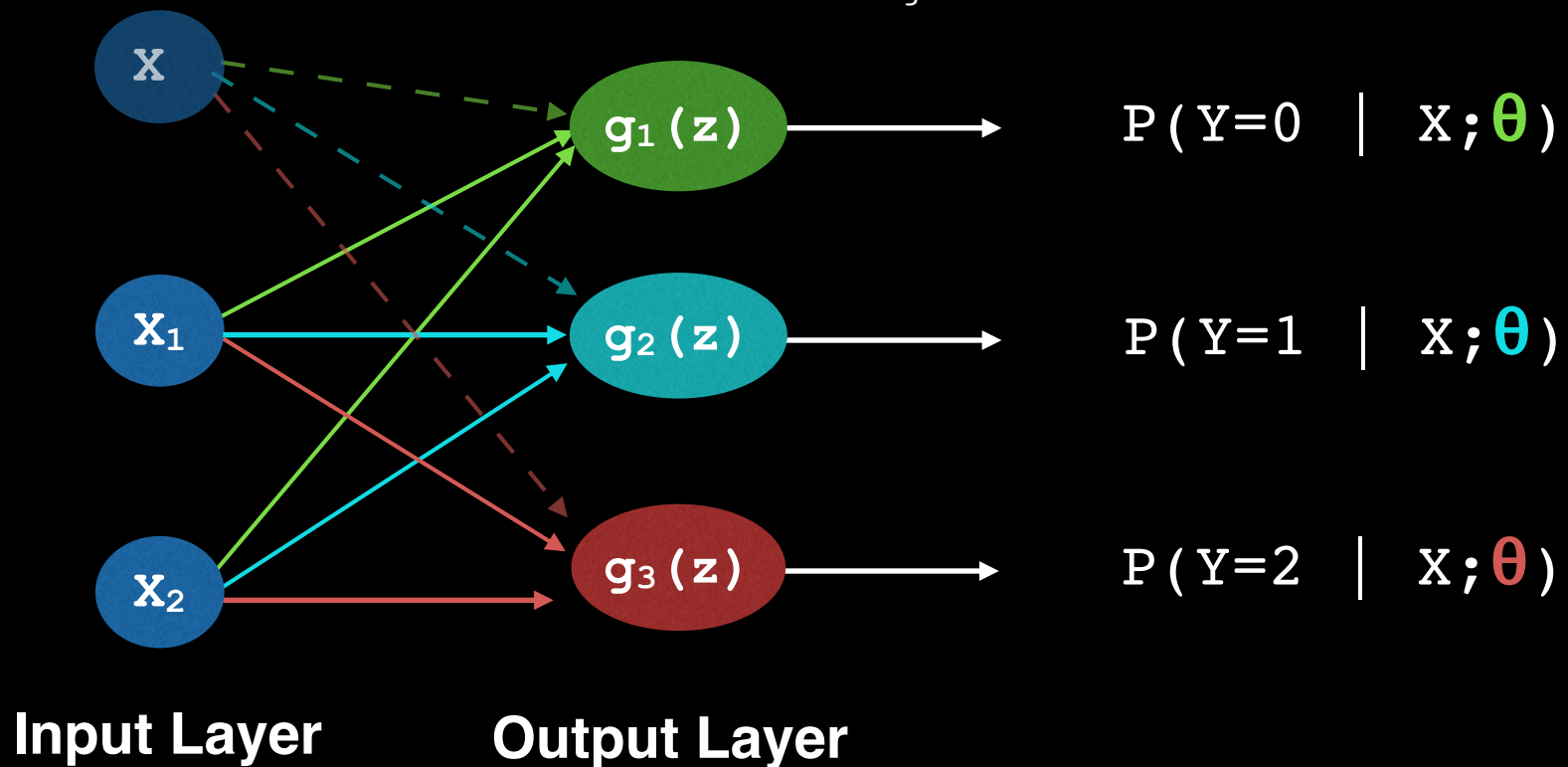
- **Continuous response (Regression)**
  - ★ 1 output node
- **Discrete response (Classification)**
  - ★ Number of unique responses

# Logistic Regression (2 classes)



# Logistic Regression (3 classes - Softmax)

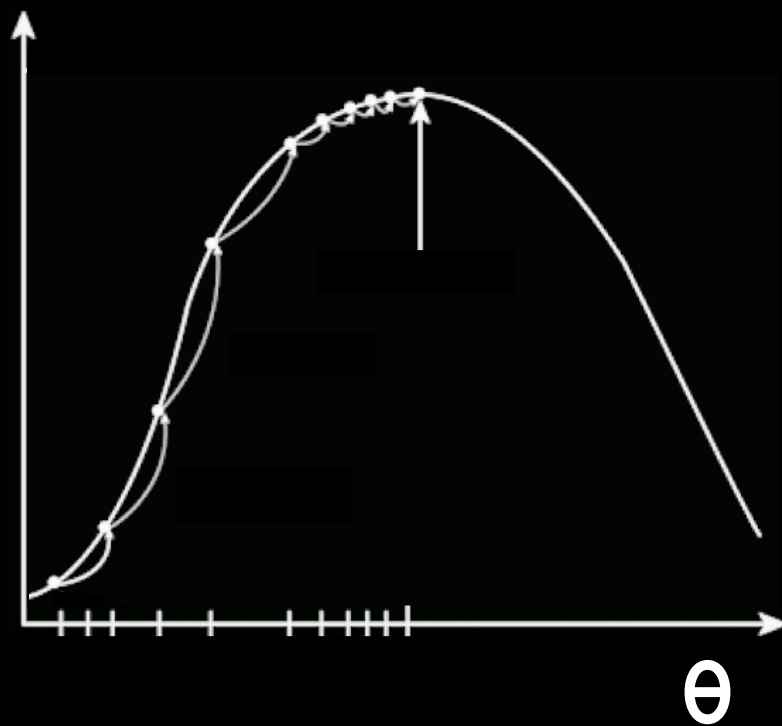
$$g_i(z) = \frac{\exp(\theta_i^T \mathbf{x})}{\sum_{j=1}^k \exp(\theta_j^T \mathbf{x})}$$



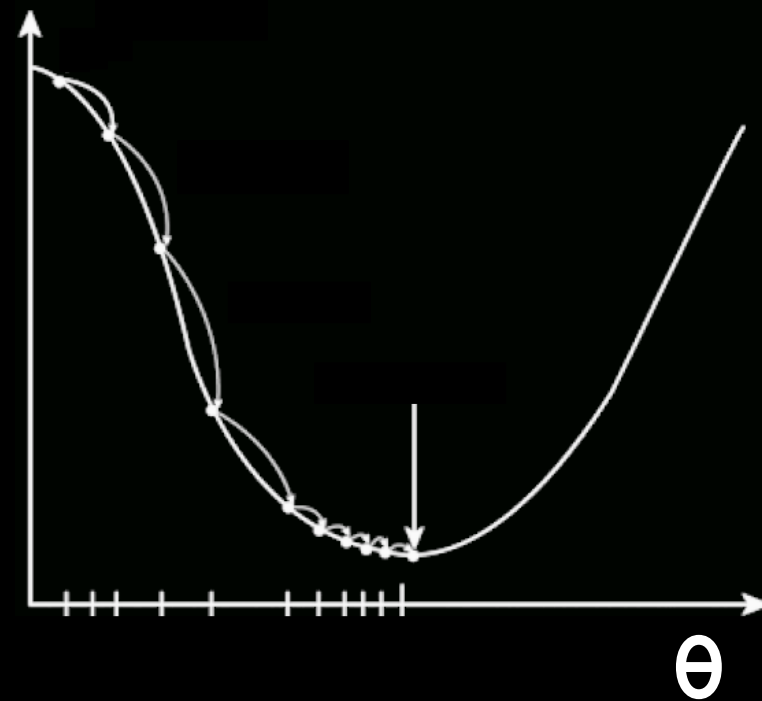
# Optimization

- Optimize  $\theta$
- Minimize loss function / Maximize Likelihood
- Gradient descent / ascent optimization
- Incrementally adjust  $\theta$  to min. loss / max. likelihood

