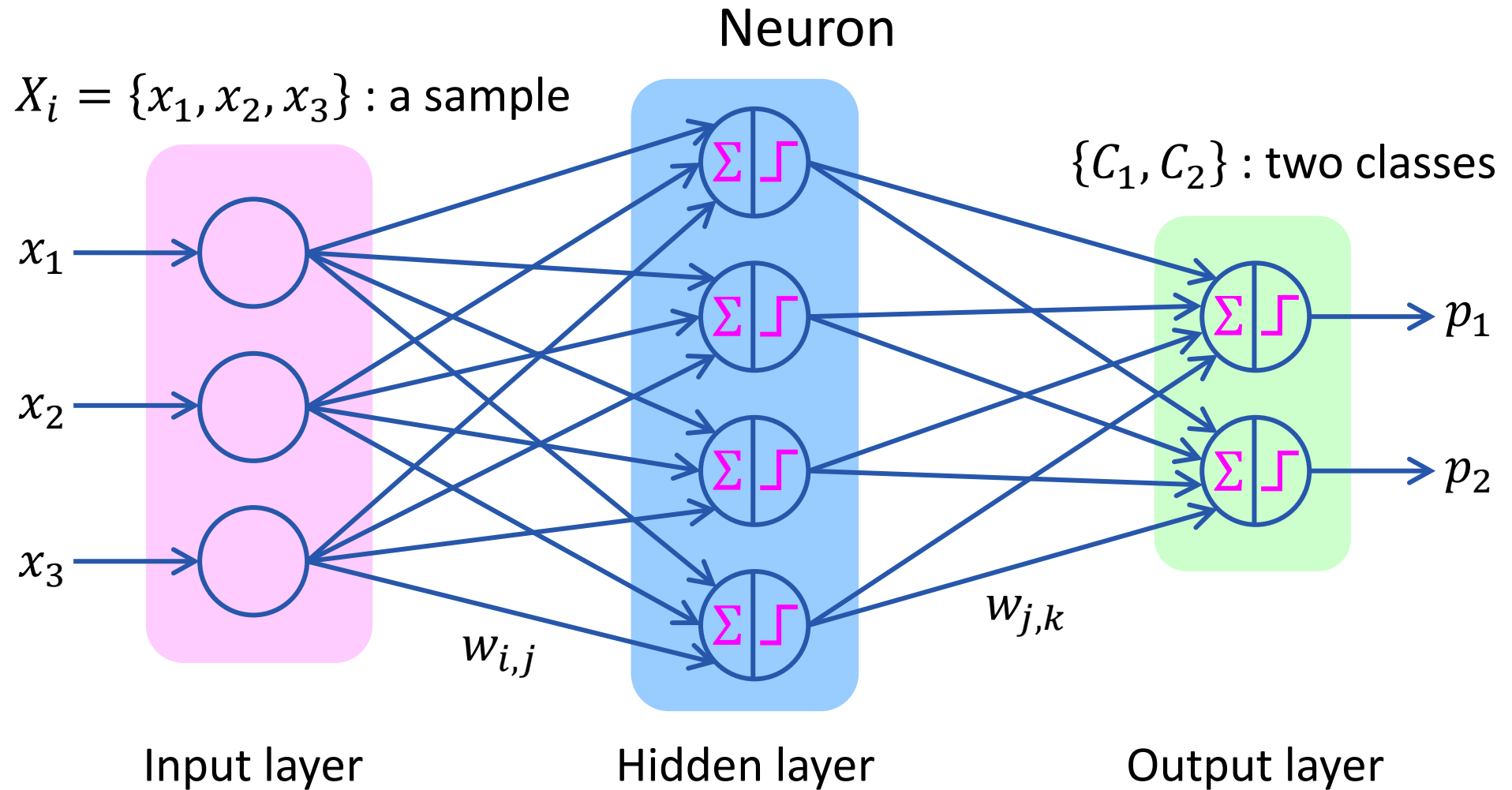


CSC 4740 / 6740 Data Mining

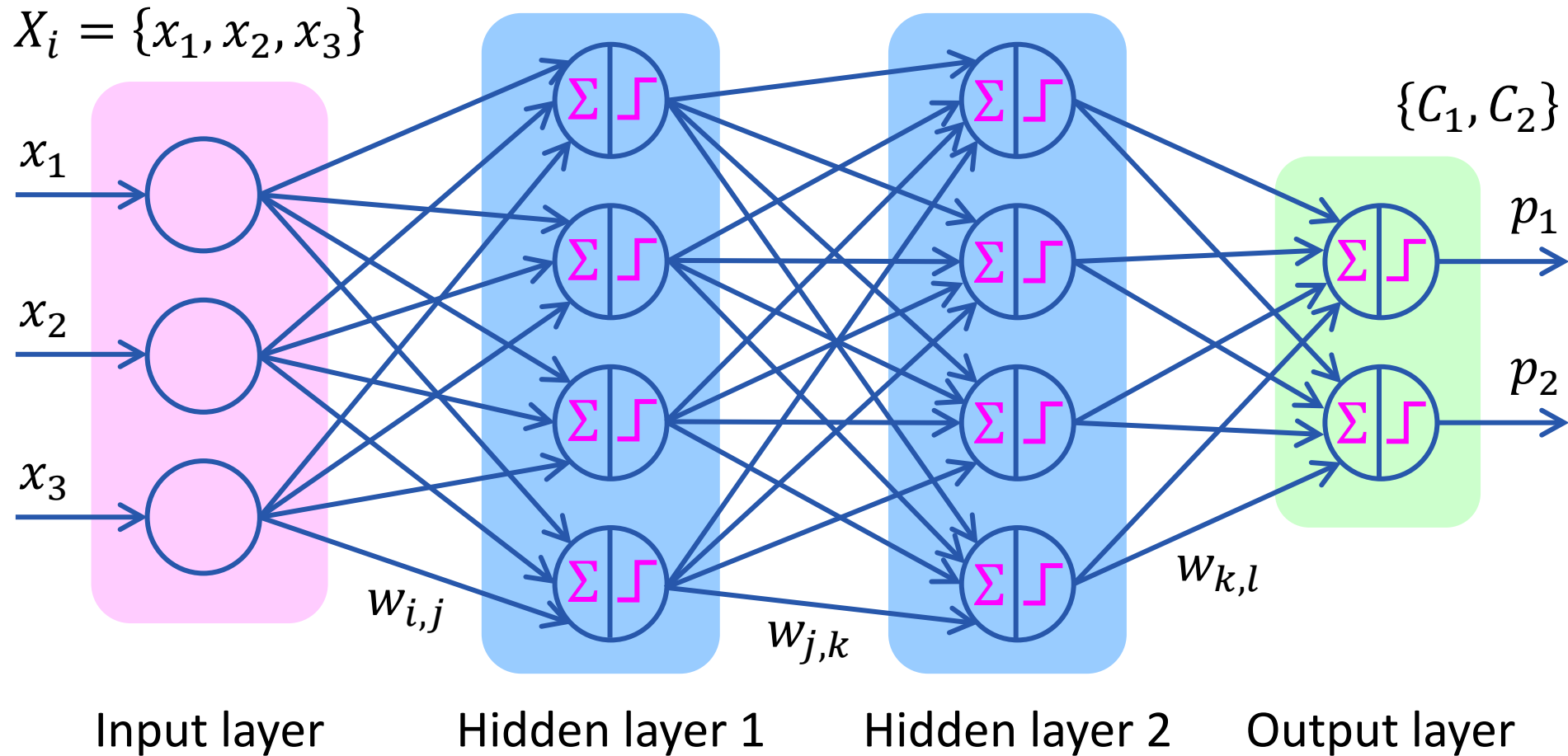
Spring 2024

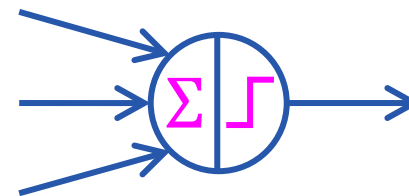
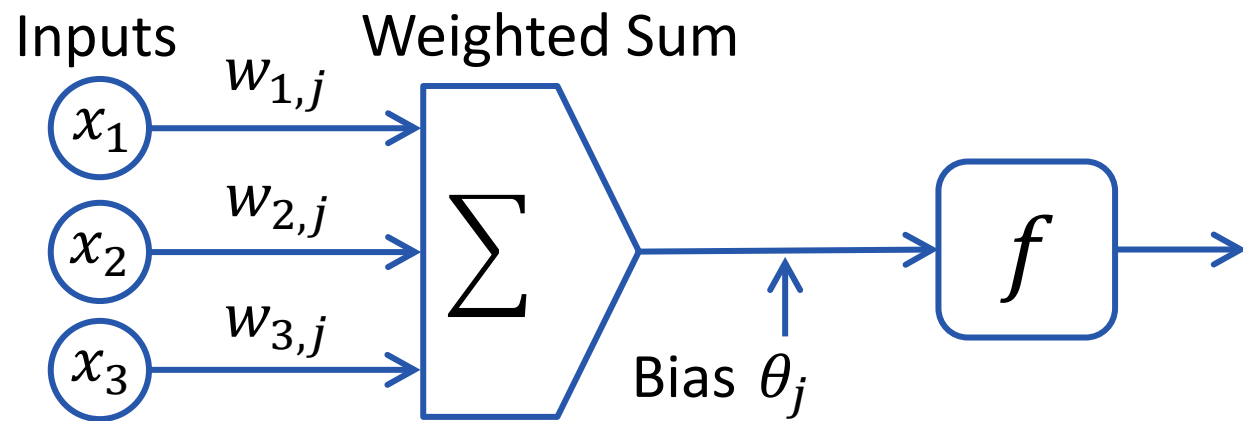
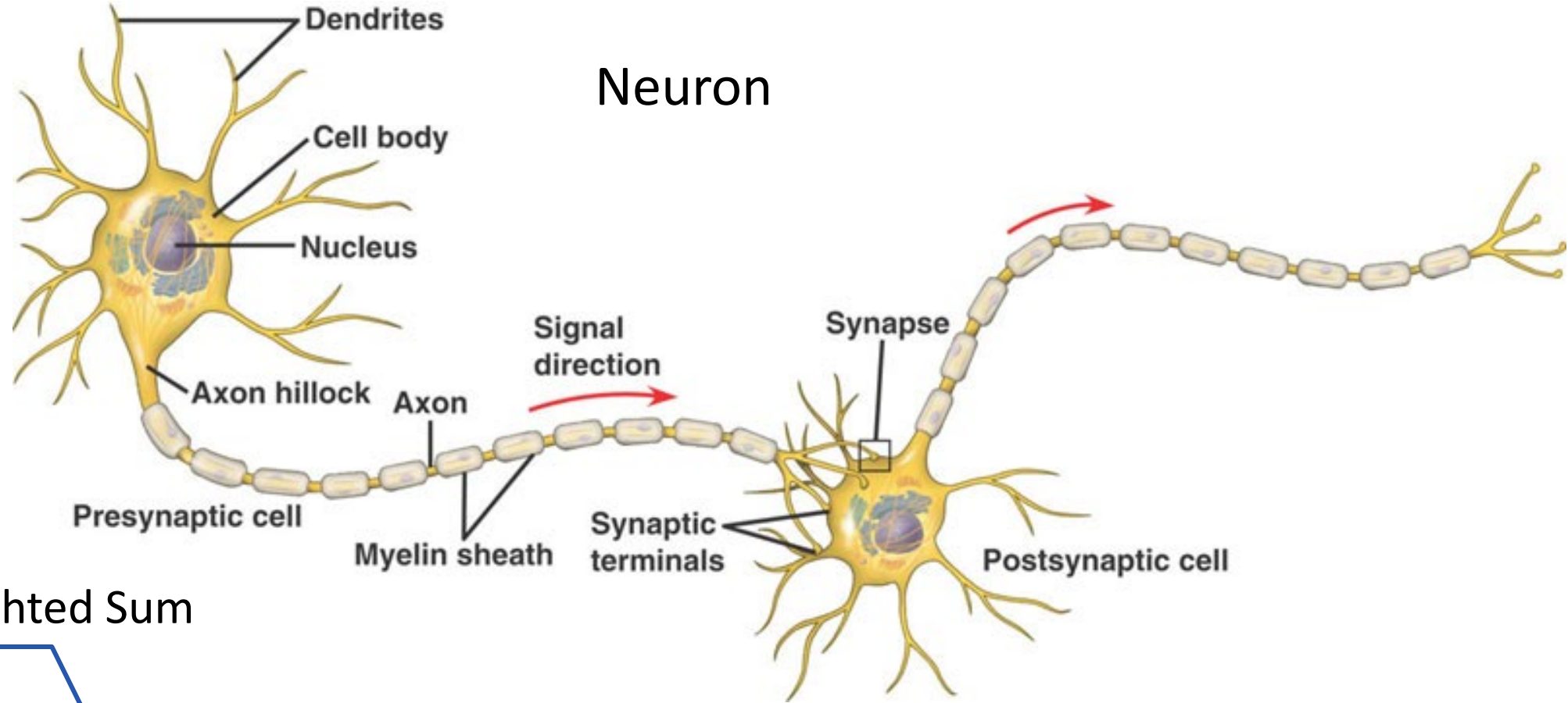
## Chapter 9 Classification: Advanced Methods

