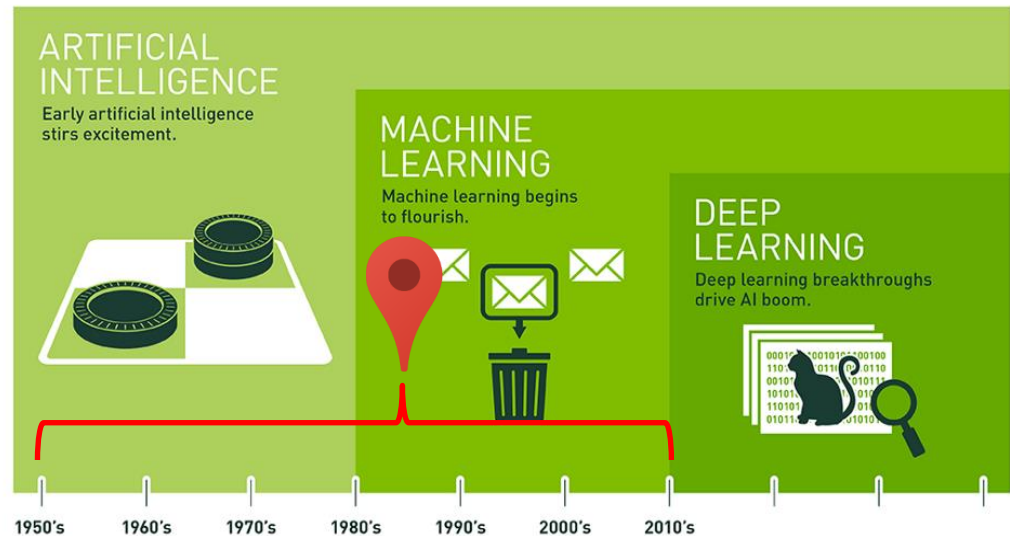# Artificial Intelligence: Introduction to Neural Networks

Perceptron, Backpropagation

# Today

- **Neural Networks**  YOU ARE HERE!
  - Perceptrons
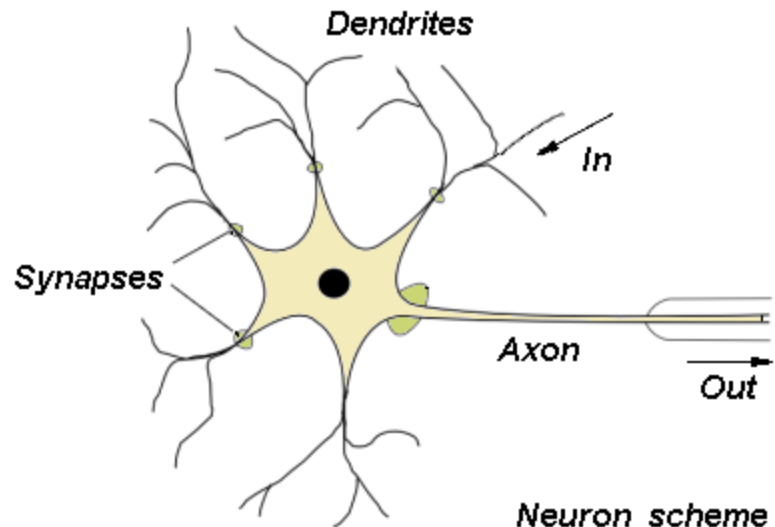  - Backpropagation

# Neural Networks

- Radically different approach to reasoning and learning
- Inspired by biology
  - the neurons in the human brain
- Set of many simple processing units (neurons) connected together
- Behavior of each neuron is very simple
  - but a collection of neurons can have sophisticated behavior and can be used for complex tasks
- In a neural network, the behavior depends on weights on the connection between the neurons
- The weights will be learned given training data

# Biological Neurons

- **Human brain =**
  - 100 billion neurons
  - each neuron may be connected to 10,000 other neurons
  - passing signals to each other via 1,000 trillion **synapses**

- **A neuron is made of:**
  - Dendrites: filaments that provide input to the neuron
  - Axon: sends an output signal
  - Synapses: connection with other neurons – releases neurotransmitters to other neurons
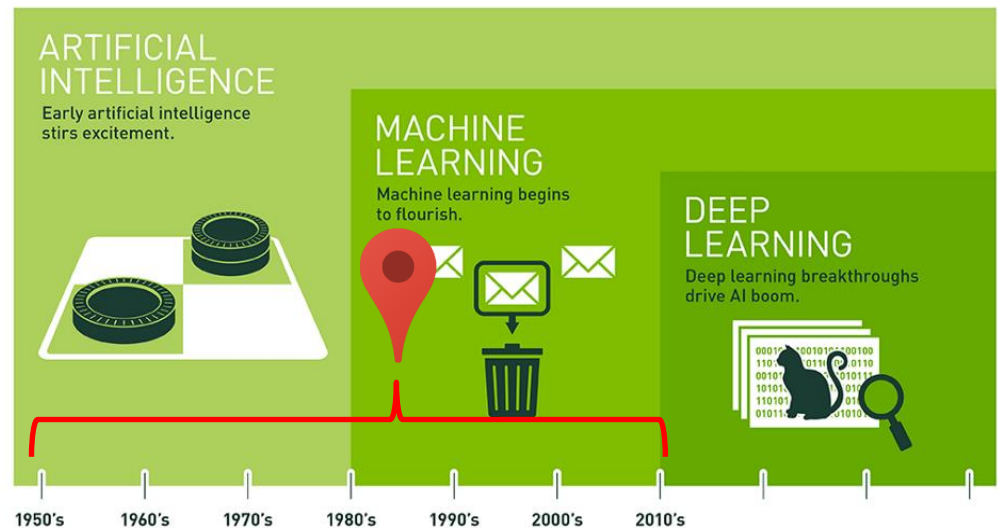
Dendrites

In

Synapses

Axon

Out

Neuron scheme

# Behavior of a Neuron

- A neuron receives inputs from its neighbors
- If enough inputs are received at the same time:
  - the neuron is activated
  - and fires an output to its neighbors
- Repeated firings across a synapse increases its sensitivity and the future likelihood of its firing
- If a particular stimulus repeatedly causes activity in a group of neurons, they become strongly associated

# Today

- ## Neural Netwo
  - Perceptrons
  - Backpropagation
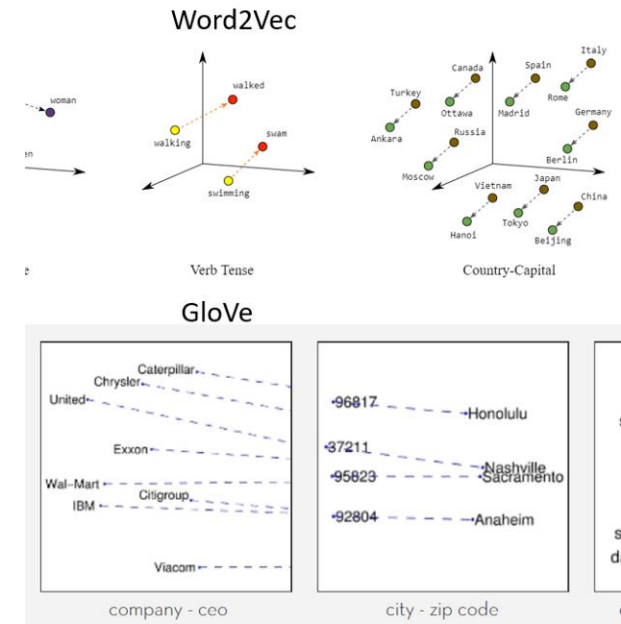
**YOU ARE HERE!**

# Feature Vector Representation

- Sources of Feature Vector x
  - Encoded image
  - Tabulated data
  - Embedded words
  - ...

