# L9: Neural Network

Shan Wang

Lingnan College, Sun Yat-sen University
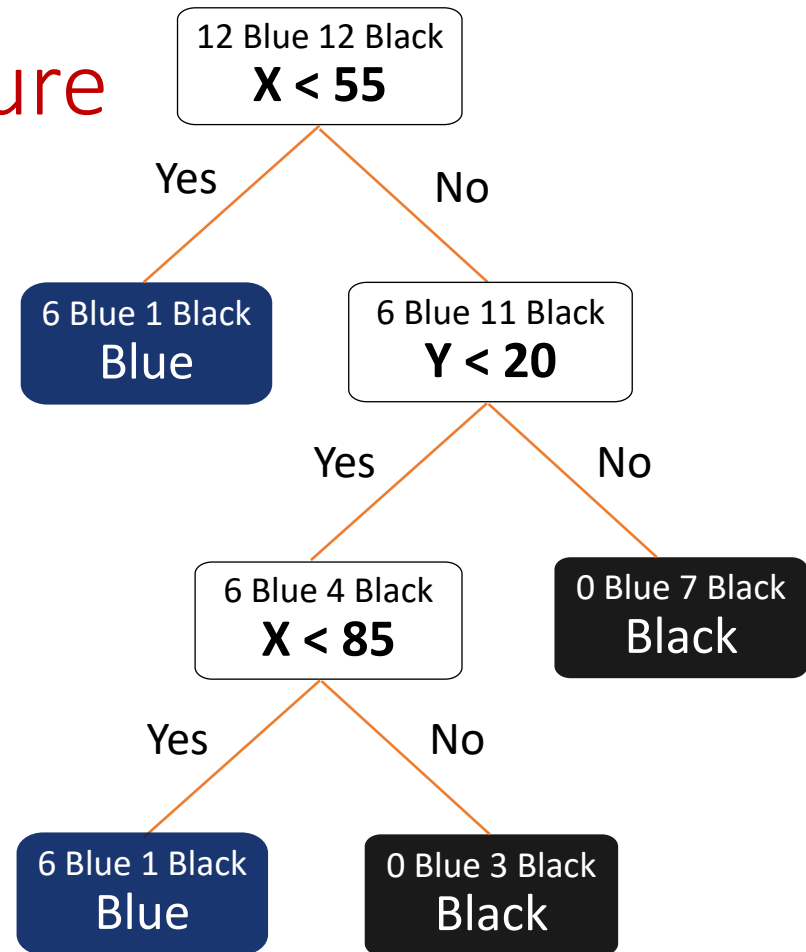
# Last lecture

- Decision Tree Model
- Strategy & Algorithm
  - ID.3: information gain
  - C4.5: information gain ratio
  - CART
    - Classification: Gini index
    - Regression: Squared error
- Regularization: pruning
- Random forest: bagging
- Application: Supreme Court Decisions

```
                    12 Blue 12 Black
                        X < 55
            Yes                    No
    6 Blue 1 Black          6 Blue 11 Black
        Blue                    Y < 20
                        Yes              No
                6 Blue 4 Black       0 Blue 7 Black
                    X < 85               Black
            Yes              No
    6 Blue 1 Black      0 Blue 3 Black
        Blue               Black
```

Information gain: with the $X$, how much the uncertainty of $Y$ decreases

# Classification Problem Revisit

$$y = \begin{cases} +1, & \text{if } f_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0 \\ -1, & \text{if } f_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0 \end{cases}$$

- Logistic regression
  - $f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \dfrac{e^{\boldsymbol{\theta}'x}}{1+e^{\boldsymbol{\theta}'x}} - h$
- SVM
  - $f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{\theta}'\phi(\boldsymbol{x}) + \theta_0$
- Decision tree
  - $f_{\boldsymbol{\theta}}(\boldsymbol{x})$: a list of rules
- What is the $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ in our brain?

# Neurons in the Brain

- The brain is composed of neurons
  - A neuron receives inputs from other neurons via synapses (突触)
  - Inputs are approximately summed
  - When the input exceeds a threshold the neuron sends an electrical spike that from the body, down the axon(轴突), to the next neuron(s) as output

# Course outline

- Supervised learning
  - Linear regression
  - Logistic regression
  - SVM and kernel
  - Tree models

- Deep learning
  - Neural networks
  - Convolutional NN
  - Recurrent NN

- Unsupervised learning
  - Clustering
  - PCA
  - EM

- Reinforcement learning
  - MDP
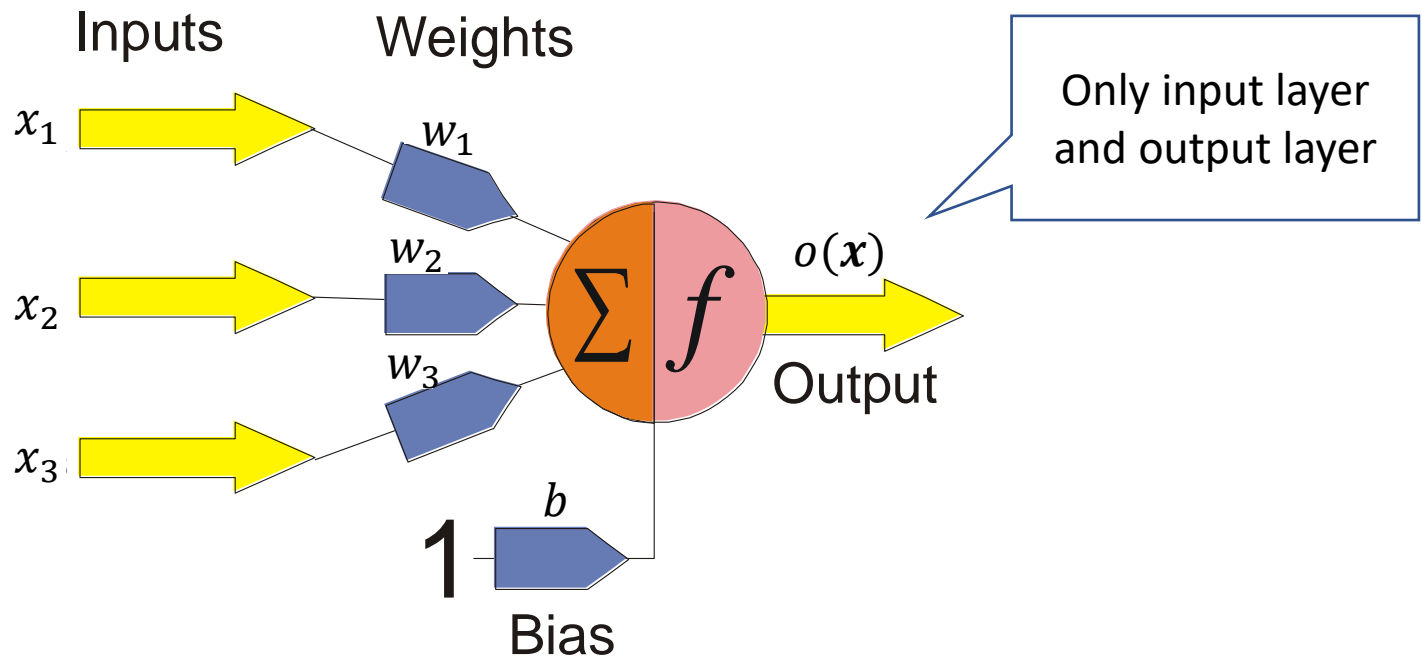  - ADP
  - Deep Q-Network

# This lecture

- One-layer Perceptron

- Multi-layer Perceptron
  - Model
  - Strategy
  - Algorithm

- Overfitting

- Application: Toyota Corolla

Reference: VE 445, Shuai LI (SJTU)

# Perceptron

One-layer Neural Network

# One-layer Perceptron

- The output is a function of the weighted sum of inputs

Inputs     Weights

$x_1$

$w_1$

$w_2$

$x_2$

$w_3$

$x_3$

$\Sigma$  $f$

$o(\boldsymbol{x})$

Output

Only input layer and output layer

**1**

$b$

Bias

$$o(\boldsymbol{x}) = f(x_1 w_1 + x_2 w_2 + x_3 w_3 + b)$$

# One-layer Perceptron (cont.)

- A one-layer perceptron is very similar to Logistic regression and SVM!