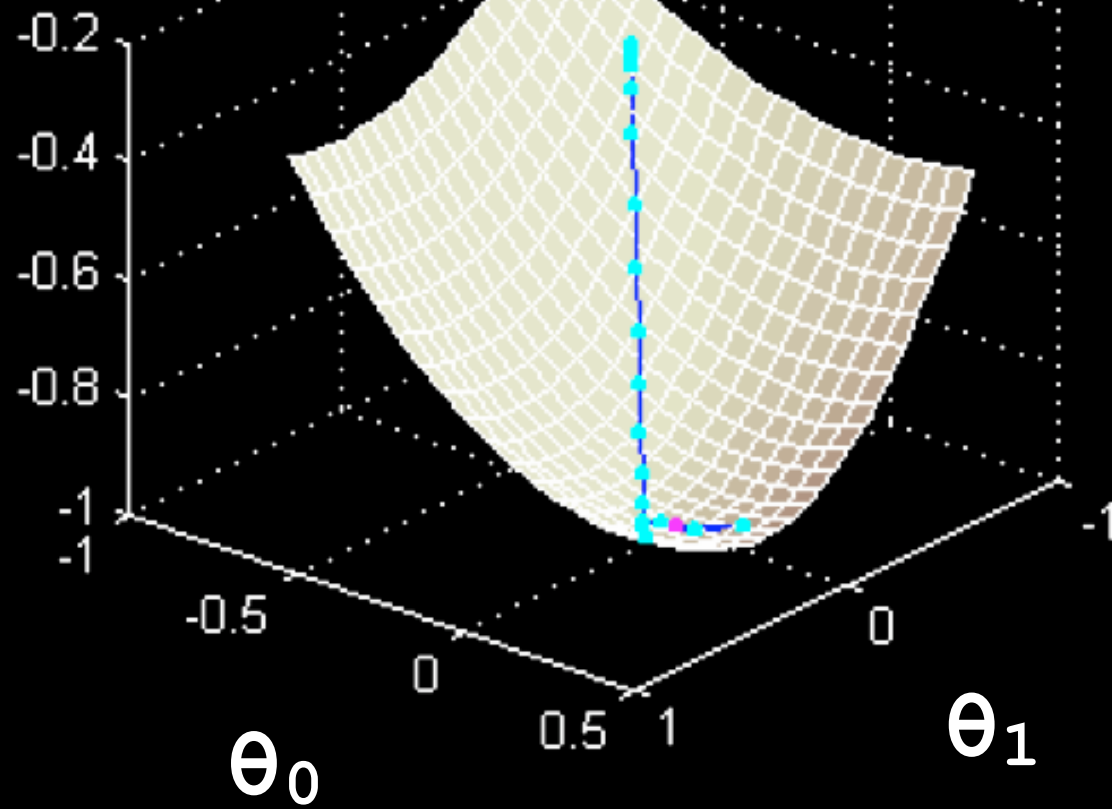
Likelihood



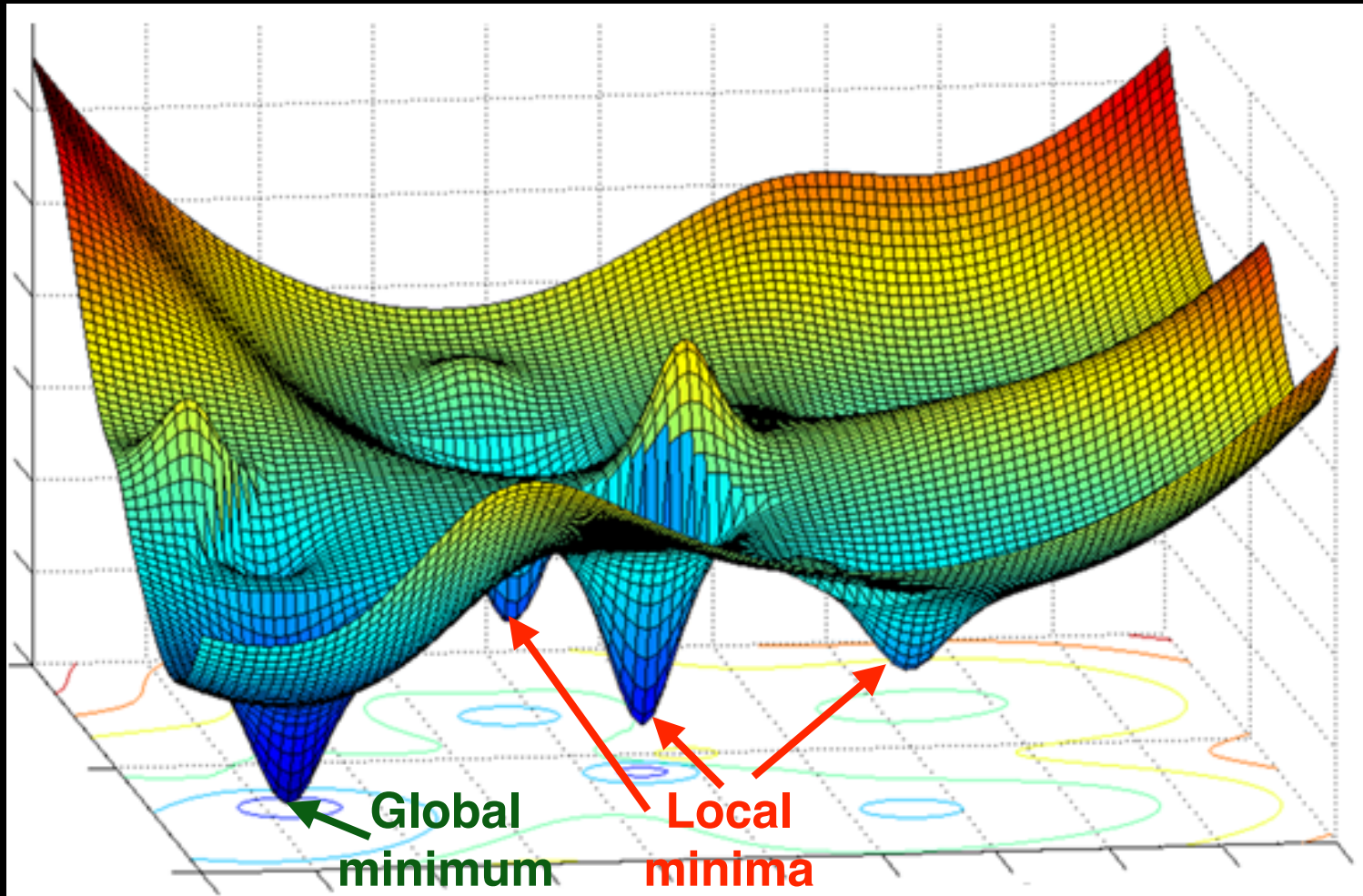
Loss



**Loss**



# Local minima





# Hyper-parameters

- **Learning Rate**

