# Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

April 15, 2015

Today:
- Artificial neural networks
- Backpropagation
- Recurrent networks
- Convolutional networks

Reading:
- Mitchell: Chapter 4
- Bishop: Chapter 5
- Quoc Le tutorial:
- Ruslan Salakhutdinov tutorial:

# Artificial Neural Networks to learn f: X → Y

- f might be non-linear function
- X (vector of) continuous and/or discrete vars
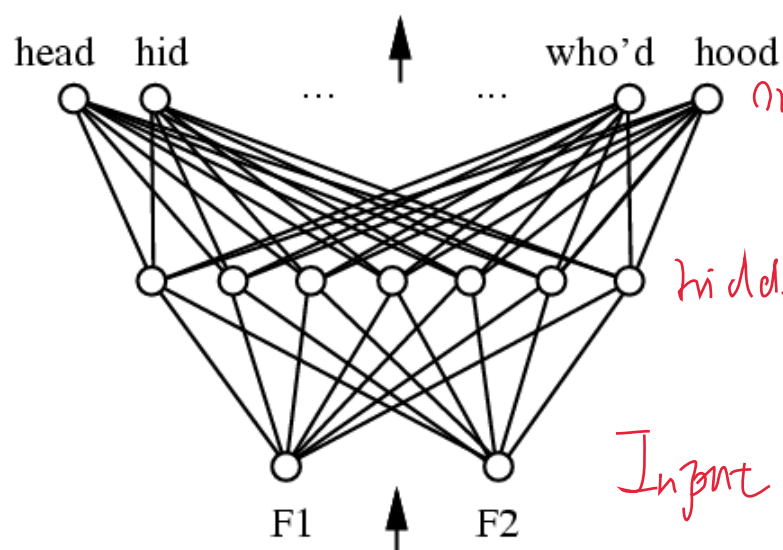- Y (vector of) continuous and/or discrete vars

- Represent f by _network_ of logistic units
- Each unit is a logistic function

$$unit\ output = \frac{1}{1 + exp(w_0 + \sum_i w_i x_i)}$$

$y = f(x) + \varepsilon, \quad \varepsilon \sim N(0, \delta^2)$

- MLE: train weights of all units to minimize sum of squared errors of predicted network outputs
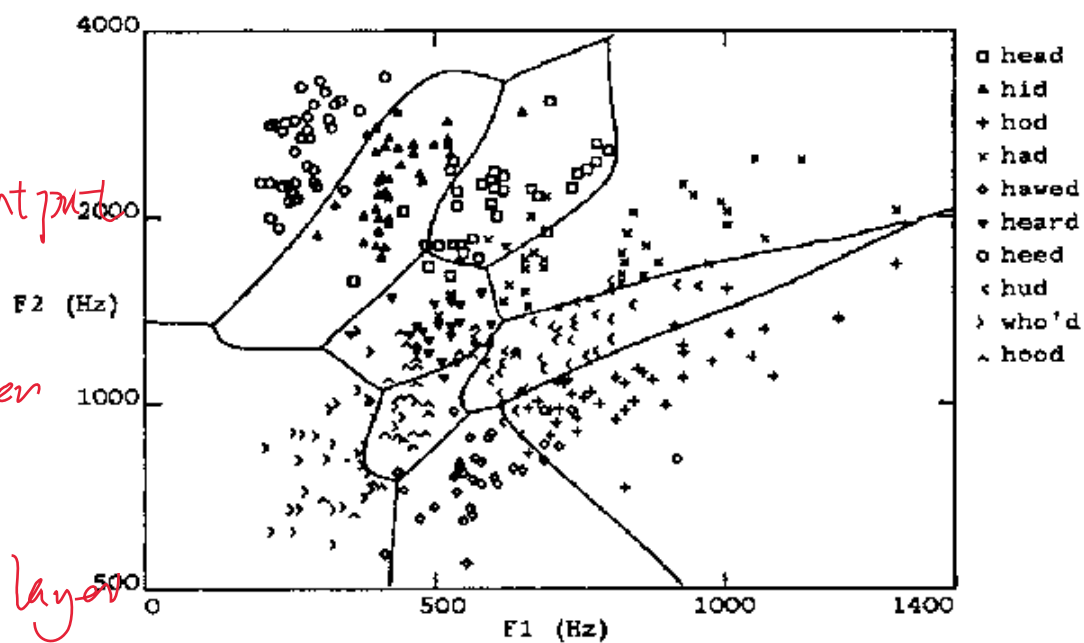- MAP: train to minimize sum of squared errors plus weight magnitudes of w.
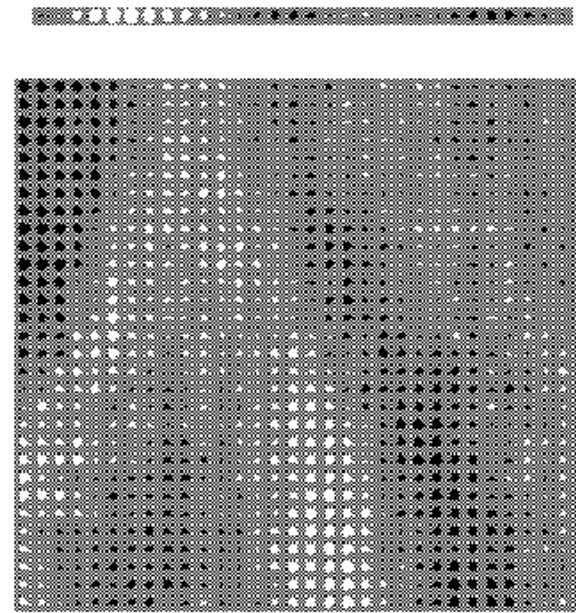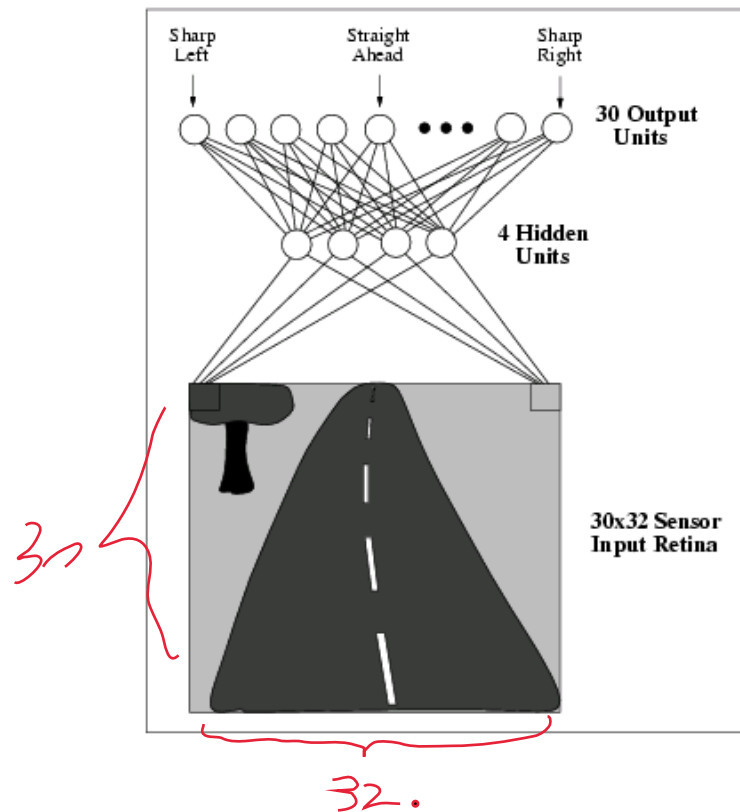
# Multilayer Networks of Sigmoid Units

ALVINN

[Pomerleau 1993]



Sharp Left    Straight Ahead    Sharp Right

30 Output Units

4 Hidden Units

30x32 Sensor Input Retina

30

32.