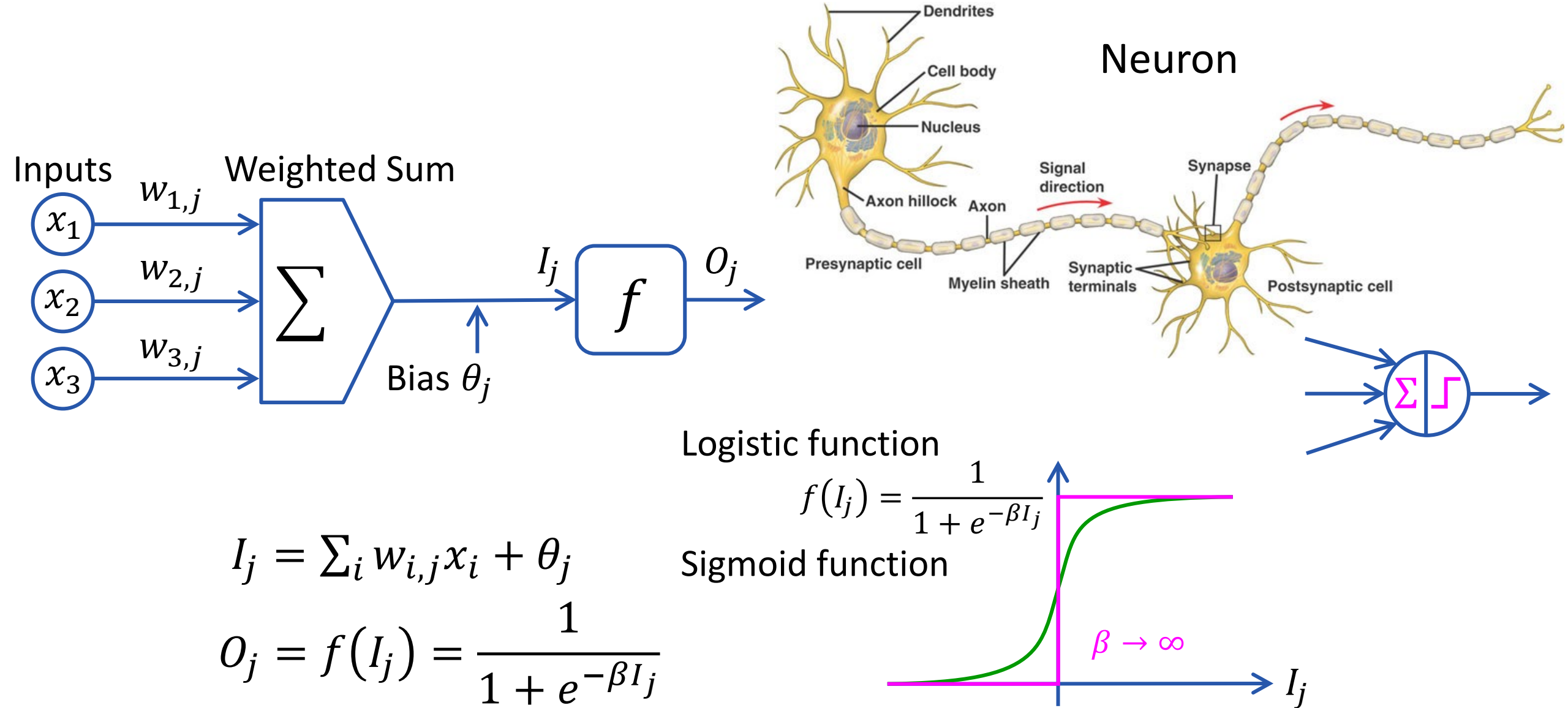
# Neural Network



# Neural Network







# Neural Network

## Topology:

- # of units in the input layer
- # of hidden layers
- # of units in each hidden layer
- # of units in the output layer

## Parameters:

- $w_{i,j}$  for each connection
- $\beta, \theta_j$  for each hidden or output unit

Neural networks can closely approximate any function

Given enough hidden units and enough training samples

# Neural Network as a Classifier

Weakness	Strength
Long learning time	High tolerance to noisy data
Some parameters are determined empirically, e.g., the network topology or “structure”	Well-suited for continuous-values inputs and outputs
Poor interpretability: Difficult to interpret the weights of the hidden units	Successful on an array of real-world data, e.g., hand-written letters
	Inherently parallel

Neural Network	SVM
Non-deterministic algorithm	Deterministic algorithm
Can easily be learned in incremental fashion	Hard to learn: learned in batch mode using quadratic programming techniques
To learn complex functions—use multilayer perceptron (nontrivial)	Using kernels can learn very complex functions



# How to Train the Neural Network?

Backpropagation Algorithm: A neural network learning algorithm

- For each training data object, the weights  $w_{i,j}$  are adjusted to minimize the mean squared error between prediction and ground truth
- Adjustments are made in the “**backwards**” direction: from the output layer , to the hidden layer

# How to Train the Neural Network?

Backpropagation Algorithm: A neural network learning algorithm

Steps:

1. Initialize weights  $w_{i,j}$  and biases  $\theta_j$  to small random numbers
2. For each training data object
  - 1) Propagate the inputs forward
  - 2) Backpropagate the error by updating weights  $w_{i,j}$  and biases  $\theta_j$
3. Terminating condition : when the error is very small, etc.

# Backpropagation Algorithm

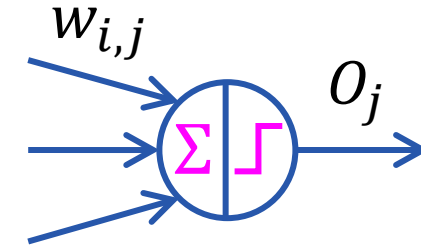
**The error:**

$$\text{Err}_j = O_j(1 - O_j)(T_j - O_j)$$

where,  $O_j$  is the actual output

$T_j$  is the known value

$O_j(1 - O_j)$  is the derivative of the logistic function



Unit  $j$  in the output layer

**Update the weights:**

$$\Delta w_{i,j} = l \cdot \text{Err}_j \cdot O_i$$

$$w_{i,j} = w_{i,j} + \Delta w_{i,j}$$

**Update the bias:**

$$\Delta \theta_j = l \cdot \text{Err}_j$$

$$\theta_j = \theta_j + \Delta \theta_j$$

$l \in (0,1)$  : learning rate

# Backpropagation Algorithm

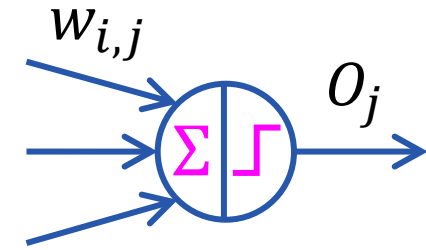
## The error:

$$\text{Err}_j = O_j(1 - O_j) \sum_k w_{j,k} \cdot \text{Err}_k$$

where,  $O_j$  is the actual output

$w_{j,k}$  is the weight of the connection from unit  $j$  to  $k$  in higher layer