- ★ Update a fraction of the gradient

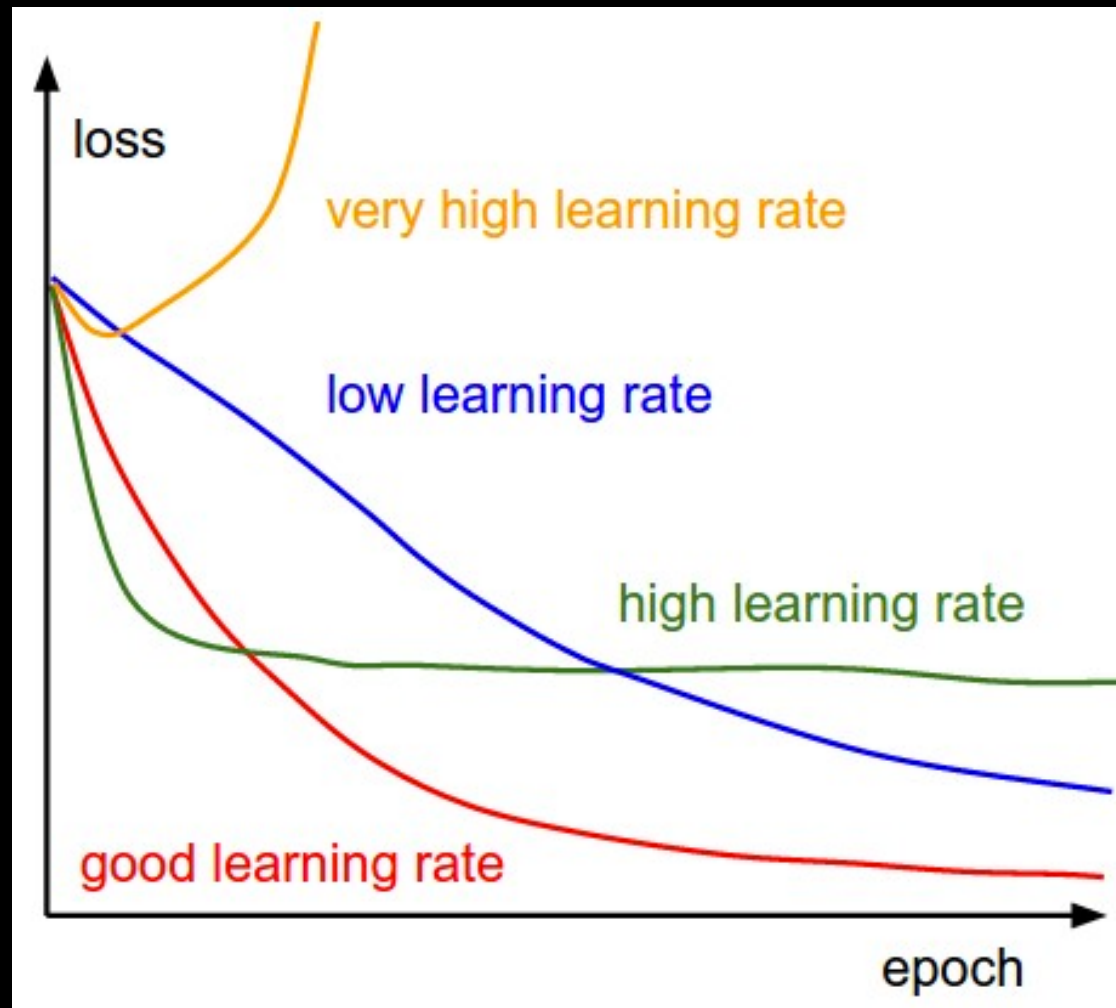
- **Momentum**

- ★ Update  $\theta$  incorporating the last “gradient”

