

# CSCI 566: Deep Learning and Its Applications

Jesse Thomason

Lecture 4: Activations, Initializations,  
Optimization, and Regularization

# CSCI 566 Resources

- **Website:** <https://csci566-spring2023.github.io/>
- **Piazza:** <https://piazza.com/usc/spring2023/csci566/info>
- **GDrive:**  
[https://drive.google.com/drive/folders/11YcbSZcJUdRSbrNBWiOs\\_HM2VkxdNKWS?usp=share\\_link](https://drive.google.com/drive/folders/11YcbSZcJUdRSbrNBWiOs_HM2VkxdNKWS?usp=share_link)
  - Includes spreadsheet with lecture slide + recording links!
- **Zoom:**  
<https://usc.zoom.us/j/95538924150?pwd=TTRBM283VCt4WD FkN1hLMmtrKzdGUT09>

# Course Project Team Registration

- Your 4-person teams should be registered now
- If you have not registered your team, fill the form ASAP!
  - [Piazza -> Resources -> Course Project -> Project Team Registration](#)
- If you haven't *found* a team, come up to the front during the break to find all the other students who haven't found a team.

# Map to the Midterm

## January 2023

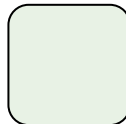
| Sunday | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday |
|--------|--------|---------|-----------|----------|--------|----------|
| 1      | 2      | 3       | 4         | 5        | 6      | 7        |
| 8      | 9      | 10      | 11        | 12       | 13     | 14       |
| 15     | 16     | 17      | 18        | 19       | 20     | 21       |
| 22     | 23     | 24      | 25        | 26       | 27     | 28       |
| 29     | 30     | 31      |           |          |        |          |

[www.a-printable-calendar.com](http://www.a-printable-calendar.com)

## February 2023

| Sunday | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday |
|--------|--------|---------|-----------|----------|--------|----------|
|        |        |         | 1         | 2        | 3      | 4        |
| 5      | 6      | 7       | 8         | 9        | 10     | 11       |
| 12     | 13     | 14      | 15        | 16       | 17     | 18       |
| 19     | 20     | 21      | 22        | 23       | 24     | 25       |
| 26     | 27     | 28      |           |          |        |          |

[www.a-printable-calendar.com](http://www.a-printable-calendar.com)



### Module 1: Neural Network Basics

# [Feb 10] A Look at February Deliverables

## February 2023

|                      |        |
|----------------------|--------|
| Project Teams Formed | Feb 3  |
| Assignment 1 Out     | Feb 10 |
| Project Proposal Due | Feb 17 |
| Assignment 1 Due     | Feb 24 |
| Midterm Exam         | Feb 24 |

| Sunday | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday |
|--------|--------|---------|-----------|----------|--------|----------|
|        |        |         | 1         | 2        | 3      | 4        |
| 5      | 6      | 7       | 8         | 9        | 10     | 11       |
| 12     | 13     | 14      | 15        | 16       | 17     | 18       |
| 19     | 20     | 21      | 22        | 23       | 24     | 25       |
| 26     | 27     | 28      |           |          |        |          |

## Next Week: Course Project Proposal Due [Feb 17]

- **Proposal:** 5-6 *slides* on problem (scope & definition) + "today" (status of literature) + challenges + directions to innovate.
  - Look at the project ideas doc populated by the TAs to get a sense of what material your proposal should cover for problem definition and background
  - Include what you hope to be your *key insight or contribution*; what will make your approach unique?
- We will share a link to collect the slides.
- You will condense these longer, more detailed slides into a single lightning pitch slide for the Project Pitch assignment.

# Overview of Today's Plan

- ~~Course organization and deliverables~~
  - Any questions before we move on?
- Activation Functions and Preprocessing
- Optimization and Regularization
- Cloud Computing Service Tutorial
- Deep Learning Software Tutorial
- Coding Assignment 1 Release Briefing

$$f: X \rightarrow Y;$$
$$f(x)=y$$

$$f: \phi_x(X) \rightarrow \phi_y(Y);$$
$$f(\mathbf{x})=\mathbf{y}$$

$$f: X \times \Theta \rightarrow Y;$$
$$f(x;\theta)=y; \theta \in \Theta \text{ for model } M$$

$$f: \phi_x(X) \times \Theta \rightarrow \phi_y(Y);$$
$$f(\mathbf{x};\theta)=\mathbf{y}$$

$$\theta^* = \min_{\theta \in \Theta} (\sum_{(x,y)} L(f(\mathbf{x};\theta), \mathbf{y})))$$

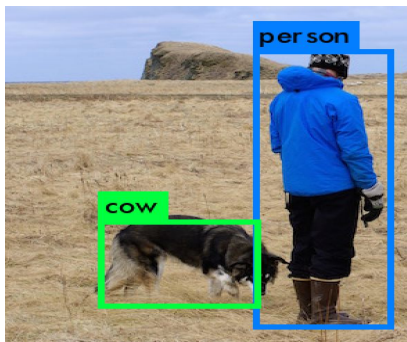
$$\theta^* = \min_{\theta \in \Theta} (\sum_{(x,y) \in D} L(f(\mathbf{x};\theta), \mathbf{y})))$$

- DL enables function approximation for  $f$  and is constrained by:

- Data representation  $\phi$
- Model architecture  $M$
- Loss function  $L$
- Optimization algorithm
- Training data  $D$



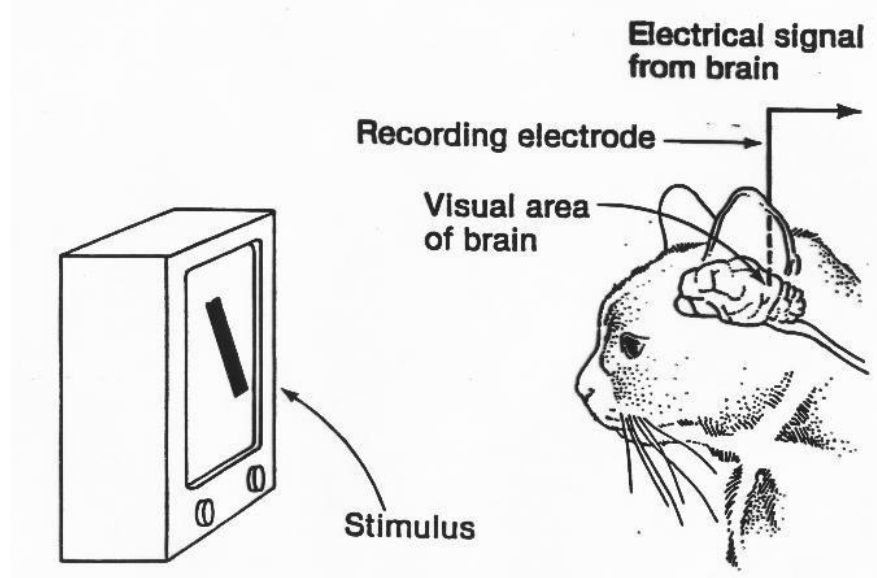
# Representations and Models



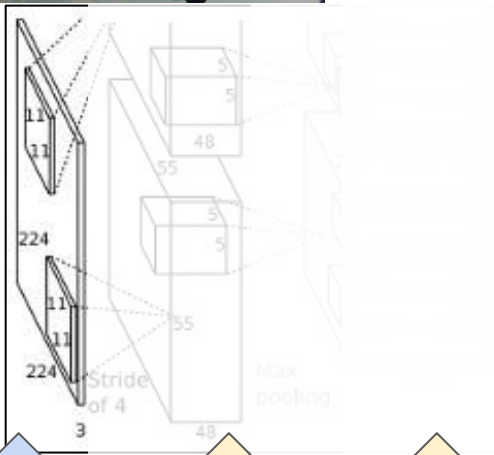
- Input space  $X$ ?
  - Pixels in the image
- Representation?
  - $\phi_x(X) = \mathbf{R}^{W,H,3}$  RGB or HSV
- Output space  $Y$ ?
  - Object classes  $C$  + bboxes
- Representation?
  - $\phi_y(Y) = C \times \mathbf{R}^{x,y,w,h} \times N$
- What model architecture might you use?
- What loss function would you want to optimize?