# Limitation of One-layer Perceptron

- XOR problem cannot be done by one-layer perceptron

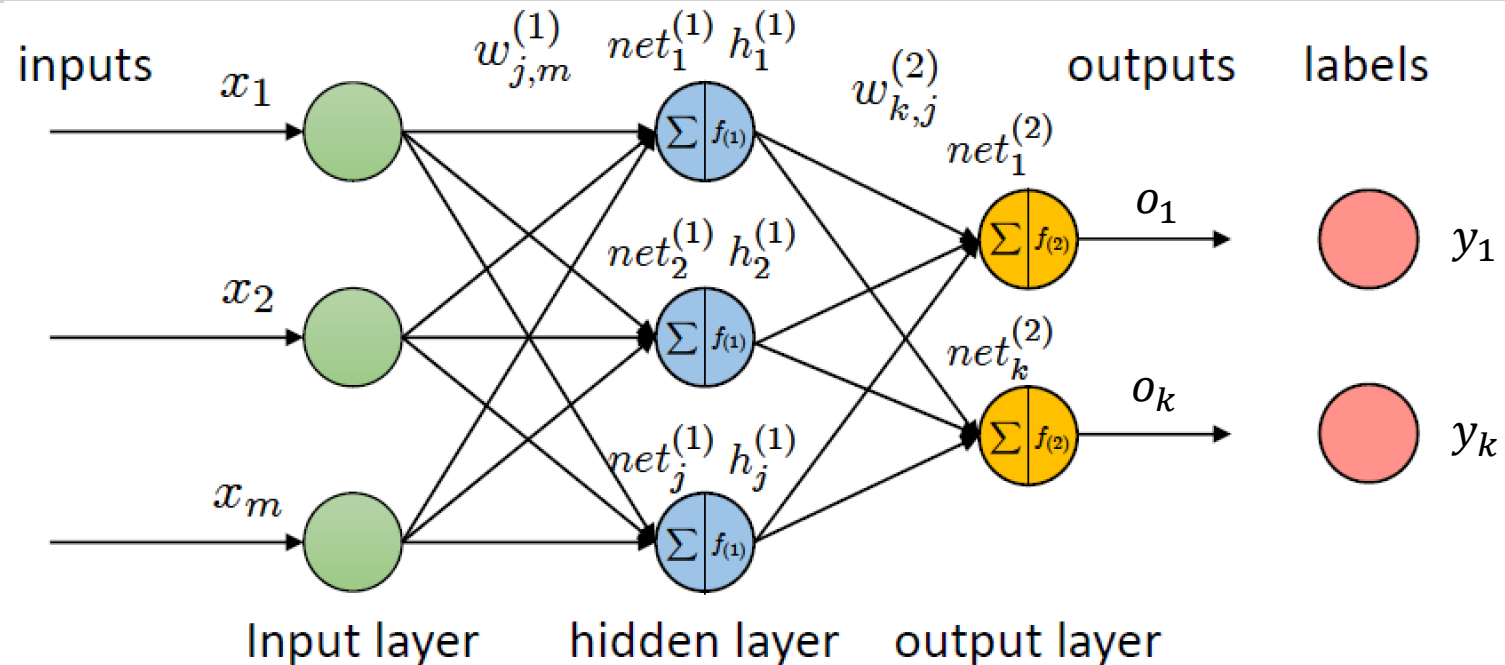| Input x | | Output y |
|:---:|:---:|:---:|
| $X_1$ | $X_2$ | $X_1$ XOR $X_2$ |
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |



Solution: add layers, between input layer and output layer!

# Multi-layer Perceptron

Model

Two-layer feedforward neural network

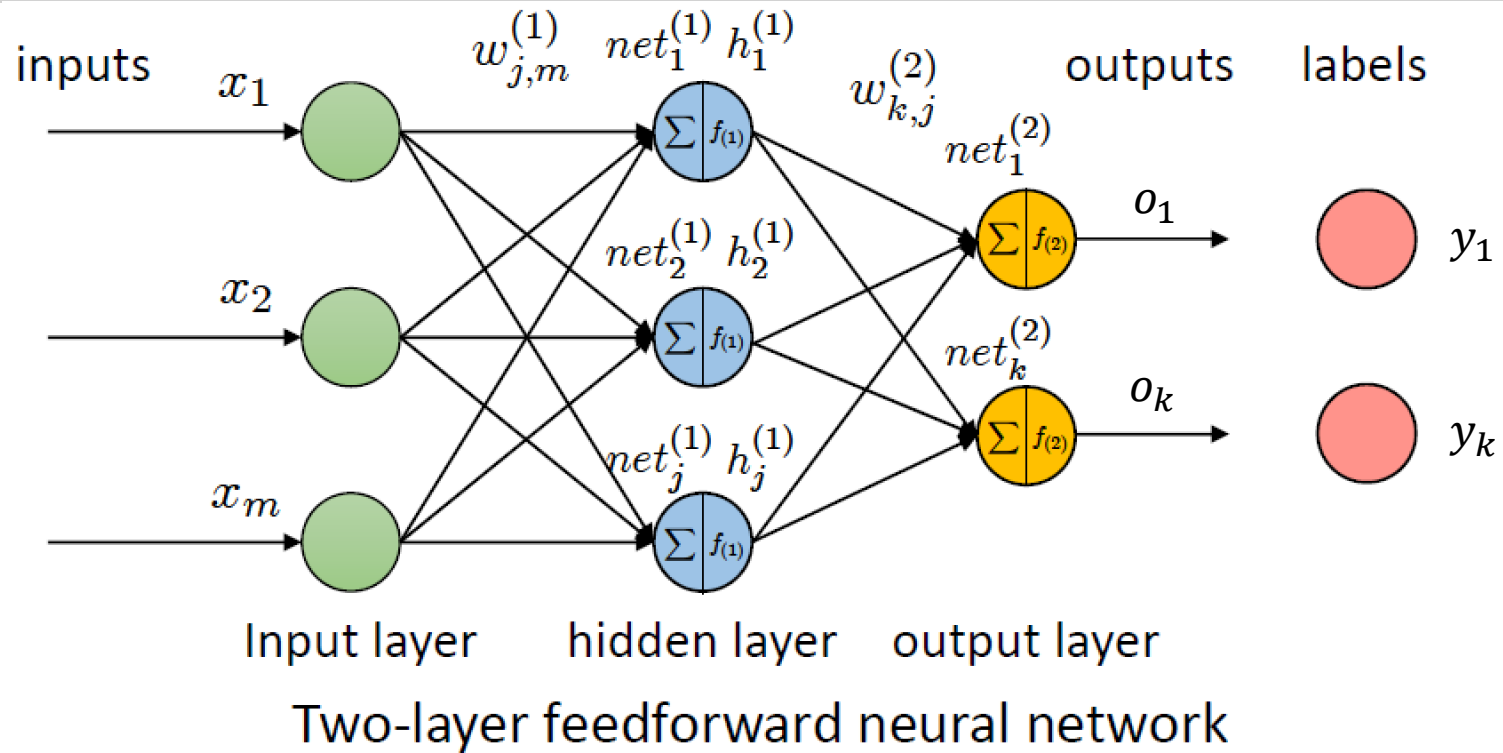- ## Multi-layer Perceptron Structure:
  - ### Input layer
    - Representing variables: $(x_1, x_2, \ldots, x_M)$
  - ### Hidden layer
  - ### Output layer
    - Classification problem: multiple neurons, representing labels: $(y_1, y_2, \ldots, y_K)$
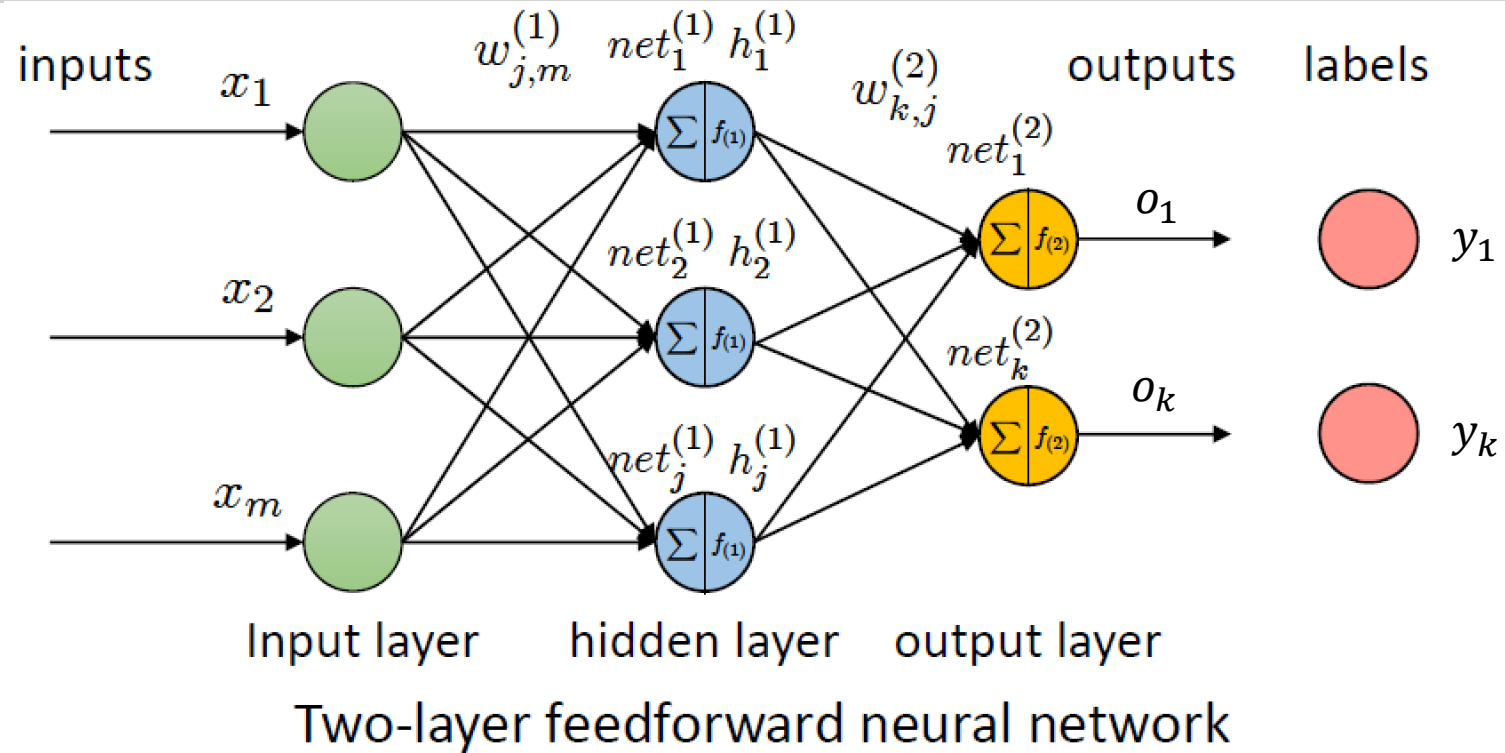    - Regression problem: single neuron

