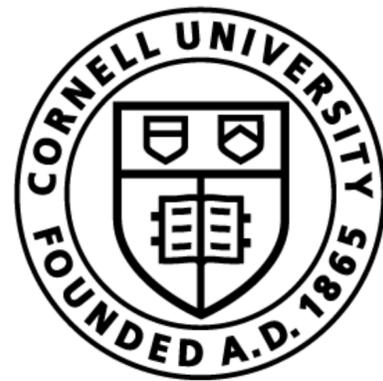# Lecture 4: FFNNs, Error Backdrop, Midterm Review.

Cornell Bowers C·IS
**Computer Science**

Tanya Goyal

CS 4740 (and crosslists): Introduction to Natural Language Processing

# Midterm

- **Closed-book**. No calculators (will not be needed)

- The test will be 2 hours long and will cover all material **up to and including error backpropagation.**

- The written questions in the HWs + practice are good representative questions for the midterm.

- I will include 2 question **very** similar to the HW questions. Make sure to understand them well!

# Today

- Recap: Feed Forward Neural Networks

- Deep Averaging Networks

- Error Backpropagation

# Feed-forward neural networks (FFNNs) for classification
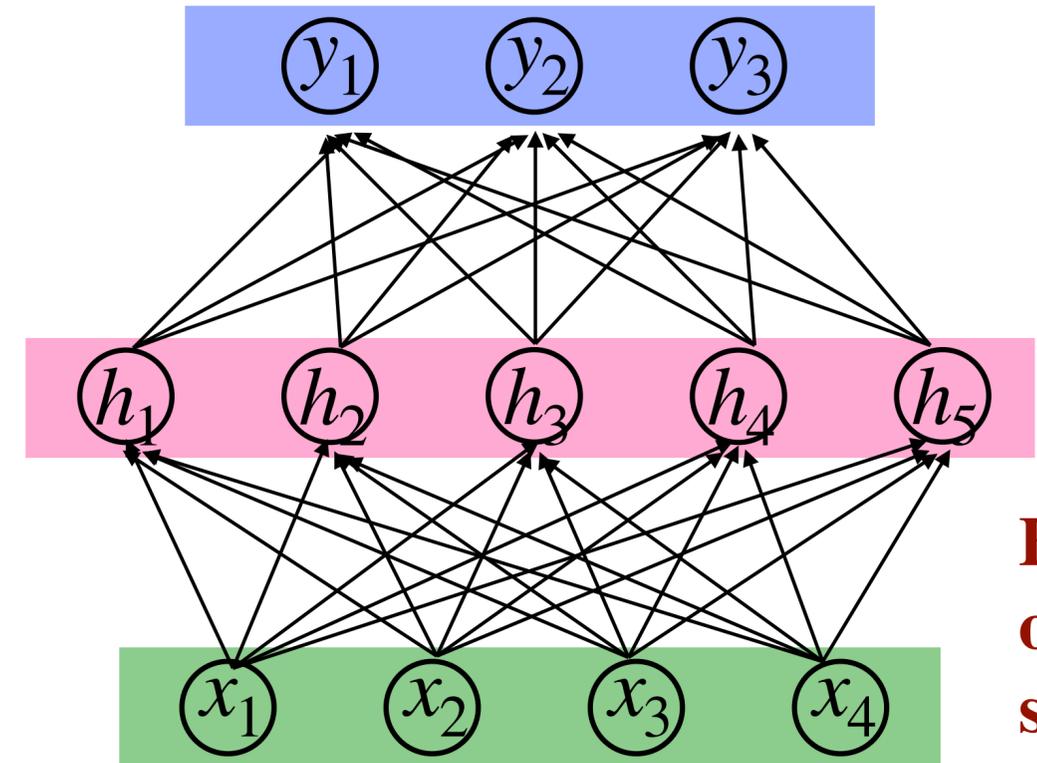
Given a input representation

$$\mathbf{x} = [x_1, \cdots x_i, \cdots x_d]$$

New vector representation

$$\mathbf{h} = [h_1, \cdots h_i, \cdots h_{d_1}]$$
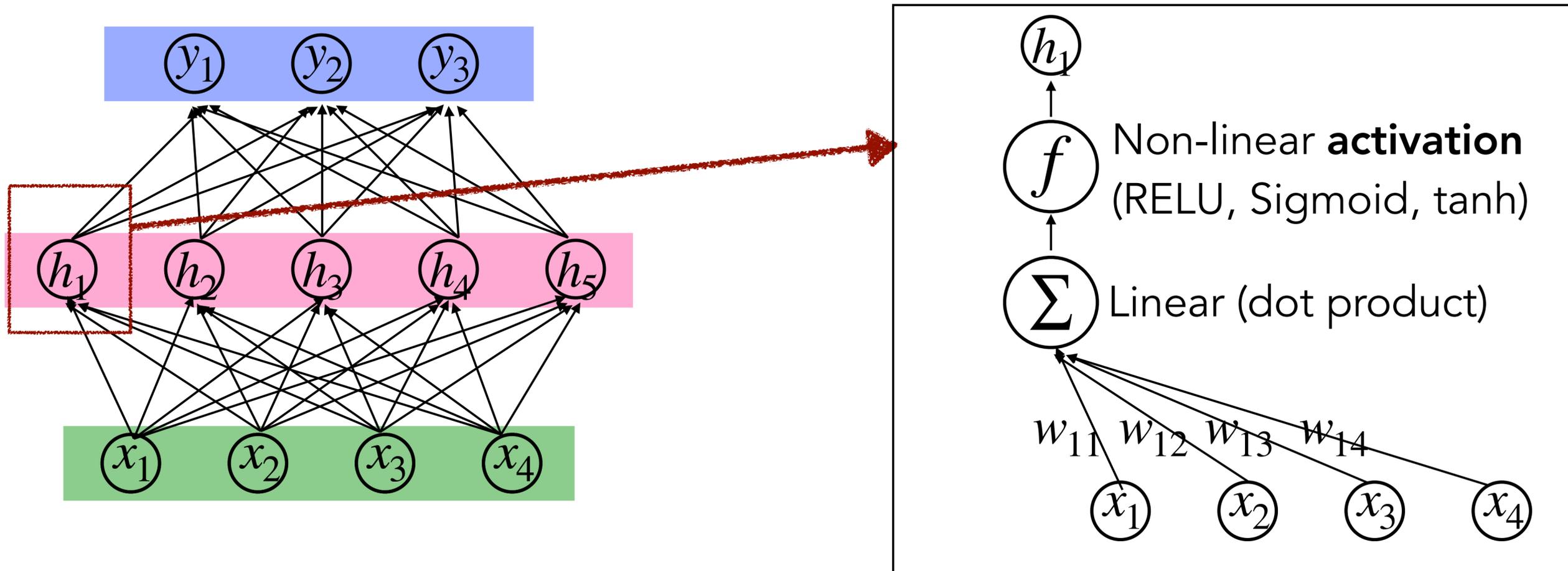
Predict: probability distribution over the labels

$$\mathbf{y} = [y_1 = P(\text{LOC}\,|\,\mathbf{x}), P(\text{PER}\,|\,\mathbf{x}), P(\text{NULL}\,|\,\mathbf{x})]$$



**Bias term omitted for simplicity, let's come back to this.**

Intuition: FFNNs learn how to better represent **x** (i.e, as **h**) so as to be able to accurately predict **y**.
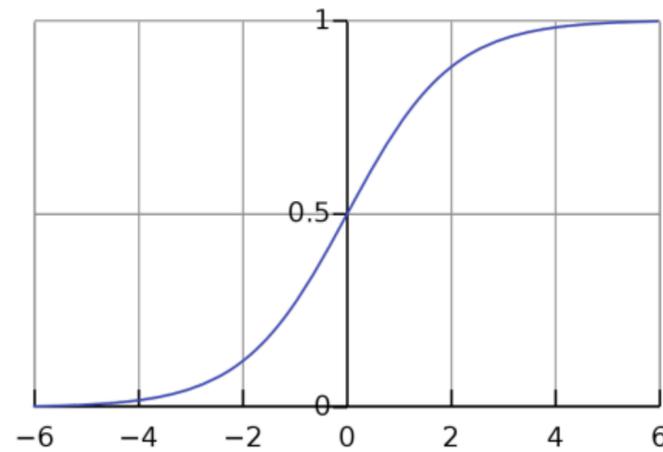
# Feed-forward neural networks (FFNNs)
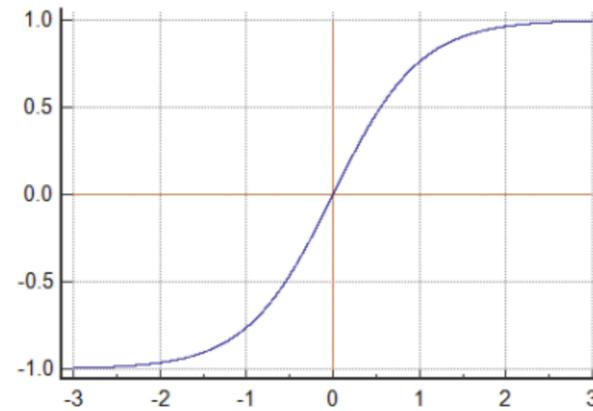
# Activation functions