# Origins of Convolutional Neural Networks

“They found that particular neurons in the visual cortexes of cats and monkeys—the areas in their brains responsible for processing visual information—didn't respond to simple points of light, but rather to lines, and in particular, lines and contours with specific orientations”



# Breaking Down AlexNet



|     |          |
|-----|----------|
| .01 | car      |
| .03 | dog      |
| .63 | person   |
| .11 | airplane |
| .02 | flower   |
| ... | ...      |

Image

Conv

Pool

Conv

Pool

Conv

Conv

Conv

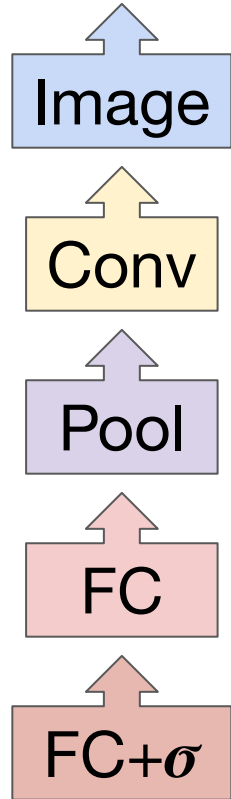
Pool

FC

FC

FC+σ

# Common Layers in a Deep Convolutional Neural Net



- $\mathbf{x}$  in our parlance; the vectorized input

- Convolution

- Pooling layer (in AlexNet, MaxPool)

- Fully-connected layer (i.e., *perceptron* layer)

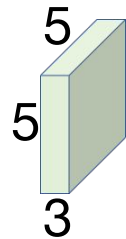
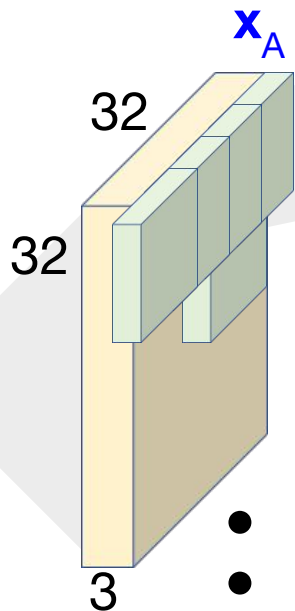
- Softmax layer (fully connected layer + exponentiated normalization)

Each layer is a differentiable function whose *input* is the *output* of the previous layer (function)

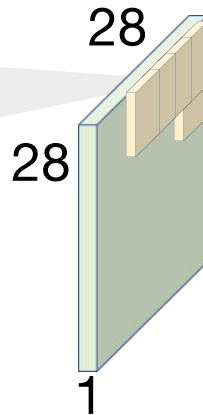
# Convolutional Layer



$$\phi_x(\text{image}) = [0, 1]^{(32 \times 32 \times 3)}$$



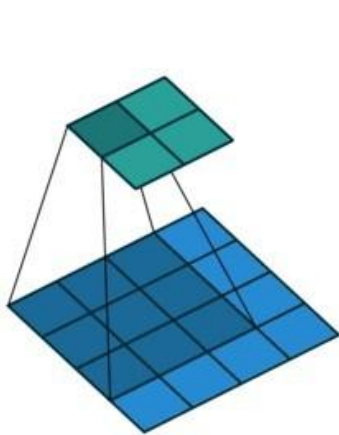
$$f_c(\mathbf{x}_A) = W\mathbf{x}_A + \mathbf{b}$$



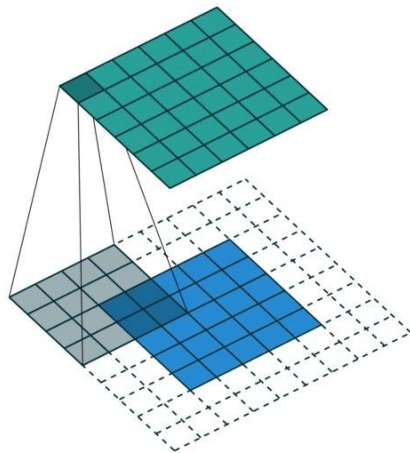
- $\mathbf{x}_A$  is a  $5 \times 5 \times 3$  subset of image tensor  $\mathbf{x}$
- Filter function  $f_c$  a learned, linear function over flattened  $\mathbf{x}_A$
- $W$  is  $1 \times (5 \times 5 \times 3)$ ,  $\mathbf{b}$  is size 1, so  $f_c$  is a single number!

# Convolutional Layer Terminology

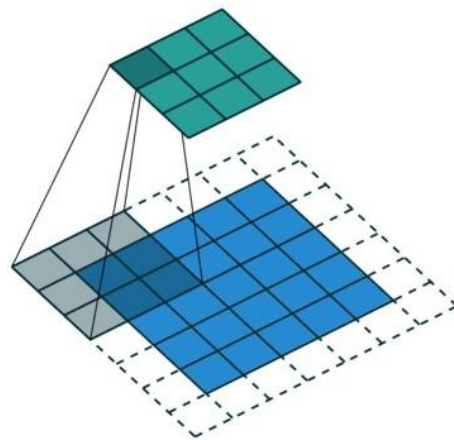
## Padding and Stride



Pad=0; Stride=1  
Kernel size?  
3x3; lose 2 w/h

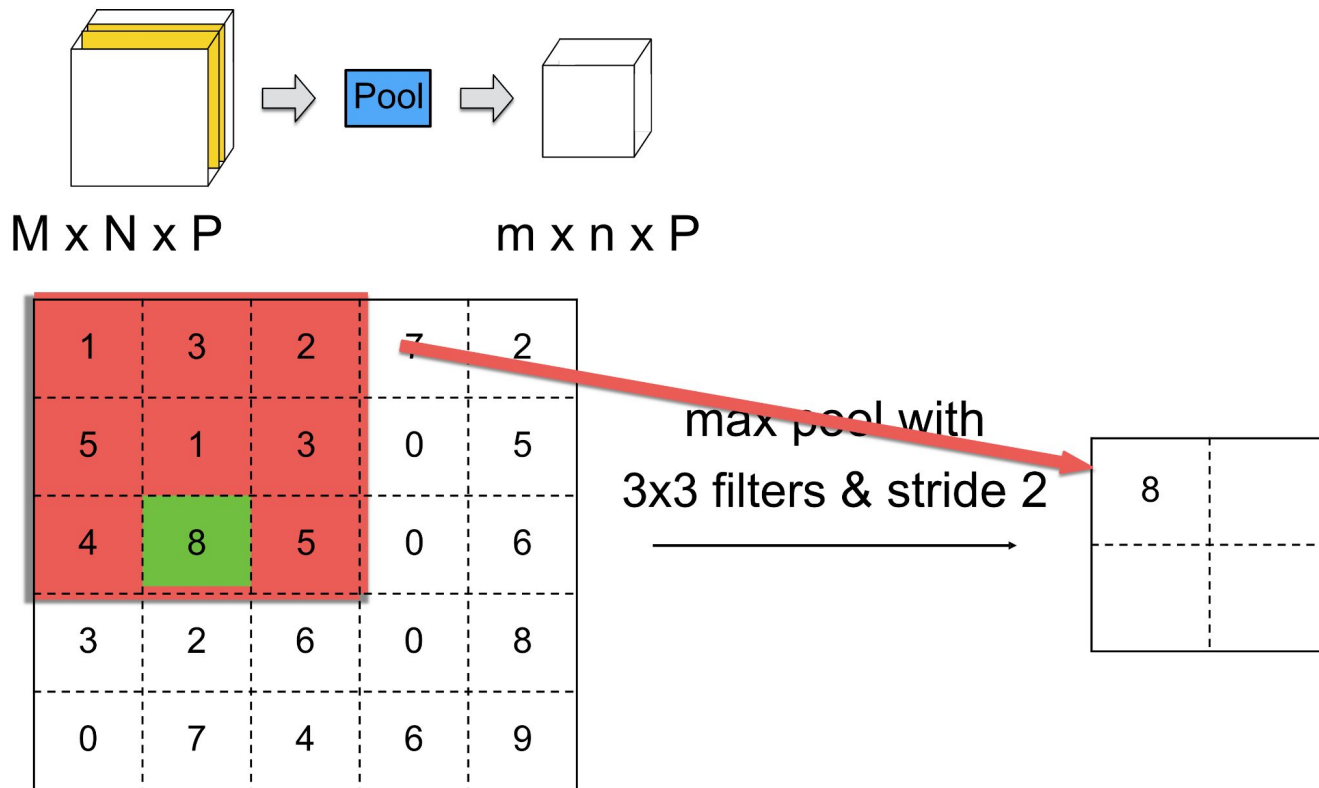


Pad=2; Stride=1  
What to pad with?  
Zeros common



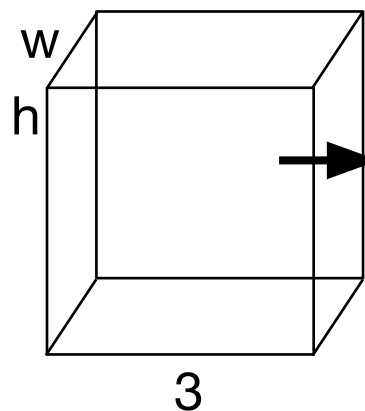
Pad=1; Stride=2  
Effect of stride?  
Sparser output

