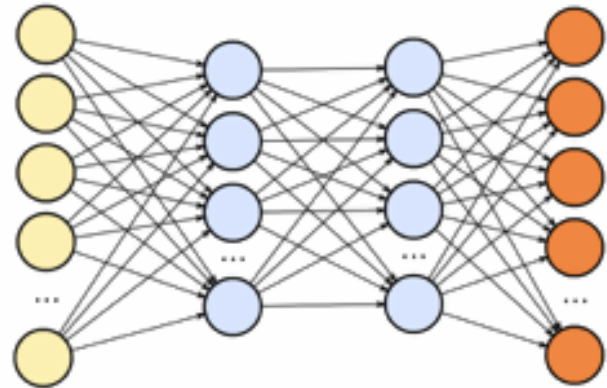


# CSE 543

Simon Du

---



# CSE543: Deep Learning

---

Instructor: Simon Du

Teaching Assistant: Ruoqi Shen, Yifang Chen

[ Course Website (contains all logistic information): <https://courses.cs.washington.edu/courses/cse543/23wi/> ]

[ Piazza: <https://piazza.com/class/lbsxy7e01whdd> ]

[ Announcements: Canvas ]

[ Homework: Canvas ]

OH: Simon: Tu 10:30-11:30

Yifang: We 2-3 PM Zoom

Ruoqi: Fr 10-11 AM Zoom

# CSE543: Deep Learning

---

## What this class is:


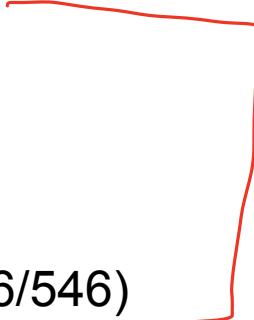

- **Fundamentals of DL:** Neural network architecture, approximation properties, optimization, generalization, generative models, representation learning
- **Preparation for further learning / research:** the field is fast-moving, you will be able to apply the fundamentals and teach yourself the latest

## What this class is not:

- An easy course: mathematically easy
- A survey course: laundry list of algorithms
- An application course: implementation of different architectures on different datasets

# Prerequisites

---

- Working knowledge of:
    - Linear algebra
    - Vector calculus
    - Probability and statistics
    - Algorithms
    - Machine learning (CSE 446/546)
  - Mathematical maturity 
  - “Can I learn these topics concurrently?”
- 
- 

# Lecture

---

- Time: Tuesday and Thursday 9:00 - 10:20AM
- MUE 153 or Zoom (see website for the schedule)
- Slides + handwritten notes (e.g., proofs)
- Please ask questions ←
- \*Recordings on Canvas
- Tentative schedule on course website

# Homework (40%)

---

- 2 homework (20%+20%)
  - Each contains both theoretical questions and will have programming
  - Related to course materials
  - Collaboration okay but must write who you collaborated with. You must write, submit, and understand your answers and code.
  - Submit on Canvas
  - Must be **typed**
  - **Two** late days
  - Tentative timeline:
    - HW 1 due: 1/27
    - HW 2 due: 2/10

# Course Project (60%)

---

- Group of 1 - 2.
- Topic: literature review (state-of-the-art) or original research. *anything related to NN*
- Some potential topics are listed on Canvas. OK to do a project ~~on~~ *not* listed.
- You can work on a project related to your research.
- Proposal (due: 1/13): 5%
  - Format: NeurIPS Latex format, ~1 - 1.5 pages
- Presentations on (3/7 and 3/9 on Zoom): **20%**
- Final report (due: 3/17): **35%**
  - Format: NeurIPS Latex format, ~8 pages
- Submit on Canvas

# Possible Topics

---

- Approximation properties
- Advanced optimization methods
- Optimization theory for deep learning
- Generalization theory for deep learning
- Deep reinforcement learning )
- Implicit regularization
- Meta-learning algorithm / theory
- Robustness
- Lottery ticket hypothesis
- Deep learning application
- ...