### sigmoid

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

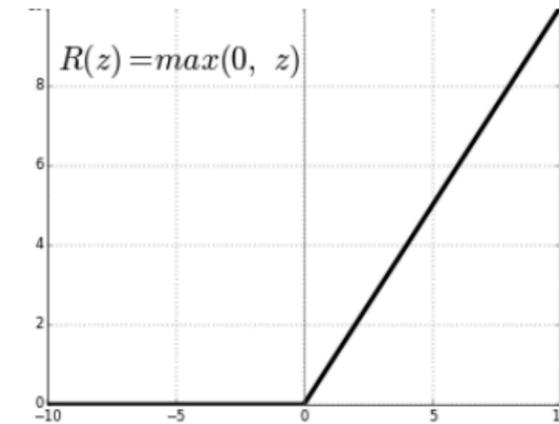$$f'(z) = f(z) \times (1 - f(z))$$

### tanh

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

$$f'(z) = 1 - f(z)^2$$

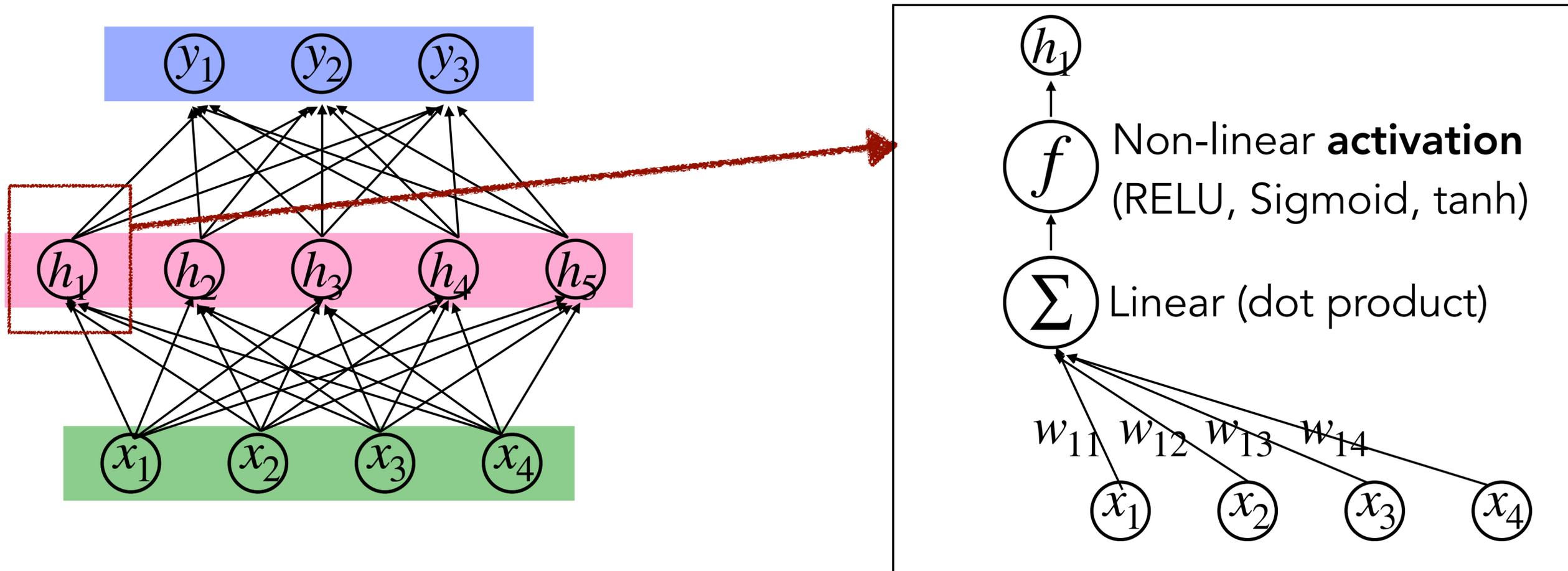### ReLU
### (rectified linear unit)

$$f(z) = \max(0, z)$$

$$f'(z) = \begin{cases} 1 & z > 0 \\ 0 & z < 0 \end{cases}$$
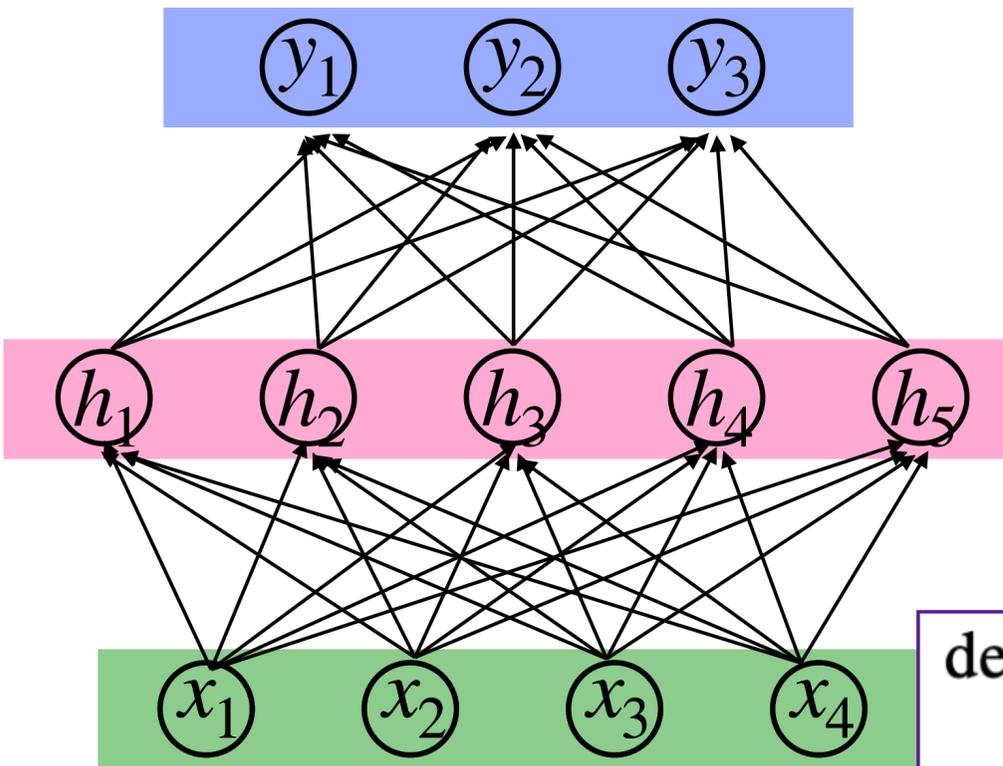
- Applied element-wise: $f([z_1, z_2, \ldots z_n]) = [f(z_1), f(z_2) \ldots f(z_n)]$

# Feed-forward neural networks (FFNNs)



- $h_i = f(\sum_j w_{ij} x_j)$

- Or, in matrix form: $x = [x_1 x_2 \ldots x_d]^T$, $w_i = [w_{i1}, w_{i2}, \ldots w_{id}]$, then $h_i = f(w_i x)$

# Feed-forward neural networks (FFNNs), Matrix version



- Input layer $x = [x_1 x_2 \ldots x_4]^T$

- Hidden Layer:

  - $h = f(Wx) \quad \in R^5$

  - $W \in \mathbb{R}^?$

defined as:

$$\text{softmax}(\mathbf{z}_i) = \frac{\exp(\mathbf{z}_i)}{\sum_{j=1}^{d} \exp(\mathbf{z}_j)} \quad 1 \le i \le d \tag{7.9}$$

Thus for example given a vector

$$\mathbf{z} = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1], \tag{7.10}$$

the softmax function will normalize it to a probability distribution (shown rounded):

$$\text{softmax}(\mathbf{z}) = [0.055, 0.090, 0.0067, 0.10, 0.74, 0.010] \tag{7.11}$$

# Feed-forward neural networks (FFNNs), matrix version



$U$ and $W$ are weights that are learnt during training.