# Max Pooling Layer

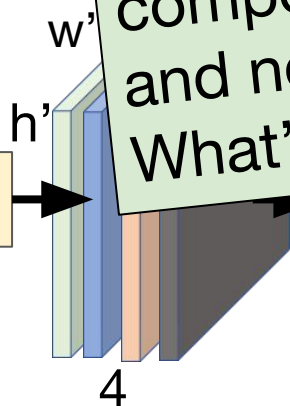


# Image Classification with a Conv and Pool Layers

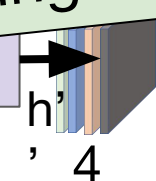
$$\phi_x(\text{image}) = [0, 1]^{(w \times h \times 3)}$$



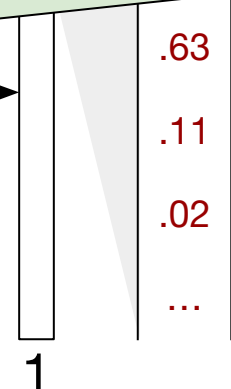
Conv



Pool



FC+σ



.63 person  
.11 airplane  
.02 flower  
...

A neural network is a composition of a composition of matrix multiplications and non-linear transformations. What's missing here?

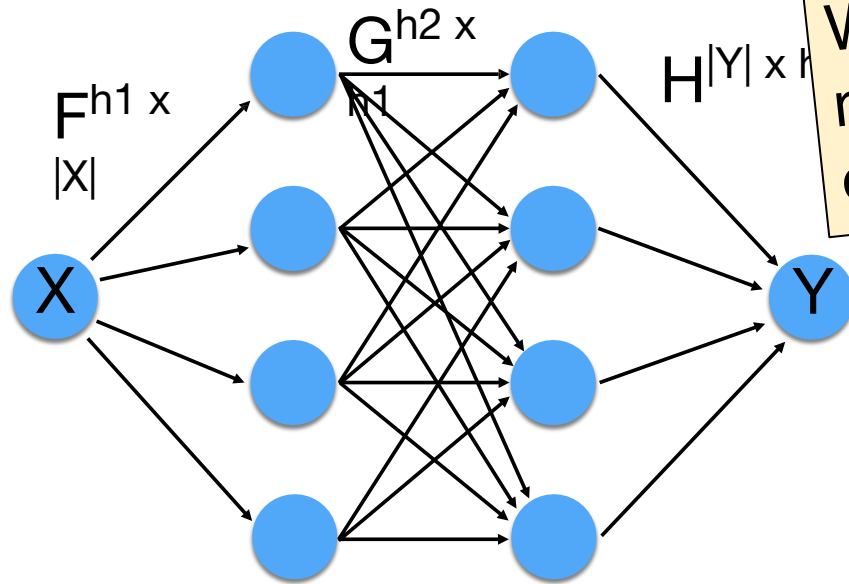
$$f_c(\mathbf{x}_A) = W_c \mathbf{x}_A + \mathbf{b}_c$$

$$f_p(f_c(\bullet)) = \max(f_c(\bullet))$$

$$FC(f_p(f_c(\bullet))) = W_{out} f_p(f_c(\bullet)) + \mathbf{b}_{out}$$

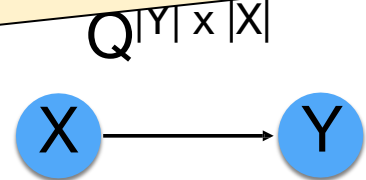
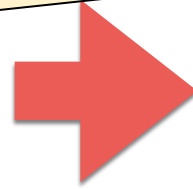


# Why Utilize Non-Linearities in a Deep Neural Network?



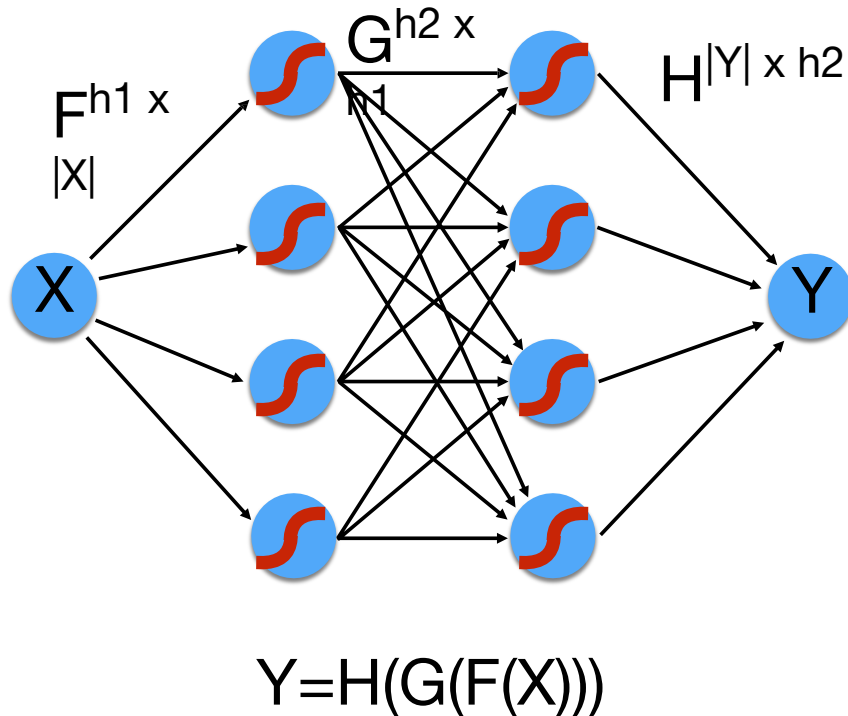
$$Y = H(G(F(X)))$$

Without nonlinear activations, multiple layers are reduced to one matrix multiplication.



$$\text{Define } Q^{|Y| \times |X|} = HGF$$
$$\text{Then } Y = QX$$

# Why Utilize Non-Linearities in a Deep Neural Network?



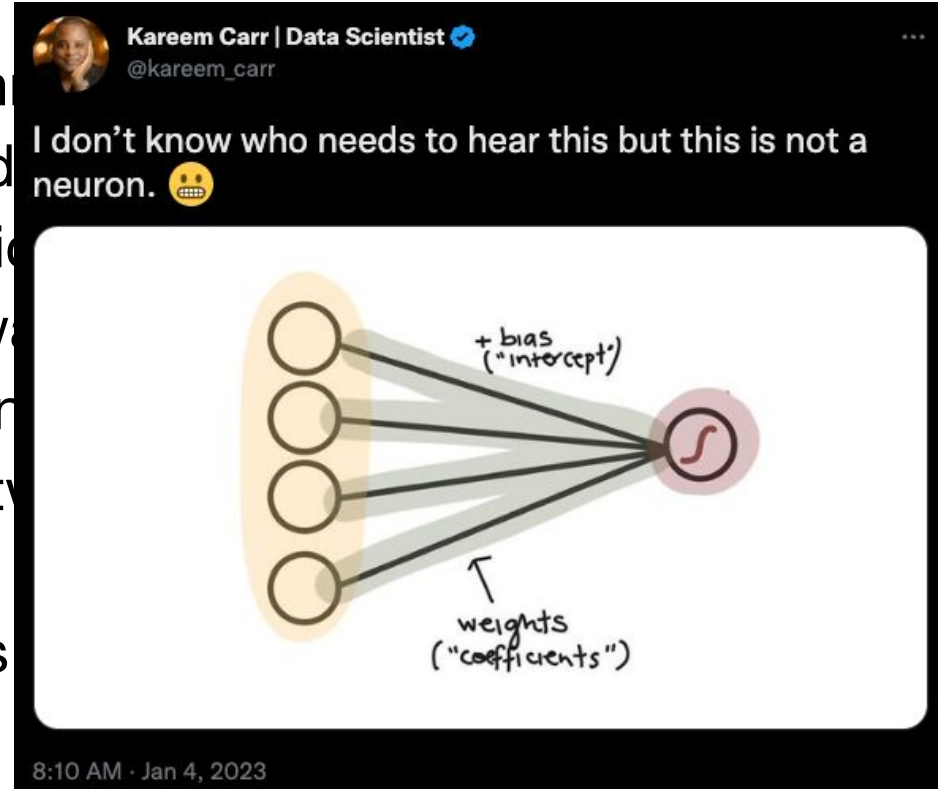
$$F(\mathbf{x}) = a(W_F(\mathbf{x}) + \mathbf{b}_F)$$

$$G(\mathbf{x}_{h1}) = a(W_G(\mathbf{x}_{h1}) + \mathbf{b}_G)$$

*Activation* function  $a$  is non-linear and not expressible as a matrix multiplication.