# Connectionist Models

Consider humans:

- Neuron switching time ˜ .001 second
- Number of neurons ˜ $10^{10}$
- Connections per neuron ˜ $10^{4-5}$
- Scene recognition time ˜ .1 second
- 100 inference steps doesn't seem like enough

$\rightarrow$ much parallel computation

Properties of artificial neural nets (ANN's):

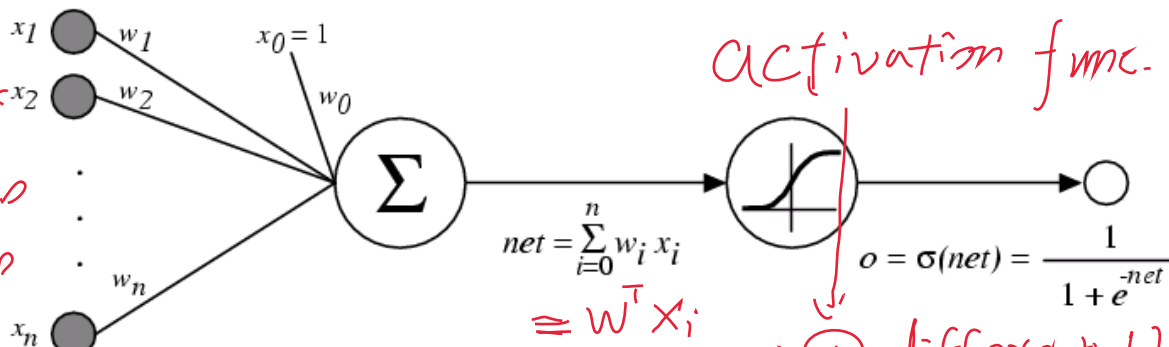- Many neuron-like threshold switching units
- Many weighted interconnections among units
- Highly parallel, distributed process

## Sigmoid Unit

$\sigma(x)$

$\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$

(Gradient vanishing)

$x_1, w_1$
$x_0 = 1, w_0$
$x_2, w_2$
$w_n$
$x_n$

activation func.

$net = \sum_{i=0}^{n} w_i x_i$

$\cong w^\top x_i$

$o = \sigma(net) = \dfrac{1}{1 + e^{-net}}$

①. differentiable (simple)

②. non-linear.

$\sigma(x)$ is the sigmoid function

$$\frac{1}{1 + e^{-x}}$$

$\log \dfrac{P(y=1 | x=x)}{P(y=0 | x=x)} = w^\top x = 0$

$w^\top x = 0$

Nice property: $\dfrac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$

We can derive gradient decent rules to train

- One sigmoid unit

Multi-Layer Perceptron (MLP)

- *Multilayer networks* of sigmoid units →
  Backpropagation

# Multi-Layer Perceptron (MLP)



$net_1$  $O_1 = \sigma(net_1)$

$net_3 = W_3^T O$

$y_1 = \sigma(net_3) = \sigma(W_3^T O)$

$y_2 = \sigma(net_4) = \sigma(W_4^T O)$

$net_2$  $O_2 = \sigma(net_2)$  $net_4$

input    hidden    output

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

$$\log \frac{P(y=1|X=x)}{P(y=0|X=x)} = net_3 = W_3^T O = W_{31} \cdot O_1 + W_{32} O_2$$

$$= W_{31} \cdot \sigma(W_1^T x) + W_{32} \cdot \sigma(W_2^T x) = 0$$

- Single-class: $p(y=1|x) = \dfrac{e^{W^T x}}{1 + e^{W^T x}}$

- multi-class: $p(y=j|x) = \dfrac{e^{W_j^T x}}{\sum_{j=1}^{K} e^{W_j^T x}}$

  (softmax)

# M(C)LE Training for Neural Networks

- Consider regression problem f:X→Y , for scalar Y

$$y = f(x) + \varepsilon \longleftarrow \quad \text{assume noise N}(0,\sigma_\varepsilon), \text{ iid}$$

deterministic

- Let's maximize the conditional data likelihood

$$W \leftarrow \arg\max_W \ \ln \prod_l P(Y^l | X^l, W)$$

$$W \leftarrow \arg\min_W \ \sum_l (y^l - \hat{f}(x^l))^2$$

Learned neural network

# MAP Training for Neural Networks

- Consider regression problem f:X→Y , for scalar Y

$$y = f(x) + \varepsilon \longleftarrow \quad \text{noise } N(0,\sigma_\varepsilon)$$
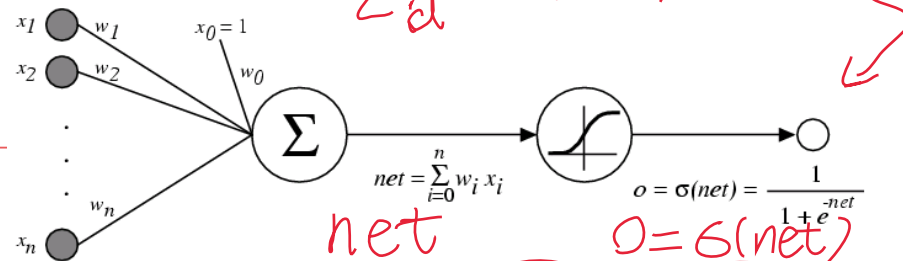
deterministic

Gaussian $P(W) = N(0,\sigma I)$

$$W \leftarrow \arg \max_W \ \ln \ P(W) \prod_l P(Y^l | X^l, W)$$

$$W \leftarrow \arg \min_W \left[ c \sum_i w_i^2 \right] + \left[ \sum_l (y^l - \hat{f}(x^l))^2 \right]$$

$$\ln P(W) \leftrightarrow c \sum_i w_i^2$$

$$f \circ g = f(g(x))$$
$$\frac{\partial f \circ g}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$$

# Error Gradient for a Sigmoid Unit



$$E = \frac{1}{2} \sum_{d} (t_d - o_d)^2$$

$net = \sum_{i=0}^{n} w_i x_i$

$o = \sigma(net) = \dfrac{1}{1 + e^{-net}}$

$net$

$O = \sigma(net)$

$$\frac{\partial E}{\partial w_i} = \sum_d \left( \frac{\partial E_d}{\partial O_d} \right) \cdot \left( \frac{\partial O_d}{\partial net_d} \right) \left( \frac{\partial net_d}{\partial w_i} \right)$$

$-(t_d - o_d)$   $O_d(1 - O_d)$   $x_{i,d}$

$$W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \in \mathbb{R}^{(n+1) \times 1}$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i}(t_d - o_d)$$

$$= \sum_d (t_d - o_d) \left( -\frac{\partial o_d}{\partial w_i} \right)$$

$$= -\sum_d (t_d - o_d) \frac{\partial o_d}{\partial net_d} \frac{\partial net_d}{\partial w_i}$$

But we know:

$$\frac{\partial o_d}{\partial net_d} = \frac{\partial \sigma(net_d)}{\partial net_d} = o_d(1 - o_d)$$

$$\frac{\partial net_d}{\partial w_i} = \frac{\partial (\vec{w} \cdot \vec{x}_d)}{\partial w_i} = x_{i,d}$$