Two-layer feedforward neural network

- Hidden layer
  - Each neuron is a function
  - The outputs are the inputs for next layer
- The link
  - $w_{j,m}^{(1)}$: weight from input $m$ to neuron $j$
  - $w_{k,j}^{(2)}$: weight from neuron $j$ to output $k$

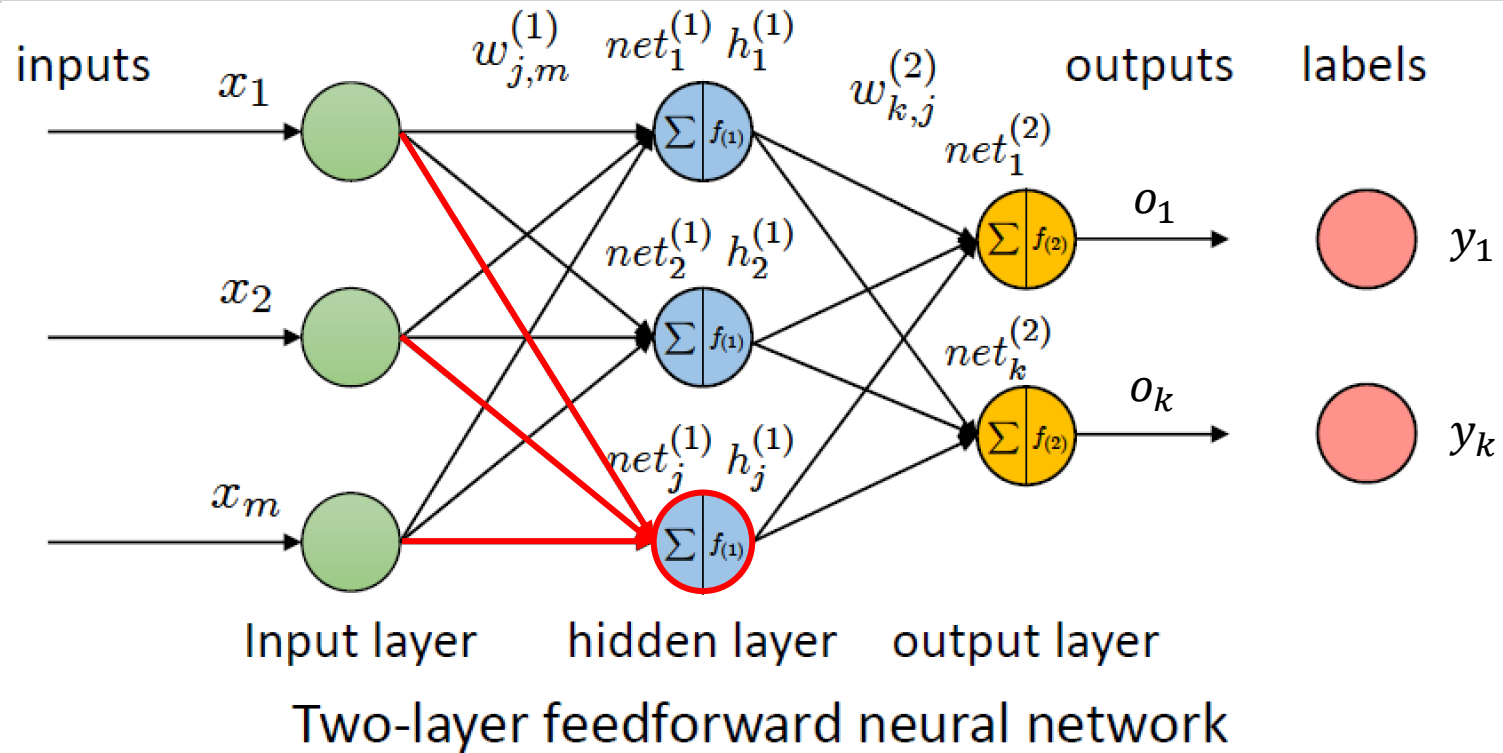**Key issue:**
**Design the structure!**
**How many layers?**
**How many nodes?**

inputs $x_1$ $x_2$ $x_m$

$w_{j,m}^{(1)}$ $net_1^{(1)}$ $h_1^{(1)}$ $w_{k,j}^{(2)}$ $net_1^{(2)}$ outputs labels

$net_2^{(1)}$ $h_2^{(1)}$ $net_k^{(2)}$ $o_1$ $y_1$

$net_j^{(1)}$ $h_j^{(1)}$ $o_k$ $y_k$

Input layer    hidden layer    output layer

Two-layer feedforward neural network

Making prediction:

$$\boldsymbol{x} = (x_1, \dots, x_M) \xrightarrow{\hspace{4cm}} h_j^{(1)} \xrightarrow{\hspace{6cm}} o_k$$

inputs $x_1$  $w_{j,m}^{(1)}$  $net_1^{(1)}$  $h_1^{(1)}$  $w_{k,j}^{(2)}$  outputs  labels