What is this **b** term and the extra input value that is always a 1?

Number of layers = Number of learnable layers (2 in this figure)

# Parameters, Hyperparameters

- What are parameters? Weights that are learned during training.

- What are hyperparameters? Set by the model developer, e.g. dimensions, learning rates, number of epochs, etc.

# Quiz 1: Dimensionality

My input $x$ has dimension $100$. I want to classify it as one of 4 labels. I designed a 2 layer FFNN.

Q1. Is this dimension information complete to construct the FFNN?

Q2: What is the dimensions of the weight matrices and the bias matrices?

Q3: Which weights are "learnable"?

# Big question: how do we learn the weights?

- In the previous slides, we saw how to derive the final output $\hat{y}$ given an input $x$ and the weights ($W$ and $b$) for all layers.

- How do we learn these weights?

- Similar strategy to how we learned weights for logistic regression:

Initialize $W$ and $b$ randomly.

Repeat

For each training datapoint $(x, y)$

        1.  Feedforward pass to compute output $\hat{y}$

        2.  Compute error, i.e. loss($y, \hat{y}$)

        3.  Update **weights in the network to reduce error.**

Until a stopping criterion is reached.

Return the learned $W$ and $b$

# Quiz detour.

# Constructing Input for different Tasks

In previous slides, we assumed we know how to represent the input text as a vector $x$.

Consider the following sentiment classification task:
- Input: $w_1 w_2 w_3 \ldots w_n$      "The move was great"
- Output: $y =$ "positive".    Y = {"postive", "negative", "neutral"}

**How should we convert the text input "The move was great" to a vector representation $x$?**

# Constructing Input for different Tasks

In previous slides, we assumed we know how to represent the input text as a vector $x$.

Consider the following sentiment classification task:
- Input: $w_1 w_2 w_3 \ldots w_n$      "The move was great"
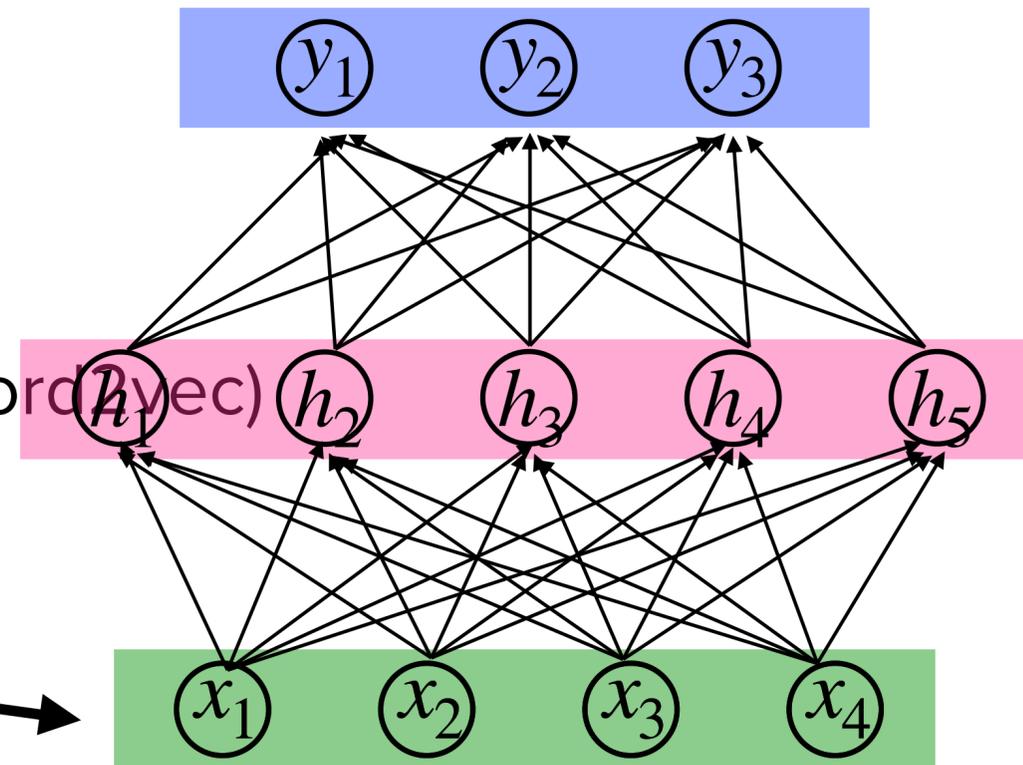- Output: $y =$ "positive".    Y = {"postive", "negative", "neutral"}
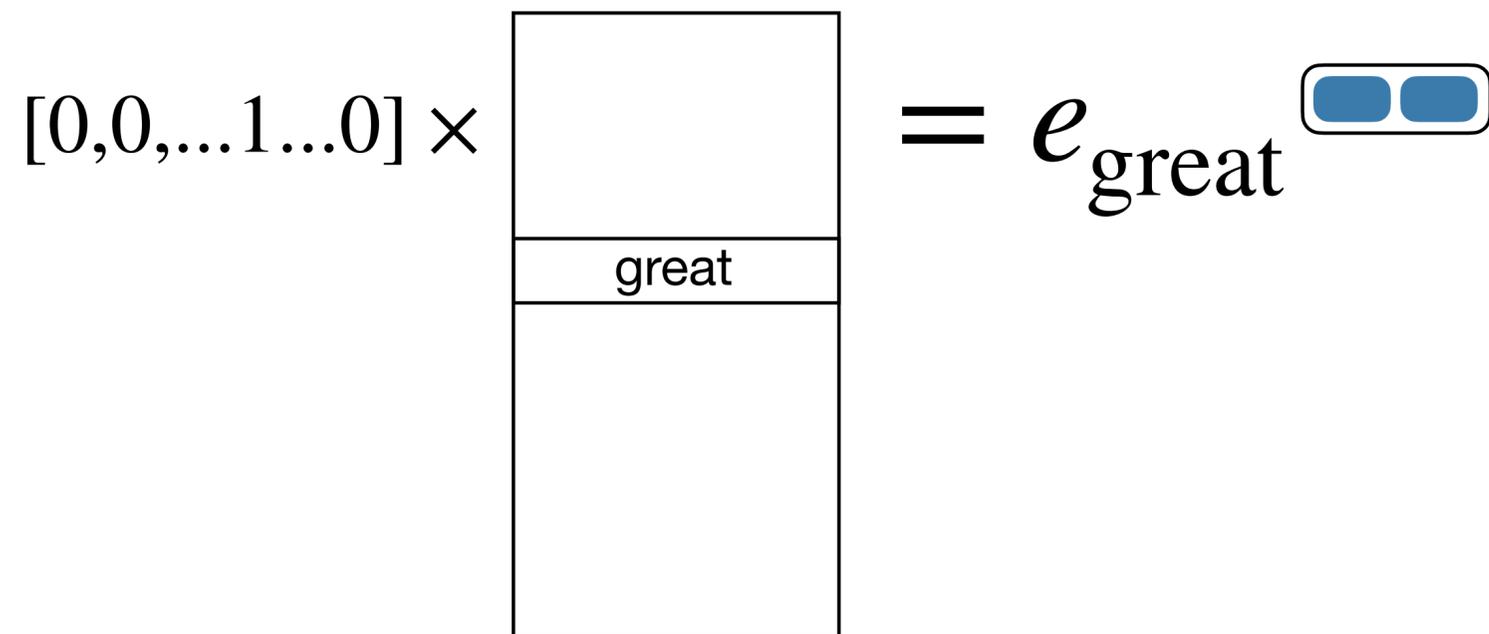
**How should we convert the text input "The move was great" to a vector representation $x$?**