**Tabular Data**

columns = attributes for those observations

| Player | Minutes | Points | Rebounds | Assists |
|--------|---------|--------|----------|---------|
| A | 41 | 20 | 6 | 5 |
| B | 30 | 29 | 7 | 6 |
| C | 22 | 7 | 7 | 2 |
| D | 26 | 3 | 3 | 9 |
| E | 20 | 19 | 8 | 0 |
| F | 9 | 6 | 14 | 14 |
| G | 14 | 22 | 8 | 3 |
| I | 22 | 36 | 0 | 9 |
| J | 34 | 8 | 1 | 3 |

Word2Vec

Verb Tense

Country-Capital

GloVe

company - ceo

city - zip code

source: Luger (2005)

# Feature Vector Representation

- Sources of Feature Vector $x$
  - Encoded image
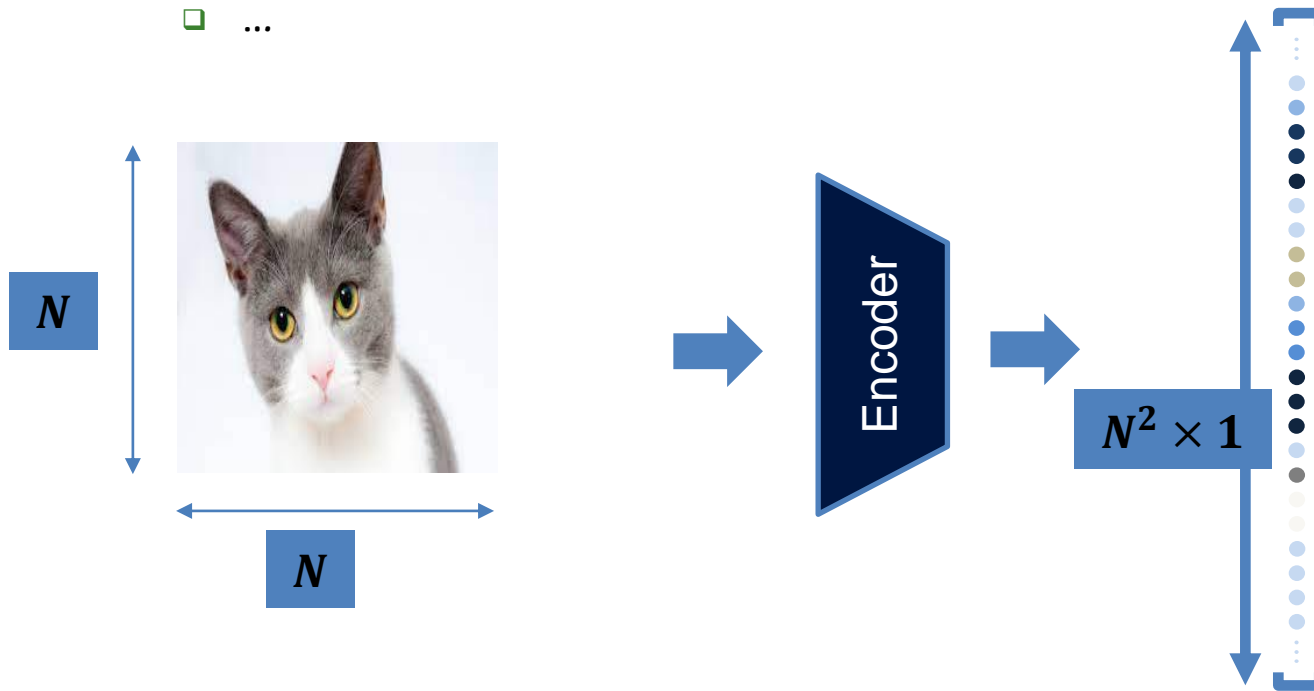  - Tabulated data
  - Embedded words
  - ...



$N$

$N$

Encoder

$N^2 \times 1$

source: Luger (2005)

# Feature Vector Representation

- Sources of Feature Vector $x$
  - Encoded image
  - Tabulated data
  - Embedded words
  - ...



$N^2 \times 1$

source: Luger (2005)

# Feature Vector Representation

- Sources of Feature Vector $x$
    - Encoded image
    - Tabulated data
    - Embedded words
    - ...