$\text{Err}_k$  is the error of unit  $k$



Unit  $j$  in the hidden layer

## Update the weights:

$$\Delta w_{i,j} = l \cdot \text{Err}_j \cdot O_i$$

$$w_{i,j} = w_{i,j} + \Delta w_{i,j}$$

## Update the bias:

$$\Delta \theta_j = l \cdot \text{Err}_j$$

$$\theta_j = \theta_j + \Delta \theta_j$$

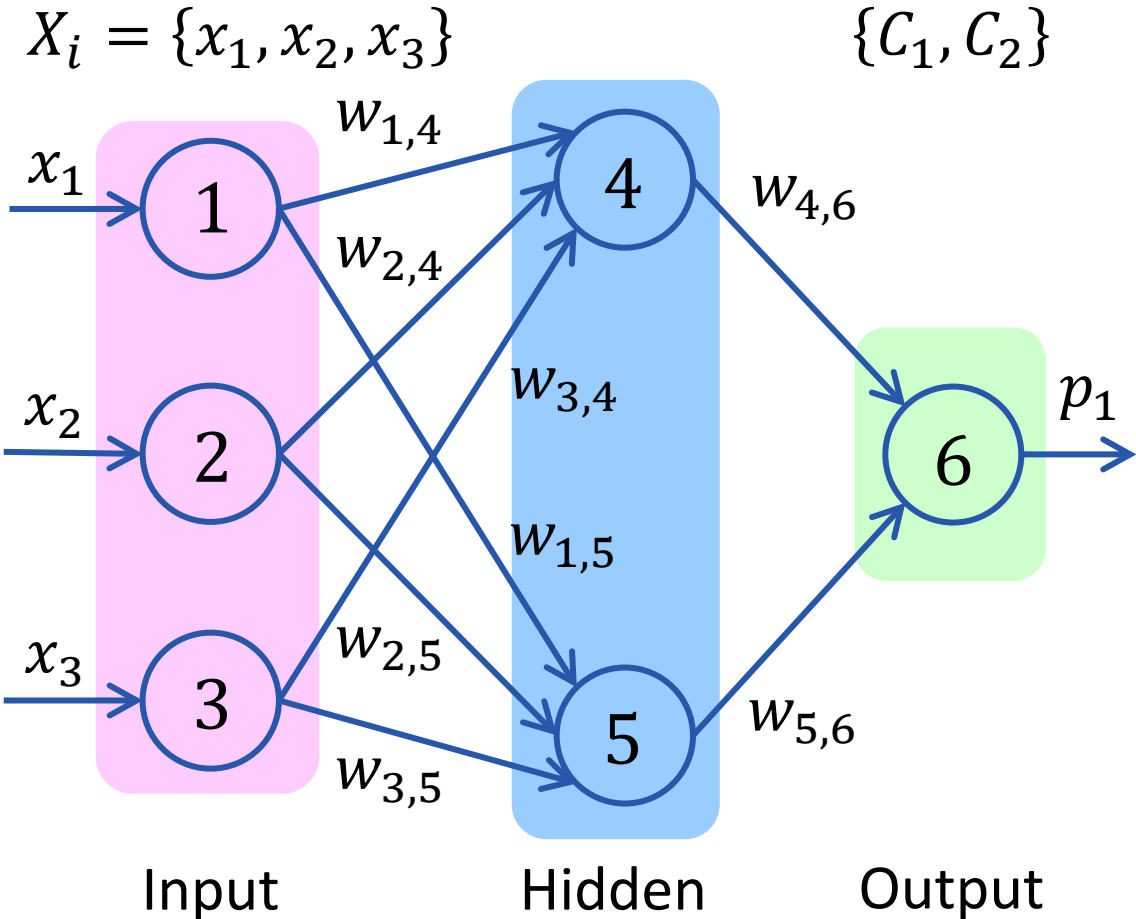
$l \in (0,1)$  : learning rate

Training data object:

$X_i = \{1, 0, 1\}$   
class label: 1

Initial values:

$w_{1,4}$	0.2
$w_{1,5}$	-0.3
$w_{2,4}$	0.4
$w_{2,5}$	0.1
$w_{3,4}$	-0.5
$w_{3,5}$	0.2
$w_{4,6}$	-0.3
$w_{5,6}$	-0.2
$\theta_4$	-0.4
$\theta_5$	0.2
$\theta_6$	0.1

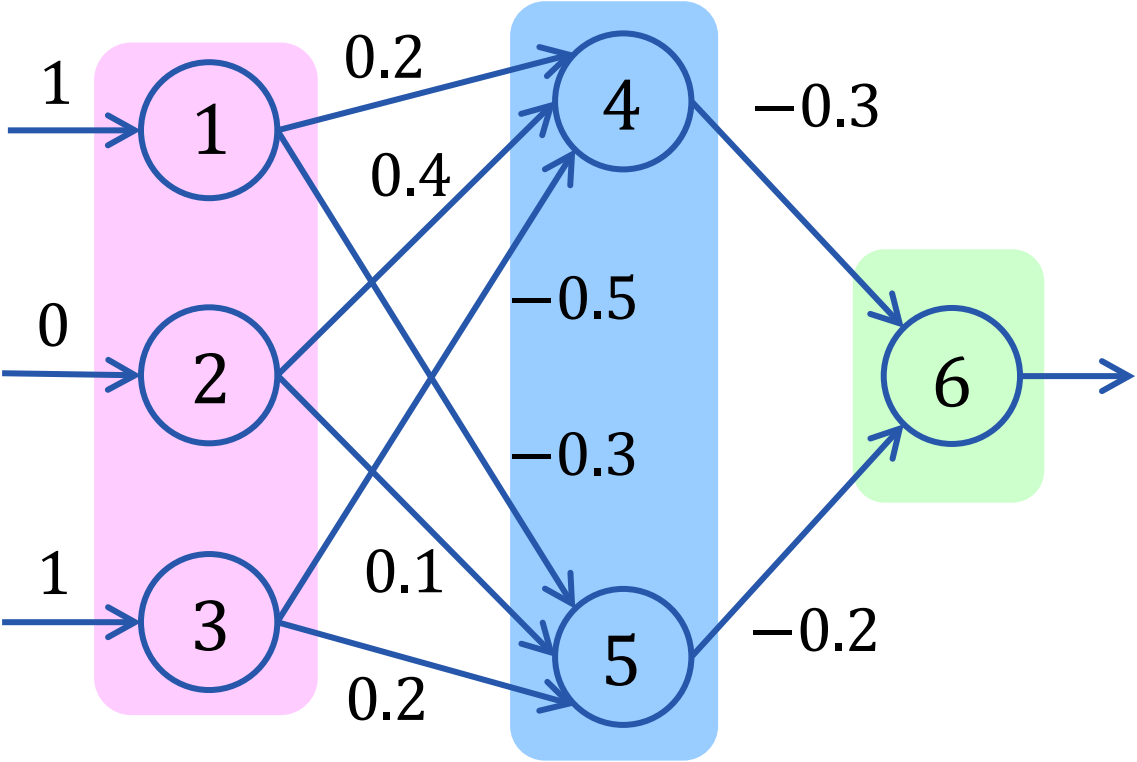


Training data object:

$X_i = \{1, 0, 1\}$   
class label: 1

Initial values:

$w_{1,4}$	0.2
$w_{1,5}$	-0.3
$w_{2,4}$	0.4
$w_{2,5}$	0.1
$w_{3,4}$	-0.5
$w_{3,5}$	0.2
$w_{4,6}$	-0.3
$w_{5,6}$	-0.2
$\theta_4$	-0.4
$\theta_5$	0.2
$\theta_6$	0.1



Unit $j$	Net Input $I_j$	Output $O_j$
4	$0.2 + 0 - 0.5 - 0.4 = -0.7$	$\frac{1}{1 + e^{0.7}} = 0.332$
5	$-0.3 + 0 + 0.2 + 0.2 = 0.1$	$\frac{1}{1 + e^{-0.1}} = 0.525$
6	$-0.3 \times 0.332 - 0.2 \times 0.525 + 0.1 = -0.105$	$\frac{1}{1 + e^{0.105}} = 0.474$

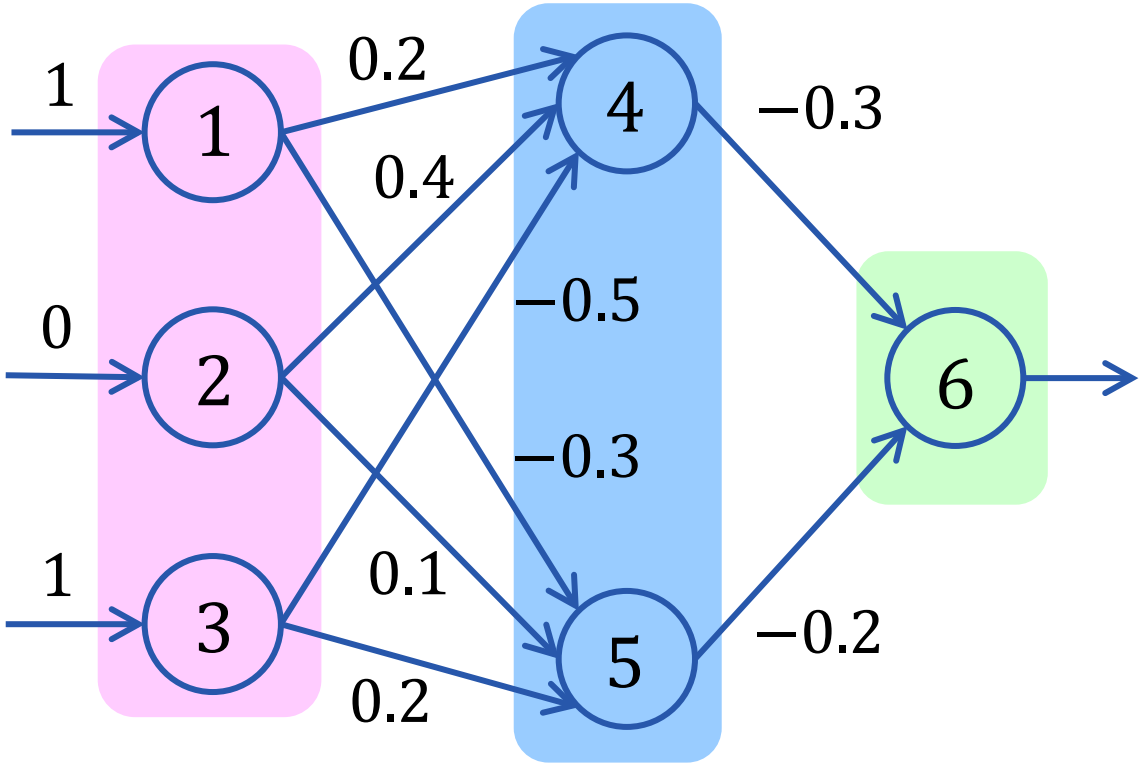
Propagate the inputs forward  $I_j = \sum_i w_{i,j}x_i + \theta_j$   $O_j = f(I_j) = \frac{1}{1 + e^{-\beta I_j}}$

Training data object:

$X_i = \{1, 0, 1\}$   
class label: 1

Initial values:

$w_{1,4}$	0.2
$w_{1,5}$	-0.3
$w_{2,4}$	0.4
$w_{2,5}$	0.1
$w_{3,4}$	-0.5
$w_{3,5}$	0.2
$w_{4,6}$	-0.3
$w_{5,6}$	-0.2
$\theta_4$	-0.4
$\theta_5$	0.2
$\theta_6$	0.1



Unit $j$	Error $Err_j$
6	$0.474 \times (1 - 0.474) \times (1 - 0.474) = 0.1311$
5	$0.525 \times (1 - 0.525) \times 0.1311 \times (-0.2) = -0.0065$
4	$0.332 \times (1 - 0.332) \times 0.1311 \times (-0.3) = -0.0087$

Backpropagate the errors

Output layer:

$$Err_j = O_j(1 - O_j)(T_j - O_j)$$

Hidden layer:

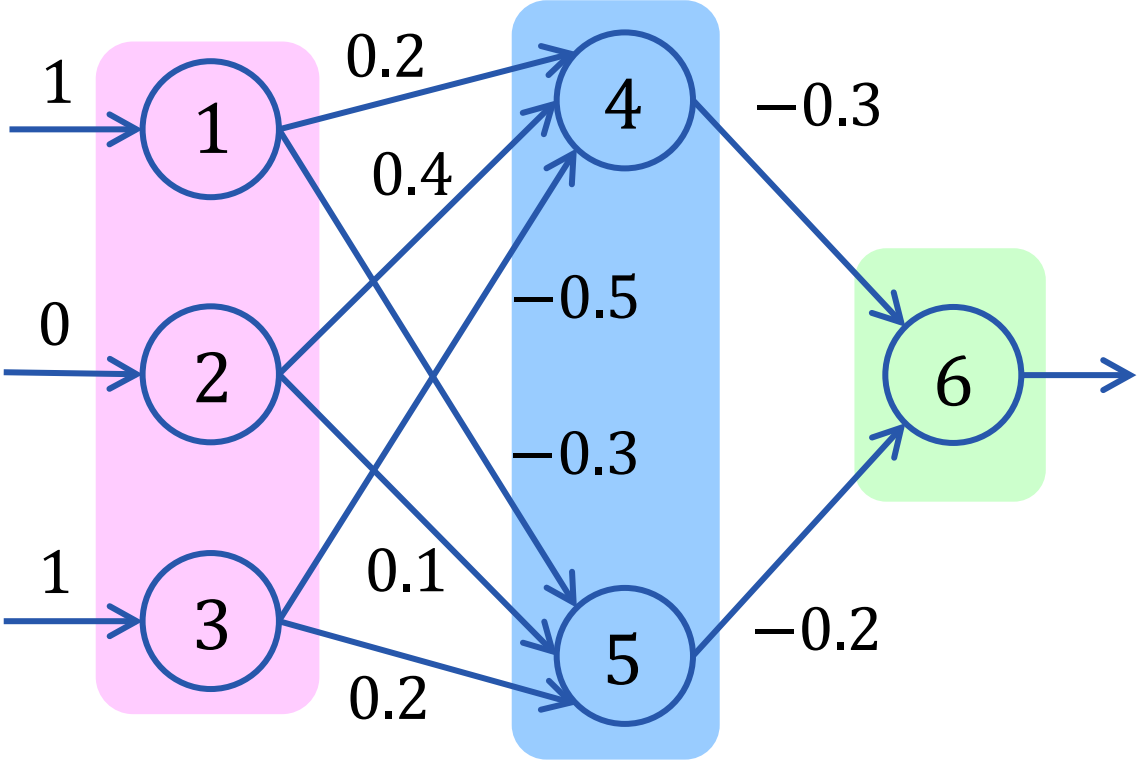
$$Err_j = O_j(1 - O_j) \sum_k w_{j,k} \cdot Err_k$$

Training data object:

$X_i = \{1, 0, 1\}$   
class label: 1

Initial values:

$w_{1,4}$	0.2
$w_{1,5}$	-0.3
$w_{2,4}$	0.4
$w_{2,5}$	0.1
$w_{3,4}$	-0.5
$w_{3,5}$	0.2
$w_{4,6}$	-0.3
$w_{5,6}$	-0.2
$\theta_4$	-0.4
$\theta_5$	0.2
$\theta_6$	0.1



weights	Error $Err_j$
$w_{4,6}$	$-0.3 + 0.9 \times 0.1311 \times 0.332 = -0.261$
$w_{5,6}$	$-0.2 + 0.9 \times 0.1311 \times 0.525 = -0.138$

Ge

Backpropagate the errors  
Learning rate  $l = 0.9$

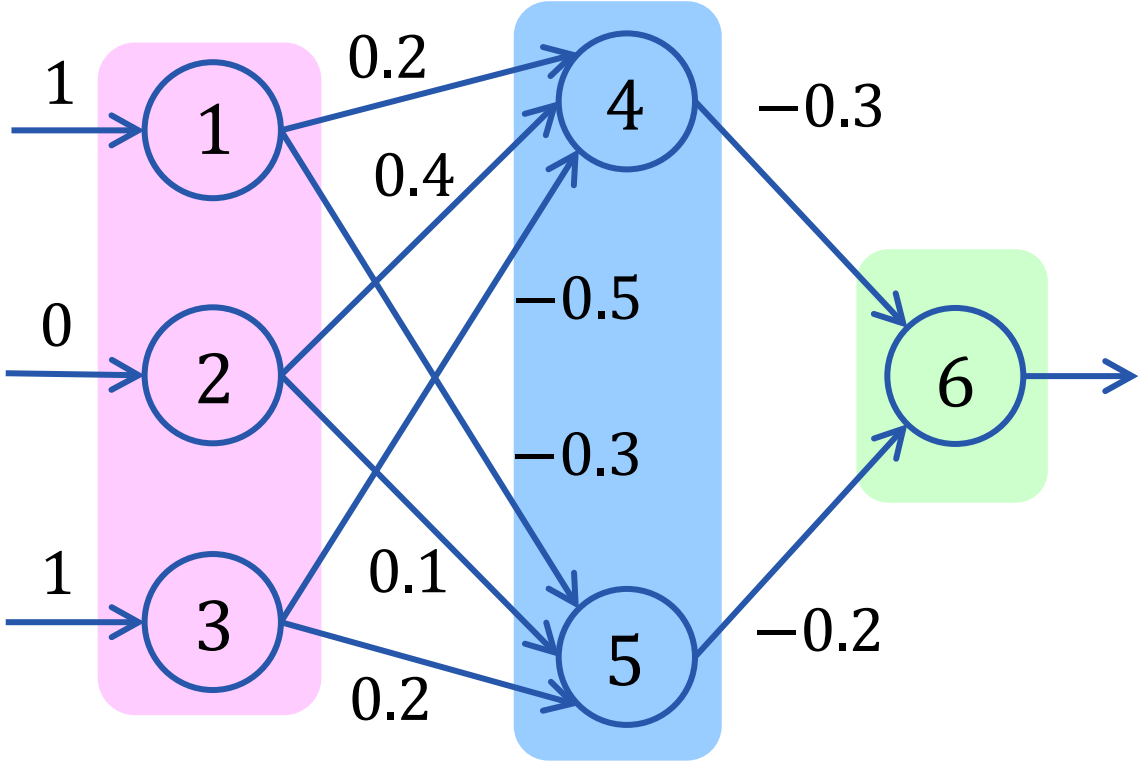
Update the weights:  $\Delta w_{i,j} = l \cdot Err_j \cdot O_i$   
 $w_{i,j} = w_{i,j} + \Delta w_{i,j}$



Training data object:  
 $X_i = \{1, 0, 1\}$   
class label: 1

Initial values:

$w_{1,4}$	0.2
$w_{1,5}$	-0.3
$w_{2,4}$	0.4
$w_{2,5}$	0.1
$w_{3,4}$	-0.5
$w_{3,5}$	0.2
$w_{4,6}$	-0.3
$w_{5,6}$	-0.2
$\theta_4$	-0.4
$\theta_5$	0.2
$\theta_6$	0.1



weights	Error Err <sub>j</sub>
$w_{1,4}$	$0.2 + 0.9 \times (-0.0087) \times 1 = 0.192$
$w_{2,4}$	$0.4 + 0.9 \times (-0.0087) \times 0 = 0.4$
$w_{3,4}$	$-0.5 + 0.9 \times (-0.0087) \times 1 = -0.508$

Backpropagate the errors  
Learning rate  $l = 0.9$

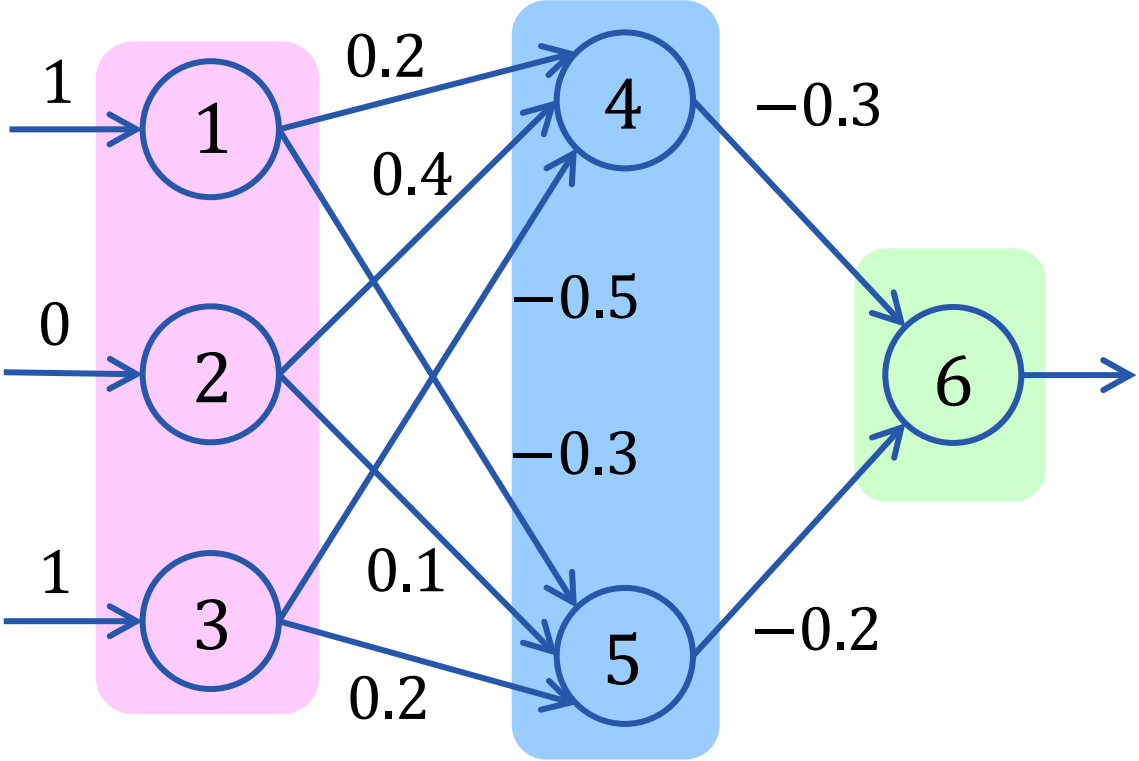
Update the weights:  $\Delta w_{i,j} = l \cdot \text{Err}_j \cdot O_i$   
 $w_{i,j} = w_{i,j} + \Delta w_{i,j}$