$net_1^{(2)}$

$o_1$

$y_1$

$net_2^{(1)}$  $h_2^{(1)}$

$x_2$

$net_k^{(2)}$

$net_j^{(1)}$  $h_j^{(1)}$

$o_k$

$x_m$

$y_k$

Input layer        hidden layer     output layer

Two-layer feedforward neural network

## Making prediction:

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right) = f_{(1)}\left(\sum_m w_{j,m}^{(1)} x_m\right)$$

$$\boldsymbol{x} = (x_1, \ldots, x_M) \xrightarrow{\hspace{5cm}} h_j^{(1)} \xrightarrow{\hspace{6cm}} o_k$$

where

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

inputs $x_1$ $w_{j,m}^{(1)}$ $net_1^{(1)}$ $h_1^{(1)}$ $w_{k,j}^{(2)}$ outputs labels

$net_1^{(2)}$

$net_2^{(1)}$ $h_2^{(1)}$

$x_2$

$net_k^{(2)}$

$net_j^{(1)}$ $h_j^{(1)}$

$x_m$

$o_1$

$y_1$

$o_k$

$y_k$

Input layer       hidden layer    output layer

Two-layer feedforward neural network

# Making prediction:

$$h_j^{(1)}=f_{(1)}\big(net_j^{(1)}\big)=f_{(1)}(\sum_m w_{j,m}^{(1)}x_m)$$

$$o_k=f_{(2)}\big(net_k^{(2)}\big)=f_{(2)}(\sum_m w_{k,j}^{(2)}h_j^{(1)})$$

$$\boldsymbol{x}=(x_1,\dots,x_M)\ \xrightarrow{\hspace{4cm}}\ h_j^{(1)}\ \xrightarrow{\hspace{4cm}}\ o_k$$

where $\qquad net_j^{(1)}=\sum_m w_{j,m}^{(1)}x_m \qquad\qquad net_k^{(2)}=\sum_m w_{k,j}^{(2)}h_j^{(1)}$

# Activation Function

**Linear activation**

$$f(net) = net$$

**Sigmoid activation**

$$f(net) = \frac{1}{1 + e^{-net}}$$

**Threshold activation**

$$f(net) = \text{sign}(net) = \begin{cases} 1, & if \quad net \geq 0, \\ -1, & if \quad net < 0. \end{cases}$$
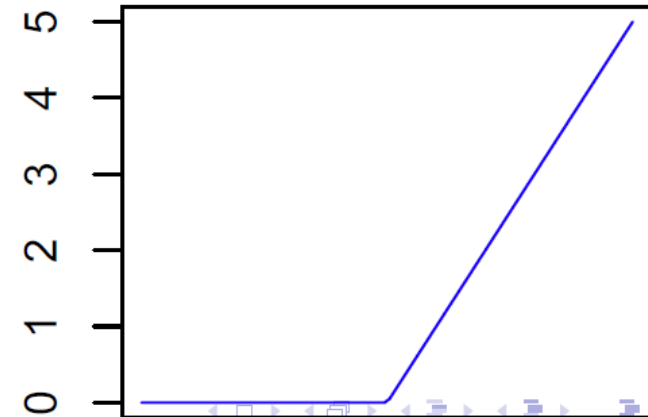
**Hyperbolic tangent activation**

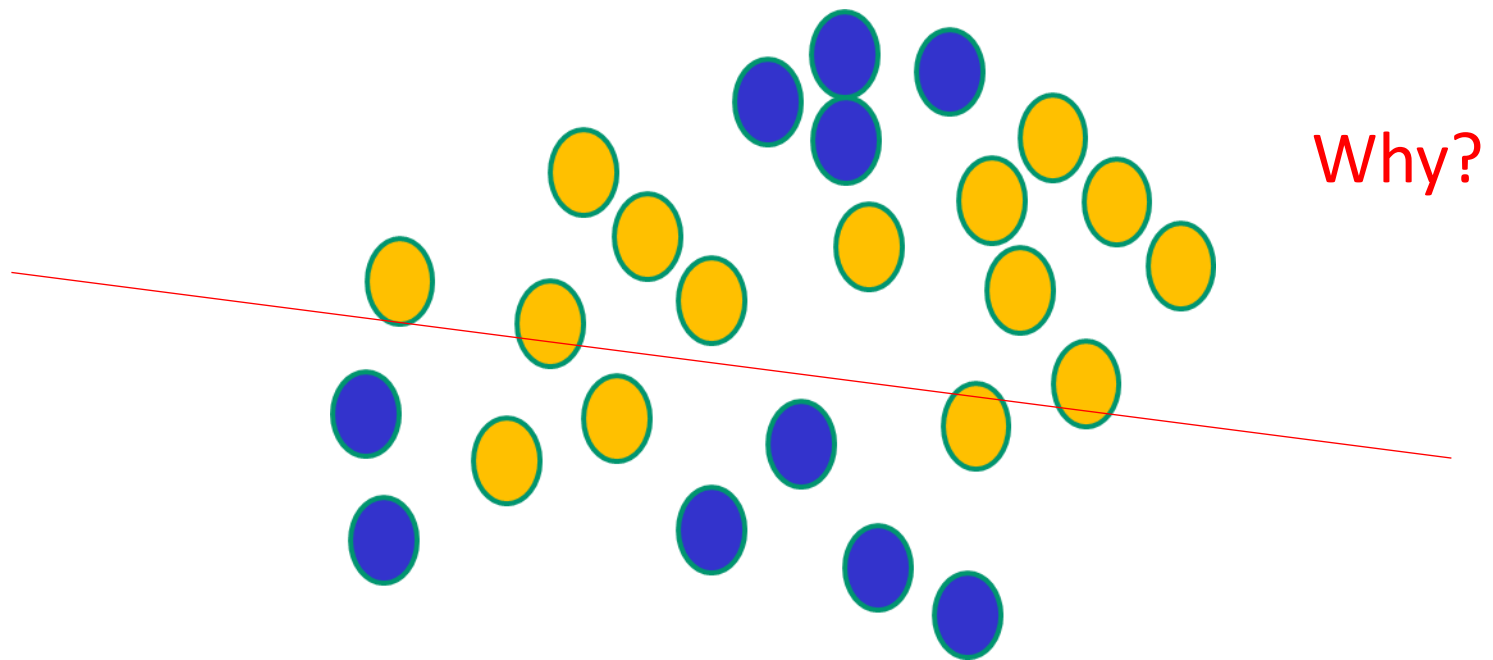$$f(net) = tanh(net) = \frac{1 - e^{-2net}}{1 + e^{-2net}}$$

rectified linear unit (ReLU)