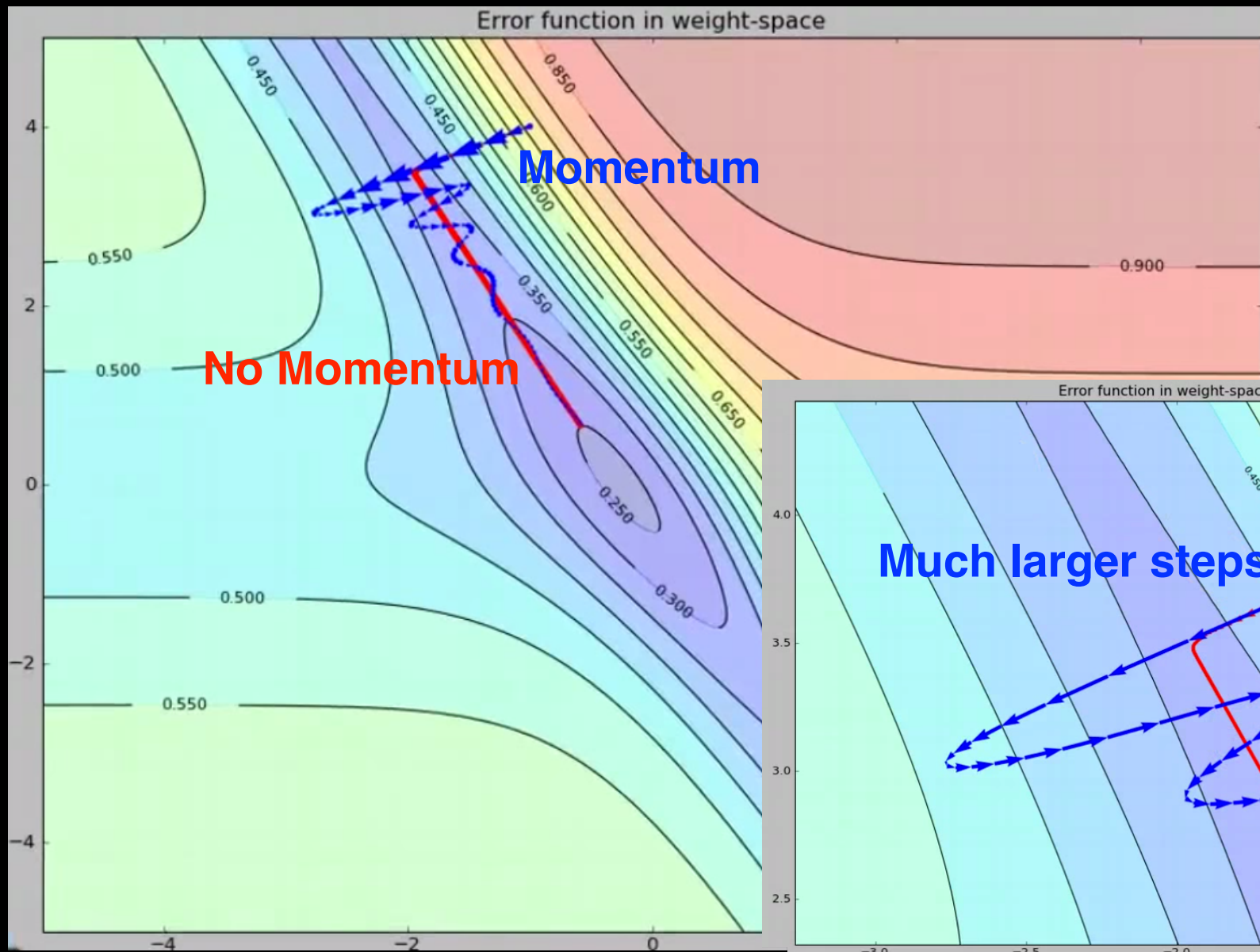
- **Epoch**

- ★ Times going through the whole set of data

# Learning Rate



# Momentum



$$V_{t+1} = \mu V_t - \alpha \nabla L(W_t)$$

Diagram illustrating the update rule for Velocity ( $V$ ) in an optimization algorithm:

- $V_{t+1}$ : Next Velocity
- $\mu$ : Momentum
- $V_t$ : Velocity
- $\alpha$ : Learning rate
- $\nabla L(W_t)$ : Gradient

$$W_{t+1} = W_t + V_{t+1}$$

Diagram illustrating the update rule for the parameters ( $W$ ):

- $W_{t+1}$ : Next  $\theta$ s
- $W_t$ : Current  $\theta$ s

# Epoch vs Iteration

- An **iteration** is an update of  $\theta$
- An **epoch** contains **iterations** through the whole data set

# Mini-batch Gradient Descent

- Compute gradient with few data points
- Shuffle data and loop through mini-batches
- Repeat until convergence
- Generally faster than batch GD  
(More frequent gradient updates)

# Review

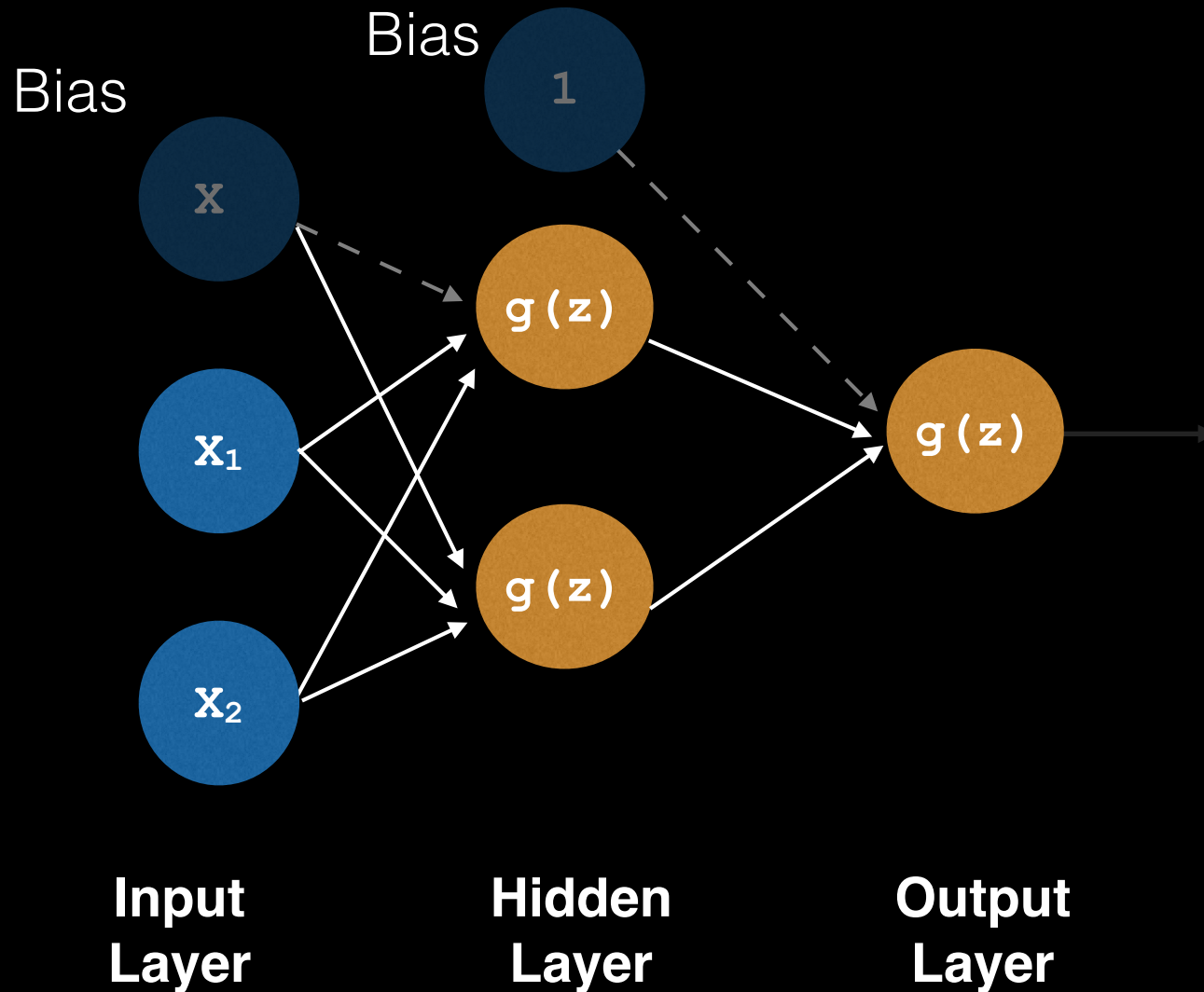
- Basic architecture of feed-forward NN
- Activation functions
- Gradient Descent Recap
- Different hyper-parameters to optimize parameters
- **Tips and Tricks**
  - ★ [research.microsoft.com/pubs/192769/tricks-2012.pdf](https://research.microsoft.com/pubs/192769/tricks-2012.pdf)
- Mini-batch GD

Break



# Multi Hidden Layer NN

# 3-layer NN



# Hidden Layer

- **Hidden:** You do not observe the output
- Can have multiple hidden layers
- 1 bias term per hidden layer

# How many hidden layers

- **Usually  $\leq 3$**
- More layers for more complex relationships
- More hidden layers = Harder to optimize
- Longer to train and tune hyper-parameters

# How many Neurons

- **Per layer:**

$$\# \text{ of outputs} \leq x \leq \# \text{ of inputs}$$

(High # of inputs)