Training data object:

$X_i = \{1, 0, 1\}$   
class label: 1

Initial values:

$w_{1,4}$	0.2
$w_{1,5}$	-0.3
$w_{2,4}$	0.4
$w_{2,5}$	0.1
$w_{3,4}$	-0.5
$w_{3,5}$	0.2
$w_{4,6}$	-0.3
$w_{5,6}$	-0.2
$\theta_4$	-0.4
$\theta_5$	0.2
$\theta_6$	0.1



weights	Error $Err_j$
$w_{1,5}$	$-0.3 + 0.9 \times (-0.0065) \times 1 = -0.306$
$w_{2,5}$	$0.1 + 0.9 \times (-0.0065) \times 0 = 0.1$
$w_{3,5}$	$0.2 + 0.9 \times (-0.0065) \times 1 = 0.194$

Ge

Backpropagate the errors  
Learning rate  $l = 0.9$

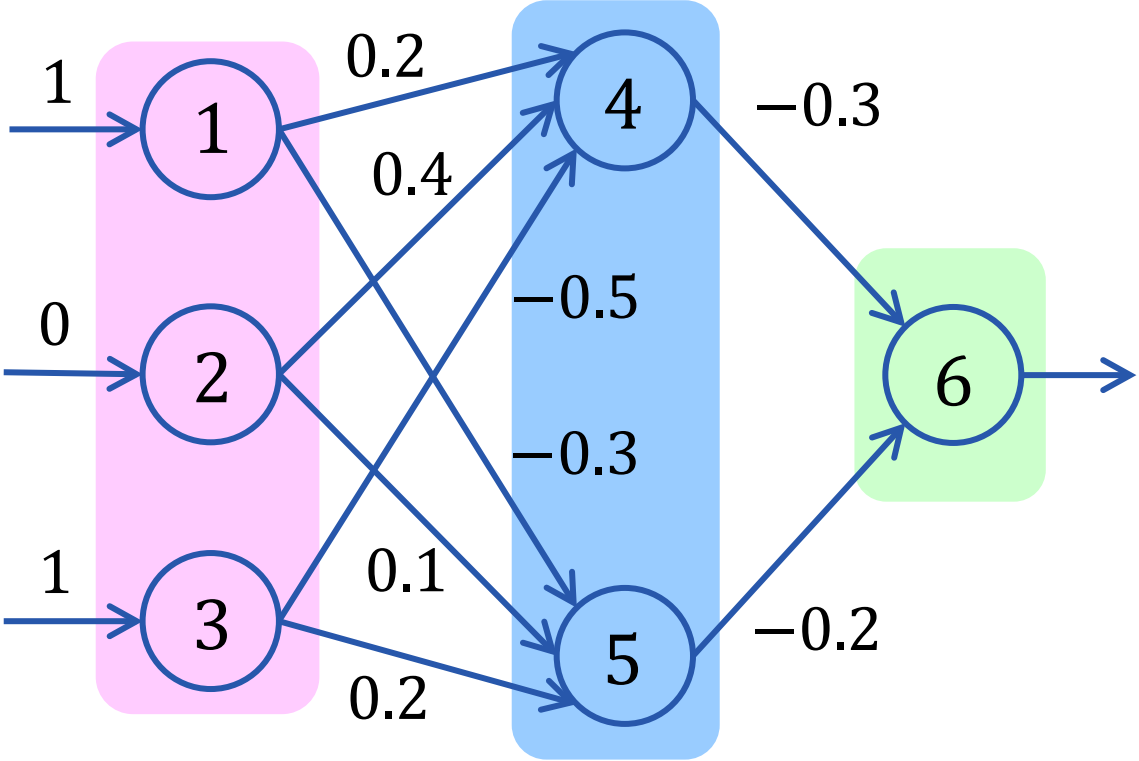
Update the weights:

$$\Delta w_{i,j} = l \cdot Err_j \cdot O_i$$
$$w_{i,j} = w_{i,j} + \Delta w_{i,j}$$

Training data object:  
 $X_i = \{1, 0, 1\}$   
class label: 1

Initial values:

$w_{1,4}$	0.2
$w_{1,5}$	-0.3
$w_{2,4}$	0.4
$w_{2,5}$	0.1
$w_{3,4}$	-0.5
$w_{3,5}$	0.2
$w_{4,6}$	-0.3
$w_{5,6}$	-0.2
$\theta_4$	-0.4
$\theta_5$	0.2
$\theta_6$	0.1



weights	Error $Err_j$
$\theta_6$	$0.1 + 0.9 \times 0.1311 = 0.218$
$\theta_5$	$0.2 + 0.9 \times (-0.0065) = 0.194$
$\theta_4$	$-0.4 + 0.9 \times (-0.0087) = -0.408$

Backpropagate the errors

Learning rate  $l = 0.9$

Update the biases:

$$\Delta\theta_j = l \cdot Err_j$$

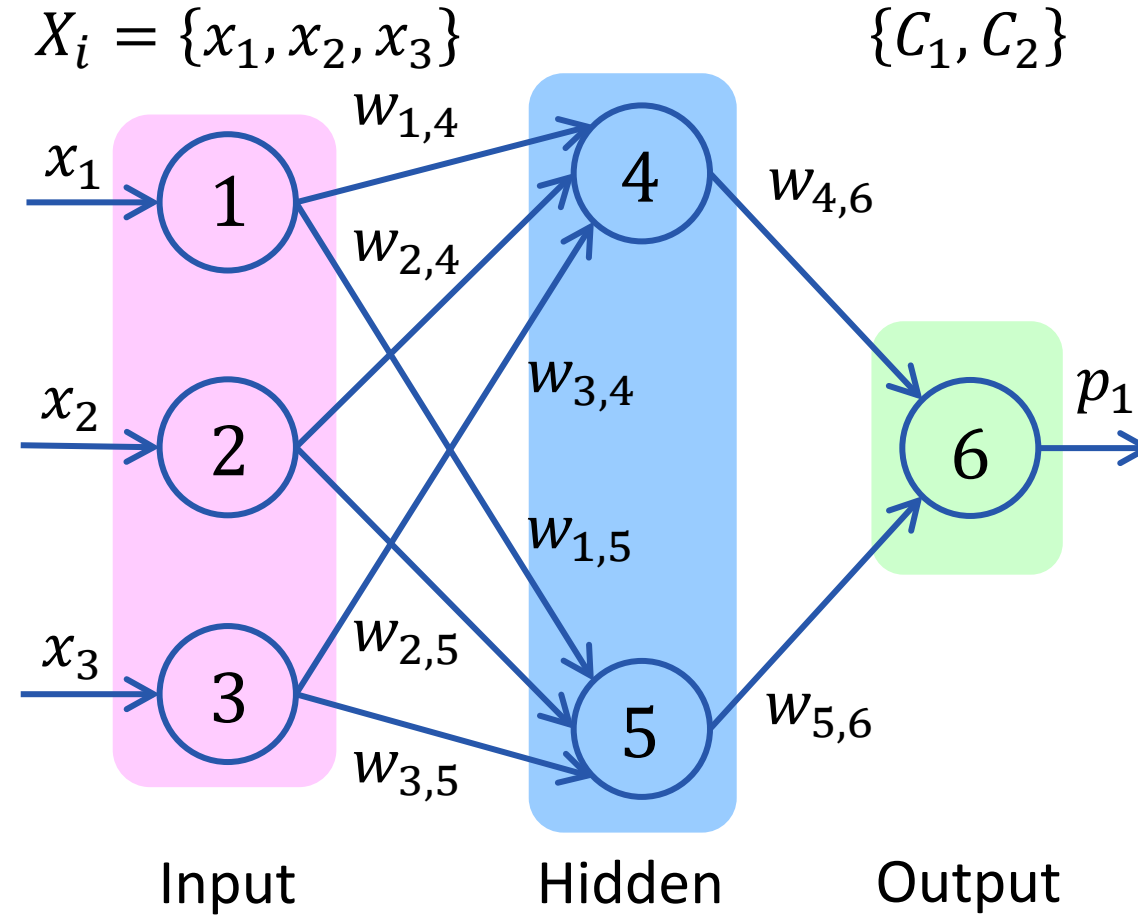
$$\theta_j = \theta_j + \Delta\theta_j$$

## Training data object:

$X_i = \{1, 0, 1\}$  label: 1

## Initial values:

weights	old	new
$w_{1,4}$	0.2	0.192
$w_{1,5}$	-0.3	-0.306
$w_{2,4}$	0.4	0.4
$w_{2,5}$	0.1	0.1
$w_{3,4}$	-0.5	-0.508
$w_{3,5}$	0.2	0.194
$w_{4,6}$	-0.3	-0.261
$w_{5,6}$	-0.2	-0.138
$\theta_4$	-0.4	-0.408
$\theta_5$	0.2	0.194
$\theta_6$	0.1	0.218