# Activation Function (cont.)

- If $f(x)$ is linear, the NN can **only** draw straight decision boundaries (even if there are many layers of units)
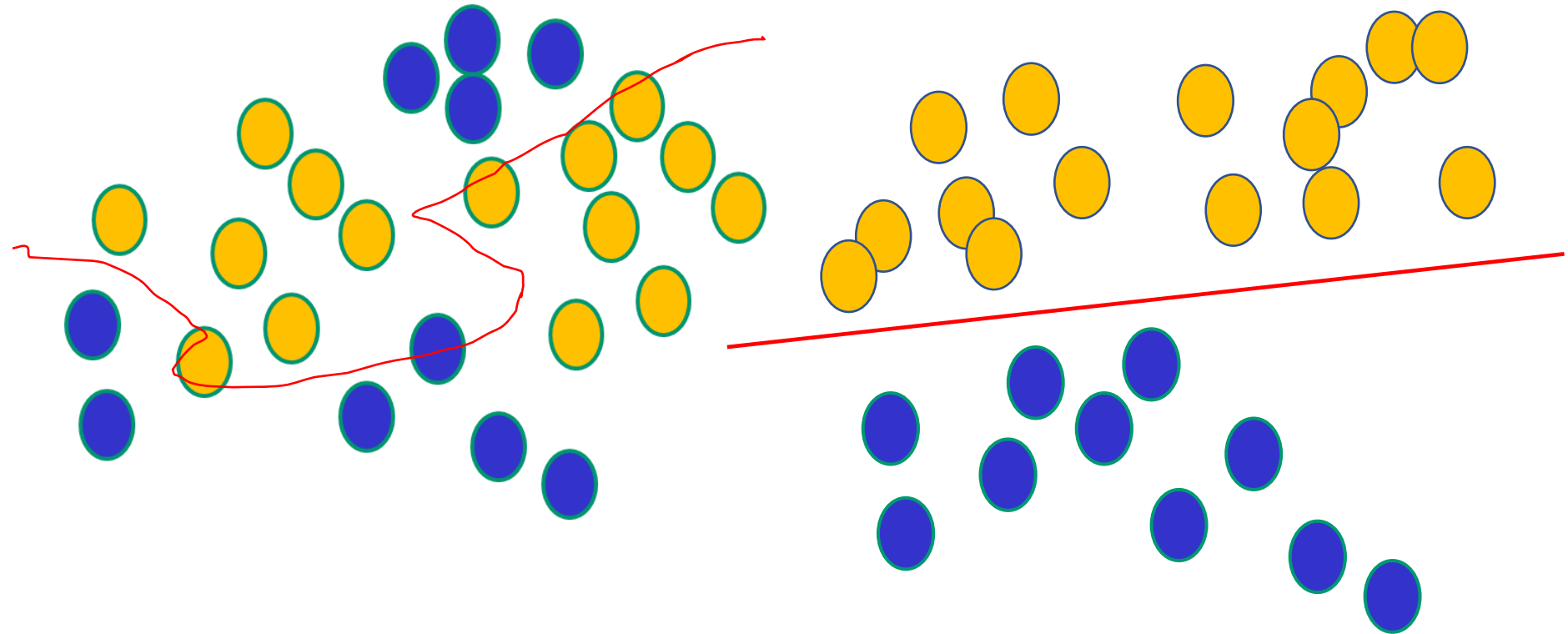
Why?

# Activation Function (cont.)

- If $f(x)$ is non-linear, a network with 1 hidden layer can, in theory, learn perfectly **any** classification problem.
  - A set of weights exists that can produce the targets from the inputs.
  - The problem is finding them.

# Compared with SVM

NNs use nonlinear $f(x)$ so they can draw complex boundaries, but keep the data unchanged

SVMs only draw straight lines, but they transform the data first in a way that makes that OK

# Multi-layer Perceptron

Strategy

# Minimize Empirical Error

- For regression problem

$$\min \frac{1}{2N} \sum_{n=1}^{N} (o^n - y^n)^2$$

  - $o^n$ is the output of the input $\boldsymbol{x}^n$ of sample $(\boldsymbol{x}^n, y^n)$
  - $y^n$ is the true value of the input $\boldsymbol{x}^n$

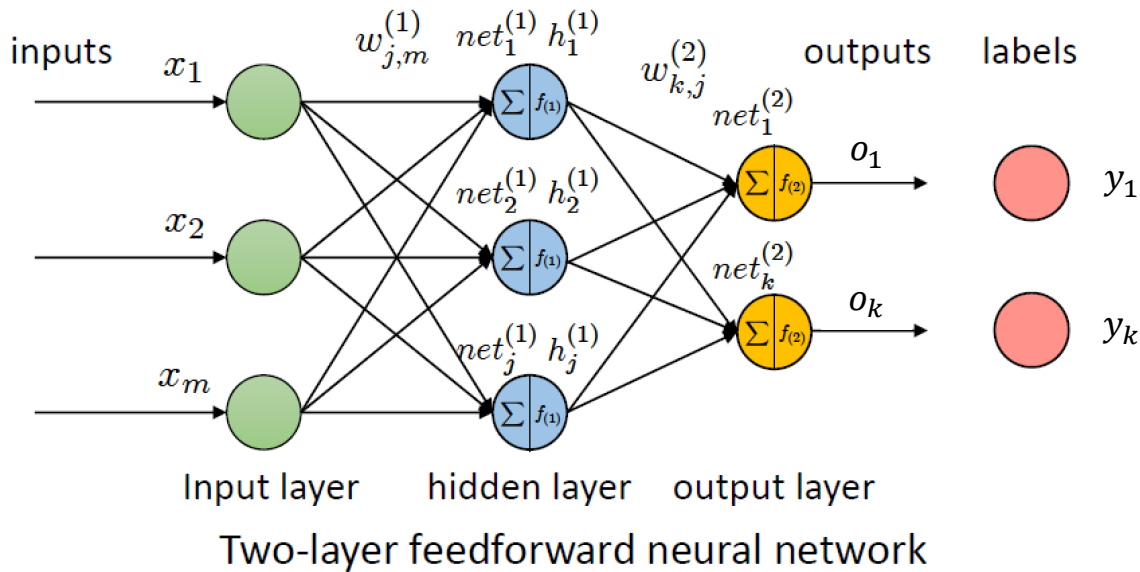- For classification problem

$$\min \frac{1}{2N} \sum_{n=1}^{N} \sum_{k=1}^{K} (o_k^n - y_k^n)^2$$

  - $o_k^n$ is the output for label $k$ of the input $\boldsymbol{x}^n$ of sample $(\boldsymbol{x}^n, \boldsymbol{y}^n)$
  - $\boldsymbol{y}^n$ is the true label of the input $\boldsymbol{x}^n$: $(0,0,1,\dots,0)$

# Multi-layer Perceptron

Algorithm

# Forward Propagation of Information



inputs

$x_1$

$w_{j,m}^{(1)}$   $net_1^{(1)}$ $h_1^{(1)}$

$w_{k,j}^{(2)}$   outputs   labels

$net_1^{(2)}$

$o_1$

$y_1$

$x_2$

$net_2^{(1)}$ $h_2^{(1)}$

$net_k^{(2)}$

$o_k$

$y_k$

$x_m$

$net_j^{(1)}$ $h_j^{(1)}$

Input layer    hidden layer   output layer

Two-layer feedforward neural network

The error:

$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

## Forward Propagation of Information:

$$\boldsymbol{x} = (x_1, \ldots, x_M) \xrightarrow{h_j^{(1)}=f_{(1)}\left(net_j^{(1)}\right)=f_{(1)}(\sum_m w_{j,m}^{(1)}x_m)} h_j^{(1)} \xrightarrow{o_k=f_{(2)}\left(net_k^{(2)}\right)=f_{(2)}(\sum_m w_{k,j}^{(2)}h_j^{(1)})} o_k$$

where

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)}x_m \qquad net_k^{(2)} = \sum_m w_{k,j}^{(2)}h_j^{(1)}$$