# Communication Channels

---

- **Announcements**
  - Canvas
- **questions about class, homework help**
  - Piazza
  - Office hours:
    - Simon Du: Tu 10:30 - 11:30 AM (in person Gates 312 and/or Zoom)
    - Ruoqi Shen:
    - Yifang Chen:
  - **Regrade requests / Personal concerns:**
    - Email to instructor or TAs

# Addcodes

---

- Email: Elle Brown ([ellean@cs.washington.edu](mailto:ellean@cs.washington.edu))  
for addcodes

# Topic 1: Review (Today)

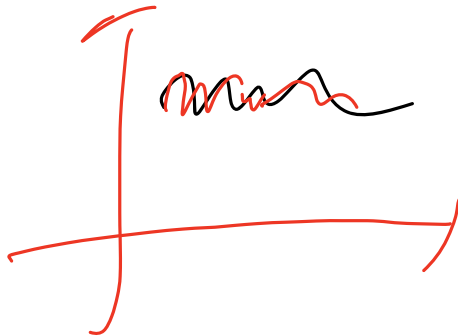
---

- ML Review: training, generalization
- Neural network basics: fully-connected neural network, gradient descent

# Topic 2: Approximation Theory

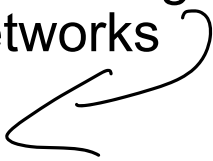
Sanity

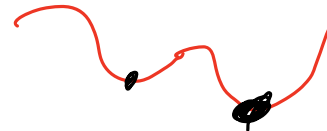
- Why neural networks can express the (regression, classification, ...) function you want?
- Construction of such desired neural networks
- Universal approximation theorem



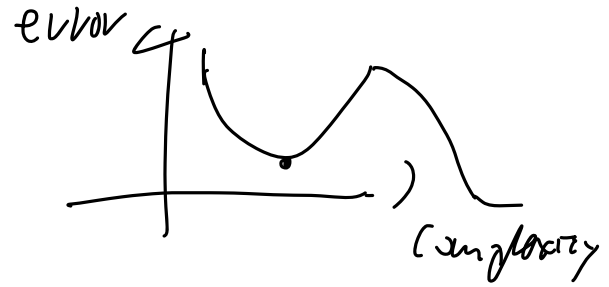
# Topic 3: Optimization

---

- Review: Back-propagation
- Auto-differentiation
- Advanced optimizers: momentum (Nesterov acceleration), adaptive method (AdaGrad, Adam)
- Techniques for improving optimization: batch-norm, layer-norm, .. *initialization*
- Theory: global convergence of gradient of over-parameterized neural networks *wide NN*
- Neural Tangent Kernel 



# Topic 4: Generalization



- Measures of generalization
- Double descent
- Techniques for improving generalization *regular*
- Generalization theory beyond VC-dimension
- Implicit regularization *algo - reg*
- Why NN outperforms kernel

$\#$  (parameters)  
 $\gg \gg \#$  data

# Topic 5: Architecture

---

- Convolutional neural network
- Recurrent neural network
  - LSTM
- Attention-based neural network
  - Transformer
- General framework

*Geometry*

# Topic 6: Representation Learning

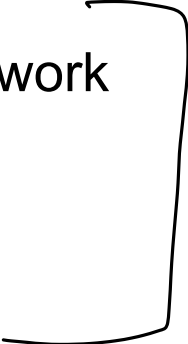
$$x \mapsto \phi(x) \rightarrow \underline{\underline{w^T \phi(x)}}$$

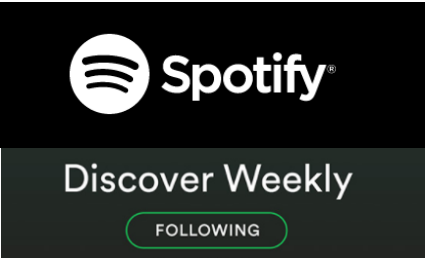
- Multi-task representation learning
- Transfer learning
- Contrastive learning
- Domain adaptation
- Meta-learning
- Theory



# Topic 7: Generative Models

---

- Generative adversarial network
  - Variational Auto-Encoder
  - Energy-based models
  - Normalizing flows
- 



ML uses past data to make predictions



# Supervised Learning Process

Collect a **dataset**

$$\{(x_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} D$$

$x_i$ : input  $\in \mathcal{R}^d$ , image, text  
 $y_i \in \{0, \dots, k\}$  classification  
 $\mathcal{R}$  regression

Decide on a **model**

$$f: \mathcal{R}^d \rightarrow \mathcal{R}$$

Find the function which fits the data best

$(f(x) - y)^2$   
 logistic

**Choose a loss function**  $l(f(x), y) \rightarrow \mathcal{R}$

**Pick the function which minimizes loss**  $\mathcal{R} \in \mathcal{R}^M$   
**on data**

$$\hat{f} \leftarrow \underset{f \in \bar{\mathcal{F}}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i) + \lambda \mathcal{R}(f)$$