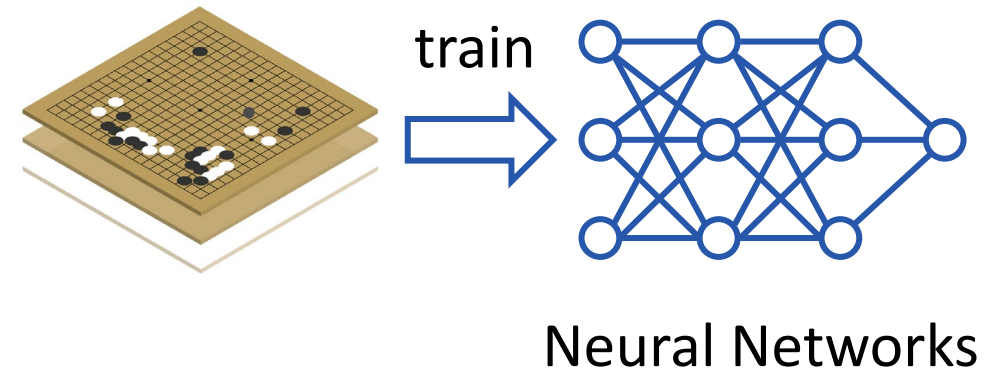
Multilayer feed-forward networks

Learning rate  $l = 0.9$

# Lee Sedol 9-dan vs AlphaGo



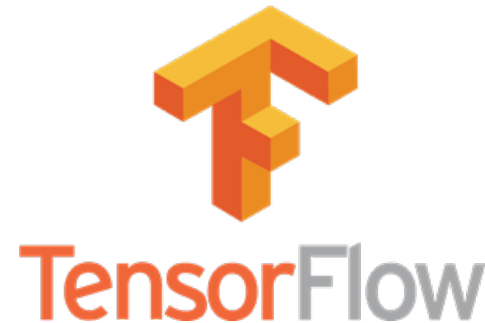
AlphaGo 4 – Lee Sedol 1



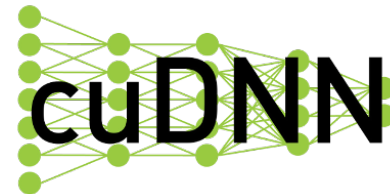
- Human expert positions
- Self-play positions

# Deep Learning Package Zoo

- Torch => PyTorch
- Caffe
- Theano (Keras, Lasagne)
- nVidia CuDNN
- Google TensorFlow
- Mxnet
- etc.



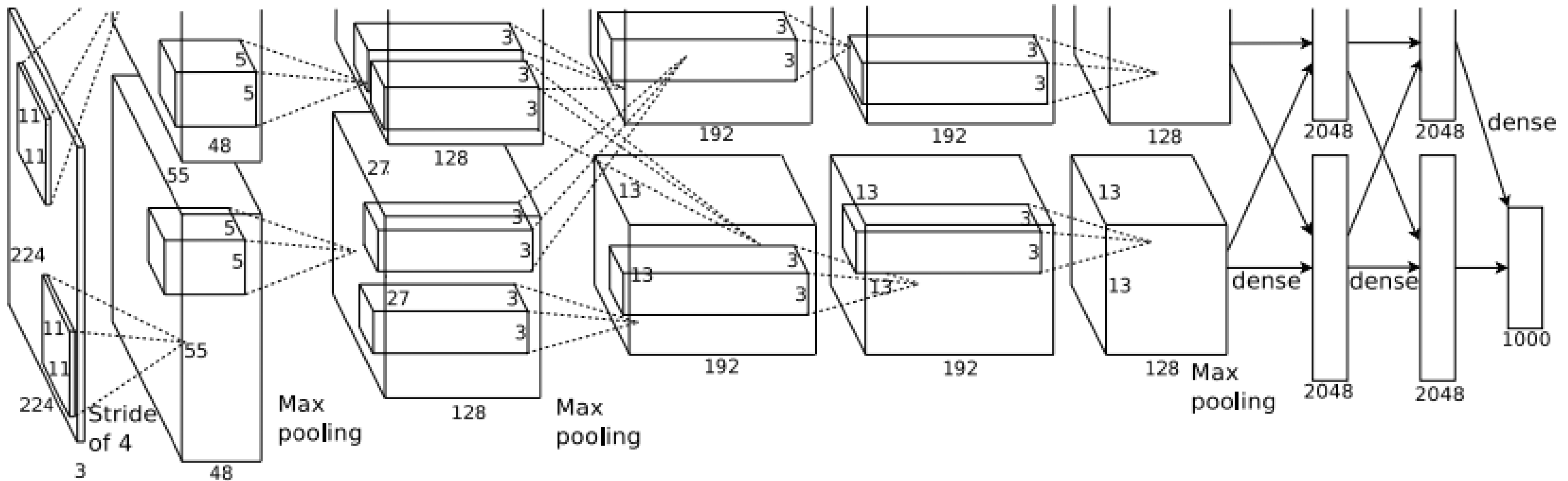
Caffe



# Pre-Trained Models

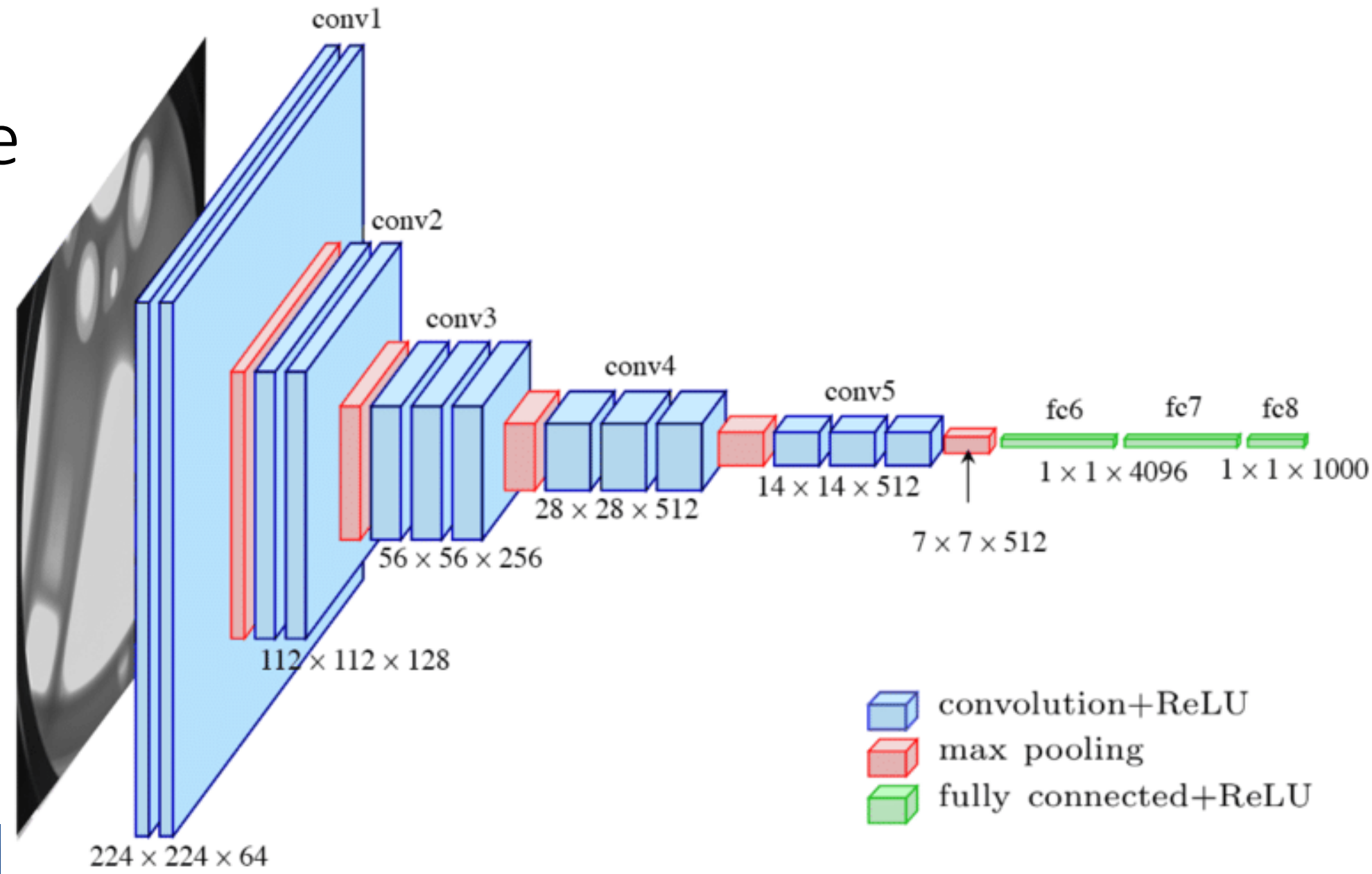
- AlexNet
- VGG-16 (Visual Geometry Group - University of Oxford)
- BERT (Bidirectional Encoder Representations from Transformers)
- GPT (Generative Pre-trained Transformer) / GPT2 / GPT3

# Architecture of AlexNet CNN





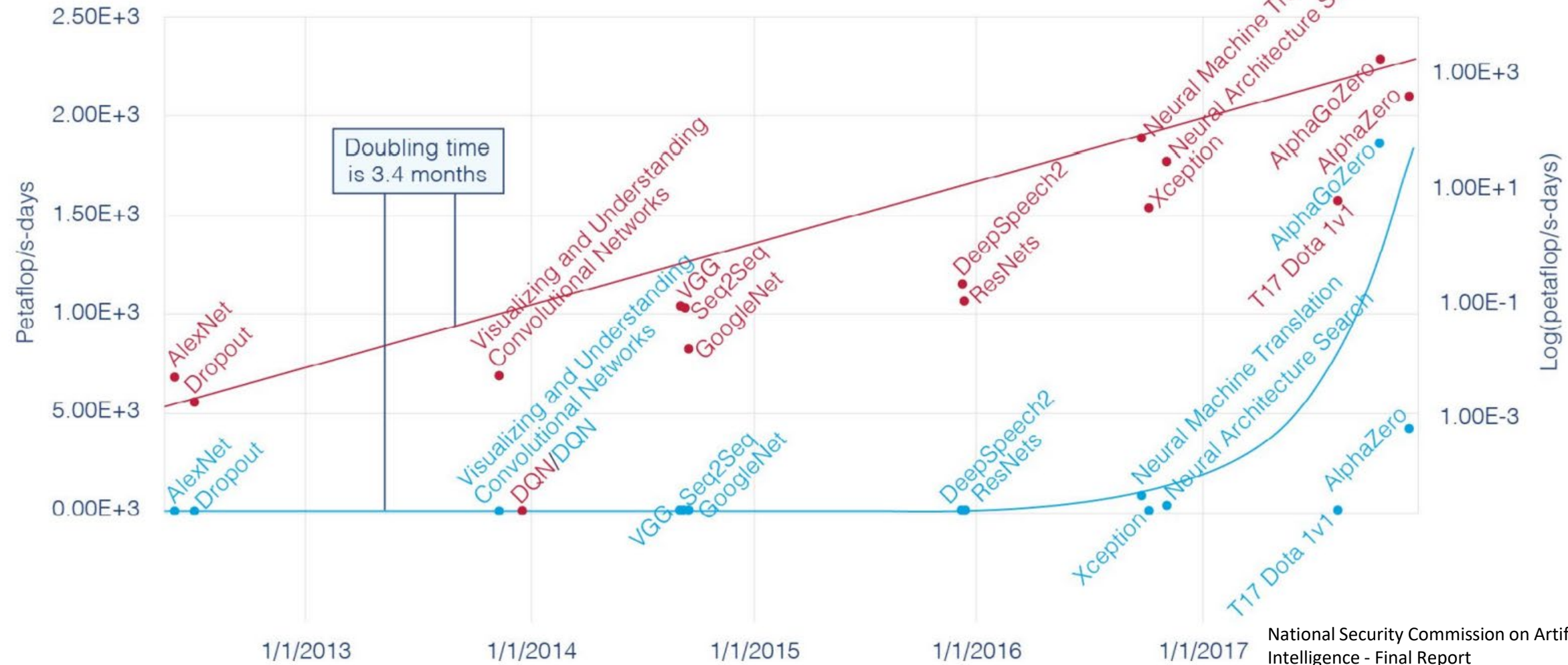
# VGG-16 Architecture



# Compute Required to Train Largest Deep Learning Models (2012-2017)

Compute Required  
to Train Largest  
Deep Learning  
Models (2012-2017)