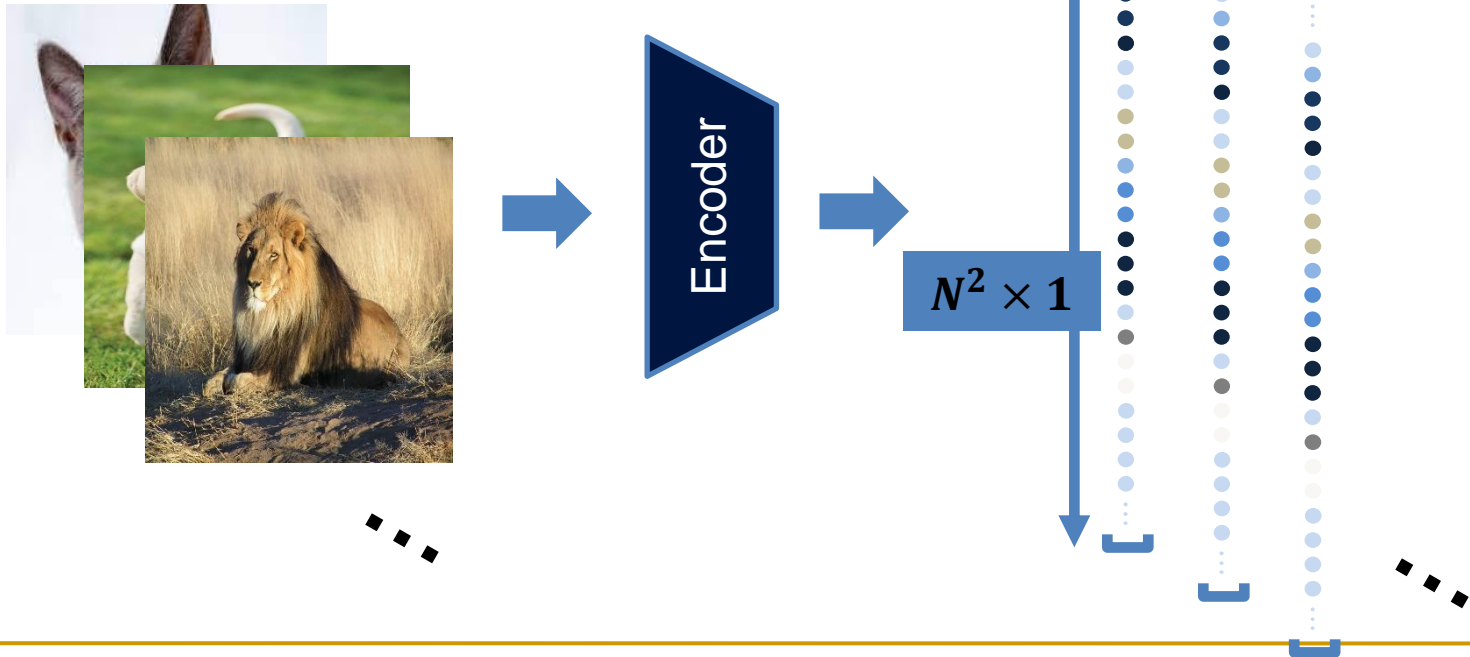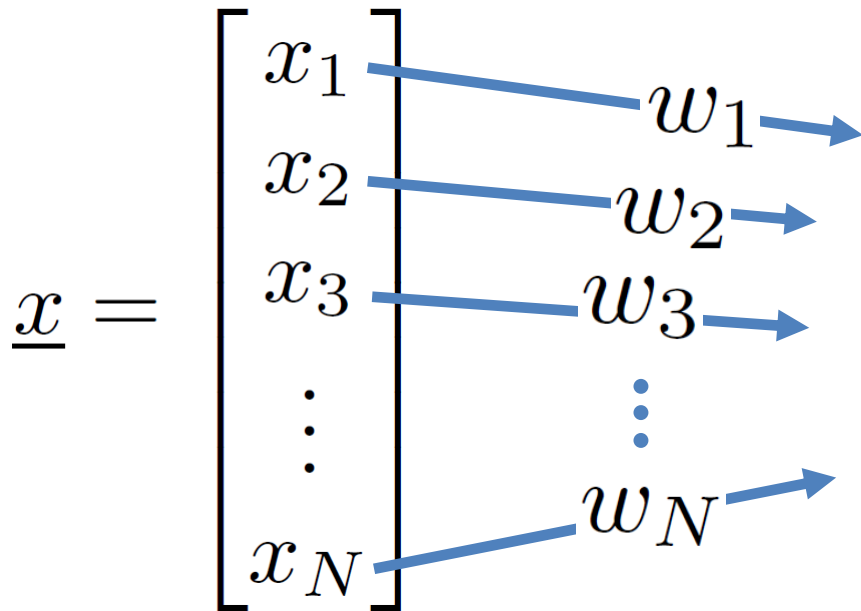$N^2 \times 1$

source: Luger (2005)

# A Perceptron Network

- **Goal:** Map *Input* Feature Vector **x** into *Output* Feature Vector **y**

$$\sum_{i=1}^{N} w_i x_i$$

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$

$w_1 \rightarrow$

$w_2 \rightarrow$

$w_3 \rightarrow$

$\vdots$

$w_N \rightarrow$

# A Perceptron Network

- **Goal:** Map *Input* Feature Vector **x** into *Output* Feature Vector **y**

$$\sum_{i=1}^{N} w_i x_i + b_j$$

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$

$w_1$

$w_2$

$w_3$

$\vdots$

$w_N$

$$\sum (\cdot) + b_j$$

# A Perceptron <u>Network</u>

- **Goal:** Map *Input* Feature Vector *x* into *Output* Feature Vector *y*

$$y_j = f\left(\sum_{i=1}^{N} w_i x_i + b_j\right)$$

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$

$w_1$

$w_2$

$w_3$

$w_N$

$$\sum(\cdot) + b_j \quad f$$

$y_j$

# A Perceptron <u>Network</u>

- **Goal:** Map *Input* Feature Vector **x** into *Output* Feature Vector **y**

$$y_j = f\left(\sum_{i=1}^{N} w_i x_i + b_j\right) = f\left(w^T x + b_j\right)$$

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} \quad \begin{array}{c} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_N \end{array} \quad \boxed{\sum(\cdot) + b_j} \,\Big|\, f \longrightarrow y_j$$

# A Perceptron <u>Network</u>

- **Goal:** Map *Input* Feature Vector ***x*** into *Output* Feature Vector ***y***

$$y_j = f\left(\sum_{i=1}^{N} w_i x_i + b_j\right) = f$$

**Activation functions**

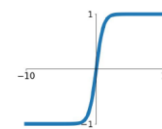**Sigmoid**
$\sigma(x) = \frac{1}{1+e^{-x}}$

**Leaky ReLU**
$\max(0.1x, x)$

**tanh**
$\tanh(x)$

**Maxout**
$\max(w_1^T x + b_1, w_2^T x + b_2)$

**ReLU**
$\max(0, x)$

**ELU**
$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$

$w_1$ $w_2$ $w_3$ $w_N$

$\sum(\cdot) + b_j$ $f$

$y_j$

# A Perceptron <u>Network</u>

- **Goal:** Map *Input* ~~Feature Vector~~ ure Vector **y**

$$y_j = f\left(\sum_{i=1}^{N}\right)$$



j'th Cell

$$\underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$

$w_1$
$w_2$
$w_3$
$\vdots$
$w_N$

$$\sum (\cdot) + b_j \quad f$$

$y_j$

?

# A Perceptron Network

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} \qquad \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_M \end{bmatrix}$$

# A Perceptron Network

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$



$$\begin{bmatrix} y_1 \\ y_2 \\ \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ \\ \\ \end{bmatrix} = f\left( \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & \cdots & w_{1,N} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} + \begin{bmatrix} b_1 \\ \\ \\ \end{bmatrix} \right)$$

# A Perceptron Network



$$
\begin{bmatrix} y_1 \\ y_2 \\ \\ \\ \\ \end{bmatrix} = f \left( \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & \cdots & w_{1,N} \\ w_{2,1} & w_{2,2} & w_{2,3} & \cdots & w_{2,N} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \\ \\ \end{bmatrix} \right)
$$

# A Perceptron Network

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_M \end{bmatrix}$$

# A Perceptron <u>Network</u>



$$\begin{bmatrix} y_1 \\ y_2 \\ \boxed{y_3} \end{bmatrix} = f\left( \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & \cdots & w_{1,N} \\ w_{2,1} & w_{2,2} & w_{2,3} & \cdots & w_{2,N} \\ \boxed{w_{3,1} \quad w_{3,2} \quad w_{3,3} \quad \cdots \quad w_{3,N}} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \boxed{b_3} \end{bmatrix} \right)$$