Use function to make prediction on new examples

$x_{\text{new}}$   
 prediction:  $\hat{f}(x_{\text{new}}) \approx y_{\text{new}}$

- 1  $f \in \mathcal{F}$
- ! model class
- ① linear  $\theta^T x$
- ② kernel. fixed  $\phi: x \rightarrow x'$   
 $\theta^T x', k(x, x')$
- ③ tree
- ④ NN

$\mathcal{R}(f)$   
 $\mathcal{R}$ : quadratic

# Framework Fix $f \in \mathcal{F}$

Goal: small test error

$$L_{te}(f) = \mathbb{E}_{(x,y) \sim D} [l(f(x), y)]$$

$$L_{tr}(f) = \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$$

$$L_{te}(f) = L_{tr}(f) + L_{te}(f) - L_{tr}(f)$$

$$= \min_{\tilde{f} \in \mathcal{F}} L_{tr}(\tilde{f})$$

$$+ L_{tr}(f) - \min_{\tilde{f} \in \mathcal{F}} L_{tr}(\tilde{f})$$

$$+ \underbrace{L_{te}(f)} - \underbrace{L_{tr}(f)}$$

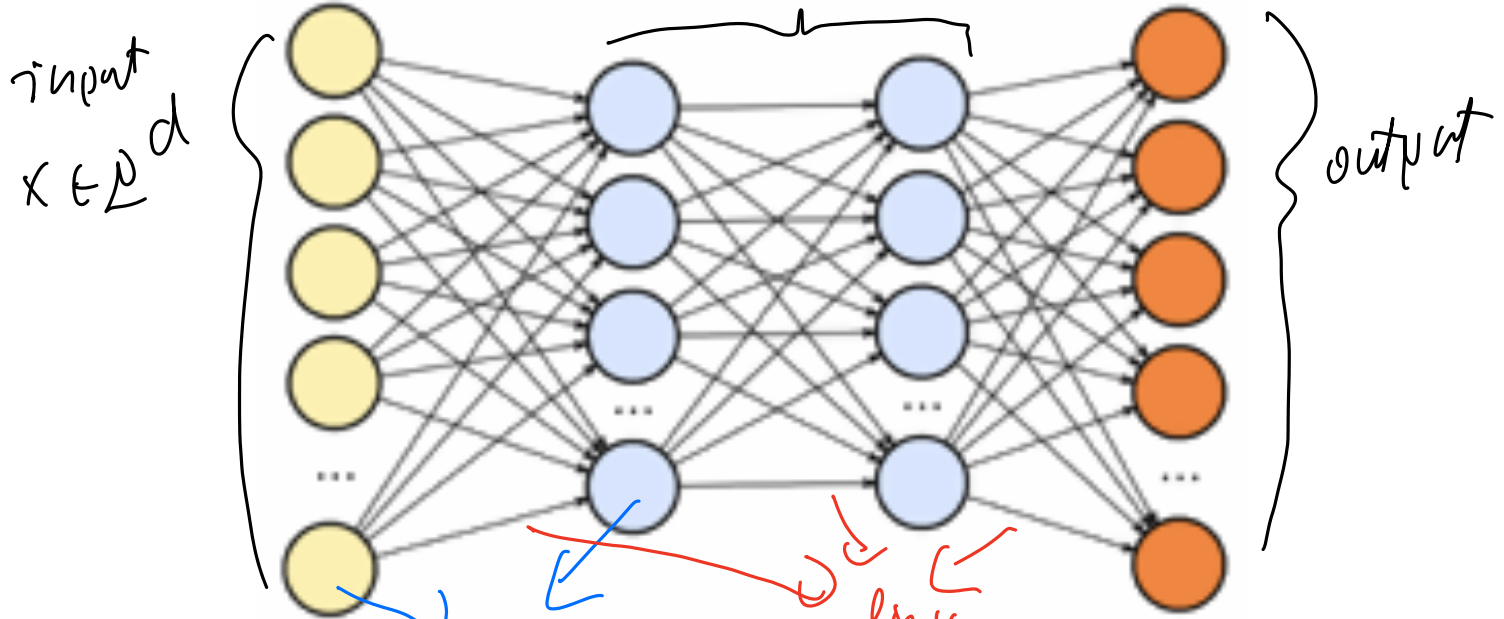
approximation  
error

opt  
error

generalization  
error

# Neural Networks

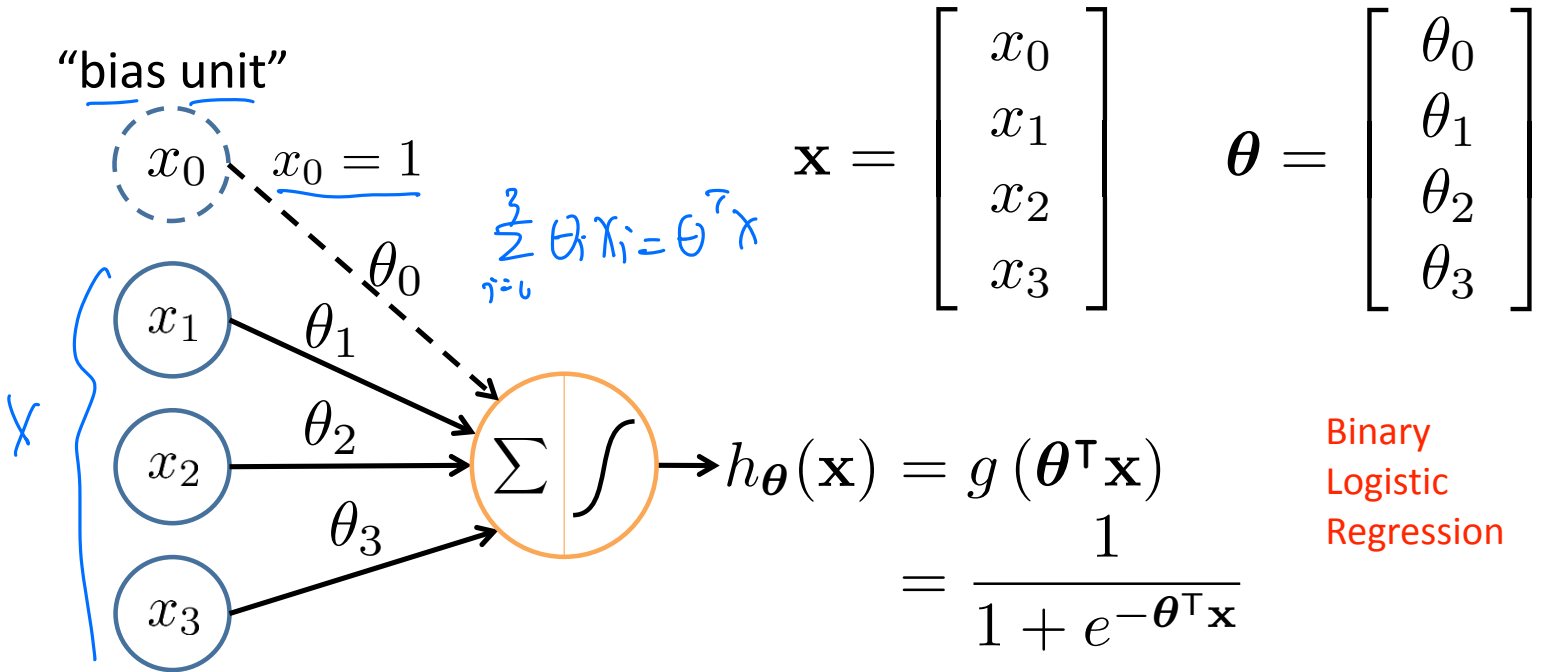
F: fully-connected NN class  
intermediate layers



- each node
- 1) input
  - 2) activation function
  - 3) output

maps the output of nodes to the input of nodes in the next layer  
each link has a weight  $\in \mathbb{R}$