So:

$$\frac{\partial E}{\partial w_i} = -\sum_{d \in D} (t_d - o_d) o_d(1 - o_d) x_{i,d}$$

$$\frac{\partial E}{\partial W} = -(t - o)^T Diag(o_d(1-o_d)) \text{✗}$$

$x_d$ = input

$t_d$ = target output

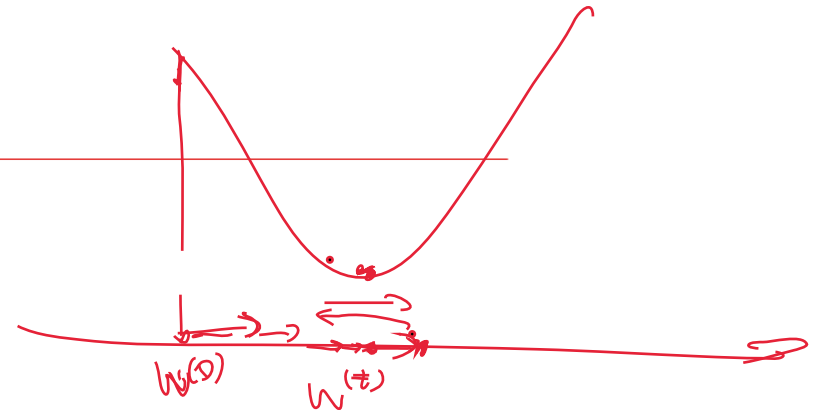$o_d$ = observed unit output

$w_i$ = weight i

$$W_i \leftarrow W_i - \eta \left( \frac{\partial E}{\partial W_i} \right)$$
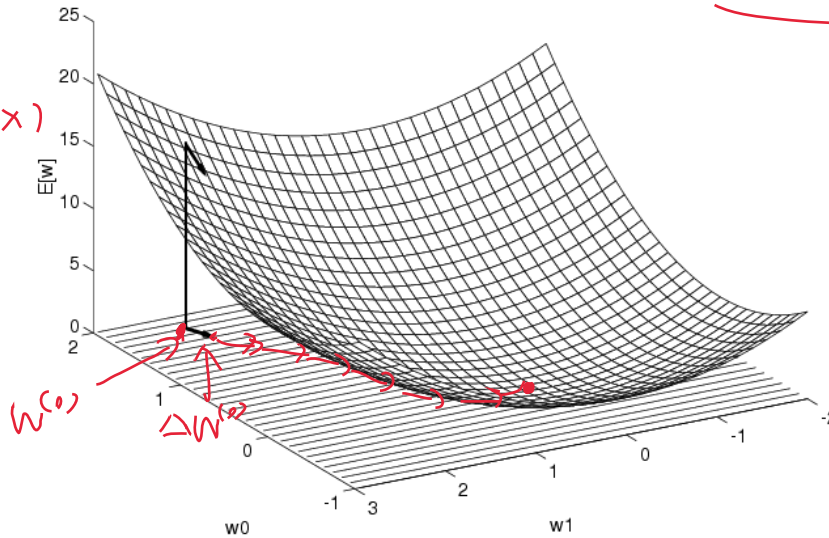
Gradient Descent (GD)

# Gradient Descent

$$E = \frac{1}{2} \sum_d (t_d - o_d)^2$$

$$w = \begin{bmatrix} w_2 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \in \mathbb{R}^{(n+1) \times 1}$$

$$\frac{\partial \cdot}{\partial} = \begin{bmatrix} \frac{\partial E}{\partial w_0} \\ \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_n} \end{bmatrix}$$

$w^{(0)}$  $\Delta w^{(0)}$

Goal: $\min\limits_{w} E[w]$

Solution: $w_i \leftarrow w_i + \Delta w_i$

$(i = 0, 1, \ldots, n)$

## Gradient

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \cdots \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

① fixed   $< 10^{-3}$

learning rate  ② line search.

i.e.,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$w^{(0)}$   $w^{(t)}$

# Incremental (Stochastic) Gradient Descent

$$E_D = \frac{1}{2} \sum_d (t_d - o_d)^2$$

$$E_d = \frac{1}{2} (t_d - o_d)^2$$

**Batch mode** Gradient Descent:    (GD)

Do until satisfied

1. Compute the gradient $\nabla E_D[\vec{w}]$
2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_D[\vec{w}]$

**Incremental mode** Gradient Descent:    Stochastic GD (SGD)

Do until satisfied

• For each training example $d$ in $D$    (Adam)

1. Compute the gradient $\nabla E_d[\vec{w}]$
2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_d[\vec{w}]$

$$E_D[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$E_d[\vec{w}] \equiv \frac{1}{2} (t_d - o_d)^2$$

*Incremental Gradient Descent* can approximate *Batch Gradient Descent* arbitrarily closely if $\eta$ made small enough

# Backpropagation Algorithm (MLE)

## (MLP)

---

Initialize all weights to small random numbers.
Until satisfied, Do

- For each training example, Do

    1. Input the training example to the network and compute the network outputs

    2. For each output unit $k$

    $$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

    3. For each hidden unit $h$

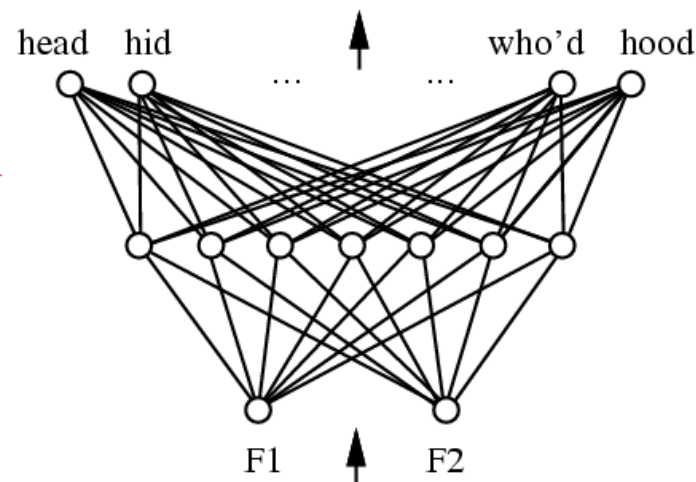    $$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in outputs} w_{h,k}\delta_k$$

    4. Update each network weight $w_{i,j}$

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

    where

    $$\Delta w_{i,j} = -\eta \frac{\partial E}{\partial w_{ij}}$$
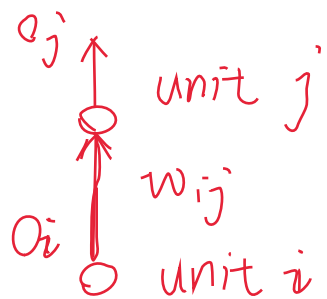
    $$\Delta w_{i,j} = \eta \delta_j x_i$$