**Simpler question? What can we do if the input is a single word.**   **Word Embeddings!**

# Constructing Input for different Tasks

**Simpler question? What can we do if the input is a single word. Word Embeddings!**

- Input: $w_1$              "great"
- Output: $y =$ "positive".   Y = {"postive", "negative"}

- Let E be an embeddings function / Embeddings matrix (e.g. word2vec)
- Represent $w_1$ as a one-hot encoding

$$[0,0,...1...0] \times \boxed{\quad\text{great}\quad} = e_{\text{great}}$$
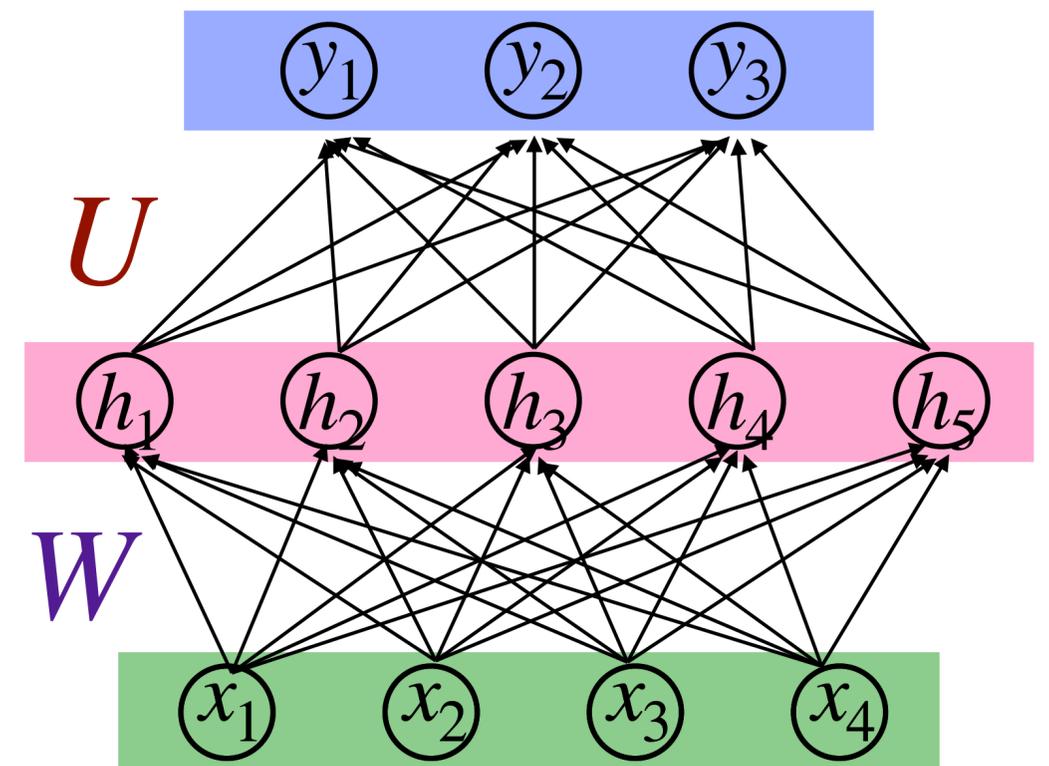
# Deep Averaging Network

**How should we convert the text input "The move was great" to a vector representation $x$?**

Consider the following sentiment classification task:
- Input: $w_1 w_2 w_3 \ldots w_n$     "The move was great"
- Output: $y =$ "positive".    Y = {"postive", "negative", "neutral"}