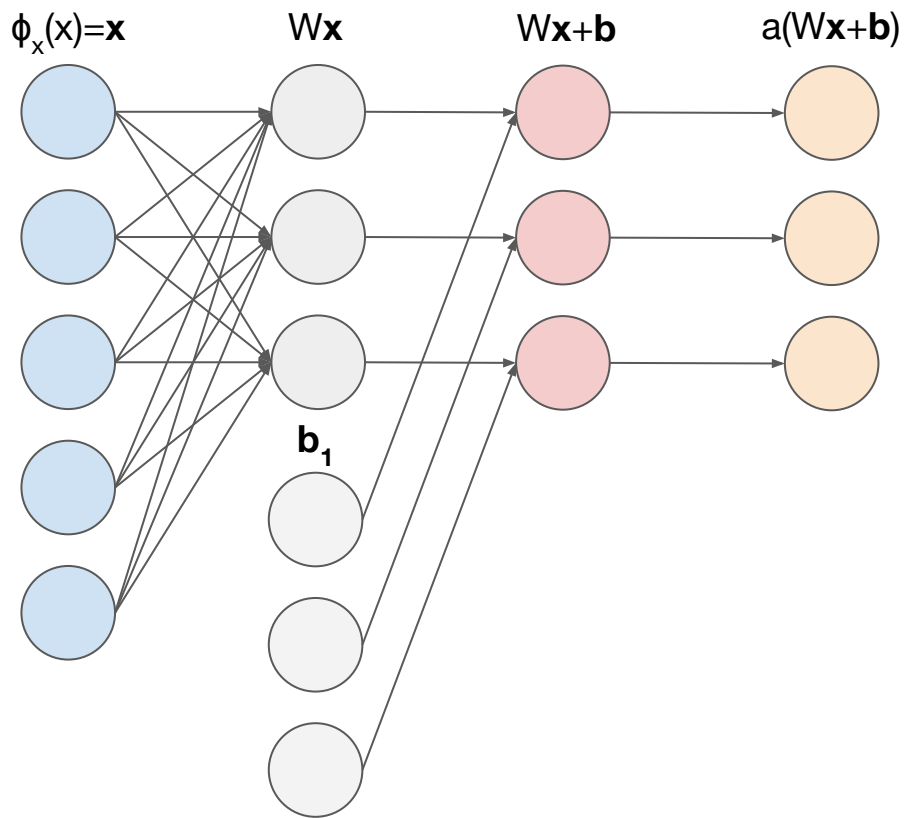
# Deep Learning Borrows\* Terminology from Neuroscience

- Brains are complex, interconnected
- Neurons are connected via dendrites through chemical and electrical synapses
- “Artificial neural networks” with input variables “**activated**” by signals
- In a “feed forward neural network” neurons are activated in order from input to output
- A “deep neural network” has multiple hidden layers



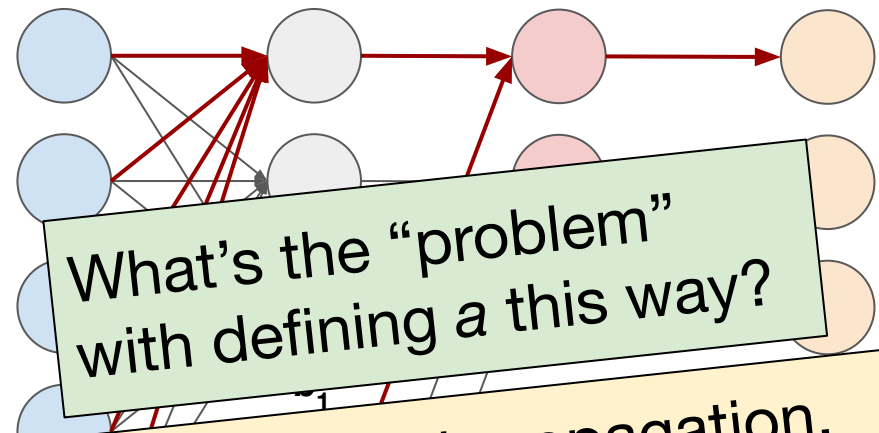
\*Abuses, one might say.

# Activation Functions



# Activation Functions

$\phi_x(x)=x$        $Wx$        $Wx+b$        $a(Wx+b)$

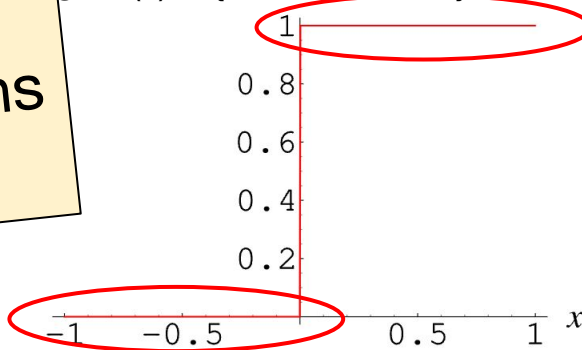


$$= a(W_{1,1}x_1 + W_{1,2}x_2 + W_{1,3}x_3 + W_{1,4}x_4 + W_{1,5}x_5 + b_1)$$

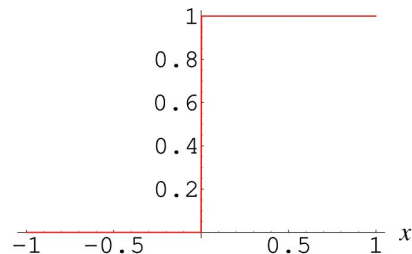
If we think about a “neuron” as inspiration, the output after a should be 1 if  $W_{1,1}x_1 + W_{1,2}x_2 + W_{1,3}x_3 + W_{1,4}x_4 + W_{1,5}x_5 + b_1$  exceeds a threshold and 0 else.

To perform backpropagation, we need differentiable functions with non-zero gradients.

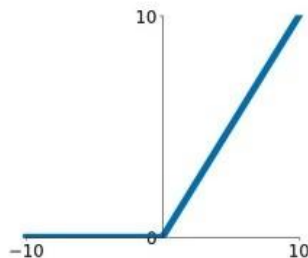
g.,  $a(x) = \{1 \text{ if } x > 0 \text{ else } 0\}$



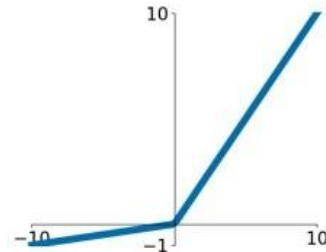
# Activation Functions



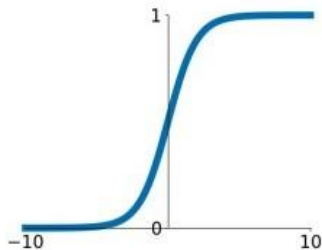
**ReLU**  
 $\max(0, x)$



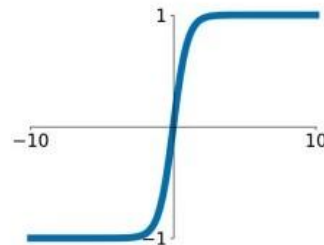
**Leaky ReLU**  
 $\max(0.1x, x)$



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



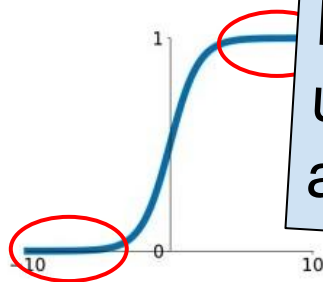
**tanh**  
 $\tanh(x)$



# Activation Functions: Sigmoid

## Sigmoid

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



Note: annoy  
uses  $\sigma$  to der  
also sigmoid

- Advantages:

- “Looks” a lot like our step function!