- **Overall:**

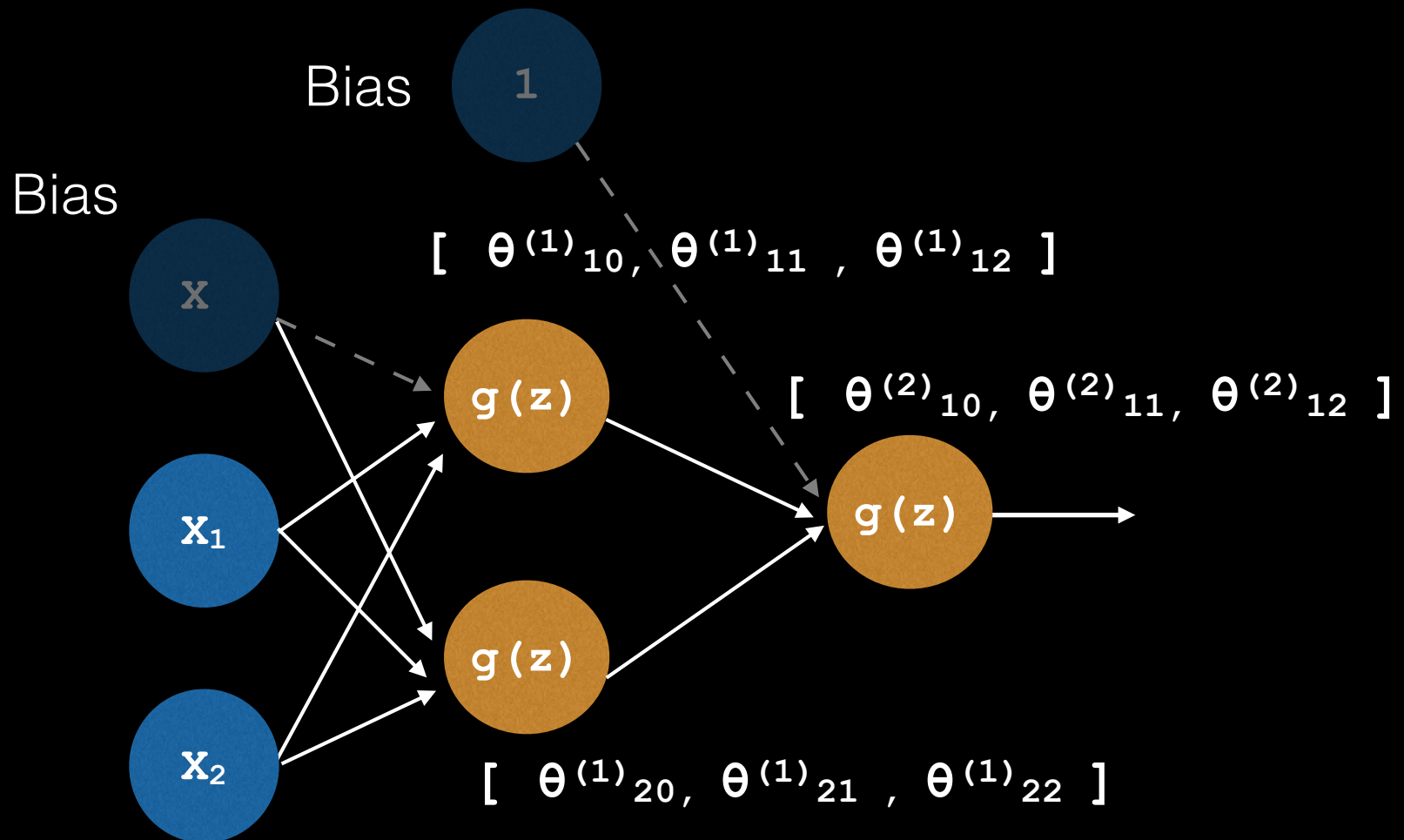
$$\# \text{ of data points} \geq \frac{\# \text{ of } \theta s}{30}$$

(Noisy data)

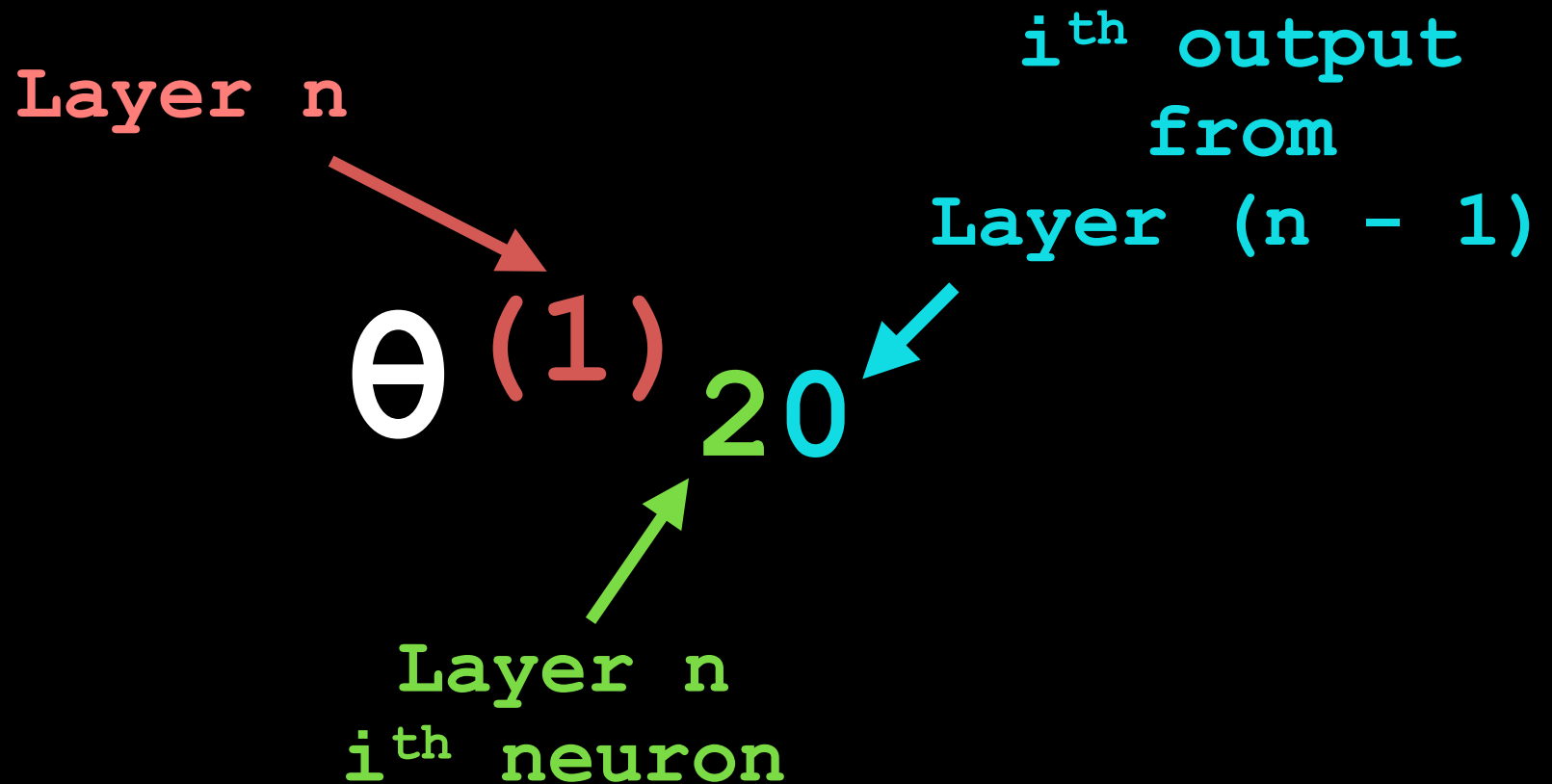
$$\# \text{ of data points} \geq \frac{\# \text{ of } \theta s}{2}$$