# A Perceptron Network

# A Perceptron <u>Network</u>

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_M \end{bmatrix} = f \left( \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & \cdots & w_{1,N} \\ w_{2,1} & w_{2,2} & w_{2,3} & \cdots & w_{2,N} \\ w_{3,1} & w_{3,2} & w_{3,3} & \cdots & w_{3,N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{M,1} & w_{M,2} & w_{M,3} & \cdots & w_{M,N} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \vdots \\ b_M \end{bmatrix} \right)$$

# A Perceptron Network



$$\underline{y} = f\left(W^T \underline{x} + \underline{b}\right)$$

$\underline{x}$: input feature vector $N \times 1$
$\underline{y}$: input feature vector $M \times 1$
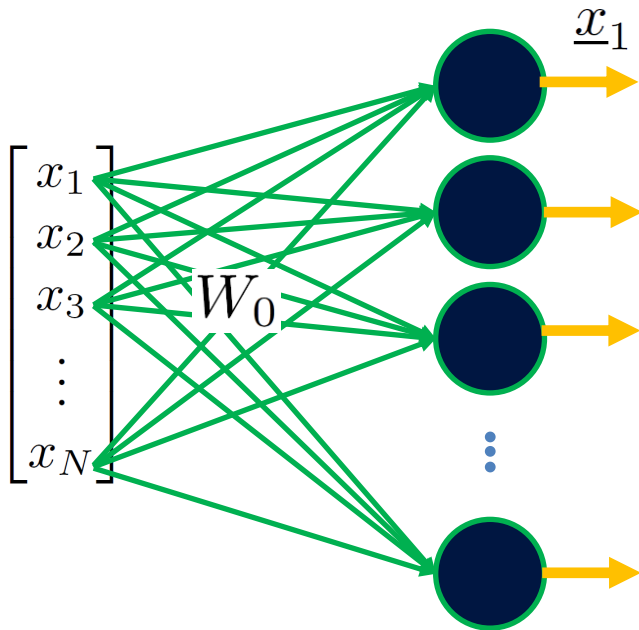$\overline{W}$: weight matrix $M \times N$
$\underline{b}$: bias vector $M \times 1$
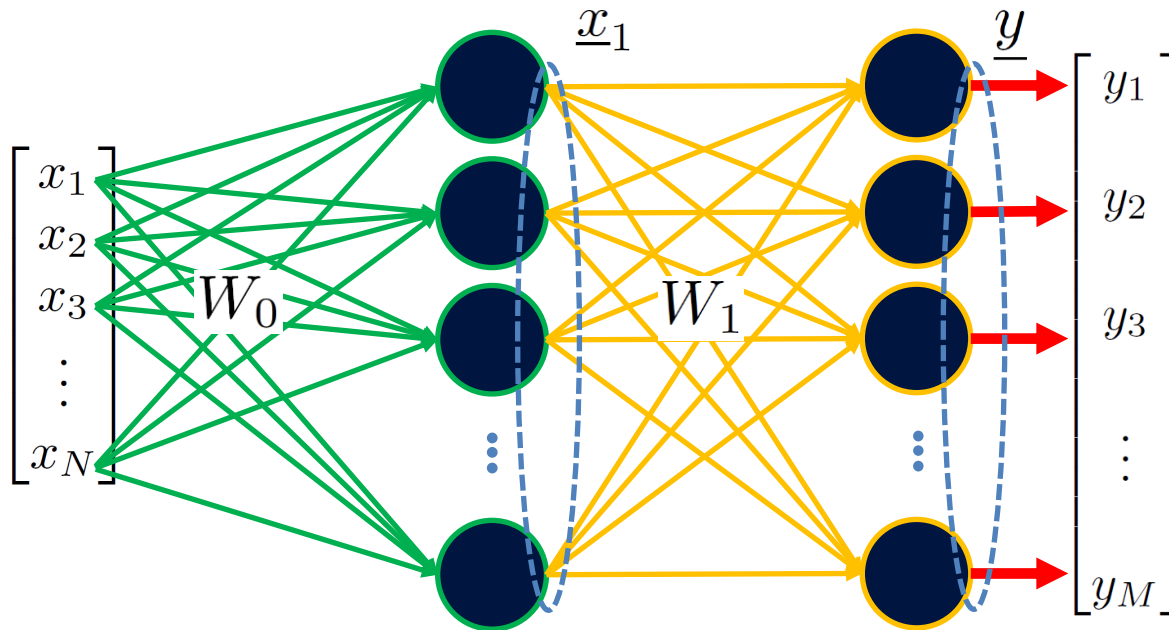$f\left(\cdot\right)$: activation function