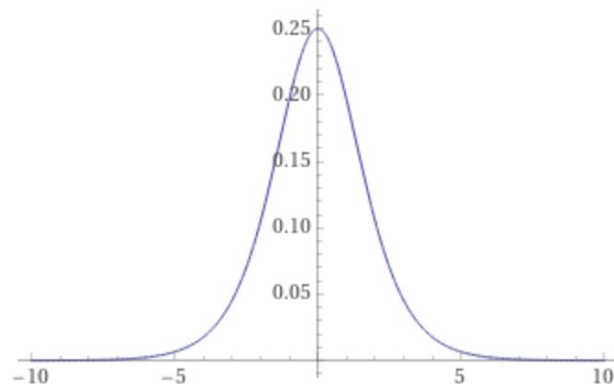
- Disadvantages:

Learning can become prohibitively slow, especially if multiple sigmoids are involved across layers!

$$\frac{d}{dx} \sigma(x) = \sigma(x)(1-\sigma(x))$$
 [\[derivation\]](#)

|      |  |                          |
|------|--|--------------------------|
| plot | $\frac{1}{1+e^{-x}} \left(1 - \frac{1}{1+e^{-x}}\right)$ | $x = -10 \text{ to } 10$ |
|------|--|--------------------------|

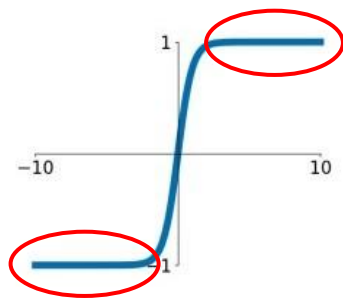
Plot



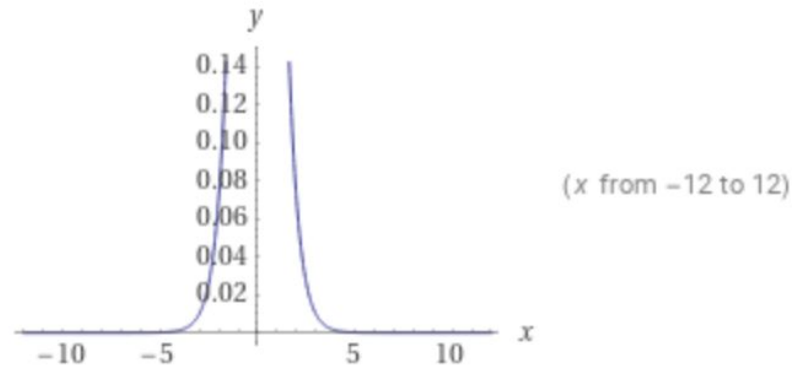
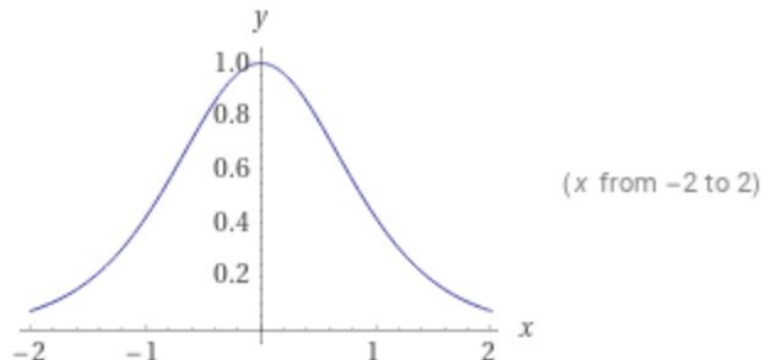
anywhere

# Activation Functions: Hyperbolic Tangent

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



- Advantages:
  - “Looks” a lot like our step fun
- Disadvantages:
  - Derivative is near 0 almost ev
  - What is  $\frac{d}{dx} \tanh(x)$ ?
    - $\frac{d}{dx} \tanh(x) = \text{sech}^2(x)$

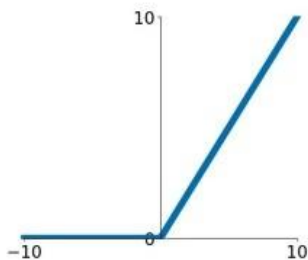




# Activation Functions: Rectified Linear Unit (ReLU)

**ReLU**

$$\max(0, x)$$



Learning can stall if activations are below zero because gradient stops flowing.

- Advantages:

- For  $x > 0$ , reliable and cheap to calculate gradient
- What is  $\frac{d}{dx} \max(0, x)$ ?
  - $\frac{d}{dx} \max(0, x) = \{1 \text{ if } x > 0; \text{ else } 0\}$

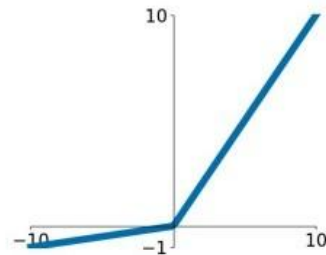
- Disadvantages:

- Derivative is exactly 0 for  $x < 0$
- Output in the positive direction is unbounded

# Activation Functions: Rectified Linear Unit (ReLU)

## Leaky ReLU

$$\max(0.1x, x)$$



- Advantages:

- Reliable and cheap to calculate gradient
- What is  $\frac{d}{dx} \max(0.1x, x)$ ?
  - $\frac{d}{dx} \max(0.1x, x) = \{1 \text{ if } x > 0; \text{ else } 0.1\}$

- Disadvantages:

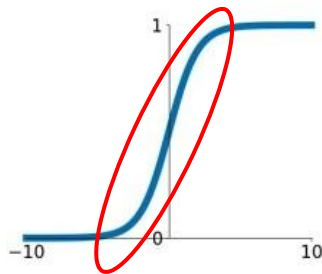
- As with ReLU, outputs are still *unbounded* and could cause problems in later layers



# What Assumption Do Our Activation Functions Make?

## Sigmoid

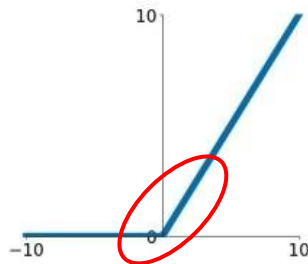
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



Our activation functions are all on their best behavior when the input  $x$  is close to zero.

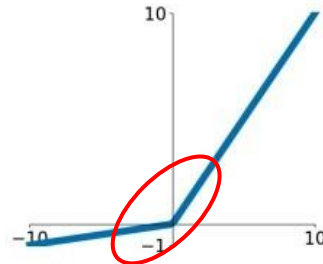
## ReLU

$$\max(0, x)$$



## Leaky ReLU

$$\max(0.1x, x)$$



# Weight Initialization: Zeros

- Since input to activation functions behaves well around zero, how about we start all our weights at zero?
- Initialize all weight matrix  $W$  and bias vector  $\mathbf{b}$  parameters are initialized to zero for all layers (e.g., fully connected, CNN filters, etc.):
  - What is the effect on the final output vector?
    - Guaranteed to be all zeros for any input!
  - Gradients will vanish; model can't learn

# Weight Initialization

- We want well-behaved, random initial parameter values that are near zero but not uniformly zero
- What's the easiest way to accomplish that goal?
  - Draws from a Gaussian distribution!

$$W \sim \mathcal{N}(\mu, \sigma^2)$$

- What's a good value for  $\mu$ ?
  - $\mu=0$
- What's a good value for  $\sigma^2$ ?
  - $0 < \sigma^2 < 1$ , e.g.,  $\sigma^2=0.01$

# Weight Initialization: Gaussian Distribution

- Draws from a Gaussian distribution

$$W \sim \mathcal{N}(\mu, \sigma^2)$$

- Advantages:
  - Easy to implement
- Disadvantages:
  - Hyperparameter we need to set:  $\sigma^2$
  - If the weights are too small, the activation signals shrink quickly during training
  - If the weights are too large, the activation signals explode during training

# Weight Initialization: Controlling Activation Variance

- Glorot and Bengio argue that when the variance of the layer outputs (and hence the downstream layer inputs) is not  $\approx 1$ , depending on the activation function, models will converge more slowly, especially when the layer output variance is  $< 1$
- “Fixes” that initialize parameters to encourage layer outputs whose variance is  $\approx 1$  depend the activation function



# Weight Initialization: Controlling Activation Variance

- Activation functions with centered average (sigmoid; tanh):
  - Set the variance to inverse of the number of input units that contribute to each logit in the layer [*Glorot Initialization*]

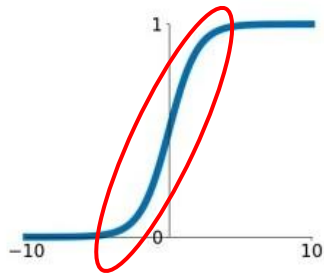
$$W \sim \mathcal{N}(\mu, \sigma^2) \longleftarrow Var(W) = \frac{1}{n_{in}}$$

- Activation functions with unbounded average (ReLU):
  - Because average activation is a higher positive number, need to adjust up to  $Var(W) = 2/n_{in}$  [*Kaiming initialization*]