# Backward Propagation of Errors



Two-layer feedforward neural network

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$
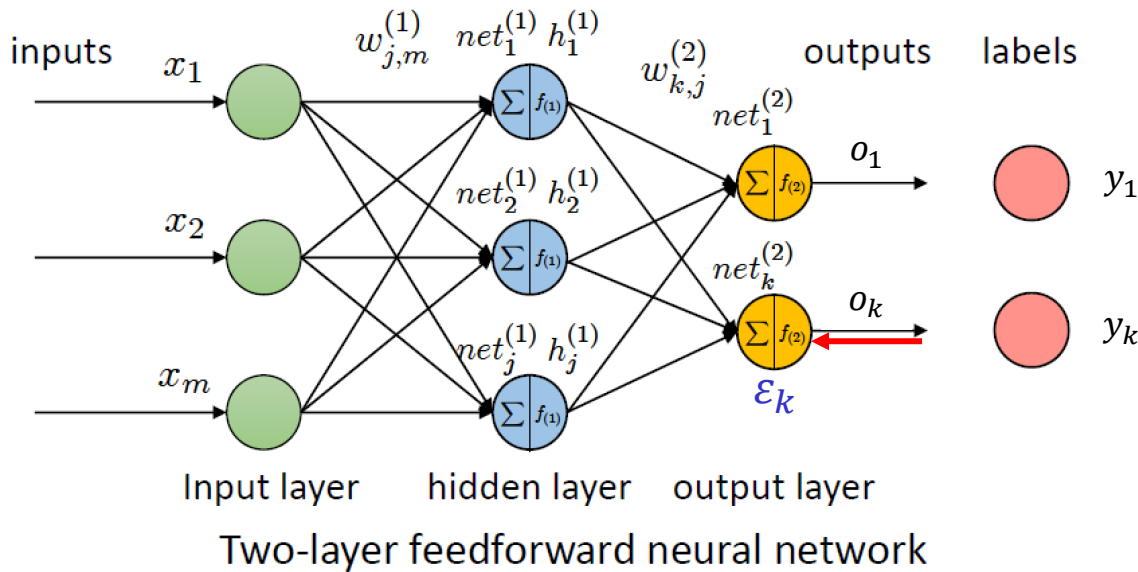
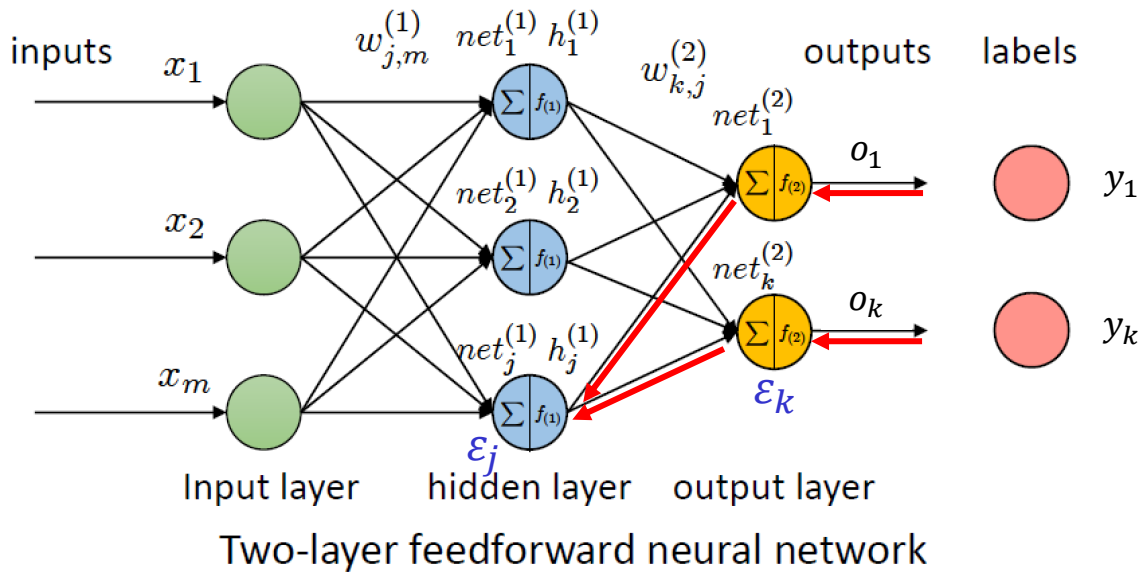$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$

$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

## Backward Propagation of Errors:

$$\varepsilon_m \longleftarrow \varepsilon_j \longleftarrow \varepsilon_k \longleftarrow \varepsilon(\boldsymbol{w})$$

# Backward Propagation of Errors



Two-layer feedforward neural network

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$
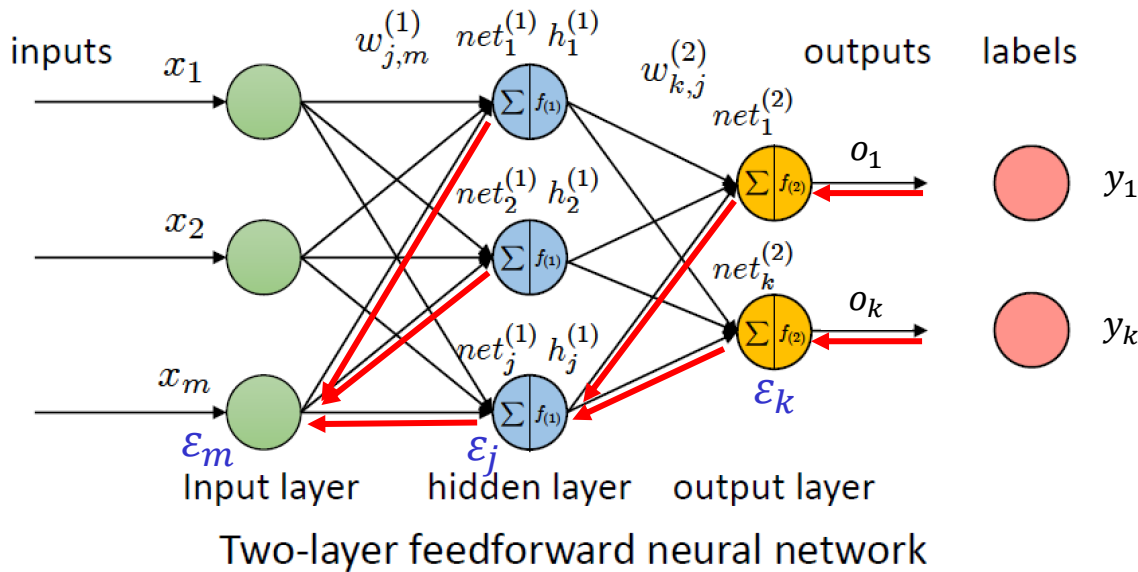
$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

## Backward Propagation of Errors:

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_k^{(2)}} = (o_k - y_k)f_{(2)}'^k$$

$$\varepsilon_m \longleftarrow \varepsilon_j \longleftarrow \varepsilon_k \longleftarrow \varepsilon(\boldsymbol{w})$$

$$\varepsilon_k$$

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_k^{(2)}} = \frac{\partial \frac{1}{2}\sum_{k=1}^{K}(f_{(2)}(net_k^{(2)}) - d_k)^2}{\partial net_k^{(2)}} = \boxed{(o_k - y_k)f_{(2)}'^k}$$

# Backward Propagation of Errors



$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$
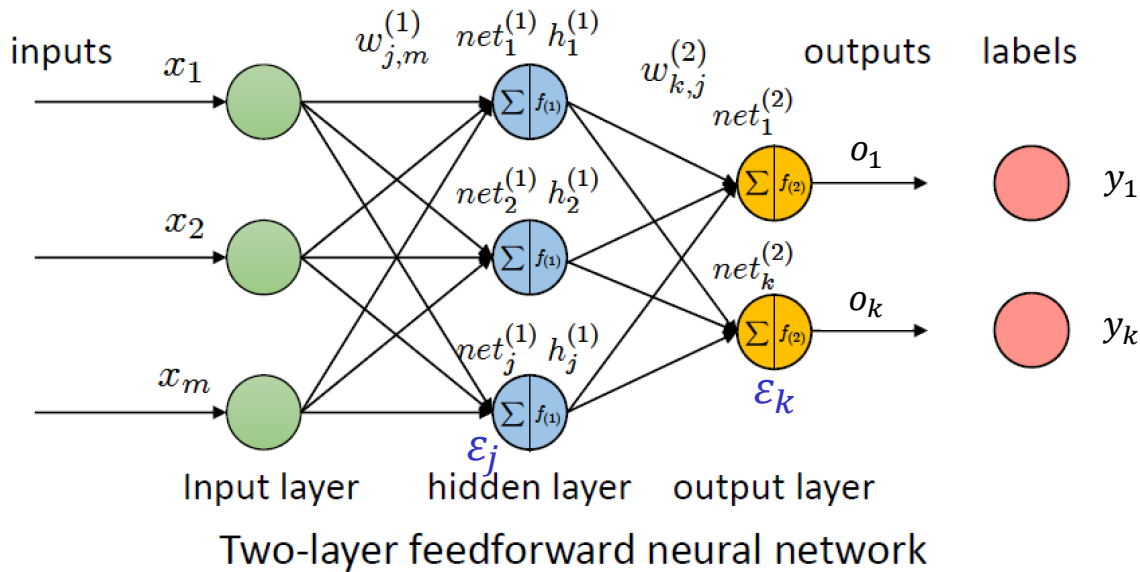
$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