# A Perceptron Network

$$\underline{x}_1 = f\left(W_0^T \underline{x} + \underline{b}_0\right)$$

# A Perceptron <u>Network</u>
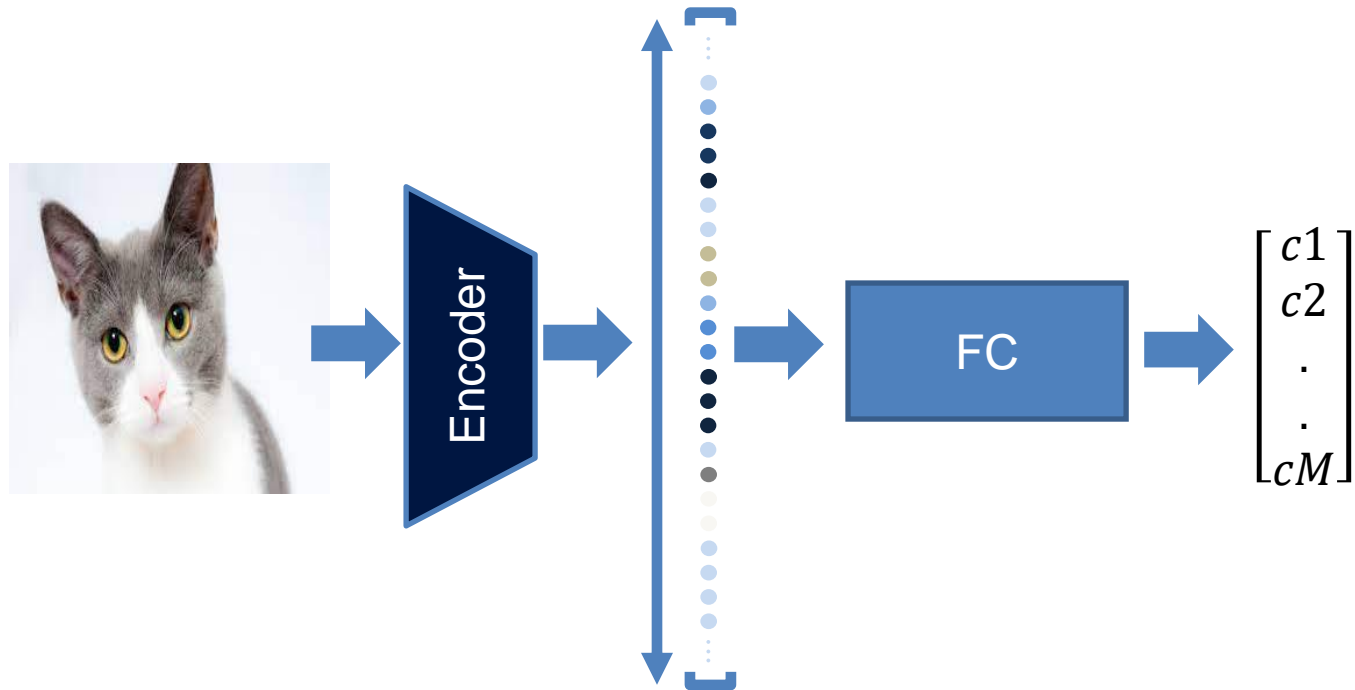


$$\underline{x}_1 = f\left(W_0^T \underline{x} + \underline{b}_0\right)$$

$$\underline{y} = f\left(W_1^T \underline{x}_1 + \underline{b}_1\right)$$

# Fully Connected (FC) Network



$$\begin{bmatrix} c1 \\ c2 \\ . \\ . \\ cM \end{bmatrix}$$

# Neural Network: FC Layer Design

- First layer design: map input feature vector into smaller vector

$[x]_{784 \times 1}$

$\left[W_0^T\right]_{784 \times 128}$

784

# Neural Network: FC Layer Design

- First layer design: activate features using non-linear activation

$$[\underline{x}_1]_{128 \times 1} = f\left(W_0^T \underline{x} + \underline{b}_0\right)$$



$[\underline{x}]_{784 \times 1}$

$[\underline{x}_1]_{128 \times 1}$

$\left[W_0^T\right]_{784 \times 128}$

784

128

# Neural Network: FC Layer Design

- Second layer design: map input feature vector into #Classes

$$[\underline{x}_1]_{128 \times 1} = f\left(W_0^T \underline{x} + \underline{b}_0\right)$$

$[\underline{x}]_{784 \times 1}$

$[\underline{x}_1]_{128 \times 1}$

784

$\left[W_0^T\right]_{784 \times 128}$

128

$\left[W_1^T\right]_{128 \times 10}$

# Neural Network: FC Layer Design

- Second layer design: no activation is required
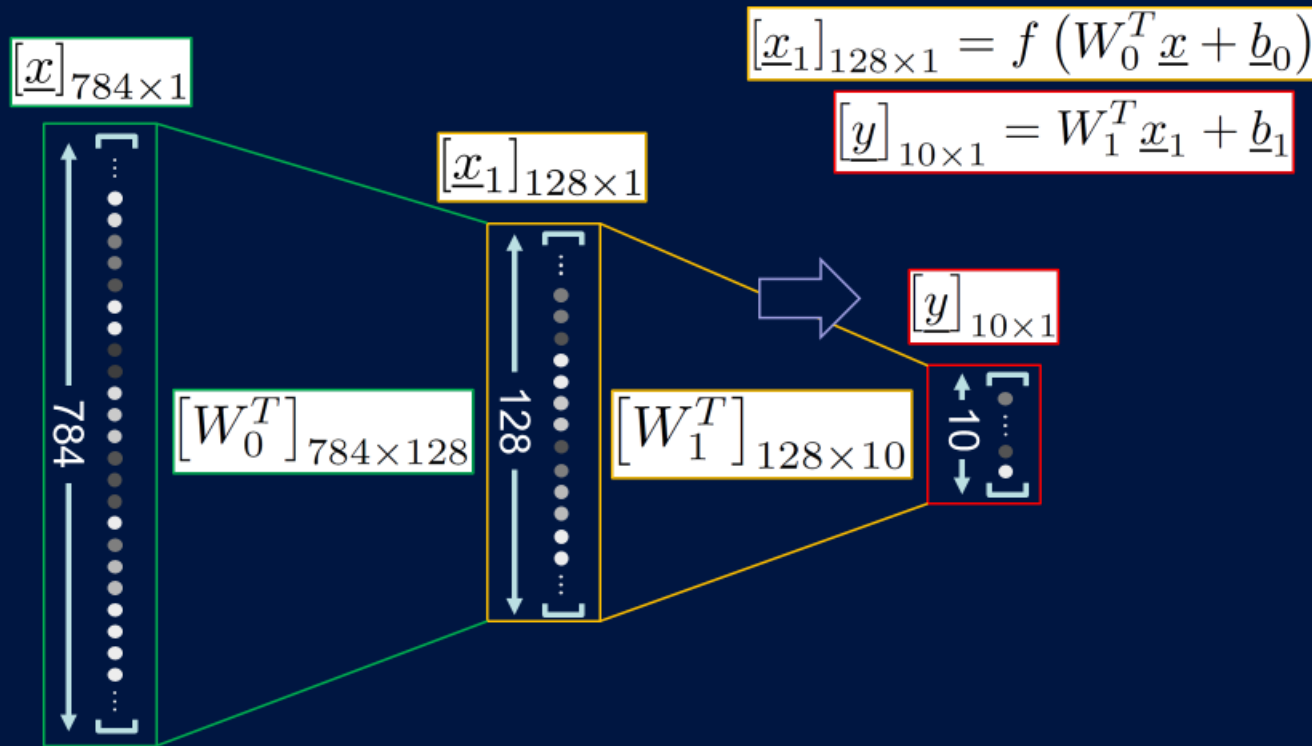


$$[\underline{x}_1]_{128\times 1} = f\left(W_0^T \underline{x} + \underline{b}_0\right)$$

$$[\underline{y}]_{10\times 1} = W_1^T \underline{x}_1 + \underline{b}_1$$

$$[\underline{x}]_{784\times 1}$$

$$[\underline{x}_1]_{128\times 1}$$

$$[\underline{y}]_{10\times 1}$$

$$\left[W_0^T\right]_{784\times 128}$$

$$\left[W_1^T\right]_{128\times 10}$$

784

128

10

# Neural Network: FC Layer Design

Quiz: What is the #Learnable Parameters?