head    hid    ...    ↑    ...    who'd    hood

F1    ↑    F2

$x_d$ = input

$t_d$ = target output

$o_d$ = observed unit output

$w_{ij}$ = wt from i to j

# Backpropagation

$O_j$

unit $j$

$w_{ij}$

$O_i$ unit $i$

$net_j = \sum_i x_{ij} w_{ij} = \sum_i O_i w_{ij}$

$x_{ij}$: input from unit $i$ to unit $j$

$(O_i = x_{ij})$

① $w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$

$$= -\eta \frac{\partial E_d}{\partial w_{ij}}$$

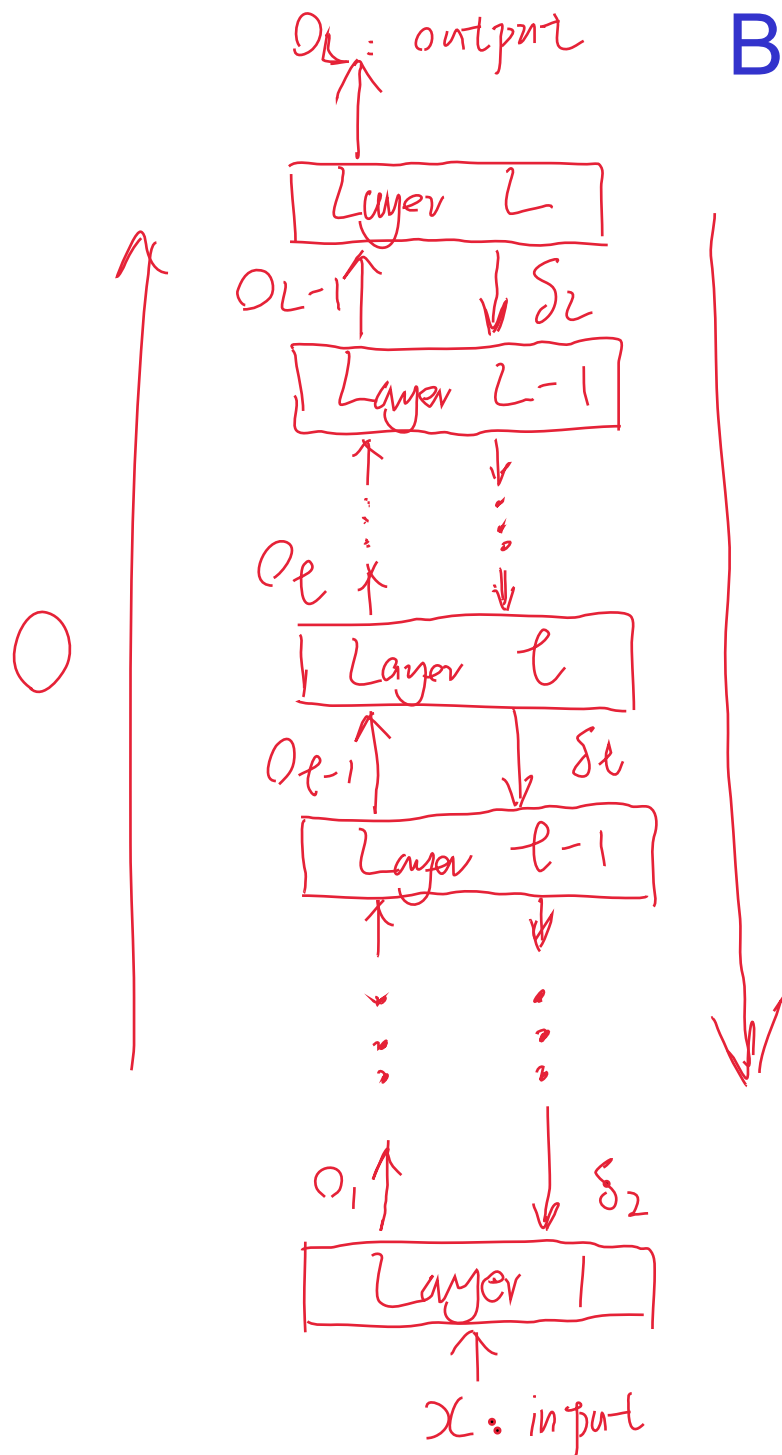② $E_d = \frac{1}{2} \sum_k (t_k - O_k)^2$

($K$: #outputs)

③ $\frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ij}}$

$x_{ij} = O_i$

$\Delta w_{ij} = -\eta \frac{\partial E_d}{\partial w_{ij}}$

$$= \eta \delta_j x_{ij} , \quad \left( \delta_j = - \frac{\partial E_d}{\partial net_j} \right)$$

• Case 1: unit $j \in$ output layer.

Output $O_j$ unit $j

$w_{ij}$

Hidden $\circ$ $\circ$ - $\circ$ $\cdots$ $\circ$

unit $i$

$$\delta_j = -\frac{\partial E_d}{\partial net_j} = -\sum_{j=1}^{K} \frac{\partial E_d}{\partial O_j} \cdot \frac{\partial O_j}{\partial net_j}$$

$$= \sum_{j=1}^{K} (t_j - O_j) O_j (1 - O_j)$$

• Case 2: unit $j \in$ Hidden Layer

unit $\ell$

Output $\bigcirc$ $\bigcirc$ $\cdots$ $\bigcirc$ $\cdots$ $\bigcirc$

Hidden: unit $j$

$w_{ij}$

Input $\bigcirc$ $\bigcirc$ $\cdots$ $\bigcirc$ - $\bigcirc$

unit $i$

$$\delta_j = -\frac{\partial E_d}{\partial net_j}$$

$$= -\sum_{\ell=1}^{K} \frac{\partial E_d}{\partial net_\ell} \cdot \frac{\partial net_\ell}{\partial O_j} \cdot \frac{\partial O_j}{\partial net_j}$$

$$= \sum_{\ell=1}^{K} \delta_\ell w_{j\ell} O_j (1 - O_j)$$

$$= O_j (1 - O_j) \sum_{\ell=1}^{K} \delta_\ell w_{j\ell}$$

# Backpropagation

$$E = \frac{1}{2}(t - o)^2$$

$O_L$ : output



**Diagram (left):**

Layer $L$

$O_{L-1}$ ↑   ↓ $\delta_L$

Layer $L-1$

$O_\ell$ ↑   ↓

Layer $\ell$

$O_{\ell-1}$ ↑   ↓ $\delta_\ell$

Layer $\ell-1$

$\vdots$

$O_1$ ↑   ↓ $\delta_2$

Layer $1$

$x$ : input

$O$ (left bracket)   $\delta$ (right bracket)

**Algorithm (right):**

1. Initialization $W^{(0)}$

2. Repeat

3. $\quad W_{ij} \leftarrow W_{ij} - \eta \cdot \delta_j \cdot O_i$

where $\delta_j = \begin{cases} (t_j - o_j) o_j (1 - o_j), & j \in \text{output} \\ o_j(1 - o_j) \sum\limits_{\ell=1}^{K} \delta_\ell W_{j\ell}, & j \in \text{hidden} \end{cases}$

4. Until convergence

$$\frac{|Q^{(t+1)} - Q^{(t)}|}{|Q^{(t)}|} < \varepsilon = 10^{-5}$$

# More on Backpropagation

- Gradient descent over entire *network* weight vector

- Easily generalized to arbitrary directed graphs

- Will find a local, not necessarily global error minimum

  - In practice, often works well (can run multiple times)

- Often include weight *momentum* $\alpha$ $\in (0,1)$

$$\Delta w_{i,j}(n) = \eta \delta_j x_{i,j} + \alpha \Delta w_{i,j}(n-1)$$

- Minimizes error over *training* examples

  - Will it generalize well to subsequent examples?

- Training can take thousands of iterations $\rightarrow$ slow!