(Completely noise free)

# How many $\theta$ s



# $\theta$ Notation





- # of layers , # of neurons are guidelines
- Much debate around

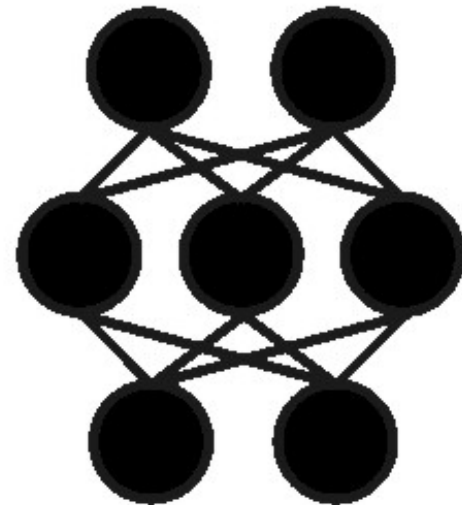
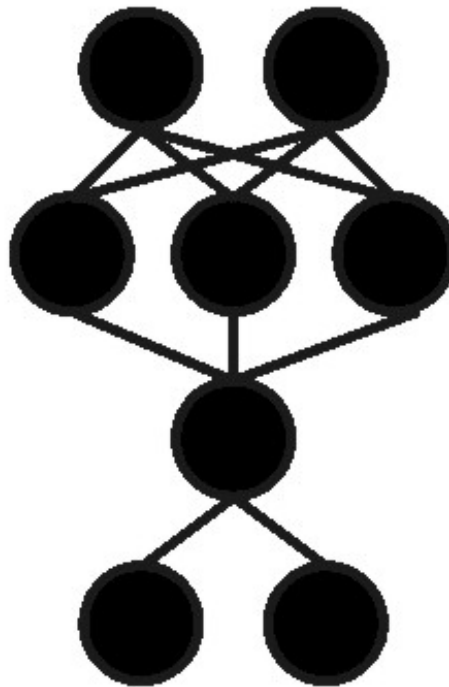
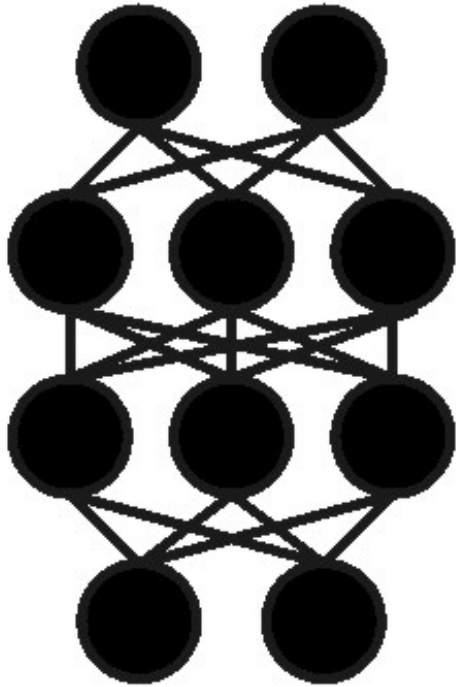
## **1. Network configuration**