# What Assumption Do Our Activation Functions Make?

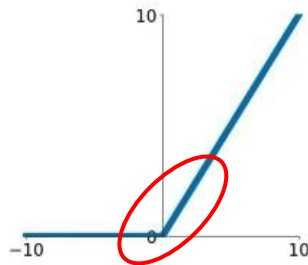
## Sigmoid

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



## ReLU

$$\max(0, x)$$

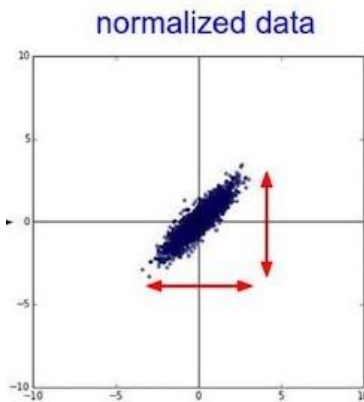
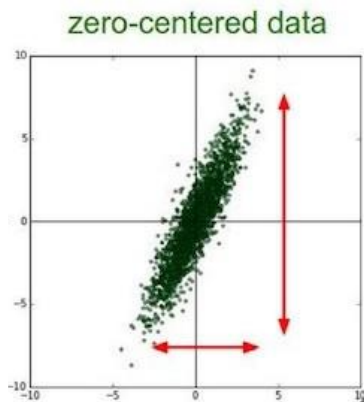
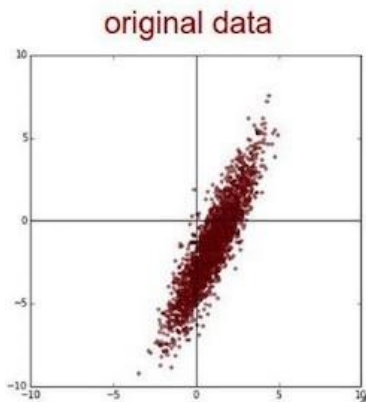


Our activation functions are all on their best behavior when the input  $x$  is close to zero.

What aspect of the path from input to output won't weight initialization alone help with?

# Input Representation and Data Preprocessing

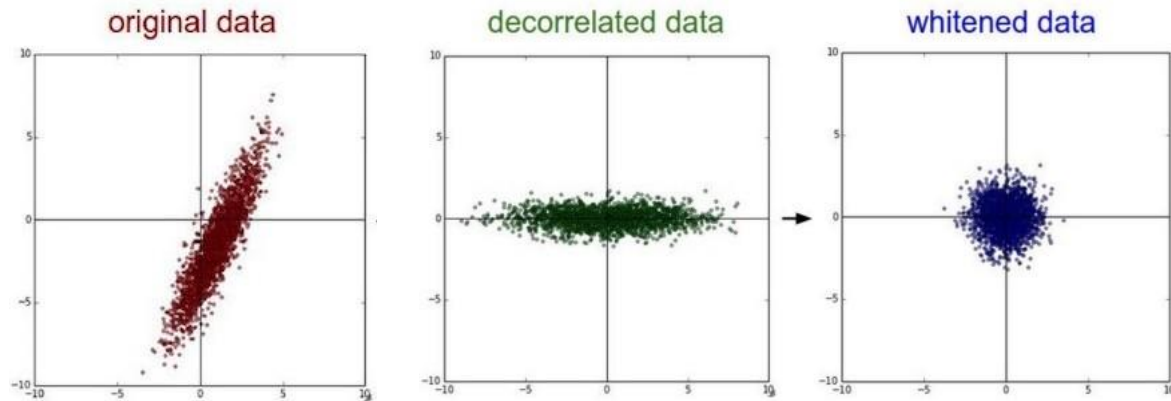
- We can make our input representation  $\phi$  a function that considers the entire training dataset to *standardize* it
- Zero-centering: shift the data to have features with mean 0
- Normalization: adjust features to have variance 1
  - (e.g., z-score normalization)



$$z = \frac{x - \mu}{\sigma}$$

# Input Representation and Data Preprocessing

- We can make our input representation  $\phi$  a function that considers the entire training dataset to *standardize* it
- Remove correlations between input features via dimensionality reduction (e.g., PCA)



# Overview of Today's Plan

- ~~Course organization and deliverables~~
- ~~Activation Functions and Preprocessing~~
  - Any questions before we move on?
- Optimization and Regularization
- Cloud Computing Service Tutorial
- Deep Learning Software Tutorial
- Coding Assignment 1 Release Briefing

# Gradient Descent

- Vanilla (true) Gradient Descent

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t)$$

- Estimate  $\theta$  over the whole dataset

- Stochastic Gradient Descent (SGD)

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t; x^{(i)})$$

- Estimate  $\theta$  at a random datapoint

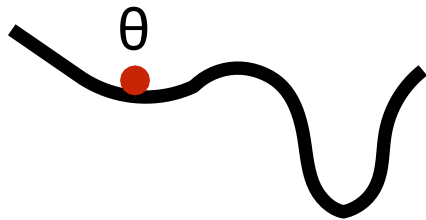
- Minibatch Gradient Descent

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t; x^{(i:i+n)})$$

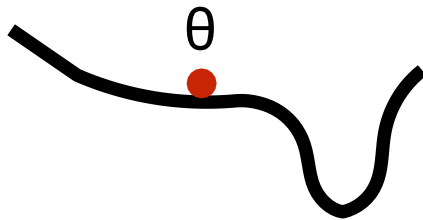
- Estimate  $\theta$  as an average over a random batch of training data
- Commonly we use “SGD” to refer to this setup

# Saddles and Local Minima

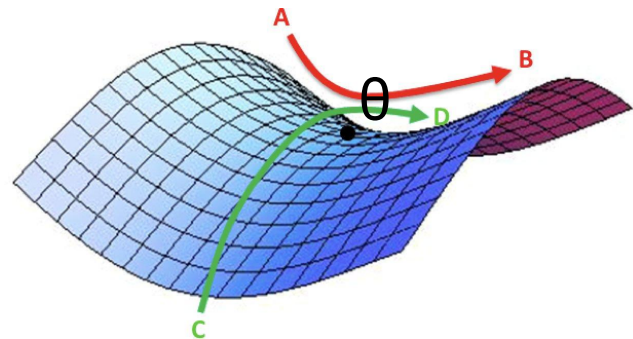
- Gradient descent drives the parameters  $\theta$  towards the closest low point in the loss space
- Gets “stuck” if the loss is locally minimal (basin) or the gradient goes in orthogonal directions across dimensions (saddle)



local minimum



saddle points

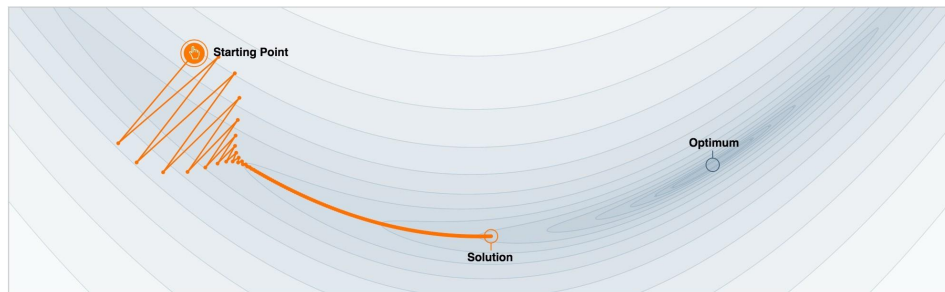


# Gradient Descent with Momentum

- Consider the history:
  - Add a fraction  $\gamma$  of the update vector of the past time step to the current update vector

$$\nu_t = \boxed{\gamma \nu_{t-1}} + \eta \nabla_{\theta_t} f(\theta_t) \quad \theta_{t+1} = \theta_t - \nu_t$$