- Using network after training is very fast

*Handwritten annotations:*

Local     Global

Local minimum.
1. momentum
2. SGD
3. different initializations

# Overfitting in ANNs



Error versus weight updates (example 1)

Training set error ◆
Validation set error +

9000

Error versus weight updates (example 2)

Training set error ◆
Validation set error +

900

overfitting

① weight decaying

② validation
(k-fold CV)

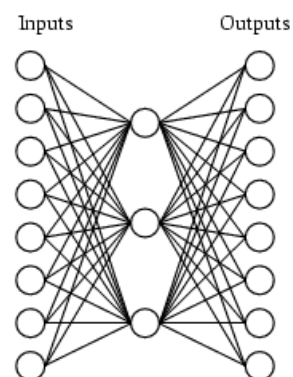# Expressive Capabilities of ANNs

$$y \leftarrow f(x)$$

Boolean functions:

- Every boolean function can be represented by network with single hidden layer

- but might require exponential (in number of inputs) hidden units

Continuous functions:

- Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer [Cybenko 1989; Hornik et al. 1989]

- Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988].

# Learning Hidden Layer Representations



Inputs          Outputs

A target function:

| Input | | Output |
|-----------|-----------|-----------|
| 10000000 | → | 10000000 |
| 01000000 | → | 01000000 |
| 00100000 | → | 00100000 |
| 00010000 | → | 00010000 |
| 00001000 | → | 00001000 |
| 00000100 | → | 00000100 |
| 00000010 | → | 00000010 |
| 00000001 | → | 00000001 |

Can this be learned??

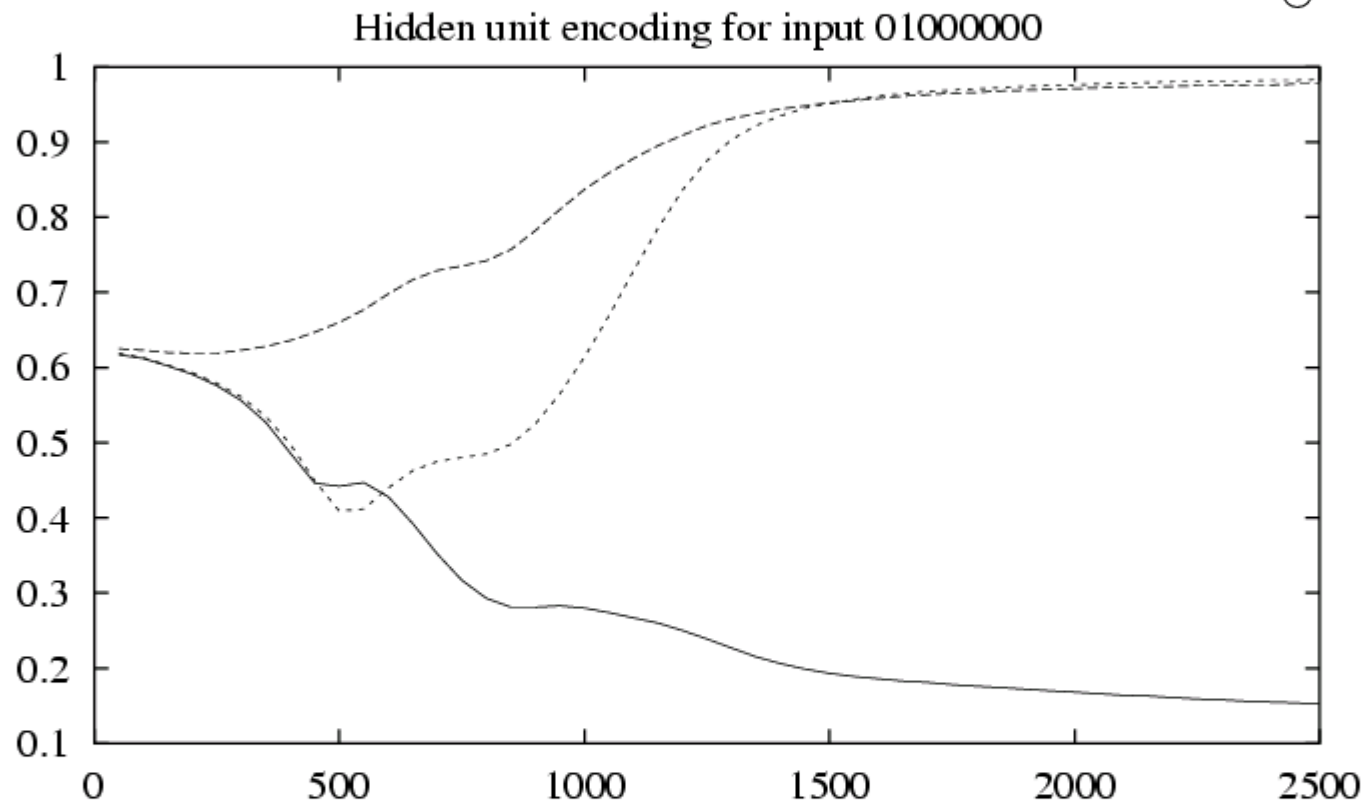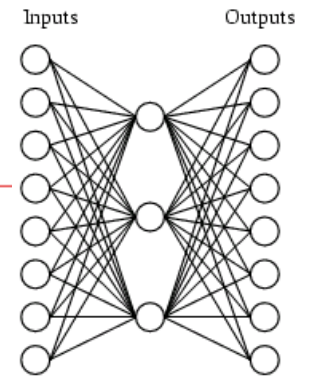# Learning Hidden Layer Representations

A network:



Learned hidden layer representation:

| Input | | Hidden Values | | | Output |
|---|---|---|---|---|---|
| 10000000 | → | .89 | .04 | .08 | → 10000000 |
| 01000000 | → | .01 | .11 | .88 | → 01000000 |
| 00100000 | → | .01 | .97 | .27 | → 00100000 |
| 00010000 | → | .99 | .97 | .71 | → 00010000 |
| 00001000 | → | .03 | .05 | .02 | → 00001000 |
| 00000100 | → | .22 | .99 | .99 | → 00000100 |
| 00000010 | → | .80 | .01 | .98 | → 00000010 |
| 00000001 | → | .60 | .94 | .01 | → 00000001 |

# Training

Sum of squared errors for each output unit

# Training



Hidden unit encoding for input 01000000

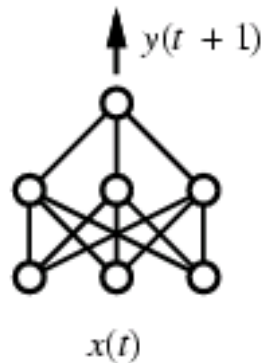# Training

Weights from inputs to one hidden unit

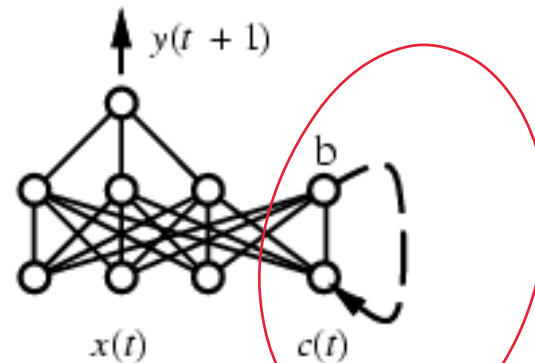# Training Networks on Time Series

- Suppose we want to predict next state of world
  - and it depends on history of unknown length
  - e.g., robot with forward-facing sensors trying to predict next sensor reading as it moves and turns

# Recurrent Networks: Time Series

- Suppose we want to predict next state of world
  - and it depends on history of unknown length
  - e.g., robot with forward-facing sensors trying to predict next sensor reading as it moves and turns