$$[\underline{x}_1]_{128 \times 1} = f\left(W_0^T \underline{x} + \underline{b}_0\right)$$

$$[\underline{y}]_{10 \times 1} = W_1^T \underline{x}_1 + \underline{b}_1$$

$[\underline{x}]_{784 \times 1}$

$[\underline{x}_1]_{128 \times 1}$

$[\underline{y}]_{10 \times 1}$

784

$\left[W_0^T\right]_{784 \times 128}$

128

$\left[W_1^T\right]_{128 \times 10}$

10

# Neural Network: FC Layer Design

#Learnable Parameters = 784x128 + 128 + 128x10 + 10 = 101,770

$$[\underline{x}_1]_{128 \times 1} = f\left(W_0^T \underline{x} + \underline{b}_0\right)$$

$$[\underline{y}]_{10 \times 1} = W_1^T \underline{x}_1 + \underline{b}_1$$

$[\underline{x}]_{784 \times 1}$

$[\underline{x}_1]_{128 \times 1}$

$[\underline{y}]_{10 \times 1}$

784

$\left[W_0^T\right]_{784 \times 128}$

128

$\left[W_1^T\right]_{128 \times 10}$

10

# Neural Network: FC Layer Design

Apply Softmax regression model to map output classes in range between [0,1]

$$\hat{y}_j = \frac{\exp y_j}{\sum_{k=1}^{10} \exp y_k}$$



$[x]_{784 \times 1}$

$[x_1]_{128 \times 1}$

$[y]_{10 \times 1}$  $[\hat{y}]_{10 \times 1}$

$[W_0^T]_{784 \times 128}$

$[W_1^T]_{128 \times 10}$

784

128

10

10

# Fully Connected NN: Error-Gradient Backpropagation

- To train NN, the gradient of error-loss is calculated with respect to each learnable parameter in the network

$$\frac{\partial \epsilon}{\partial w_{0,i,j}}, \ \frac{\partial \epsilon}{\partial b_{0,j}}, \ \frac{\partial \epsilon}{\partial w_{1,i,j}}, \ \frac{\partial \epsilon}{\partial b_{1,j}}$$

# Fully Connected NN: Error-Gradient Backpropagation

- To train NN, the gradient of error-loss is calculated with respect to each learnable parameter in the network

$$\frac{\partial \epsilon}{\partial w_{0,i,j}}, \frac{\partial \epsilon}{\partial b_{0,j}}, \frac{\partial \epsilon}{\partial w_{1,i,j}}, \frac{\partial \epsilon}{\partial b_{1,j}}$$

Rumelhart et al., 1986 introduced Backpropagation for training NN

- Error gradients are used to update the network parameters in iterative minimization using gradient descent

$$\text{e.g. } w_{0,i,j}^{k+1} \leftarrow w_{0,i,j}^{k} - \eta \frac{\partial \epsilon}{\partial w_{0,i,j}^{k+1}}$$
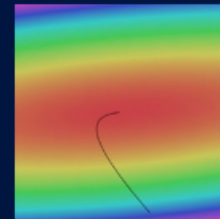
# Fully Connected NN: Error-Gradient Backpropagation

- To train NN, the gradient of error-loss is calculated with respect to each learnable parameter in the network

$$\frac{\partial \epsilon}{\partial w_{0,i,j}}, \frac{\partial \epsilon}{\partial b_{0,j}}, \frac{\partial \epsilon}{\partial w_{1,i,j}}, \frac{\partial \epsilon}{\partial b_{1,j}}$$

- Error gradients are used to update the network parameters in iterative minimization using gradient descent

$$\text{e.g. } w_{0,i,j}^{k+1} \leftarrow w_{0,i,j}^{k} - \eta \frac{\partial \epsilon}{\partial w_{0,i,j}^{k+1}}$$

# Fully Connected NN: Error-Gradient Backpropagation

- To train NN, the gradient of error-loss is calculated with respect to each learnable parameter in the network

$$\frac{\partial \epsilon}{\partial w_{0,i,j}}, \quad \frac{\partial}{\partial b_0}$$
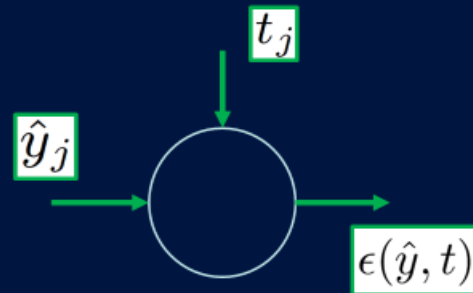
- Error gradien
  iterative mini

$$e.g. \; w_{0,i,j}^{k+1} \leftarrow$$

```
90          iteration_p = 0
91      for p in group['params']:
92          if p.grad is None:
93              iteration_p += 1
94              continue
95          d_p = p.grad.data
96          if weight_decay != 0:
97              d_p.add_(weight_decay, p.data)
98          if momentum != 0:
99              param_state = self.state[p]
100             if 'momentum_buffer' not in param_state:
101                 buf = param_state['momentum_buffer'] = torch.clone(d_p).detach()
102             else:
103                 buf = param_state['momentum_buffer']
104                 buf.mul_(momentum).add_(1 - dampening, d_p)
105             if nesterov:
106                 d_p = d_p.add(momentum, buf)
107             else:
108                 d_p = buf

110         p.data.add_(-group['lr'], d_p)
111         iteration_p += 1
112
113     return loss
```

# Fully Connected NN: Calculating Error-Gradients using Chain-Rule

- Classification sub-module

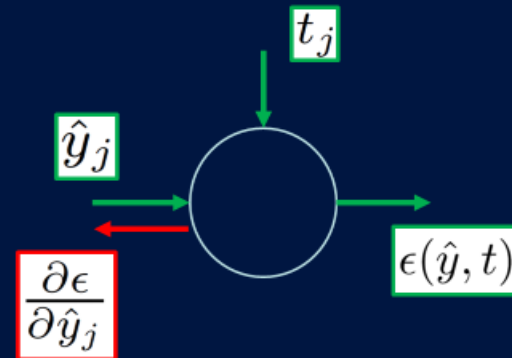$$\epsilon(\hat{y}, t) = -\sum_{k=1}^{10} t_k \ln \hat{y}_k = -\underline{t}^T \ln \underline{\hat{y}}$$

# Fully Connected NN: Calculating Error-Gradients using Chain-Rule

- Classification sub-module

$t_j$

$\hat{y}_j$

$\epsilon(\hat{y}, t)$

$\dfrac{\partial \epsilon}{\partial \hat{y}_j}$

$$\epsilon(\hat{y}, t) = -\sum_{k=1}^{10} t_k \ln \hat{y}_k = -\underline{t}^T \ln \underline{\hat{y}}$$

$$\frac{\partial \epsilon}{\partial \hat{y}_j} = -\frac{t_j}{\hat{y}_j}$$

# Fully Connected NN: Calculating Error-Gradients using Chain-Rule

- Probability sub-module

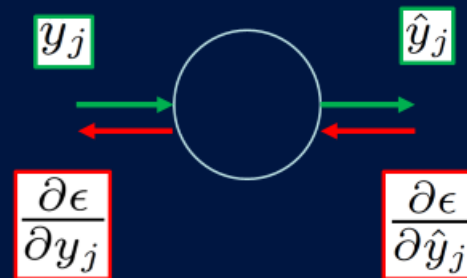$$\hat{y}_j = \frac{\exp y_j}{\sum_{k=1}^{10} \exp y_k}$$