$[0,0,1...0...0] \times$   $= e_{\text{great}}$

$[0,1,...0...0] \times$   | great |   $= e_{\text{The}}$

$[0,0,...1...0] \times$   $= e_{\text{movie}}$

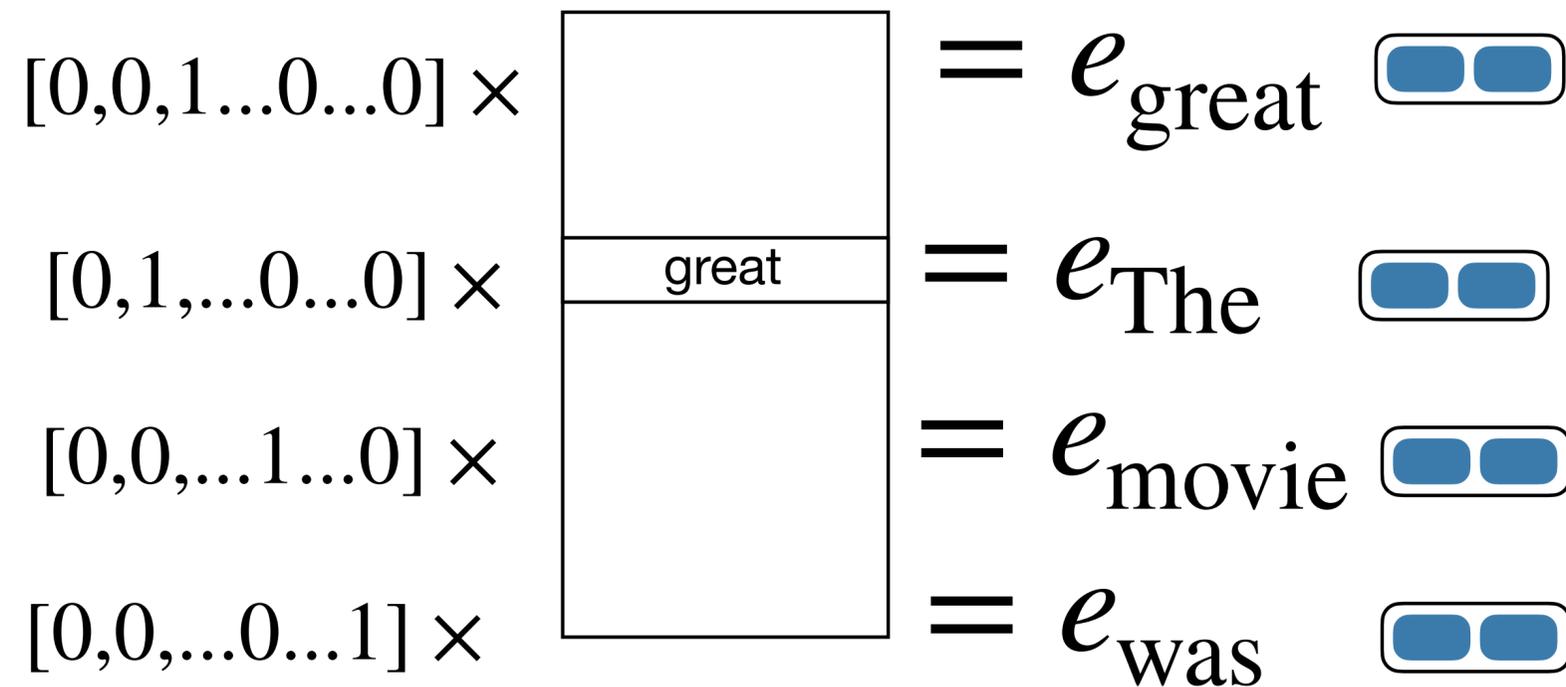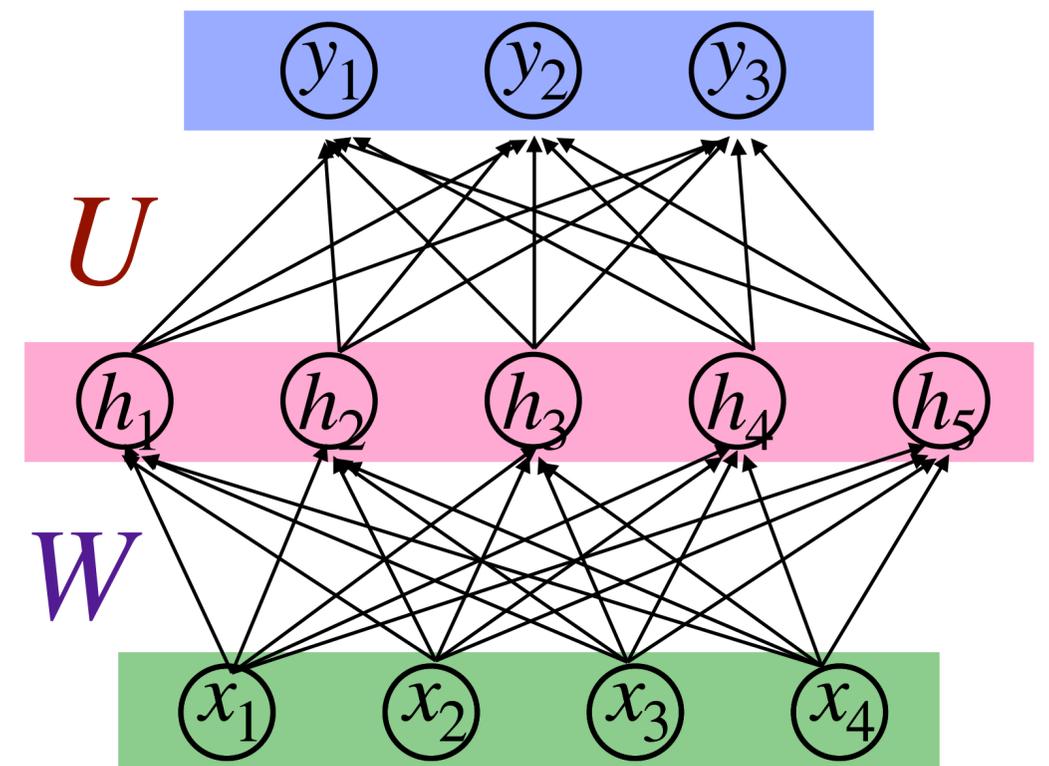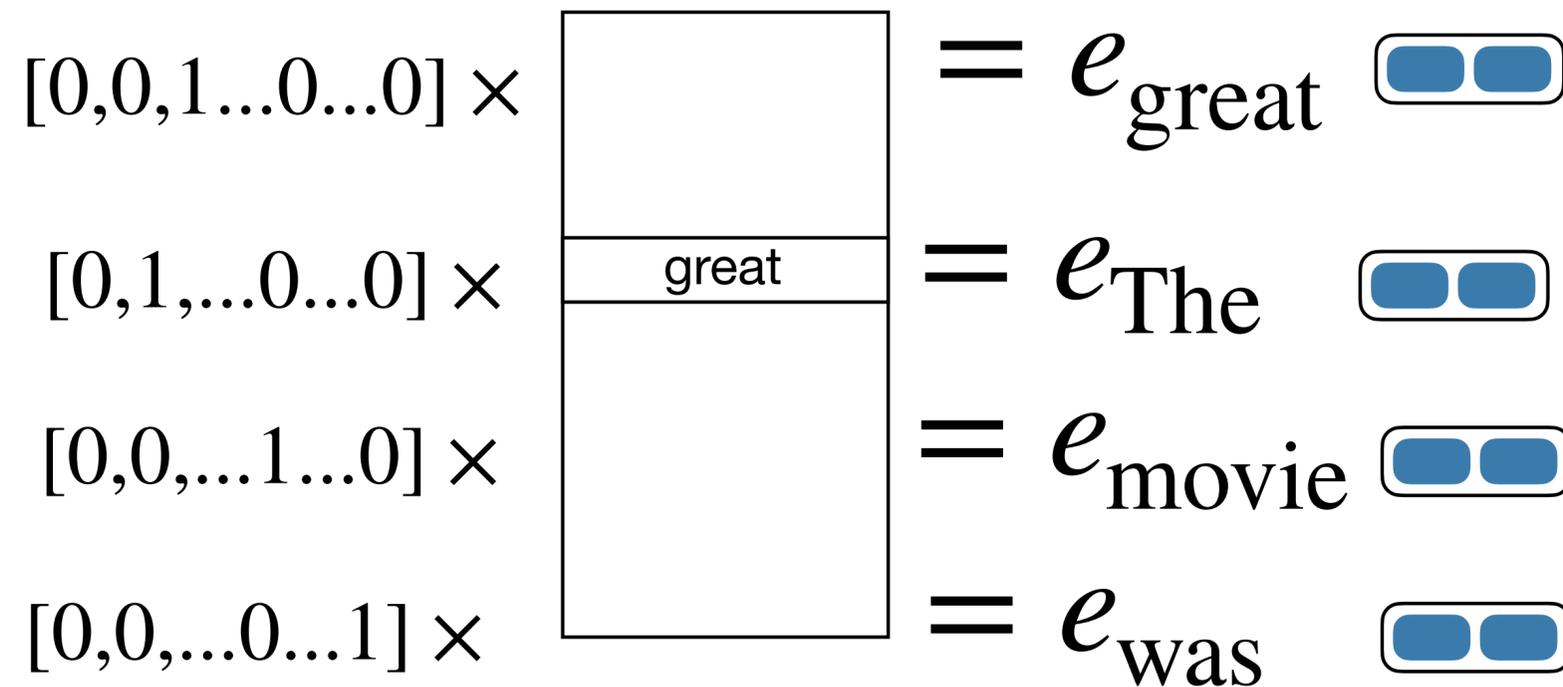$[0,0,...0...1] \times$   $= e_{\text{was}}$

$U$

$W$

# Deep Averaging Network

**How should we convert the text input "The move was great" to a vector representation $x$?**

Consider the following sentiment classification task:
- Input: $w_1 w_2 w_3 \ldots w_n$     "The move was great"
- Output: $y =$ "positive".   Y = {"postive", "negative", "neutral"}

$U$



$[0,0,1...0...0] \times$ | great | $= e_{\text{great}}$

$[0,1,...0...0] \times$ | $= e_{\text{The}}$

$[0,0,...1...0] \times$ | $= e_{\text{movie}}$

$[0,0,...0...1] \times$ | $= e_{\text{was}}$

$W$

$$x = \frac{1}{n} \sum_{i=1:n} e_{w_i}$$

# Quiz: Neural N-gram models

Consider constructing a feedforward neural network (FFNN) for the task of trigram language modeling.

(a) Provide a diagram of your network.

(b) What is purpose of the input layer of your network? What is its dimension (i.e., size)?

(c) What is purpose of the output layer of your network? What is its dimension (i.e., size)?

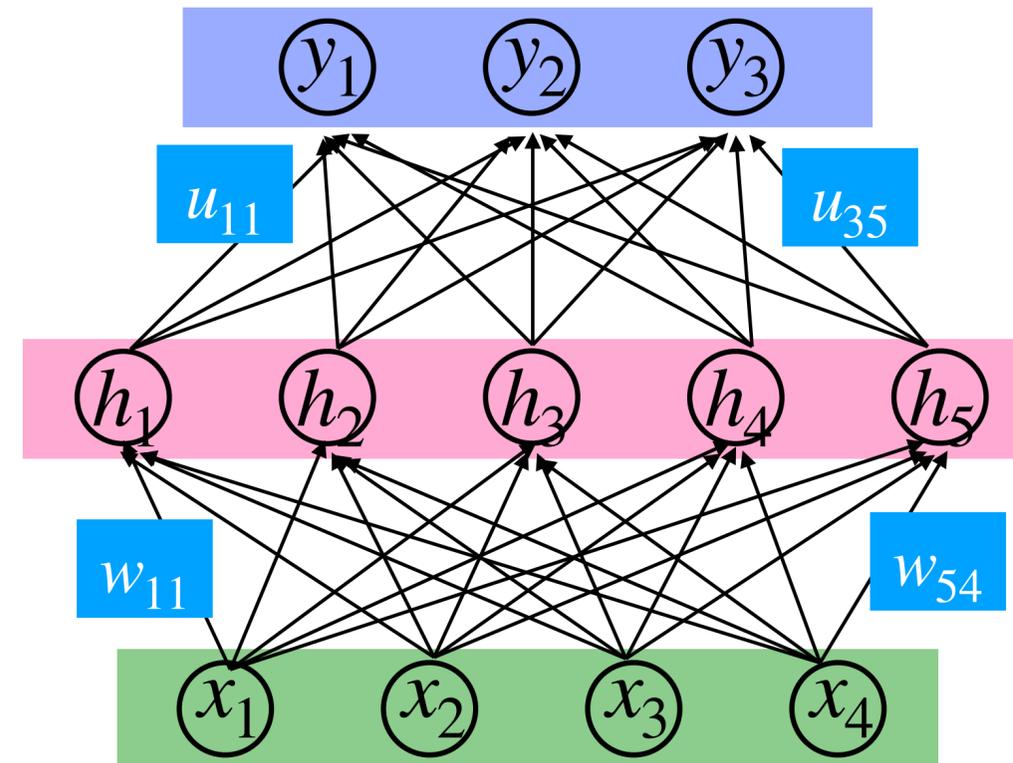# Big question: how do we learn the weights?