## **2. Hyper-parameter tuning**

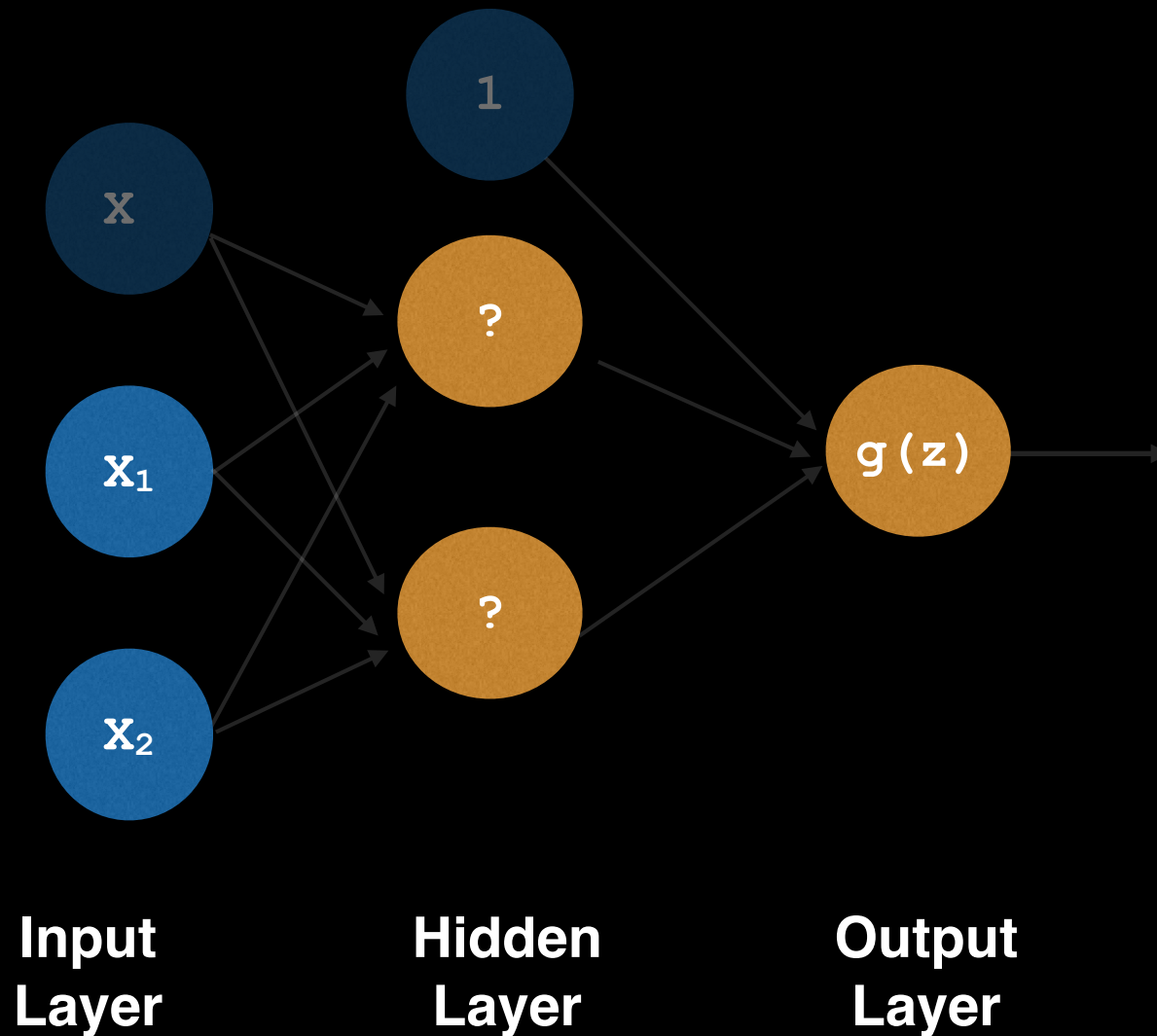
- Mostly trial-and-error

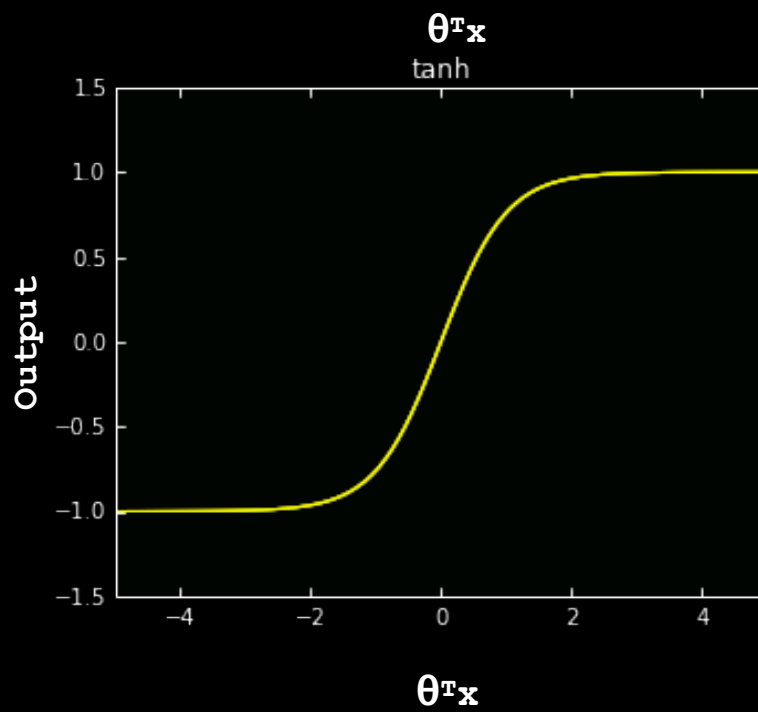
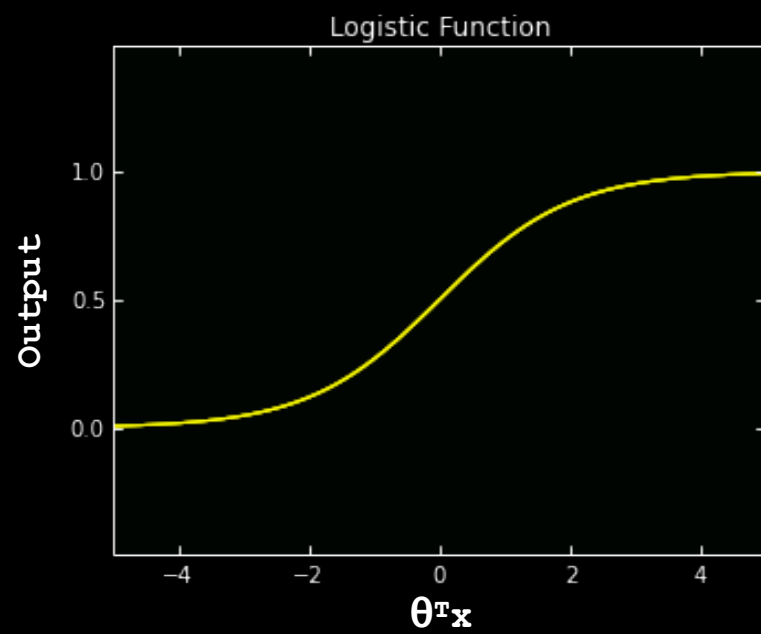
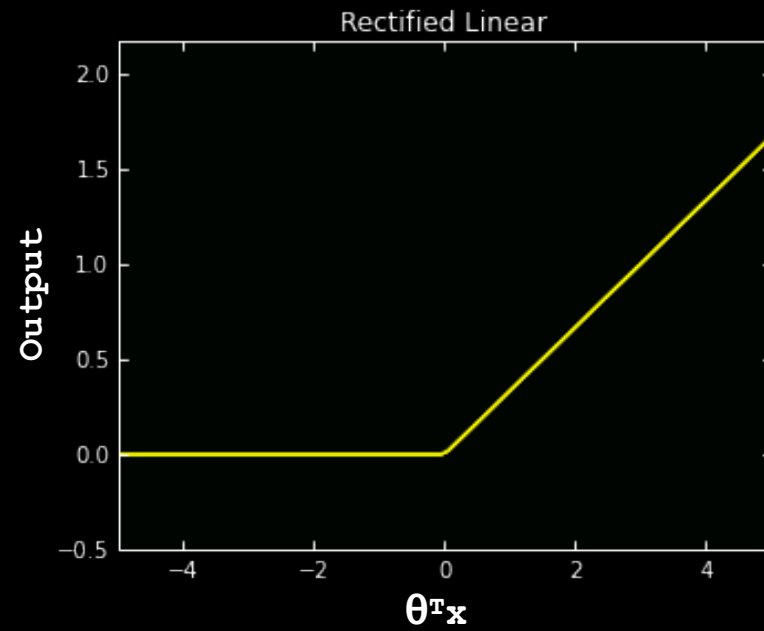


# Network Topology



# Activation Function (Hidden Layer)





# Rectify Linear (ReL/ ReLU)

- Most common activation function
- Quicker to optimize (Fewer epochs)

# Rectify Linear Cont.

- Zero out data where  $\theta^T \mathbf{x} < 0$
- Emphasize data where  $\theta^T \mathbf{x} > 0$  (Without bounds)
- Sparsity / Overfitting