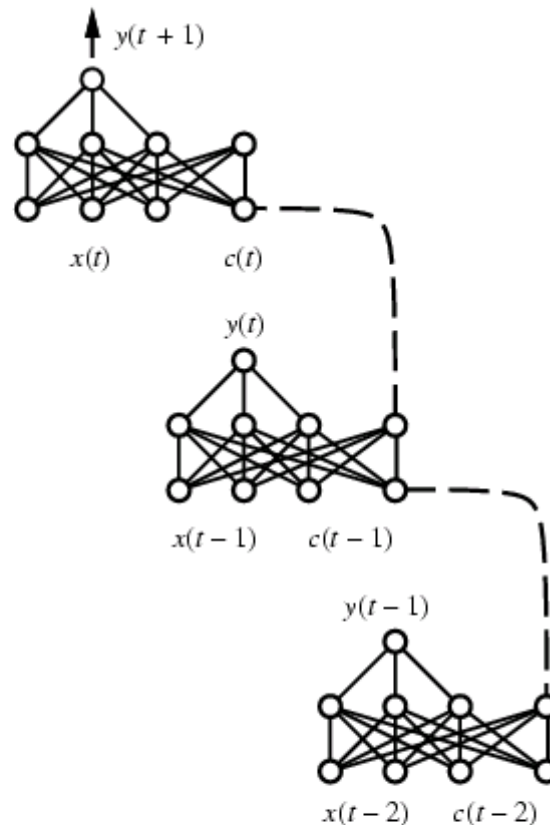- Idea: use hidden layer in network to capture state history
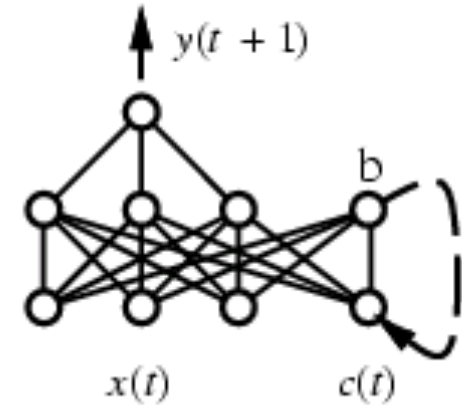


(a) Feedforward network     (b) Recurrent network

# Recurrent Networks on Time Series
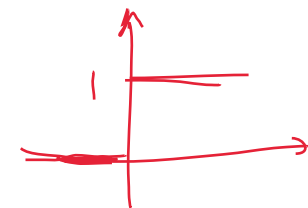
How can we train recurrent net??



$y(t + 1)$

$x(t)$  $c(t)$

( parameter sharing )

$y(t + 1)$

$x(t)$  $c(t)$

$y(t)$

$x(t - 1)$  $c(t - 1)$

$y(t - 1)$

$x(t - 2)$  $c(t - 2)$

(c) Recurrent network
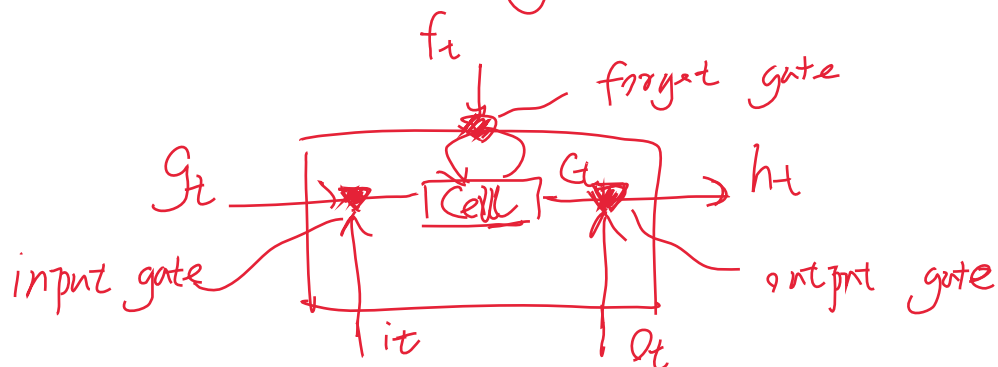unfolded in time

$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$

$ReLU$

# Recurrent Networks on Time Series

Gradient vanishing / exploding

$$\frac{\partial C_{t'}}{\partial C_t} = \prod_{k=1}^{t'-t} f_{t+k}$$

* LSTM (Long-Short Term Memory)



$f_t$ forget gate

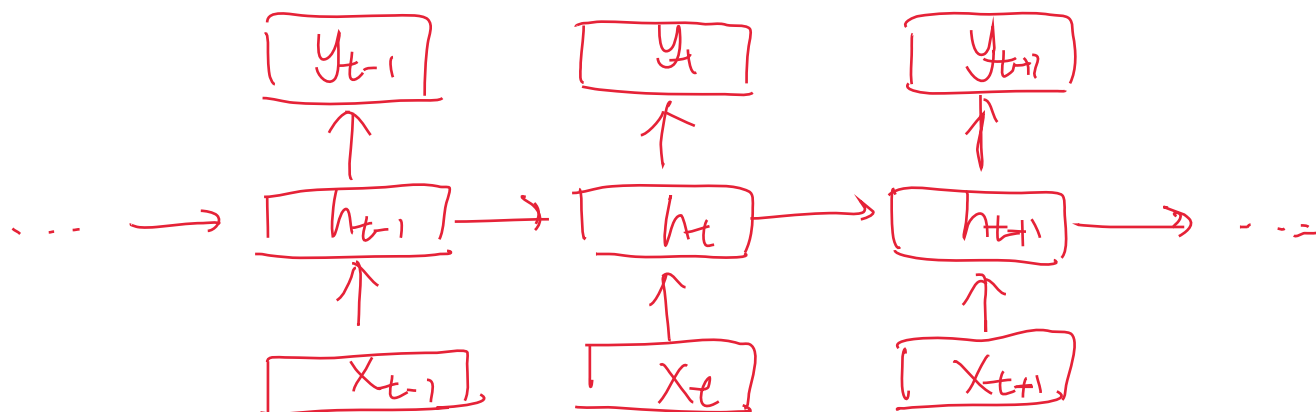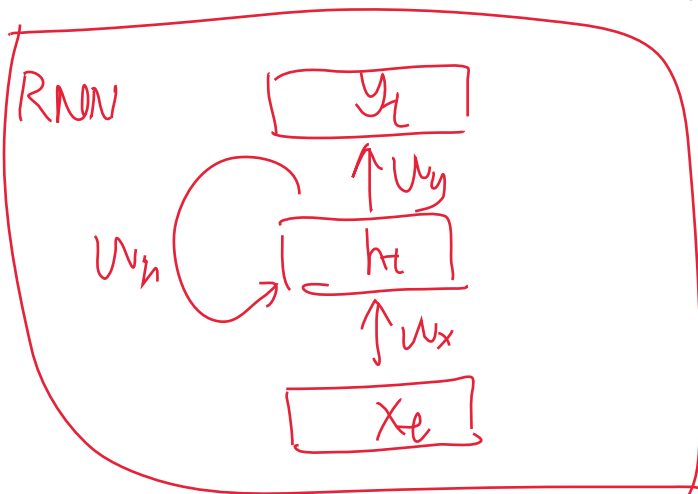$g_t$ input gate