# Single Node

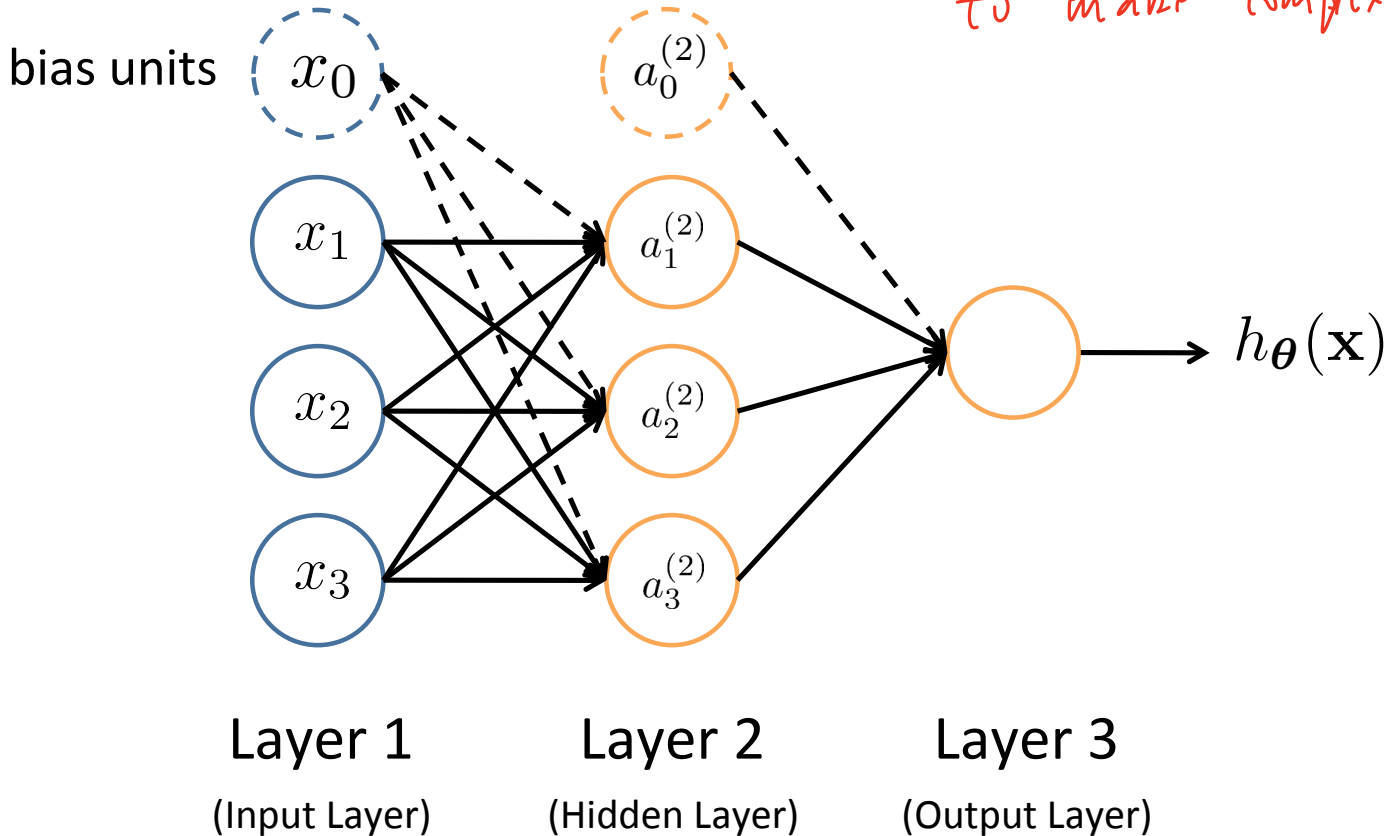


Sigmoid (logistic) activation function:  $g(z) = \frac{1}{1 + e^{-z}}$

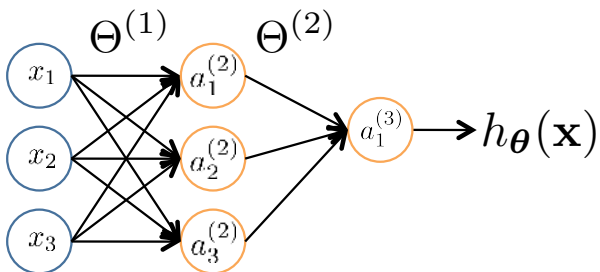
other! ReLU  $\max(z, 0)$

# Neural Network

Key:  
compose single func  
to make complex func



m neurons per layer  
 $m^2$



$a_i^{(j)}$  = "activation" of unit  $i$  in layer  $j$

$\Theta^{(j)}$  = weight matrix stores parameters from layer  $j$  to layer  $j + 1$