# Fully Connected NN: Calculating Error-Gradients using Chain-Rule
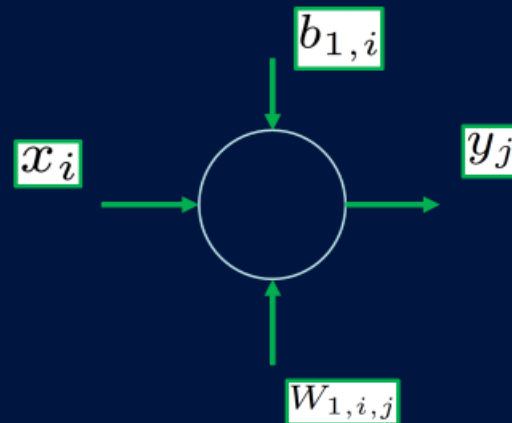
- Probability sub-module

$$\hat{y}_j = \frac{\exp y_j}{\sum_{k=1}^{10} \exp y_k}$$

$$\frac{\partial \epsilon}{\partial y_j} = \frac{\partial \epsilon}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial y_j} = \frac{\partial \epsilon}{\partial \hat{y}_j} \cdot \hat{y}_j \cdot (1 - \hat{y}_j)$$

The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
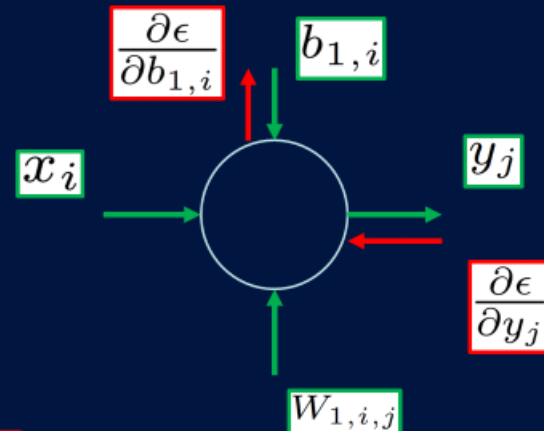UNIVERSITY OF TORONTO

FACULTY OF APPLIED SCIENCE & ENGINEERING

# Fully Connected NN: Calculating Error-Gradients using Chain-Rule

- FC Layer sub-module

$$y_j = \sum_{k=1}^{128} W_{1,k,j} x_k + b_{1,j}$$

$b_{1,i}$

$x_i$

$y_j$

$W_{1,i,j}$

# Fully Connected NN: Calculating Error-Gradients using Chain-Rule

- FC Layer sub-module

$$\boxed{\frac{\partial \epsilon}{\partial b_{1,i}}} \quad \boxed{b_{1,i}}$$

$$\boxed{x_i} \qquad \boxed{y_j}$$

$$\boxed{\frac{\partial \epsilon}{\partial y_j}}$$

$$\boxed{W_{1,i,j}}$$

$$\boxed{y_j = \sum_{k=1}^{128} W_{1,k,j} x_k + b_{1,j}}$$

$$\boxed{\frac{\partial \epsilon}{\partial b_{1,i}} = \frac{\partial \epsilon}{\partial y_j} \cdot \frac{\partial y_j}{\partial b_{1,i}} = \frac{\partial \epsilon}{\partial y_j} \cdot 1}$$
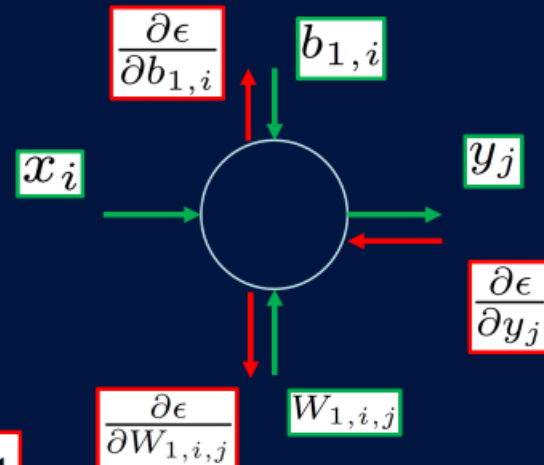
# Fully Connected NN: Calculating Error-Gradients using Chain-Rule

- FC Layer sub-module

$$y_j = \sum_{k=1}^{128} W_{1,k,j} x_k + b_{1,j}$$

$$\frac{\partial \epsilon}{\partial b_{1,i}} = \frac{\partial \epsilon}{\partial y_j} \cdot \frac{\partial y_j}{\partial b_{1,i}} = \frac{\partial \epsilon}{\partial y_j} \cdot 1$$

$$\frac{\partial \epsilon}{\partial W_{1,i,j}} = \frac{\partial \epsilon}{\partial y_j} \cdot \frac{\partial y_j}{\partial W_{1,i,j}} = \frac{\partial \epsilon}{\partial y_j} \cdot x_i$$

$\dfrac{\partial \epsilon}{\partial b_{1,i}}$  $b_{1,i}$

$x_i$  $y_j$

$\dfrac{\partial \epsilon}{\partial y_j}$

$\dfrac{\partial \epsilon}{\partial W_{1,i,j}}$  $W_{1,i,j}$

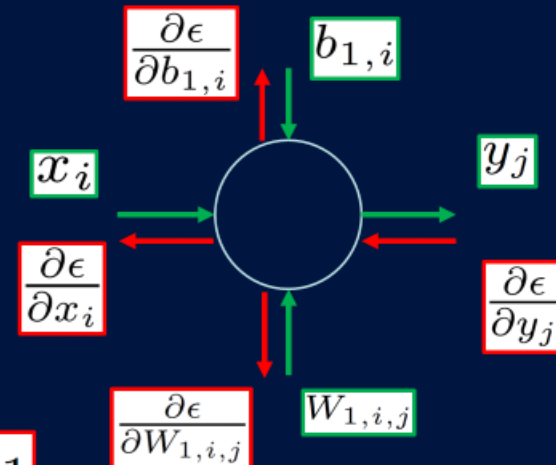# Fully Connected NN: Calculating Error-Gradients using Chain-Rule