Cell $C_t$ → $h_t$

output gate

$it$ $o_t$

$$\begin{cases} C_t = f_t \odot C_{t-1} + i_t \odot g_t \\ h_t = o_t \odot \phi(C_t) \end{cases}$$
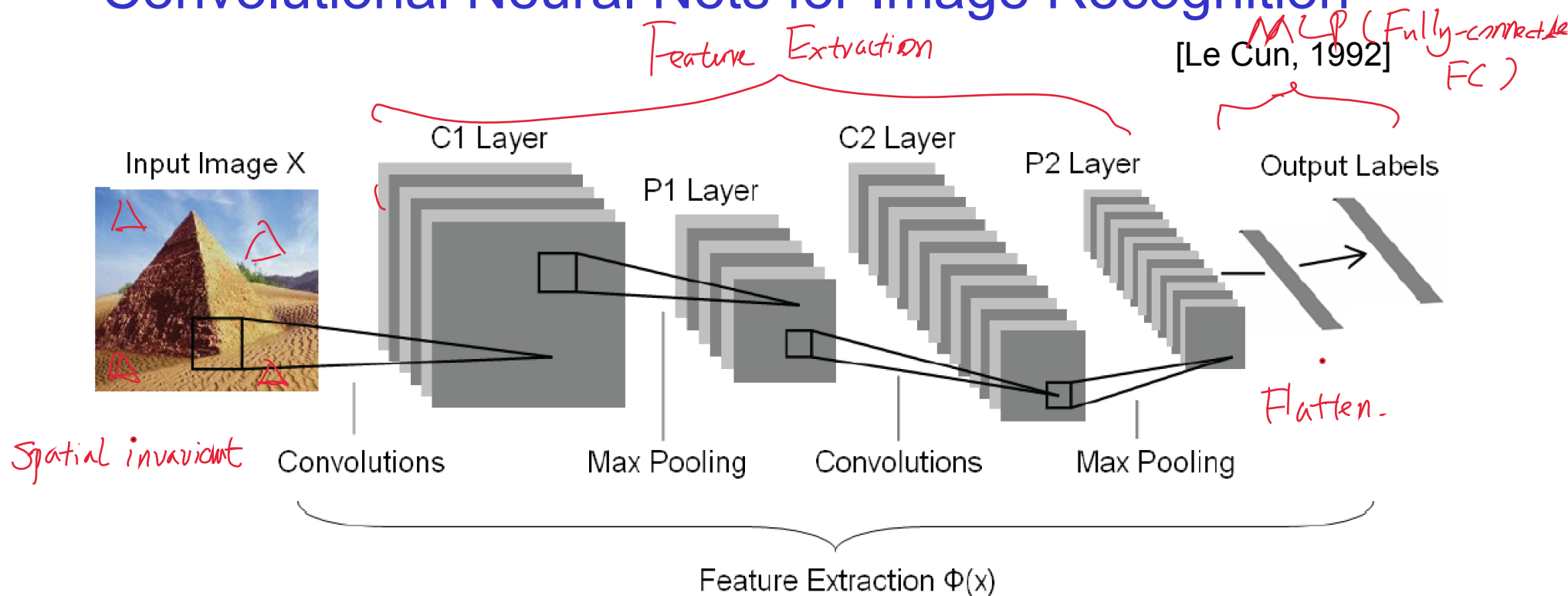
$\odot$ : element-wise product.

* GRU + Attention

(Gate Recurrent Unit)

$$\begin{bmatrix} a \\ b \\ c \end{bmatrix} \odot \begin{bmatrix} d \\ e \\ f \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \rightarrow \begin{matrix} ad \\ be \\ cf \end{matrix}$$
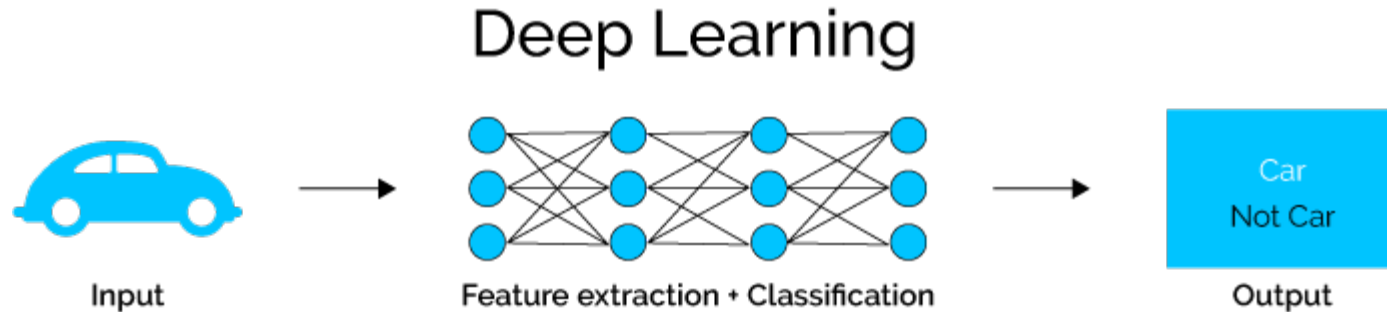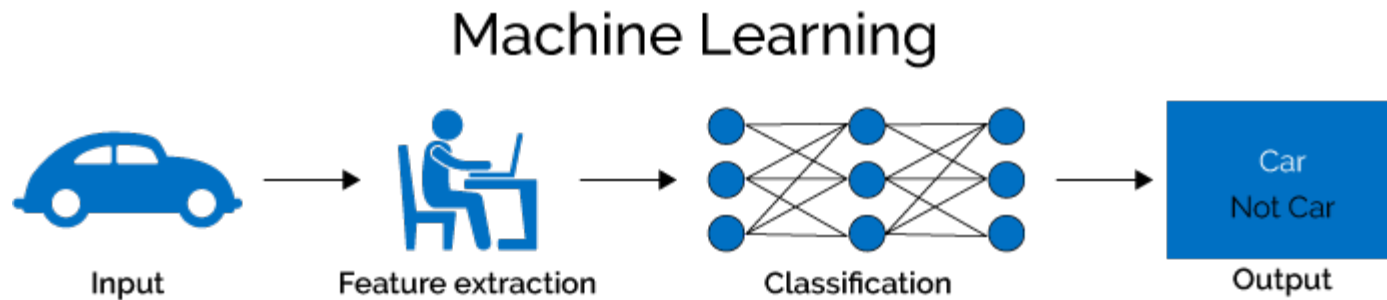
RNN



$y_t$

$\uparrow U_y$

$U_h$ → $h_t$

$\uparrow U_x$

$x_t$

$\cdots \longrightarrow$

$y_{t-1}$ → $h_{t-1}$ → $x_{t-1}$

$y_t$ → $h_t$ → $x_t$

$y_{t+1}$ → $h_{t+1}$ → $x_{t+1}$ $\longrightarrow \cdots$

# Convolutional Neural Nets for Image Recognition



Feature Extraction

[Le Cun, 1992]

MLP (Fully-connected FC)

Input Image X

C1 Layer

P1 Layer

C2 Layer

P2 Layer

Output Labels

Spatial invariant

Convolutions   Max Pooling   Convolutions   Max Pooling

Flatten.

Feature Extraction Φ(x)

- specialized architecture: mix different types of units, not completely connected, motivated by primate visual cortex

- many shared parameters, stochastic gradient training

- very successful!  now many specialized architectures for vision, speech, translation, …

# CNNs - Convolution Layer(s)

In CNNs and deep learning in general, the **features are learned** rather than manually selected during the data preprocessing phase.
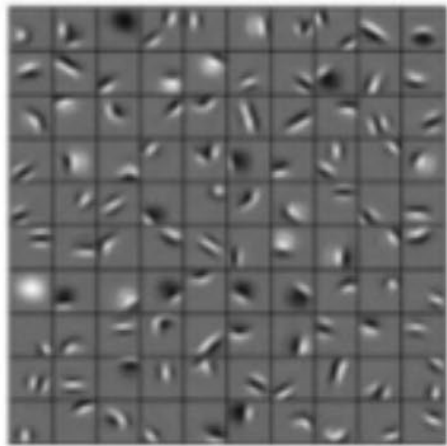
# CNNs - Convolution Layer(s)

There can be multiple convolution layers where earlier layers captures low-level features such as edges or colours while added layers capture high-level features such as facial parts.