Initialize $W$ and $b$ randomly.

Repeat

    For each training datapoint $(x, y)$

        1.  Feedforward pass to compute output $\hat{y}$

        2.  Compute error, i.e. loss$(y, \hat{y})$

        3.  Update **weights in the network to reduce error.**

Until a stopping criterion is reached.

Return the learned $W$ and $b$

- We want to learn weight updates to all "parameters" like $w_{11}, \ldots w_{54} \ldots u_{11} \ldots u_{35}, b_1 \ldots b_n$

- We need to know $\dfrac{\partial L}{\partial w_{11}}$, etc. to update the weights.

**Challenge:** We need to compute this derivative wrt weight parameters that appear in very early layers of the network, even though loss is computed at the very end.

# Solution: Backpropagation

- Computationally efficient algorithm (automatic differentiation)

- Implemented using **computation graphs.**

  - A representation of the process of computing a mathematical expression.

  - Each operation is modeled as a node in a graph.

  - Useful to think of computation graphs as modeling data dependencies
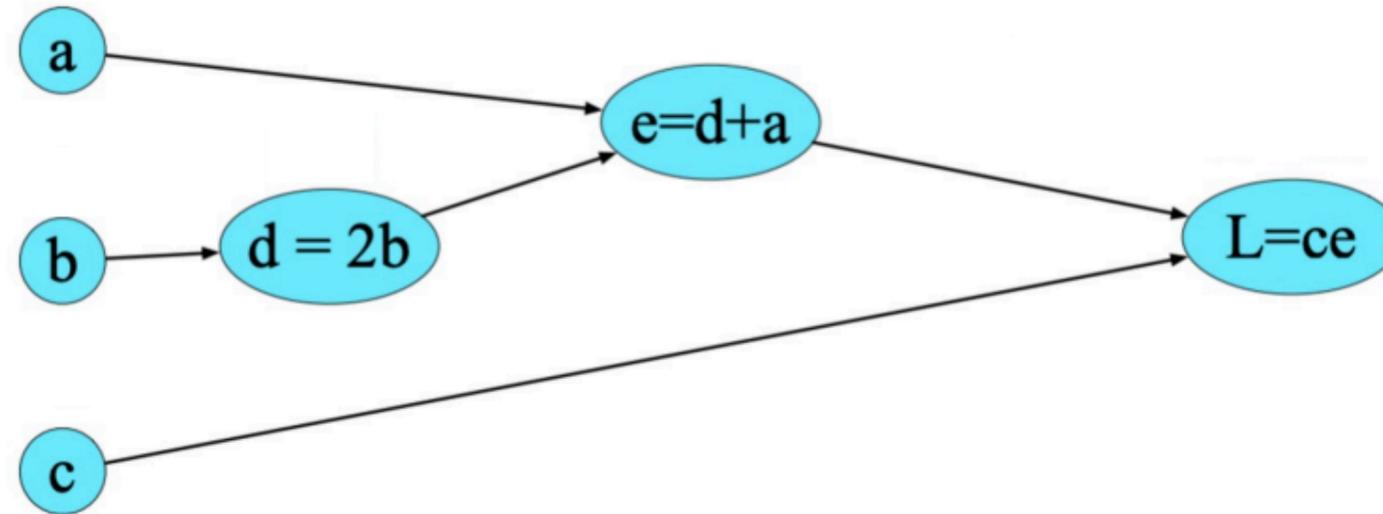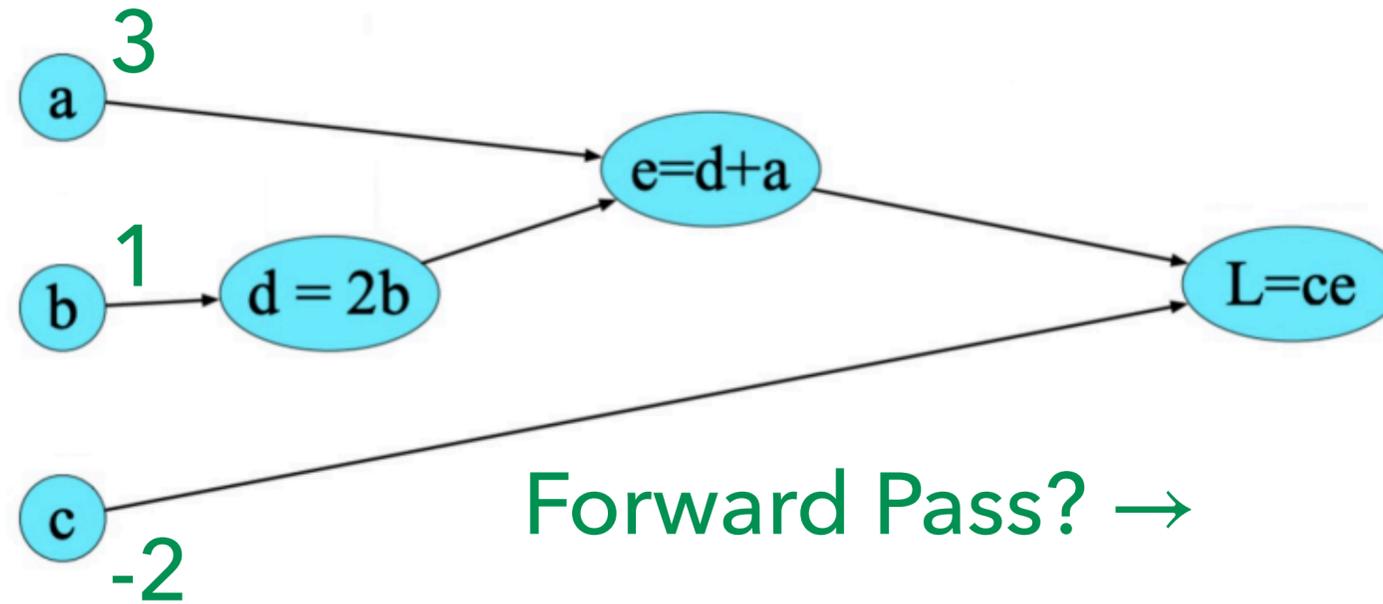
# Computational Graphs

## Computational Graph Construction

$$L(a, b, c) = c(a + 2b)$$
$$d = 2 * b$$
$$e = a + d$$
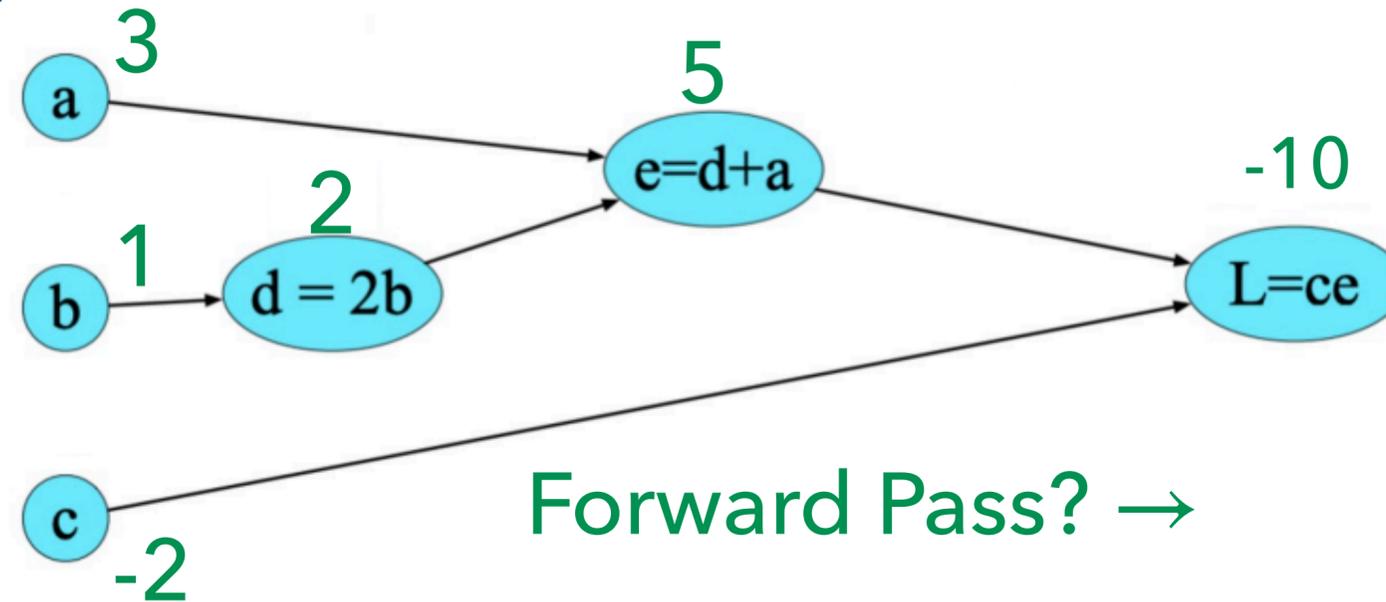$$L = c * e$$