$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

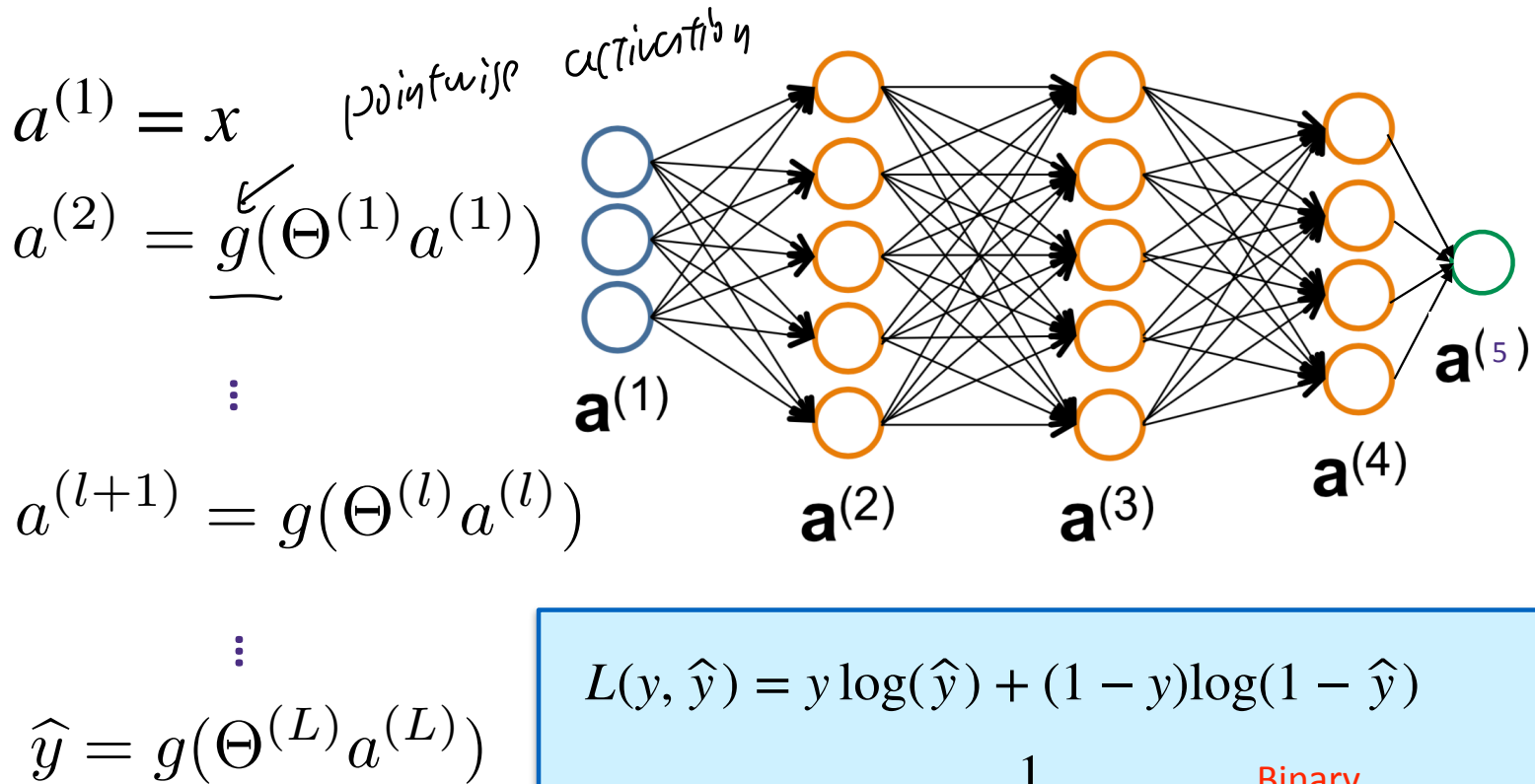
$$h_{\Theta}(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

If network has  $s_j$  units in layer  $j$  and  $s_{j+1}$  units in layer  $j+1$ , then  $\Theta^{(j)}$  has dimension  $s_{j+1} \times (s_j+1)$ .

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4} \quad \Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$