- Linear Scale
- Log Scale
- Exponential growth on linear scale
- Exponential growth on logarithmic scale



## AlexNet to AlphaGo Zero: A 300,000x Increase in Compute

Log Scale

Linear Scale

Petaflop/s-days

1e+4

1e+3

1e+2

1e+1

1e+0

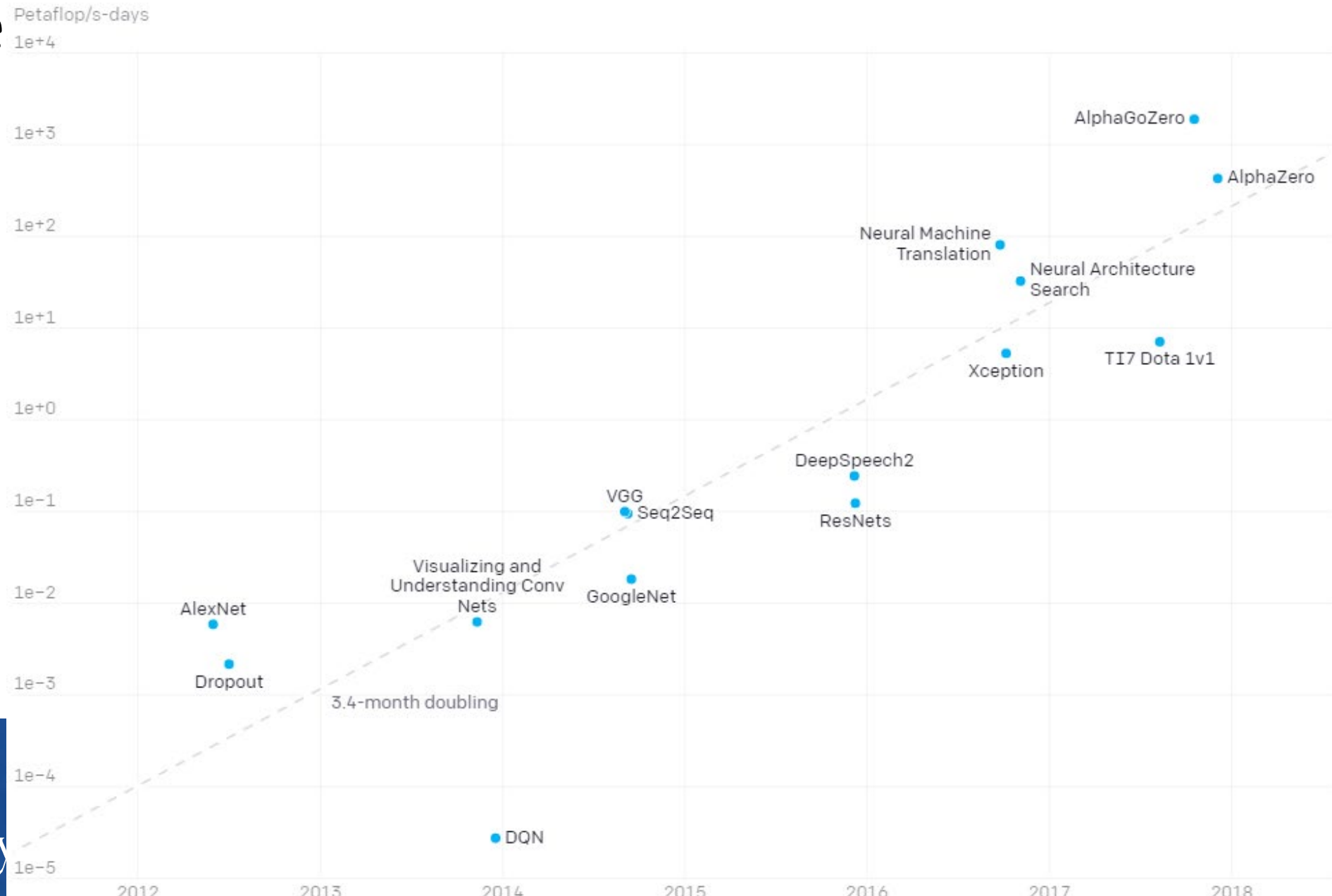
1e-1

1e-2

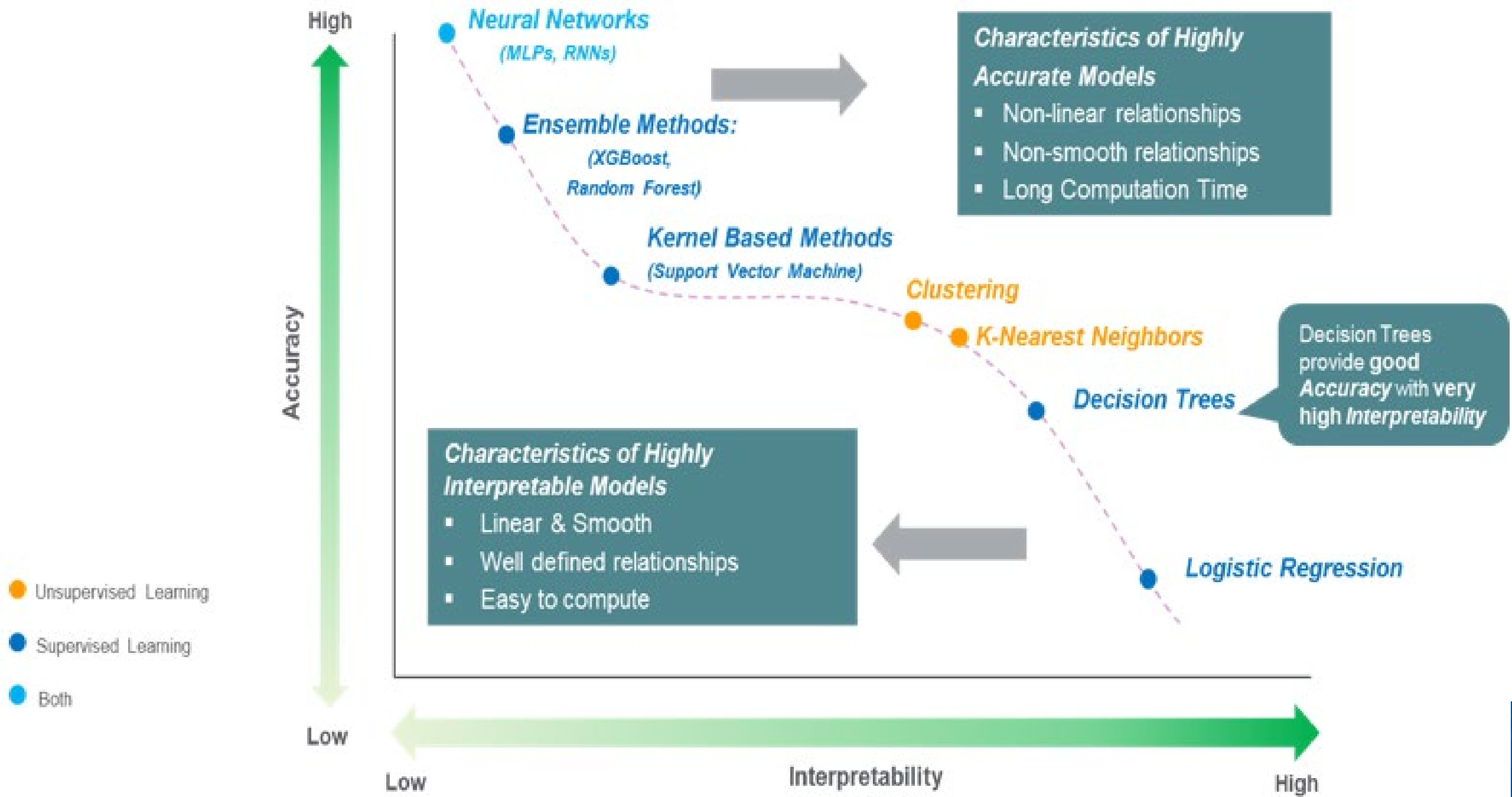
1e-3

1e-4

1e-5



AI and Compute - OpenAI  
<https://openai.com/blog/ai-and-compute/>



# $k$ -Nearest Neighbor ( $k$ -NN) Classifier

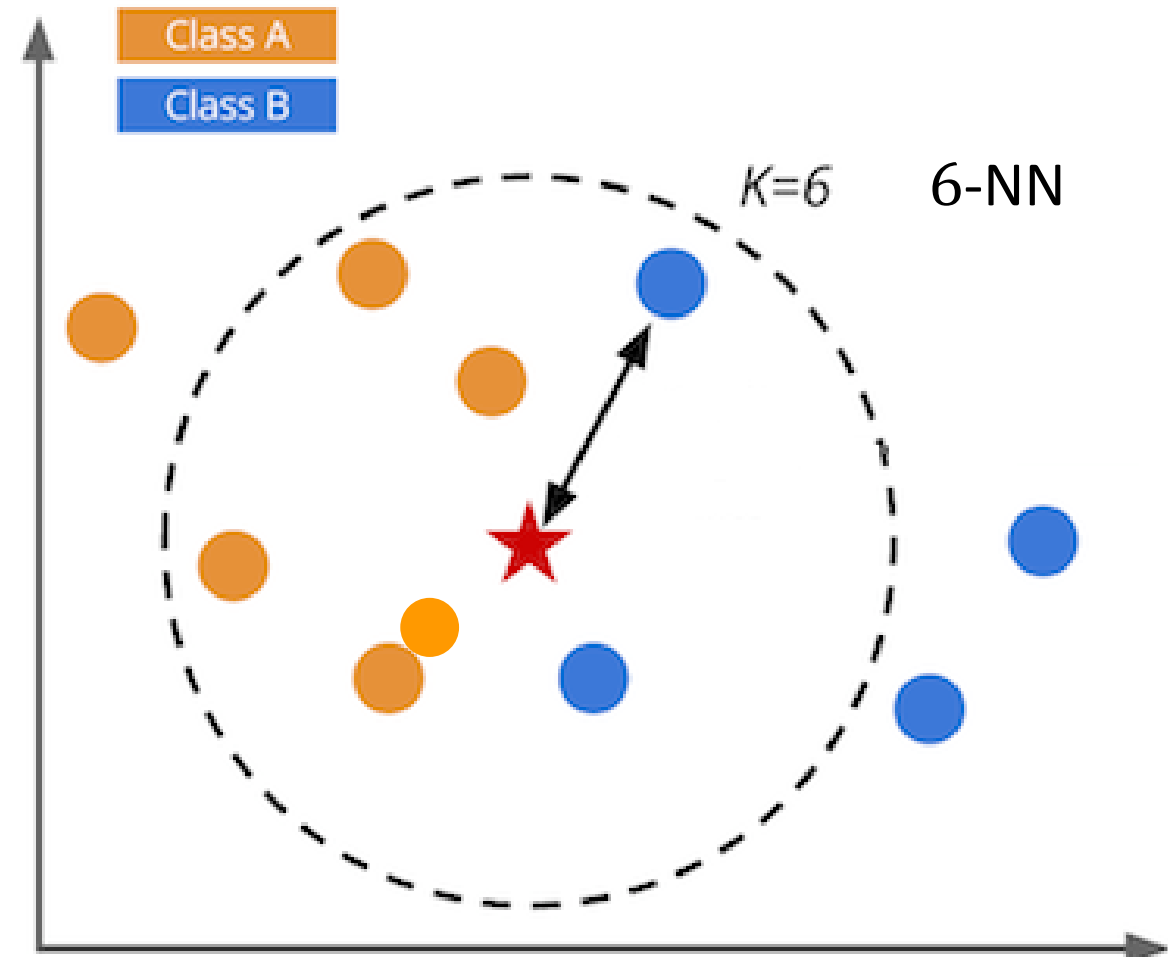
1. Compute the Distance/Similarity