inputs · Input layer · hidden layer · output layer · outputs · labels

Two-layer feedforward neural network

## Backward Propagation of Errors:

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_j^{(1)}} = \sum_{k=1}^{K} \varepsilon_k w_{k,j}^{(2)} f_{(1)}'^{j} \qquad \frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_k^{(2)}} = (o_k - y_k) f_{(2)}'^{k}$$

$$\varepsilon_m \longleftarrow \varepsilon_j \longleftarrow \varepsilon_k \longleftarrow \varepsilon(\boldsymbol{w})$$

$$\varepsilon_j$$

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_j^{(1)}} = \sum_{k=1}^{K} \frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_k^{(2)}} \frac{\partial net_k^{(2)}}{\partial net_j^{(1)}} = \sum_{k=1}^{K} \varepsilon_k \frac{\partial \sum_j w_{k,j}^{(2)} f_{(1)}\left(net_j^{(1)}\right)}{\partial net_j^{(1)}} = \sum_{k=1}^{K} \varepsilon_k w_{k,j}^{(2)} f_{(1)}'^{j}$$

# Backward Propagation of Errors



Two-layer feedforward neural network

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$
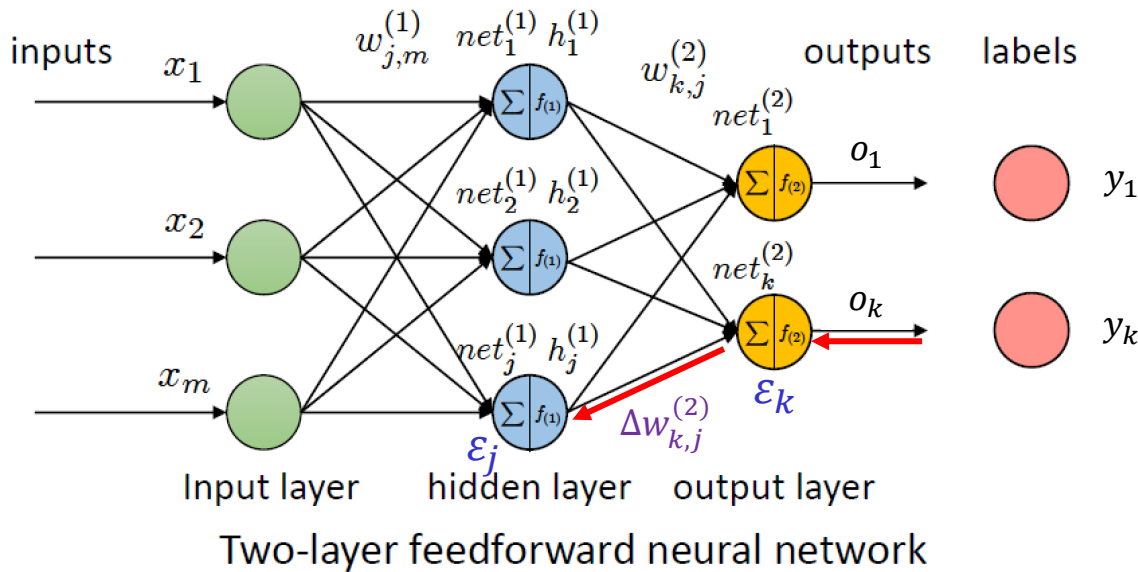
$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$

$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

## Backward Propagation of Errors:

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial x_m} = \sum_{j=1}^{J} \varepsilon_j w_{j,m}^{(1)} \qquad \frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_j^{(1)}} = \sum_{k=1}^{K} \varepsilon_k w_{k,j}^{(2)} f_{(1)}'^j \qquad \frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_k^{(2)}} = (o_k - y_k) f_{(2)}'^k$$

$$\varepsilon_m \longleftarrow \varepsilon_j \longleftarrow \varepsilon_k \longleftarrow \varepsilon(\boldsymbol{w})$$

$$\boxed{\varepsilon_m}$$

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial x_m} = \sum_{j=1}^{J} \frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_j^{(1)}} \frac{\partial net_j^{(1)}}{\partial x_m} = \sum_{j=1}^{J} \varepsilon_j \frac{\partial \sum_m w_{j,m}^{(1)} x_m}{\partial x_m} = \boxed{\sum_{j=1}^{J} \varepsilon_j w_{j,m}^{(1)}}$$

# Backward Propagation of Errors



Two-layer feedforward neural network

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$
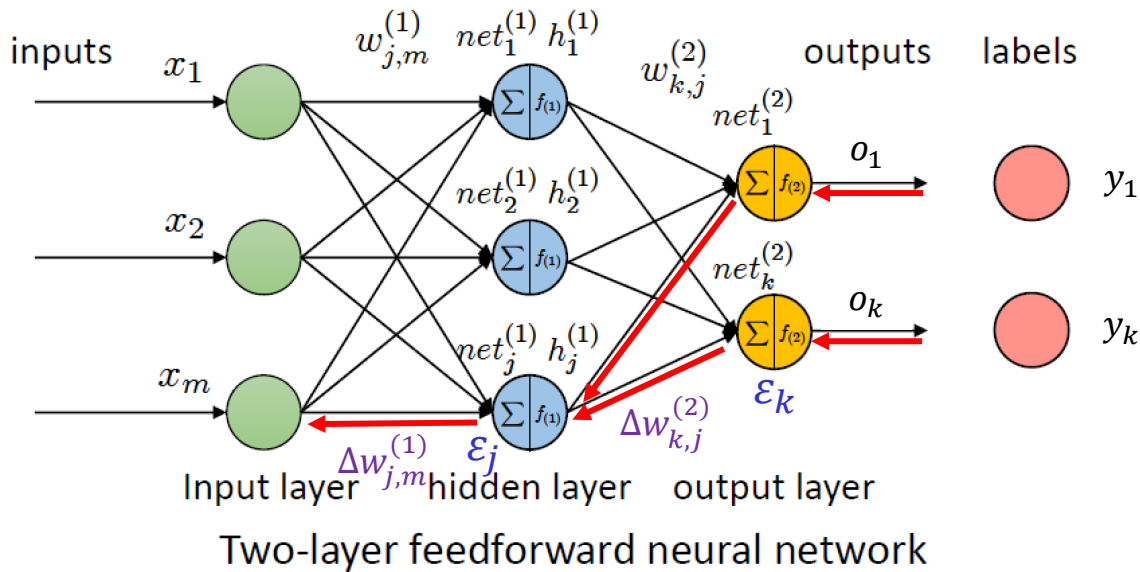
$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$

$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

## Backprop to learn the parameters:

$$w_{j,m}^{(1)} \longleftarrow w_{k,j}^{(2)} \longleftarrow \varepsilon(\boldsymbol{w})$$

# Backward Propagation of Errors



Two-layer feedforward neural network

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$

$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K} (o_k - y_k)^2$$

## Backprop to learn the parameters:

$$w_{j,m}^{(1)} \longleftarrow w_{k,j}^{(2)} \xleftarrow{\frac{\partial \varepsilon(\boldsymbol{w})}{\partial w_{k,j}^{(2)}} = \varepsilon_k h_j^{(1)}} \varepsilon(\boldsymbol{w})$$

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial w_{k,j}^{(2)}} = \frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_k^{(2)}} \frac{\partial net_k^{(2)}}{\partial w_{k,j}^{(2)}} = \boxed{\varepsilon_k h_j^{(1)}} \longleftarrow \boxed{\Delta w_{k,j}^{(2)}}$$