# Multi-layer Neural Network - Binary Classification



$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

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

Binary  
Logistic  
Regression

# Multi-layer Neural Network - Binary Classification

$$a^{(1)} = x \quad x \in \mathcal{L} \cup \mathcal{U}$$

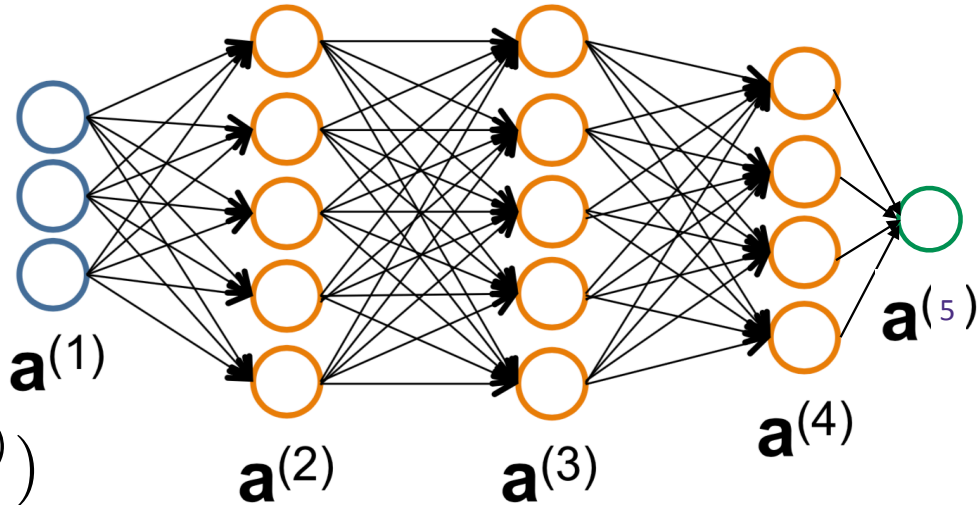
$$a^{(2)} = \underline{\sigma}(\Theta^{(1)} a^{(1)})$$

⋮

$$a^{(l+1)} = \sigma(\Theta^{(l)} a^{(l)})$$

⋮

$$\hat{y} = \underbrace{g(\Theta^{(L)} a^{(L)})}_{\text{logistic}}$$



$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

$$\sigma(z) = \max\{0, z\} \quad g(z) = \frac{1}{1 + e^{-z}}$$

Binary  
Logistic  
Regression

# Multiple Output Units: One-vs-Rest



Pedestrian



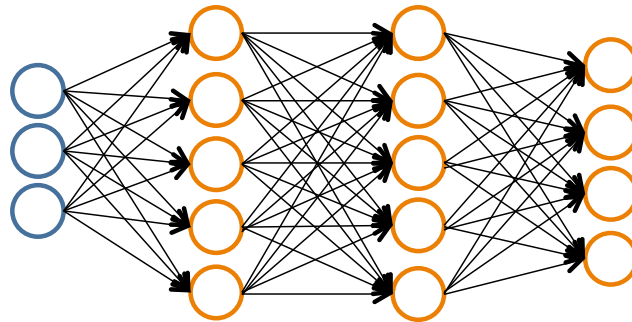
Car



Motorcycle



Truck



*(cross-entropy)*

$$h_{\Theta}(\mathbf{x}) \in \mathbb{R}^K$$

Multi-class  
Logistic  
Regression

We want:

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

when pedestrian

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

when car

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

when motorcycle

$$h_{\Theta}(\mathbf{x}) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

when truck

# Multi-layer Neural Network - Regression

$$a^{(1)} = x$$

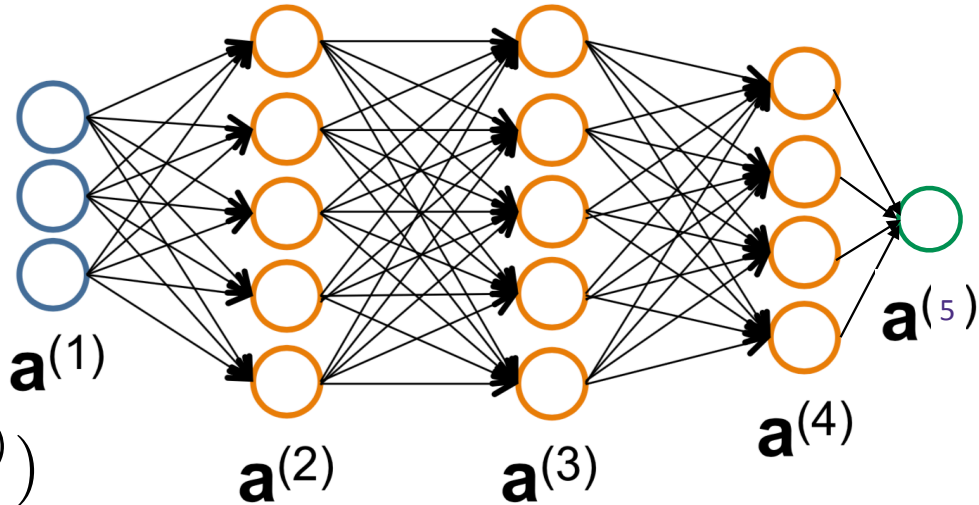
$$a^{(2)} = \sigma(\Theta^{(1)} a^{(1)})$$