Below are examples of **feature maps** created from a convolution layer.
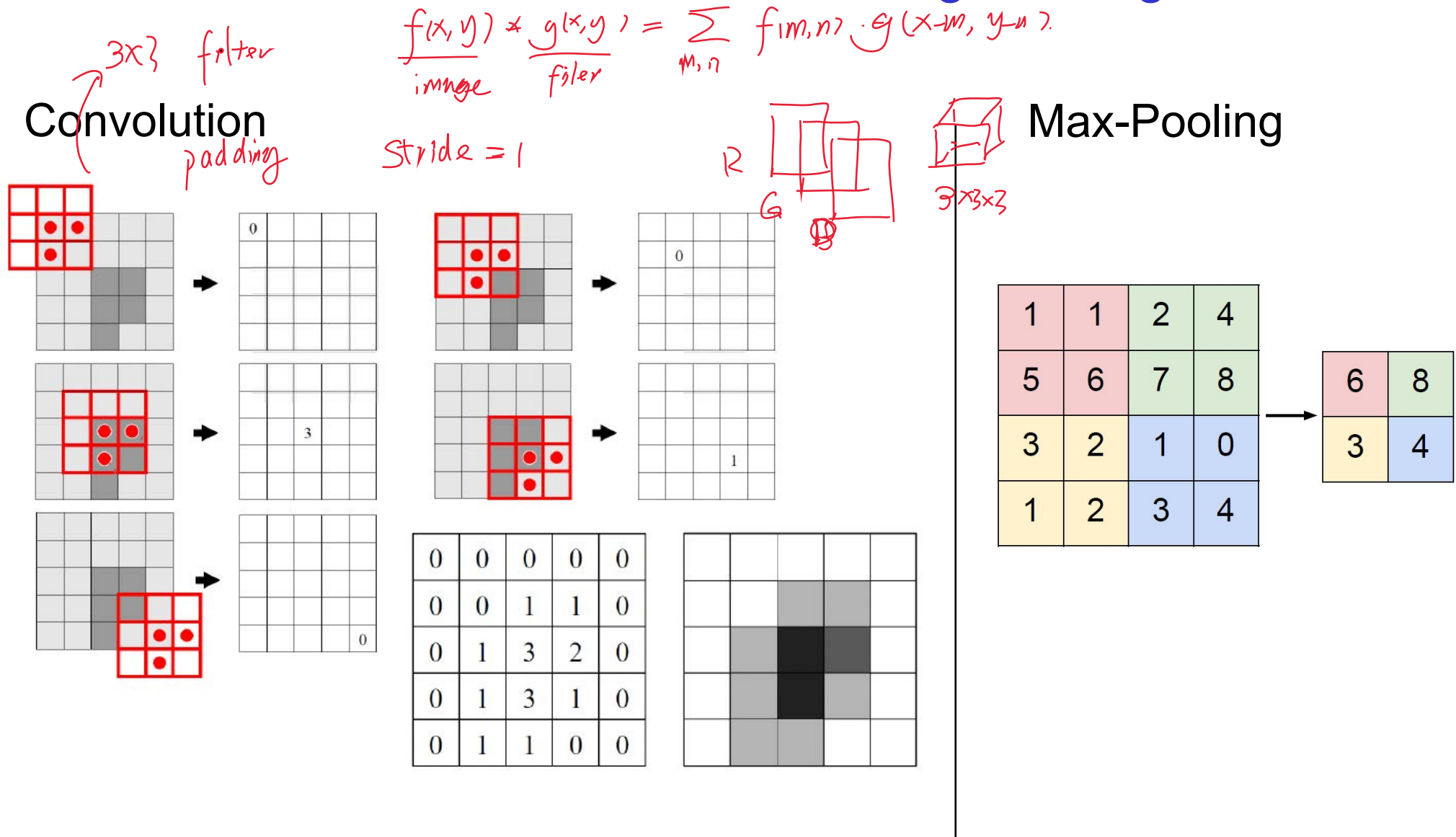
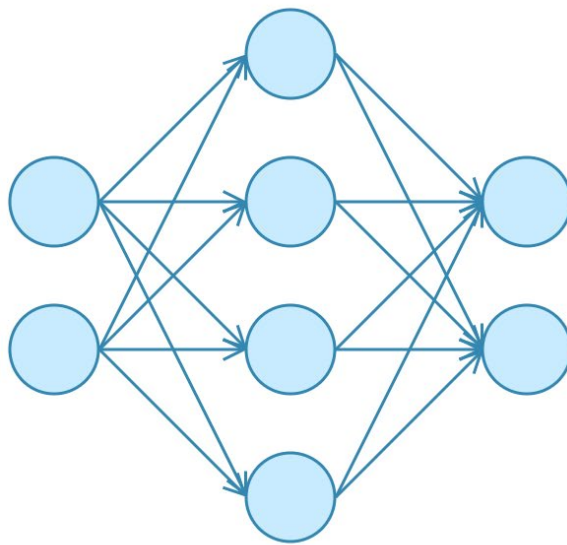Low-level feature     Mid-level feature     High-level feature

# Convolutional Neural Nets for Image Recognition

$$f(x,y) * g(x,y) = \sum_{m,n} f(m,n) \cdot g(x-m, y-n)$$

image — $f(x,y)$, filter — $g(x,y)$

3x3 filter

## Convolution

padding

Stride = 1

R

G

B

3x3x3

## Max-Pooling



| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

| 6 | 8 |
|---|---|
| 3 | 4 |

| 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 3 | 2 | 0 |
| 0 | 1 | 3 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

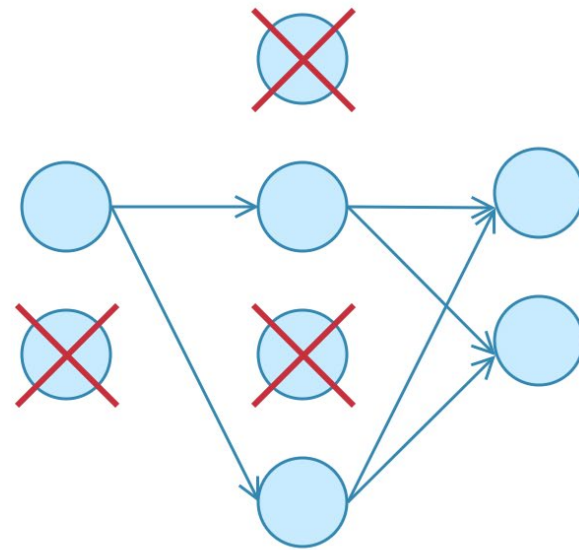http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html

# Dropout - Regularization Technique

This is basically the idea of dropout - **disabling neurons with probability p** so that the network isn't dependent on one node.

Dropout can be applied to input or hidden layers, but not output.



No Dropout                    With Dropout

# Artificial Neural Networks: Summary

- Highly non-linear regression/classification
- Hidden layers learn intermediate representations
- Potentially millions of parameters to estimate
- Stochastic gradient descent, local minima problems

- Deep networks have produced real progress in many fields
  - computer vision
  - speech recognition
  - mapping images to text
  - recommender systems
  - …
- They learn very useful non-linear representations