# Computational Graphs

## Computational Graph Construction

$$L(a, b, c) = c(a + 2b)$$
$$d = 2 * b$$
$$e = a + d$$
$$L = c * e$$



Forward Pass? →

# Computational Graphs

Computational Graph Construction

$$L(a, b, c) = c(a + 2b)$$
$$d = 2 * b$$
$$e = a + d$$
$$L = c * e$$



Forward Pass? →

How do we compute derivative of $L$ wrt $a, b, c$?

Main Idea: use chain rule.    If $f(x) = u(v(x))$, then $\dfrac{\partial f}{\partial x} = \dfrac{\partial u}{\partial v} \dfrac{\partial v}{\partial x}$

$$\frac{\partial L}{\partial c} = e$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial a}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}$$
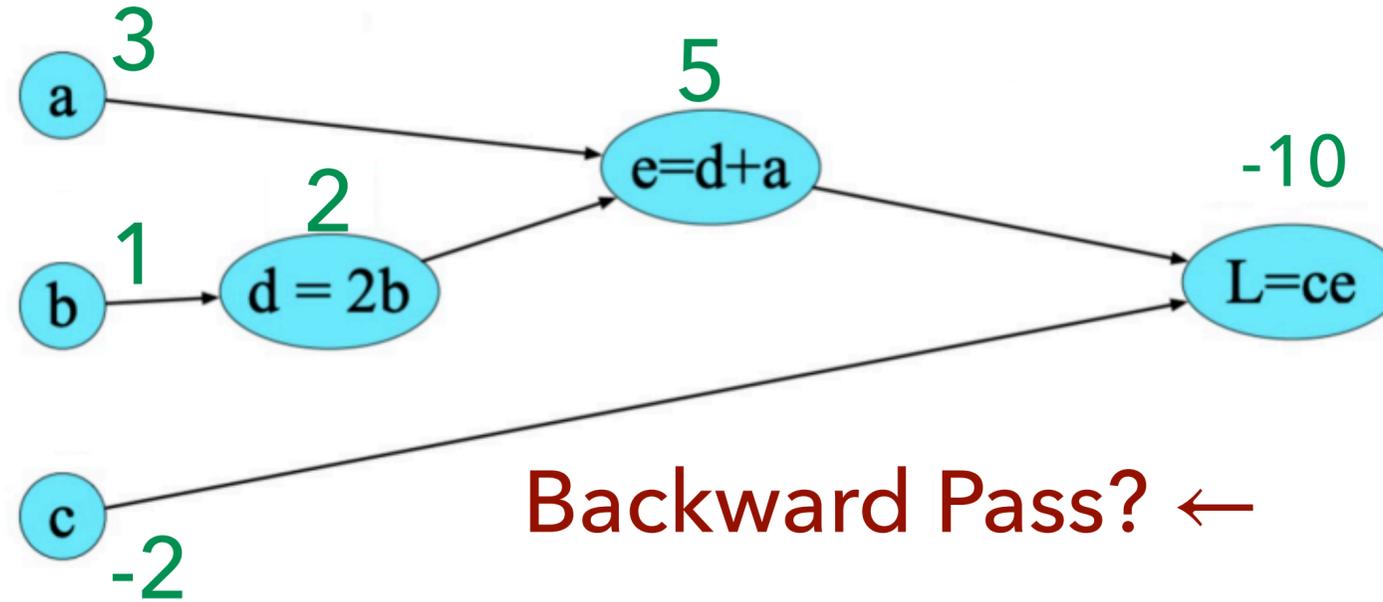
# Computational Graphs

## Computational Graph Construction

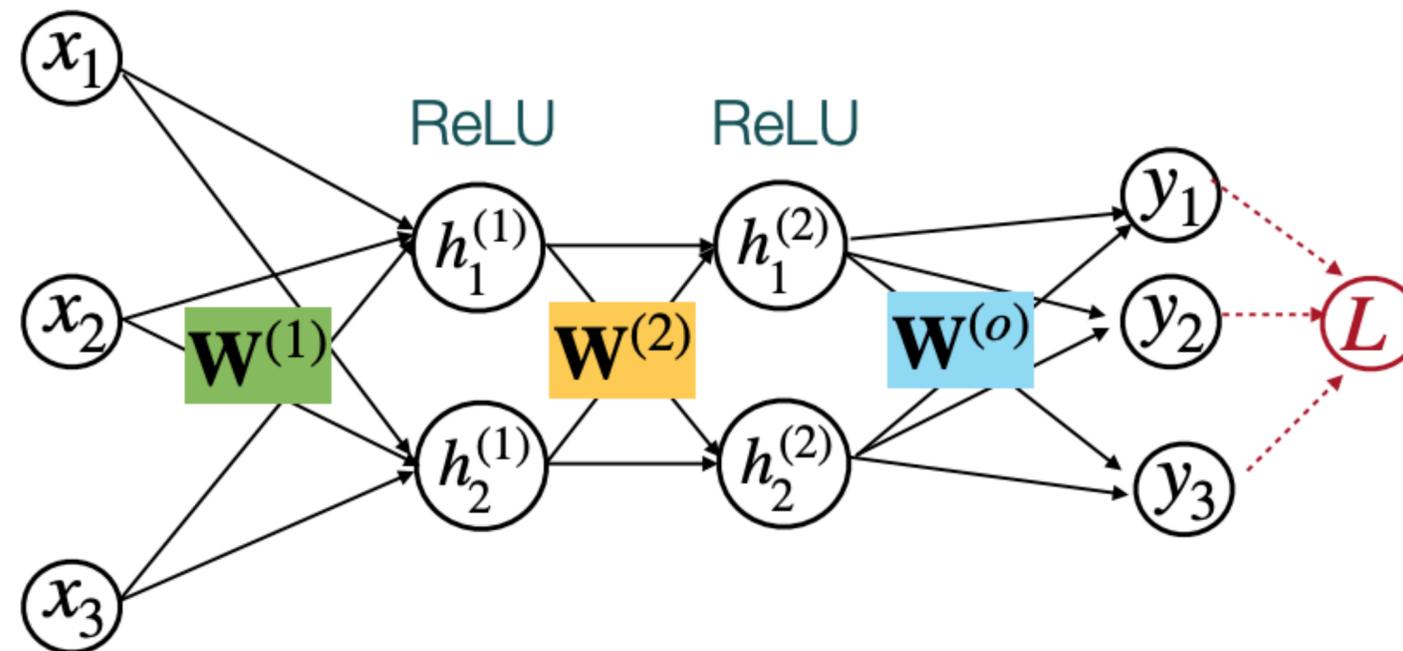$$L(a,b,c) = c(a+2b)$$
$$d = 2*b$$
$$e = a+d$$
$$L = c*e$$