⋮

$$a^{(l+1)} = \sigma(\Theta^{(l)} a^{(l)})$$

⋮

$$\hat{y} = \Theta^{(L)} a^{(L)}$$



$$L(y, \hat{y}) = (y - \hat{y})^2$$

$$\sigma(z) = \max\{0, z\}$$

Regression

$$\mathcal{F} = \{f: \mathcal{R}^n \rightarrow \mathcal{R}\}$$

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)})$$

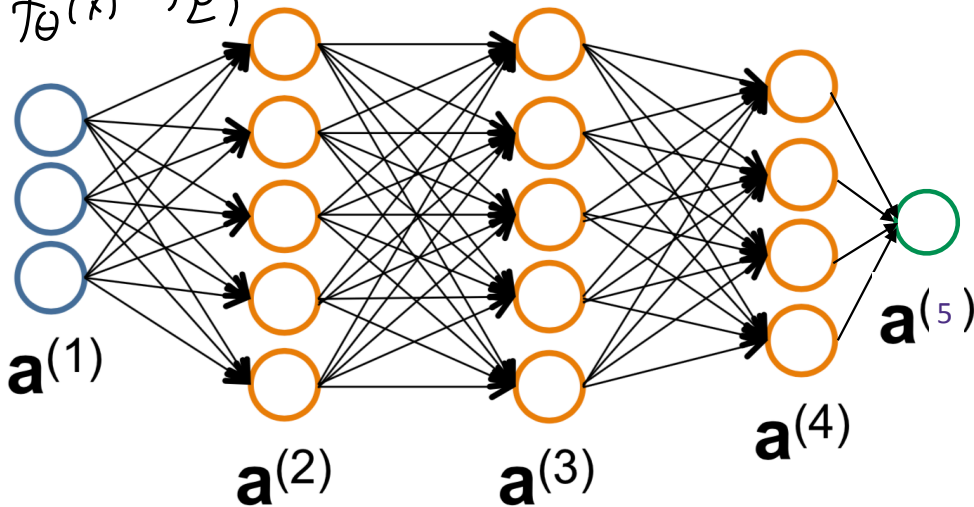
$$\vdots$$

$$z^{(l+1)} = \Theta^{(l)} a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

$$\vdots$$

$$\hat{y} = g(\Theta^{(L)} a^{(L)})$$



$$L(y, \hat{y}) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$$

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

$\Theta^{(0)}$   $\eta$ : step size / learning rate iteration

Gradient Descent:  $\Theta^{(l)} \leftarrow \Theta^{(l)} - \eta \cdot \nabla_{\Theta^{(l)}} L(y, \hat{y}) \quad \forall \underline{l}$

Gradient Descent:  $\Theta^{(l)} \leftarrow \Theta^{(l)} - \eta \nabla_{\Theta^{(l)}} L(y, \hat{y}) \quad \forall l$

Seems simple enough, why are packages like PyTorch, Tensorflow, Theano, Cafe, MxNet synonymous with deep learning?

1. Automatic differentiation

2. Convenient libraries

3. GPU support

## Gradient Descent:

Seems simple enough,  
Theano, Cafe, MxNet s

1. Automatic differ

2. Convenient libra

```
class Net(nn.Module):

    def __init__(self):
        super(Net, self).__init__()
        # 1 input image channel, 6 output channels, 3x3 square convolution
        # kernel
        self.conv1 = nn.Conv2d(1, 6, 3)
        self.conv2 = nn.Conv2d(6, 16, 3)
        # an affine operation: y = Wx + b
        self.fc1 = nn.Linear(16 * 6 * 6, 120) # 6*6 from image dimension
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        # Max pooling over a (2, 2) window
        x = F.max_pool2d(F.relu(self.conv1(x)), (2, 2))
        # If the size is a square you can only specify a single number
        x = F.max_pool2d(F.relu(self.conv2(x)), 2)
        x = x.view(-1, self.num_flat_features(x))
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x
```

```
# create your optimizer
optimizer = optim.SGD(net.parameters(), lr=0.01)

# in your training loop:
optimizer.zero_grad() # zero the gradient buffers
output = net(input)
loss = criterion(output, target)
loss.backward()
optimizer.step() # Does the update
```