GD *without* momentum



GD *with* momentum



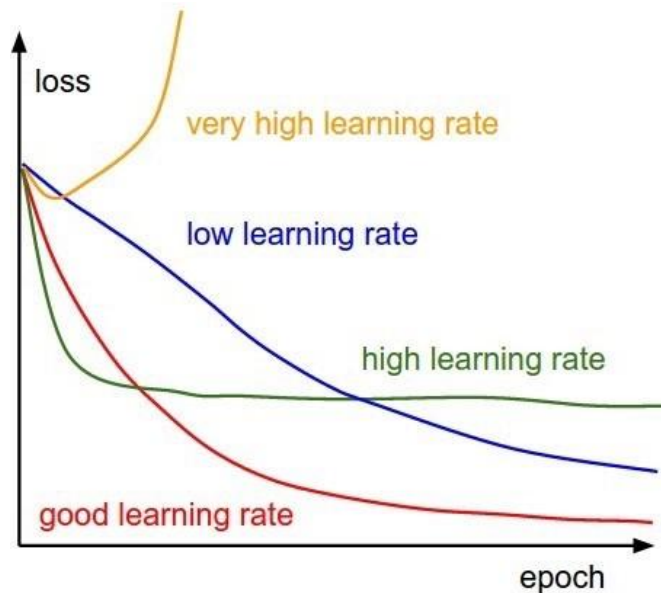


# Optimizer Parameters: Learning Rate

- How fast or slow we update the model parameters

$$\theta_{t+1} = \theta_t - \eta_t \nabla f(\theta_t)$$

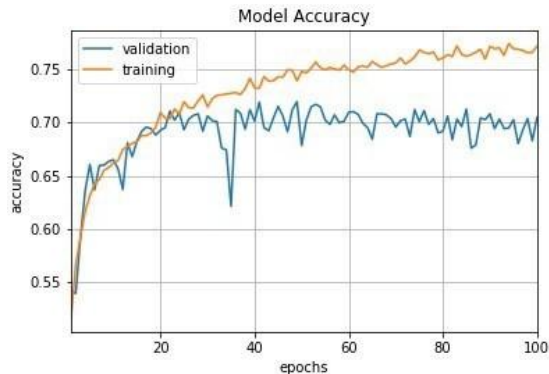
- If the learning rate is too large, the model parameters “explode” (move too fast; gain in magnitude)
- If the learning rate is too small, the model parameters move slowly and may be unable to escape local minima



# Optimizer Parameters: Learning Rate

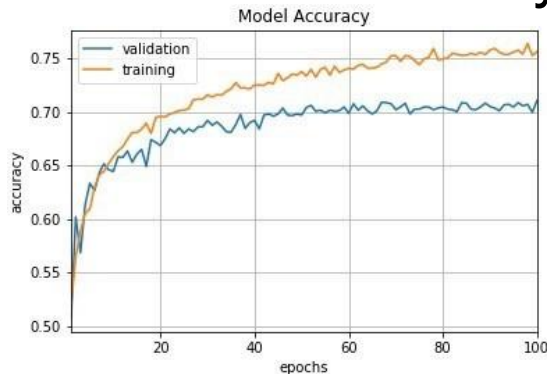
- Learning rate *schedulers* change the LR during training

## Constant



$$lr = lr_0$$

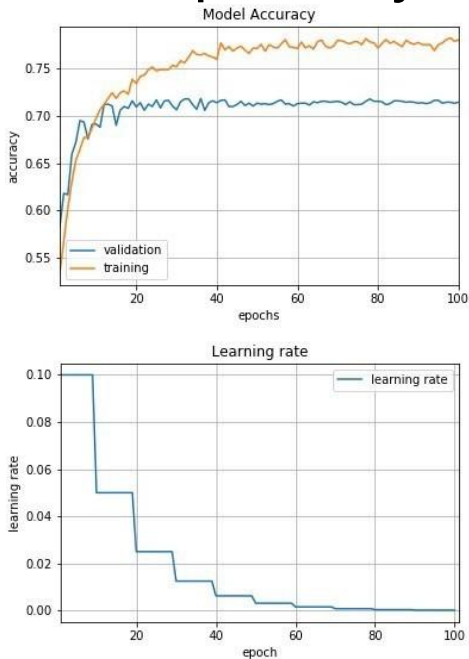
## Time-based decay



$$lr = \frac{lr_0}{1 + k * iter}$$

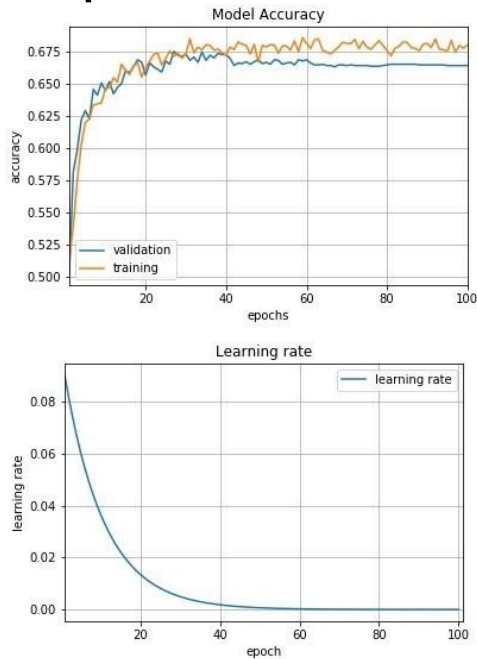
# Optimizer Parameters: Learning Rate

## Step decay



$$lr = lr_0 * floor(\frac{epoch}{epoch\_drops})$$

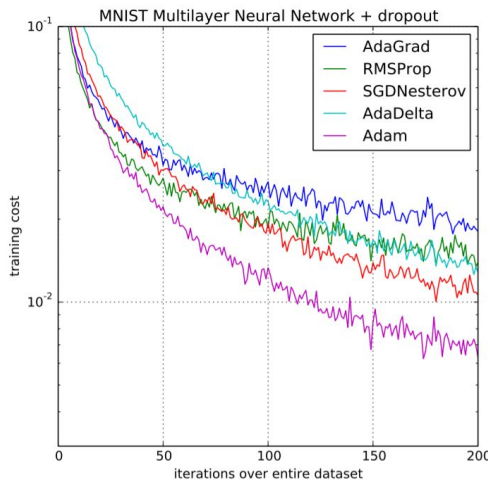
## Exponential decay



$$lr = \frac{lr_0}{\exp^{k*iter}}$$

# Optimizer Selection In Practice

- Aside from tutorial code, you probably haven't seen “SGD” as the selection for an optimizer
- SGD with momentum is the backbone of most optimizers, with adaptive learning rate schedulers cooked in
- So which optimizer is best?
  - It depends!
- In practice, the optimizer and its parameters are *hyperparameters* you'll jitter around when training a network



# Parameter Regularization During Training

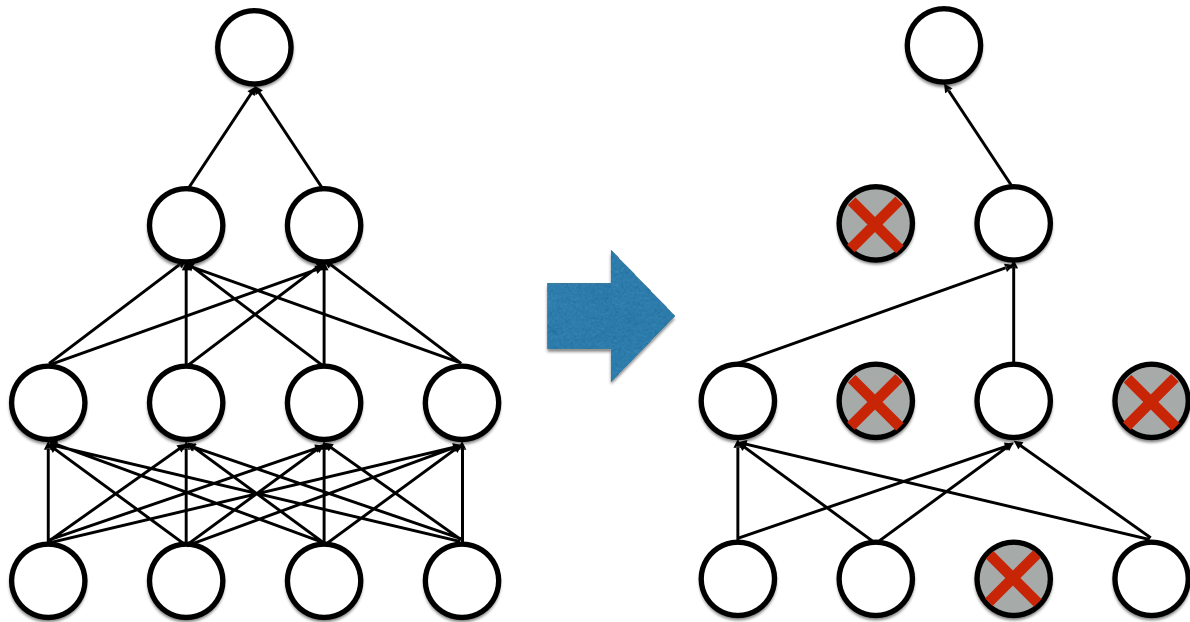
- We can add a regularization “penalty” to our loss  $L$

$$L = L_{original} + \lambda R(\theta)$$

- L1 parameter regularization  $R(\theta) = \sum_i \sum_j |\theta_{i,j}|$
- L2 parameter regularization  $R(\theta) = \sum_i \sum_j \theta_{i,j}^2$
- L1+L2 parameter regularization  $R(\theta) = \sum_i \sum_j |\theta_{i,j}| + \alpha \theta_{i,j}^2$