$$\frac{\partial L}{\partial c} = e$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial e}\frac{\partial e}{\partial a}$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial e}\frac{\partial e}{\partial d}\frac{\partial d}{\partial b}$$
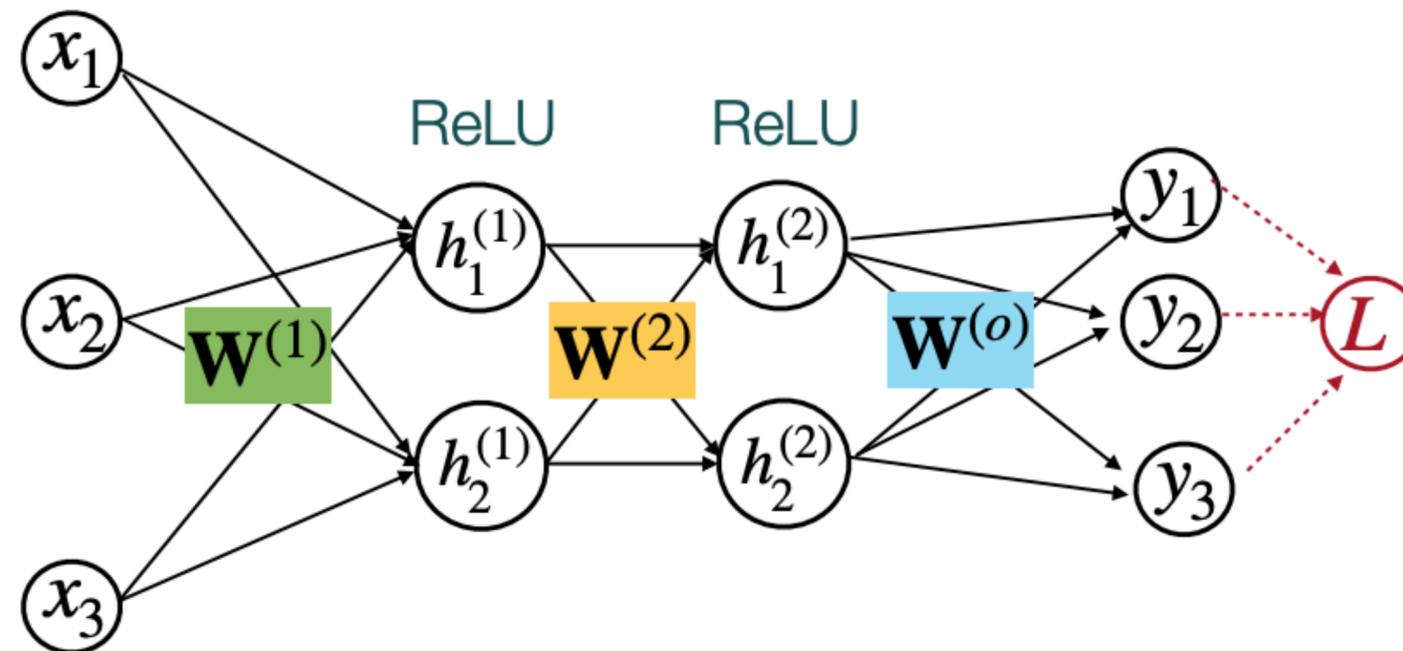
# Backpropagation for a 3 layer NN

# Backpropagation for a 3 layer NN



Given: $x_1, x_2, x_3$ and the class label $y$ (a single training example)

Goal:
$$\frac{\partial L}{\partial W^{(1)}},$$
$$\frac{\partial L}{\partial W^{(2)}},$$
$$\frac{\partial L}{\partial W^{(o)}}$$

$\mathbf{W}^{(1)}$    $\mathbf{W}^{(2)}$    $\mathbf{W}^{(o)}$

ReLU    ReLU

$x_1$  $x_2$  $x_3$  $h_1^{(1)}$  $h_2^{(1)}$  $h_1^{(2)}$  $h_2^{(2)}$  $y_1$  $y_2$  $y_3$  $L$

# Backpropagation for a 3 layer NN

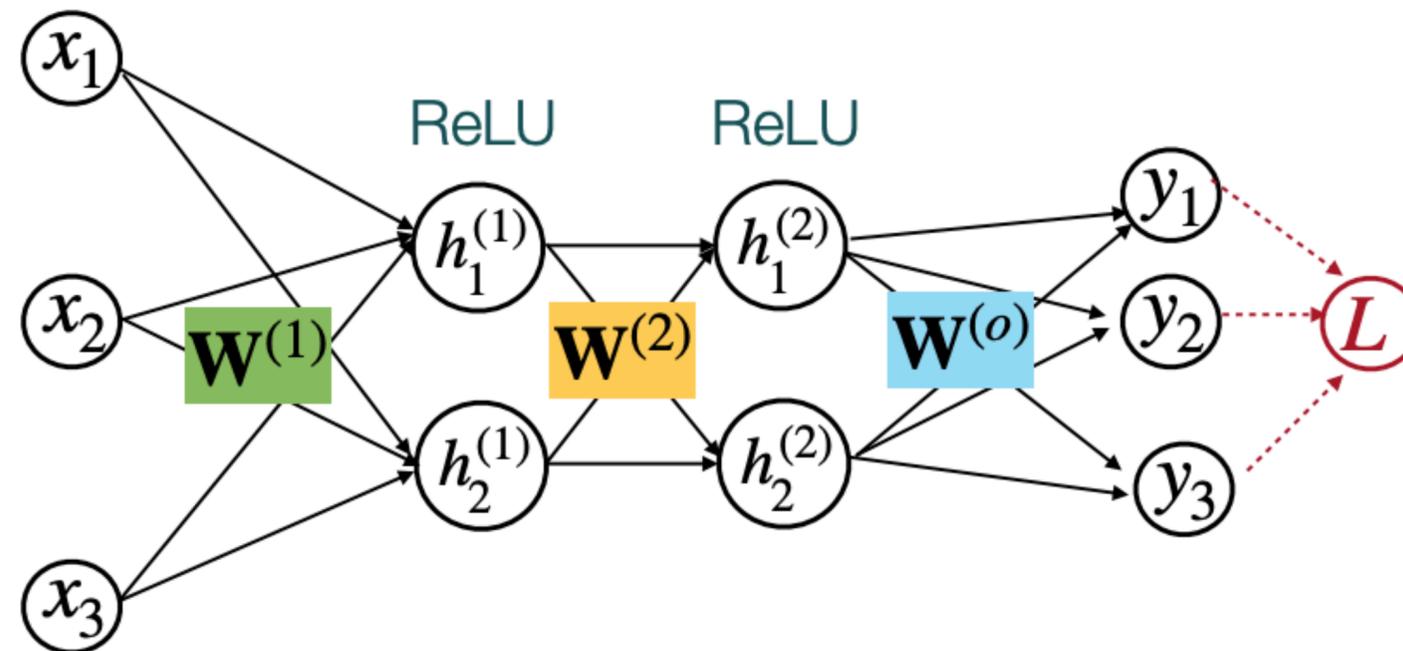

Given: $x_1, x_2, x_3$ and the class label $y$ (a single training example)

Goal:
$$\frac{\partial L}{\partial W^{(1)}},$$
$$\frac{\partial L}{\partial W^{(2)}},$$
$$\frac{\partial L}{\partial W^{(o)}}$$

Slide Credit: Alane Suhr

# Backpropagation for a 3 layer NN

# Backpropagation in PyTorch

```python
1   import torch.nn as nn
2   import torch.nn.functional as F
3
4   class Net(nn.Module):
5       def __init__(self):
6           super().__init__()
7           self.fc1 = nn.Linear(784, 128)
8           self.fc2 = nn.Linear(128, 64)
9           self.fc3 = nn.Linear(64, 10)
10
11      def forward(self, x):
12          x = F.relu(self.fc1(x))
13          x = F.relu(self.fc2(x))
14          x = self.fc3(x)
15          return x
16
```

```python
1   import torch.optim as optim
2
3   net = Net()
4   criterion = nn.CrossEntropyLoss()
5   optimizer = optim.SGD(net.parameters(), lr=0.001, momentum=0.9)
```

```python
1   outputs = net(inputs)
2   loss = criterion(outputs, labels)
3   loss.backward()
4   optimizer.step()
```

*PyTorch did back-propagation for you in this one line of code!*

# Slide Acknowledgements

▸ Earlier versions of this course offerings including materials from Claire Cardie, Marten van Schijndel, Lillian Lee

▸ CS 288 course by Alane Suhr (UC Berkeley).