- FC Layer sub-module

$$y_j = \sum_{k=1}^{128} W_{1,k,j} x_k + b_{1,j}$$

$$\frac{\partial \epsilon}{\partial b_{1,i}} = \frac{\partial \epsilon}{\partial y_j} \cdot \frac{\partial y_j}{\partial b_{1,i}} = \frac{\partial \epsilon}{\partial y_j} \cdot 1$$

$$\frac{\partial \epsilon}{\partial W_{1,i,j}} = \frac{\partial \epsilon}{\partial y_j} \cdot \frac{\partial y_j}{\partial W_{1,i,j}} = \frac{\partial \epsilon}{\partial y_j} \cdot x_i$$

$$\frac{\partial \epsilon}{\partial x_i} = \frac{\partial \epsilon}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_i} = \frac{\partial \epsilon}{\partial y_j} \cdot W_{1,i,j}$$
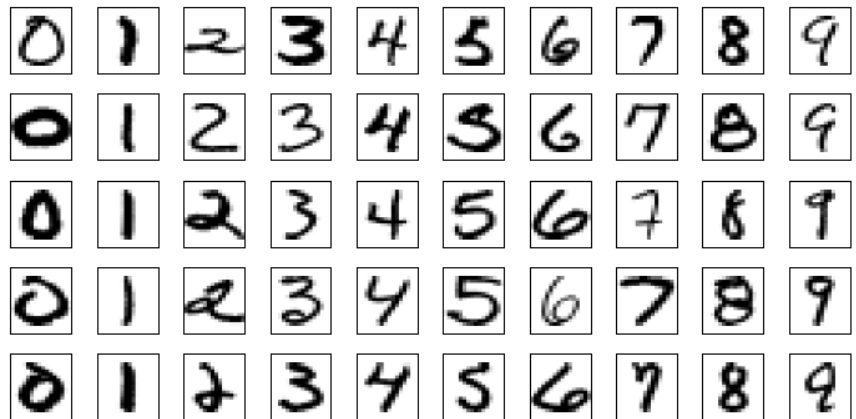
$\frac{\partial \epsilon}{\partial b_{1,i}}$ $b_{1,i}$

$x_i$ $y_j$

$\frac{\partial \epsilon}{\partial x_i}$ $\frac{\partial \epsilon}{\partial y_j}$

$\frac{\partial \epsilon}{\partial W_{1,i,j}}$ $W_{1,i,j}$

The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
UNIVERSITY OF TORONTO

FACULTY
OF APPLIED
SCIENCE &
ENGINEERING

46

# Applications of Neural Networks

- Handwritten digit recognition
  - Training set = set of handwritten digits (0…9)
  - Task: given a bitmap, determine what digit it represents
  - Input: 1 feature for each pixel of the bitmap
  - Output: 1 output unit for each possible character (only 1 should be activated)
  - After training, network should work for fonts (handwriting) never encountered

- Related pattern recognition applications:
  - recognize postal codes
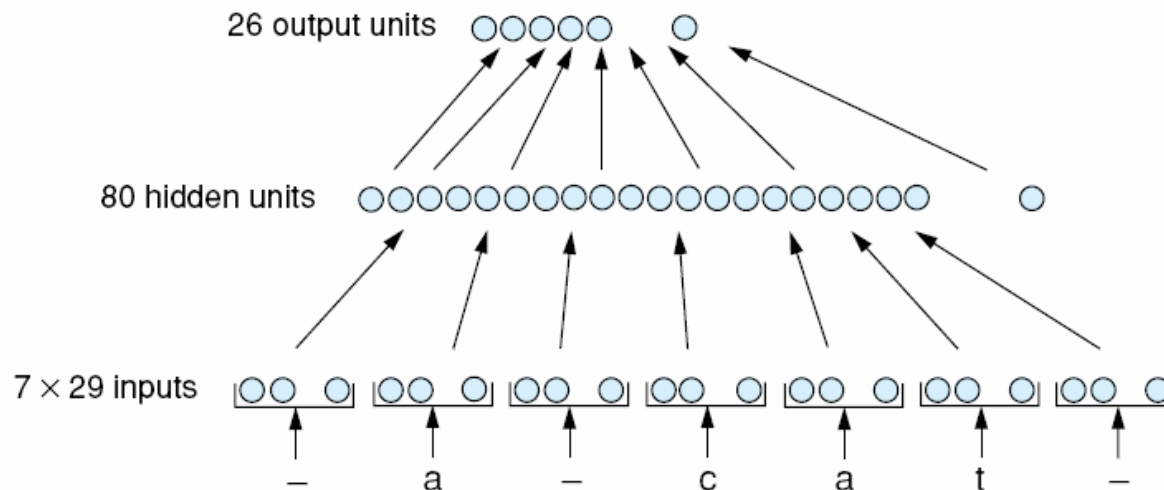  - recognize signatures
  - …

# Applications of Neural Networks

- Speech synthesis
  - Learning to pronounce English words
  - Difficult task for a rule-based system because English pronunciation is highly irregular
  - Examples:
    - letter "c" can be pronounced [k] (*cat*) or [s] (*cents*)
    - *Woman* vs *Women*
  - NETtalk:
    - uses the context and the letters around a letter to learn how to pronounce a letter
    - Input: letter and its surrounding letters
    - Output: phoneme

# NETtalk Architecture



Ex:  *a cat*  →   *c is pronounced K*

- Network is made of 3 layers of units
- input unit corresponds to a 7 character window in the text
- each position in the window is represented by 29 input units (26 letters + 3 for punctuation and spaces)
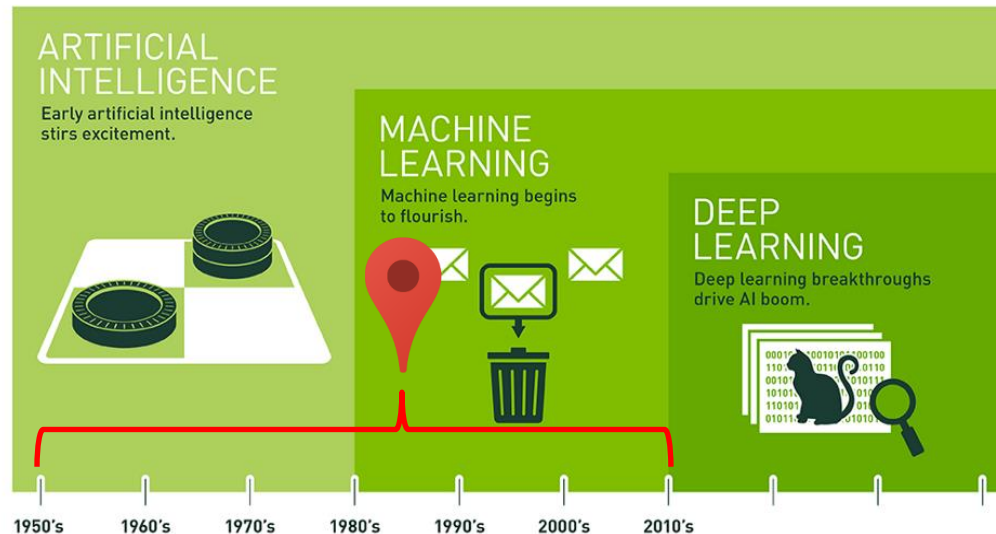- 26 output units – one for each possible phoneme

Listen to the output through iterations: https://www.youtube.com/watch?v=gakJlr3GecE

# Neural Networks

- Disadvantage:
    - result is not easy to understand by humans (set of weights compared to decision tree)… it is a black box

- Advantage:
    - robust to noise in the input (small changes in input do not normally cause a change in output) and *graceful* degradation

# Today

- **Introduction to Neural Networks**
  - Perceptrons
  - Backpropagation