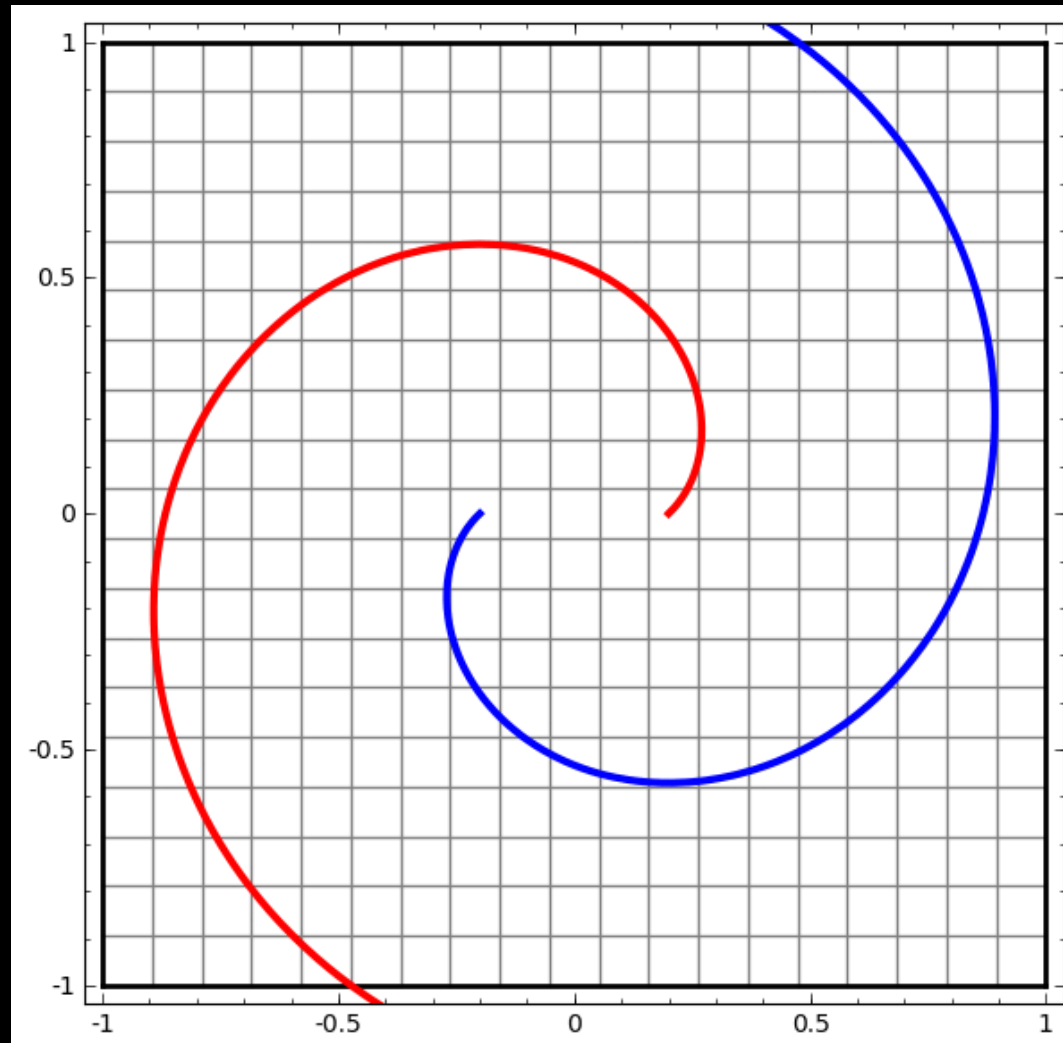
# Tanh

- Similar to sigmoid, except over range  $(-1, 1)$
- Used if you features are continuous and has -ve values
- Quicker to optimize than sigmoid

# Sigmoid / Logistic

- Over the range  $(0, 1)$
- Used if you features are binary

# Activation Function Transformation

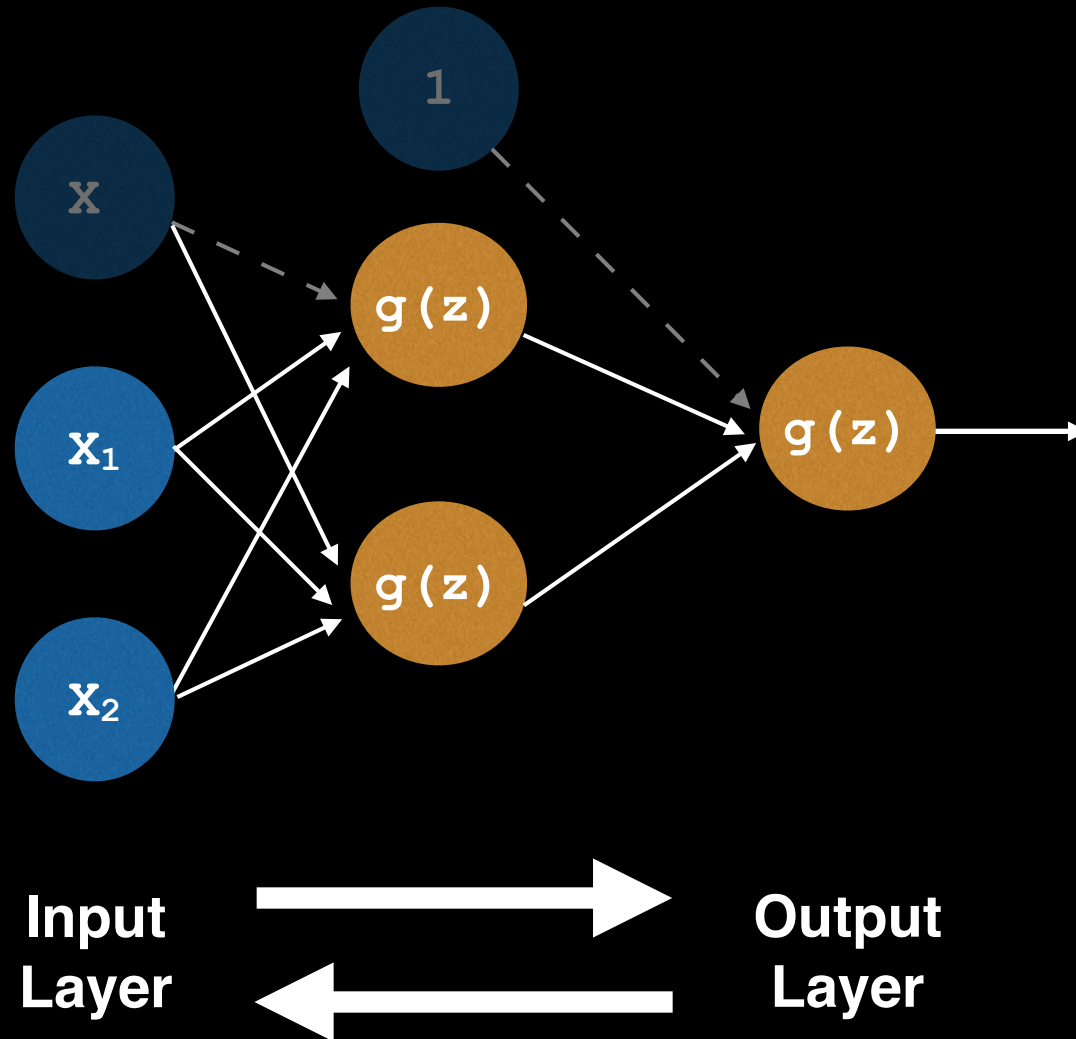




Almost there ...



# Forward / Backward Propagation



# Next Steps

## **Convolutional Neural Network**

- NN with a convolution layer(s) as the first layer(s)

## **Recurrent Neural Network**

- NN where information moves in a directed cycle

# Next Step Resources

- **CS231n: Convolutional Neural Networks for Visual Recognition**

- ★ [cs231n.github.io/](https://cs231n.github.io/)

- **More Mathematical Treatment**

- ★ [deeplearning.net/tutorial/lenet.html](https://deeplearning.net/tutorial/lenet.html)



**TRAIN A NEURAL NET, THEY  
SAID**

**IT WILL BE PERFECT, THEY  
SAID**

memegenerator.net