$$d_{i,j} = \sqrt{(x_{i,1} - x_{j,1})^2 + (x_{i,2} - x_{j,2})^2}$$

2. Rank the remaining nodes

3. Select the Top- $k$  nodes

4. Majority Vote

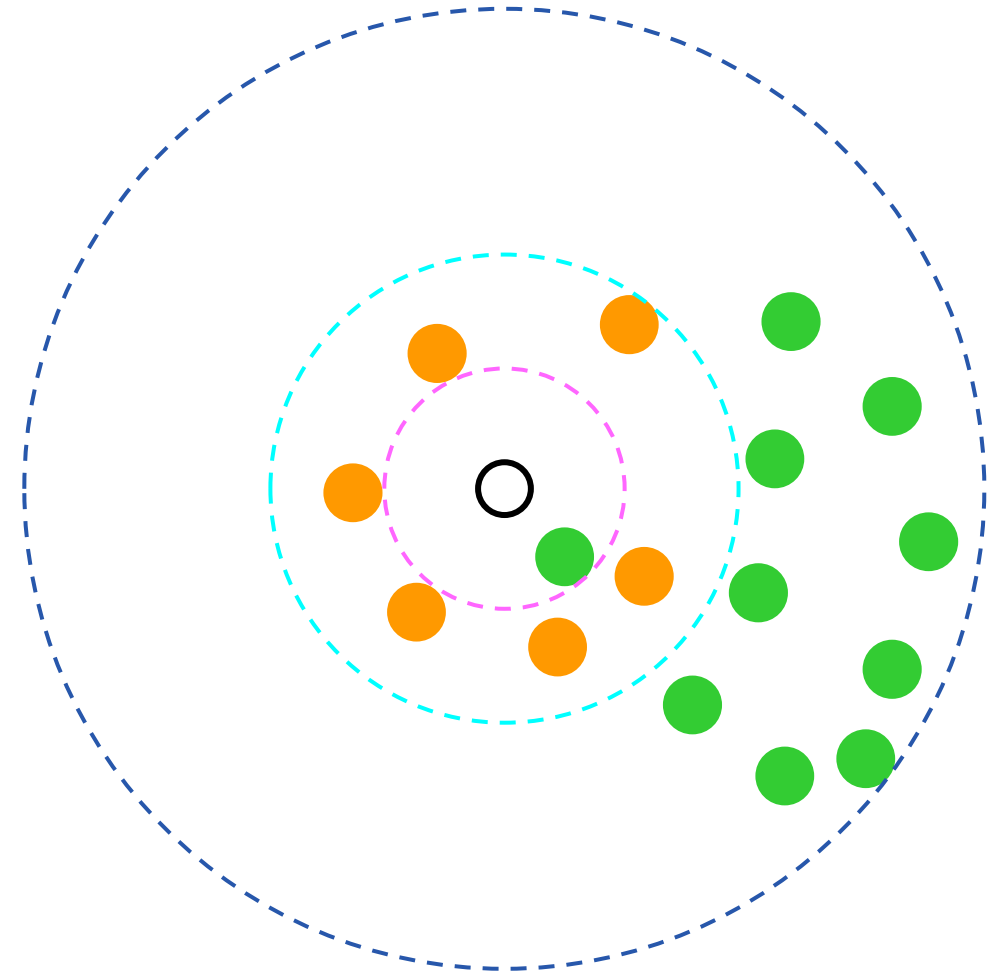
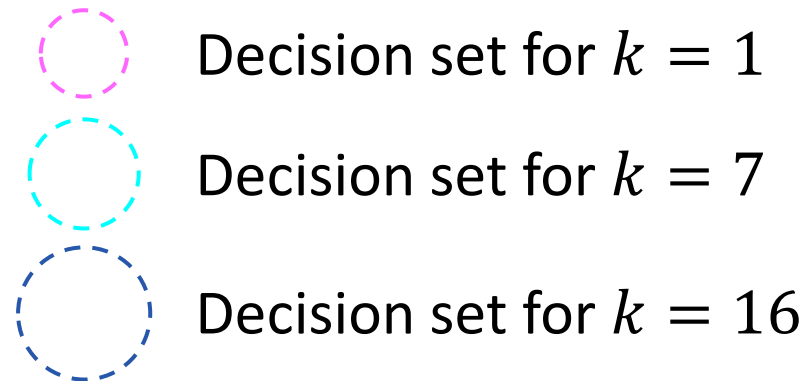


# How to choose parameter $k$ ?

“too small”  $k$ : sensitive to outliers

“too large”  $k$ : many objects from other classes

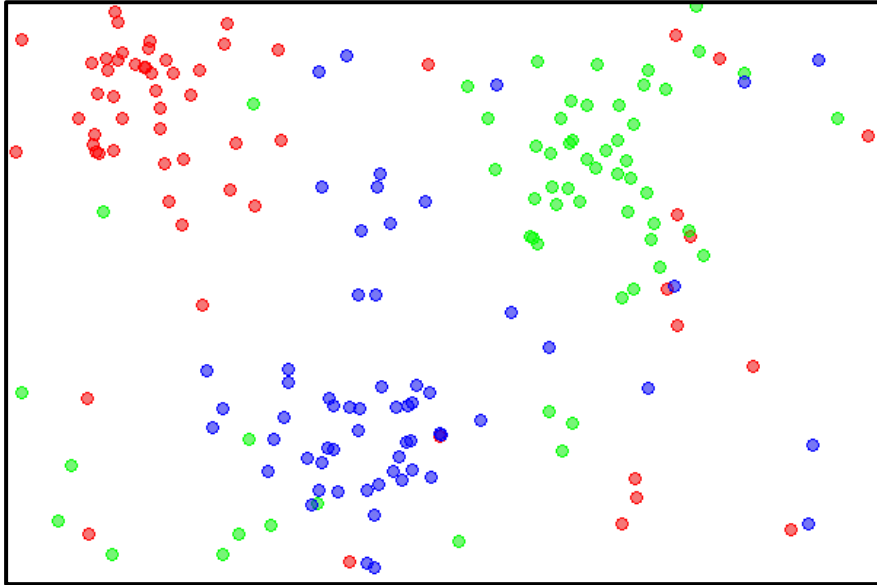
medium  $k$ : highest accuracy



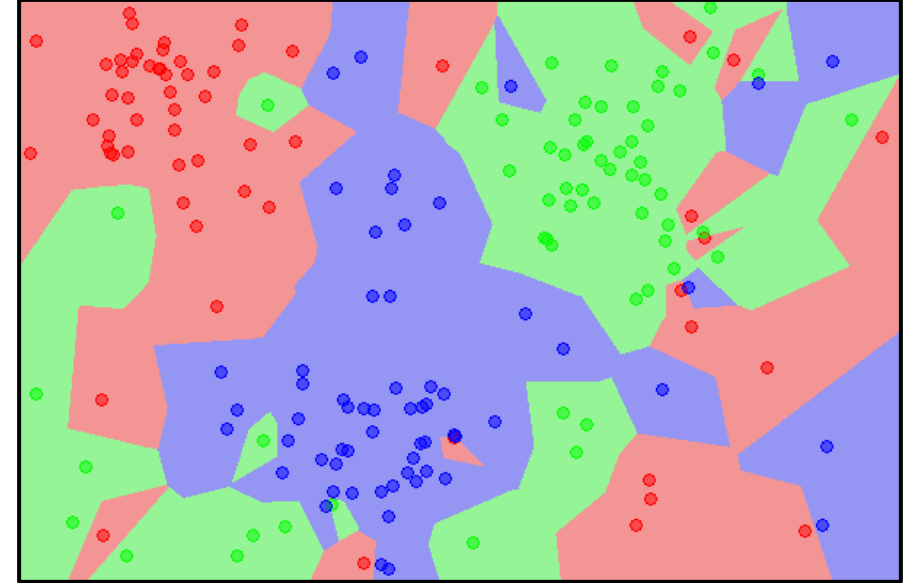


# Difference Between 1-NN and 5-NN

Original dataset



1-NN



5-NN