# Implicit Ensemble Models: Dropout

- At each iteration, randomly shutdown some neurons with a certain probability (usually 0.5)



# Implicit Ensemble Models: Dropout

- At each iteration, randomly shutdown some neurons with a certain probability (usually 0.5)
- Intuition:
  - Forces the network to have a redundant representation
    - Same pieces of information are “stored” in multiple different neurons
  - Training a bunch of small models that share parameters
  - Each binary mask is one such “small model”
- At inference time, these models are “ensembled”

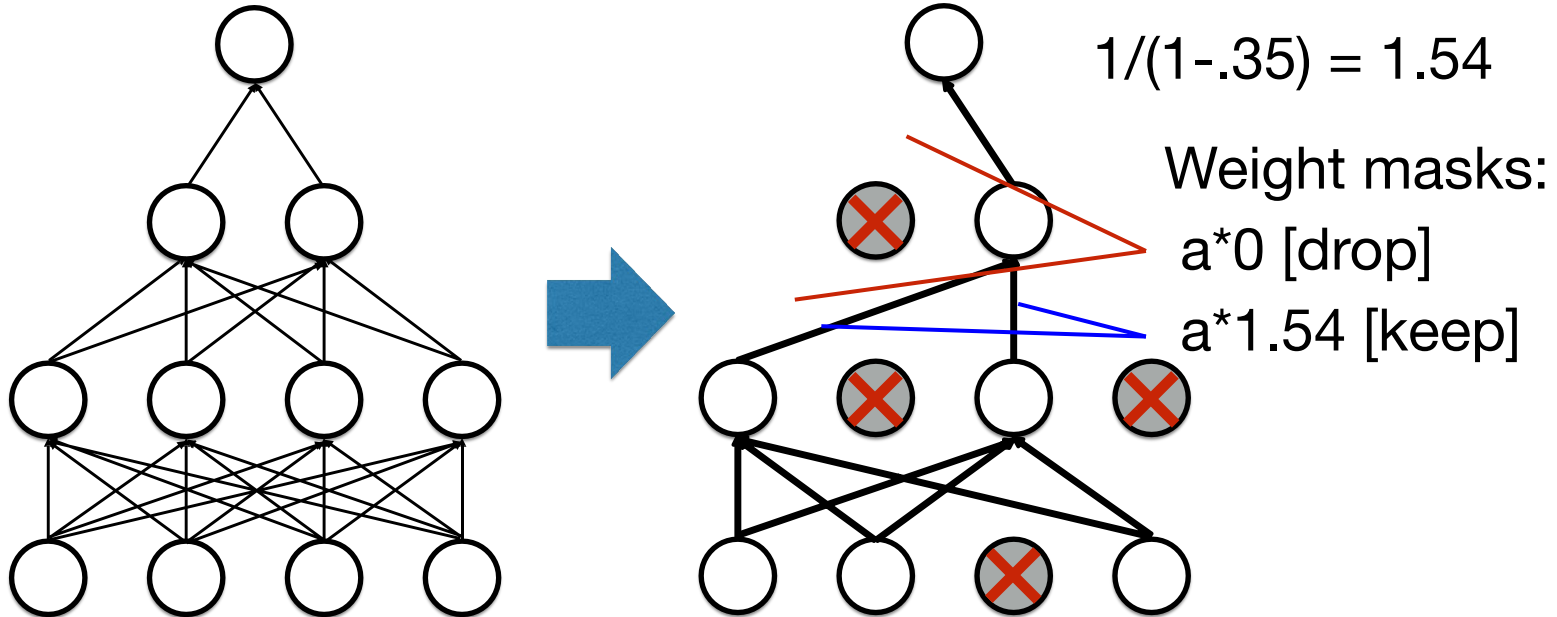
# Training Versus Inference Considerations for Dropout

- If we randomly drop 0.5 nodes during training time but use all during test time, what might go wrong?
  - Test time inputs at each layer will be, on average, x2 what was seen at training time!
- To compensate, during training time we apply an *inverted dropout* correction
  - At training time, for every node not dropped out, increase its output by a magnifying factor of  $1/(1-d)$  for  $d$  probability of dropping each unit; e.g., x10 for dropout 0.9



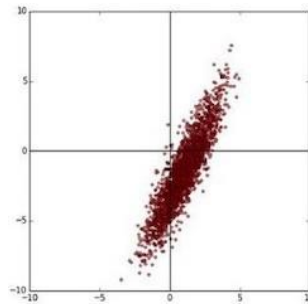
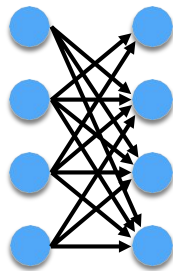
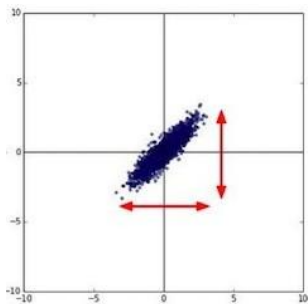
# Training Versus Inference Considerations for Dropout

- To compensate, during training time we apply an *inverted dropout* correction; e.g., for dropout value 0.35:

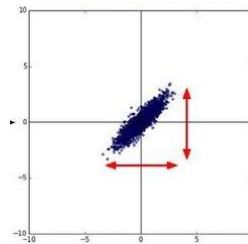


# Normalization Between Layers

- After each layer, the distribution of activation signals changes (Internal covariate shift)
- As a network becomes deeper, distribution shifts more
- Just as we preprocess the input data with our representation function to make features well-behaved, we can consider re-normalizing outputs per layer



$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$



# BatchNorm Layer

- At training time, compute the estimated mean and variance of the activations of a layer across the current batch

$$\hat{x}_i^{(k)} = \frac{x_i^{(k)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$
$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

- How could we make this transformation a little more flexible?
- Introduce learnable parameters to the BatchNorm Layer that enable *scaling* and *shifting* individual input indices

# BatchNorm Layer

At inference time, we still need to calculate mean/var of the activations. Should we do it over the batch?

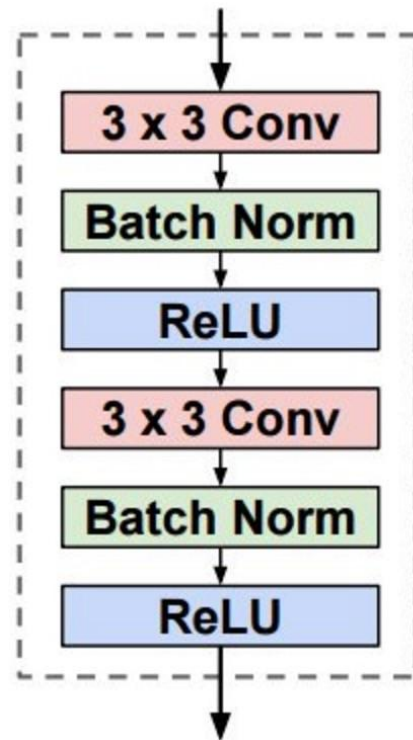
No! Might be out of expected distribution. Use empirical averages from training data!

$$y_i^{(k)} = \gamma^{(k)} \hat{x}_i^{(k)} + \beta^{(k)}$$

- Gamma and Beta are learnable vectors
- These params allow the transformed input to have flexible, learned mean & var
- BatchNorm layer is differentiable

# BatchNorm Layer

- Usually added after convolutional layers or fully connected layers but before non-linearities
- Why?
  - Point of BatchNorm is to help data get into a space where the activation function behaves well



# Overview of Today's Plan

- ~~Course organization and deliverables~~
- ~~Activation Functions and Preprocessing~~
- ~~Optimization and Regularization~~
  - Any questions before we move on?
- Cloud Computing Service Tutorial
- Deep Learning Software Tutorial
- Coding Assignment 1 Release Briefing



## Action Items for You

- If you have not registered with a team, do so **ASAP**; the registration is currently a week late
- Start firming up your project proposal ideas and drafting your slides; these proposal slides are due to us in a week
- Coding Assignment 1 will be discussed by TAs and is due two weeks from today
- Keep in mind also that the midterm exam will be in two weeks; the exam will be mostly multiple choice and submitted via a google form; it must be taken during the class period



# CSCI 566: Deep Learning and Its Applications

Jesse Thomason

Lecture 4: Activations, Initializations,  
Optimization, and Regularization