# Backward Propagation of Errors



$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$

$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

inputs

outputs   labels

Input layer   hidden layer   output layer

Two-layer feedforward neural network

## Backprop to learn the parameters:

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial w_{j,m}^{(1)}} = \varepsilon_j x_m \qquad \frac{\partial \varepsilon(\boldsymbol{w})}{\partial w_{k,j}^{(2)}} = \varepsilon_k h_j^{(1)}$$

$$w_{j,m}^{(1)} \longleftarrow w_{k,j}^{(2)} \longleftarrow \varepsilon(\boldsymbol{w})$$

$$\frac{\partial \varepsilon(\boldsymbol{w})}{\partial w_{j,m}^{(1)}} = \frac{\partial \varepsilon(\boldsymbol{w})}{\partial net_j^{(1)}} \frac{\partial net_j^{(1)}}{\partial w_{j,m}^{(1)}} = \boxed{\varepsilon_j x_m} \longleftarrow \boxed{\Delta w_{j,m}^{(1)}}$$

# Backpropagation (BP) Algorithm



inputs $x_1$, $x_2$, $x_m$

$w_{j,m}^{(1)}$ $net_1^{(1)}$ $h_1^{(1)}$

$w_{k,j}^{(2)}$ outputs labels

$net_1^{(2)}$

$net_2^{(1)}$ $h_2^{(1)}$ $o_1$ $y_1$

$net_k^{(2)}$

$net_j^{(1)}$ $h_j^{(1)}$ $o_k$ $y_k$

$\varepsilon_k$

$\Delta w_{j,m}^{(1)}$ $\varepsilon_j$ $\Delta w_{k,j}^{(2)}$

Input layer hidden layer output layer

Two-layer feedforward neural network

$$net_j^{(1)} = \sum_m w_{j,m}^{(1)} x_m$$

$$h_j^{(1)} = f_{(1)}\left(net_j^{(1)}\right)$$

$$net_k^{(2)} = \sum_j w_{k,j}^{(2)} h_j^{(1)}$$

$$o_k = f_{(2)}\left(net_k^{(2)}\right)$$

$$\varepsilon(\boldsymbol{w}) = \frac{1}{2}\sum_{k=1}^{K}(o_k - y_k)^2$$

## Update the parameters:

Error of $k$

$$w_{k,j}^{(2)\prime} = w_{k,j}^{(2)} - \eta \Delta w_{k,j}^{(2)} = w_{k,j}^{(2)} - \eta \varepsilon_k h_j^{(1)} \longleftarrow \text{Output of } j$$

Learning rate

$$w_{j,m}^{(1)\prime} = w_{j,m}^{(1)} - \eta \Delta w_{j,m}^{(1)} = w_{j,m}^{(1)} - \eta \varepsilon_j x_m \longleftarrow \text{Output of } m$$

Error of $j$

# Chain Rule

- Update the weight on the link from node $j$ to $i$

$$w_{i,j}{'} = w_{i,j} - \eta \varepsilon_i h_j$$

  - $h_j$ is the output of node $j$

- If node $i$ is in output layer

$$\varepsilon_i = (o_i - y_i)f'^i$$

- If node $i$ is in hidden layer

$$\varepsilon_i = \sum_{k=1}^{K} \varepsilon_k \, w_{k,i} f'^i$$



$j \quad w_{i,j} \qquad w_{k,i} \quad k$

$i$

Repeating updates, after a number of iterations, we will get in the (local) optima

# Overfitting

# Avoiding Overfitting

- A typical curve showing performance during training
- But here is performance on unseen data, **not** in the training set

# Avoiding Overfitting

- Use as much training data as possible, and ensure as much diversity as possible in the data
  - This cuts down on the potential existence of features that might be discriminative in the training data, but are otherwise spurious

- Jittering (adding noise)
  - During training, every time an input pattern is presented, it is randomly perturbed
  - The idea of this is that spurious features will be `washed out' by the noise, but valid discriminatory features will remain
  - The problem with this approach is how to correctly choose the level of noise

# Avoiding Overfitting (cont.)

- Early stopping
    - During training, keep track of the network's performance on a separate *validation set* of data
    - At the point where error continues to improve on the training set, but starts to get worse on the validation set, that is when training should be stopped, since it is starting to overfit on the training data
    - The problem here is that this point is far from always clear cut

# Avoiding Overfitting (cont.)

- Regularization
  - Limit the summation of the weights
- Dropout
  - In training, for each sample, dropout a neuron with some probability



(a) Standard Neural Net          (b) After applying dropout.

# Dropout on an image classification task

# NN Playground

- [playground.tensorflow.org](playground.tensorflow.org)

# Application

Toyota Corolla

# Toyota Corolla

- A line of subcompact and compact cars

- Introduced in 1966

- Best selling nameplate in the world

# Used Car Valuation

- If you have a car to sell or want to buy a used car

- Knowing how to calculate the estimated value of the car is paramount in making sure that you get the best value for your money

# Data

- On used Toyota Corolla

- 1436 record

- 38 attributes
  - Price
  - Age
  - KM
  - HP
  - …

Goal:
Predict the price of a used Toyota Corolla
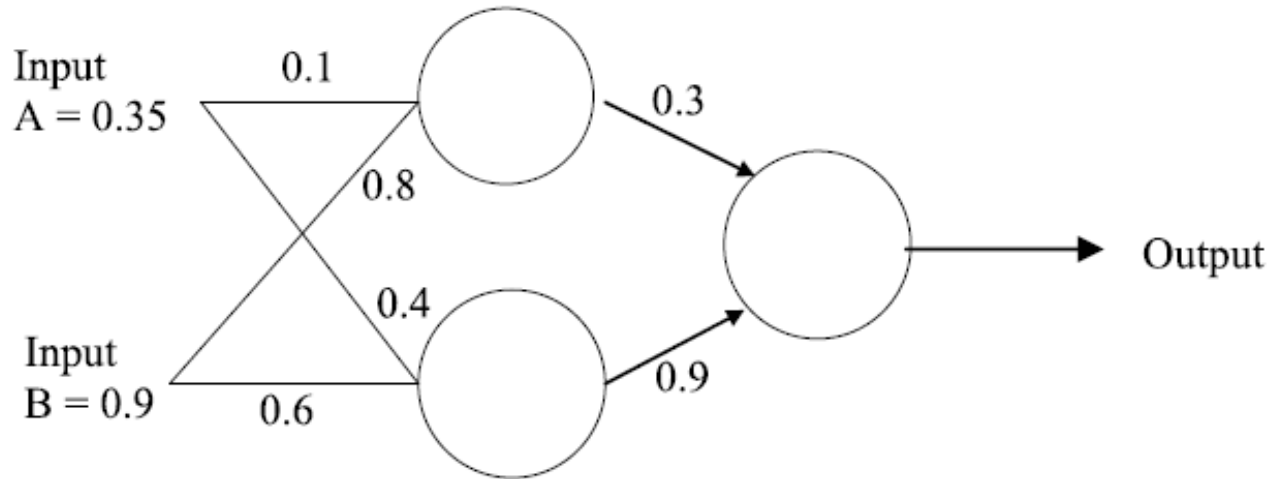
Method:
Neural Network

# Let's get our hands dirty!

# Lecture 9 wrap-up

✓One-layer Perceptron

✓Multi-layer Perceptron

   ✓Model

   ✓Strategy

   ✓Algorithm

✓Overfitting

✓Application: Toyota Corolla

# Assignment 9

- Consider the simple network below:



- Assume that the neurons have a Sigmoid activation function and
  - 1. Perform a forward pass on the network
  - 2. Perform a reverse pass (training) once (target = 0.5)
  - 3. Perform a further forward pass and comment on the result

- i.e.,
  - 1. Calculate the input and output for each neuron
  - 2. Calculate the new weights for each layer after once training
  - 3. Calculate the new error and compare

- Send your answer to TA (any form, e.g., word, pdf, photo …)

- Due: June 5, 11pm

- TA: Mr. Xiong, xiongym3@mail2.sysu.edu.cn

# Next lecture

- Supervised learning
  - Linear regression
  - Logistic regression
  - SVM and kernel
  - Tree models

- Deep learning
  - Neural networks
  - Convolutional NN
  - Recurrent NN

- Unsupervised learning
  - Clustering
  - PCA
  - EM

- Reinforcement learning
  - MDP
  - ADP
  - Deep Q-Network

# Questions?

Shan Wang (王杉)

https://wangshan731.github.io/