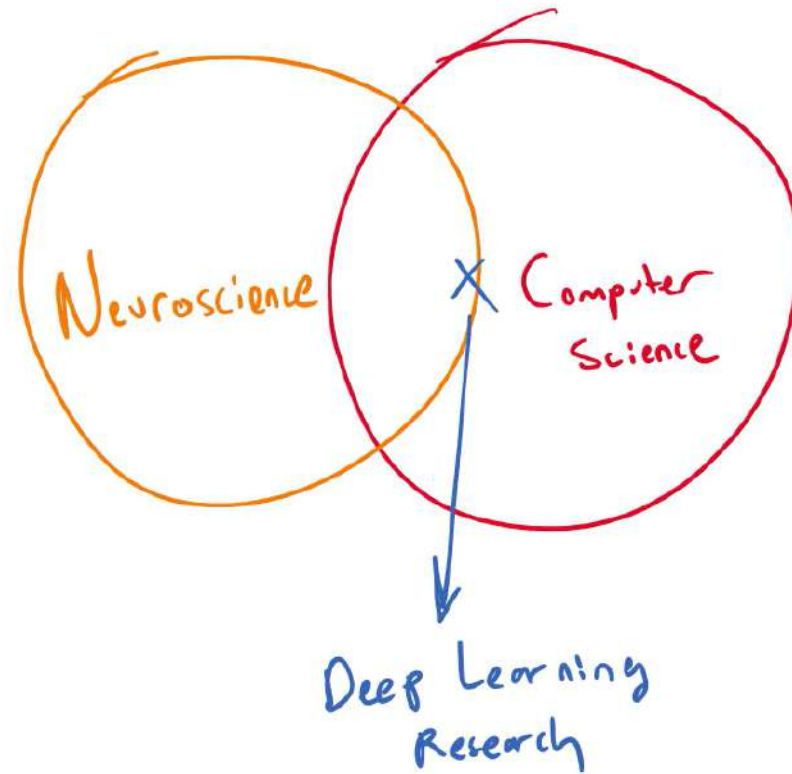
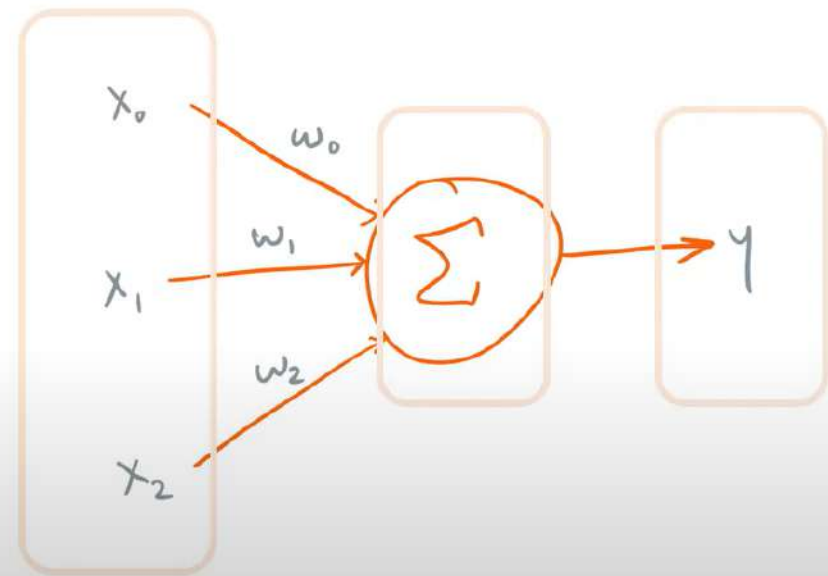
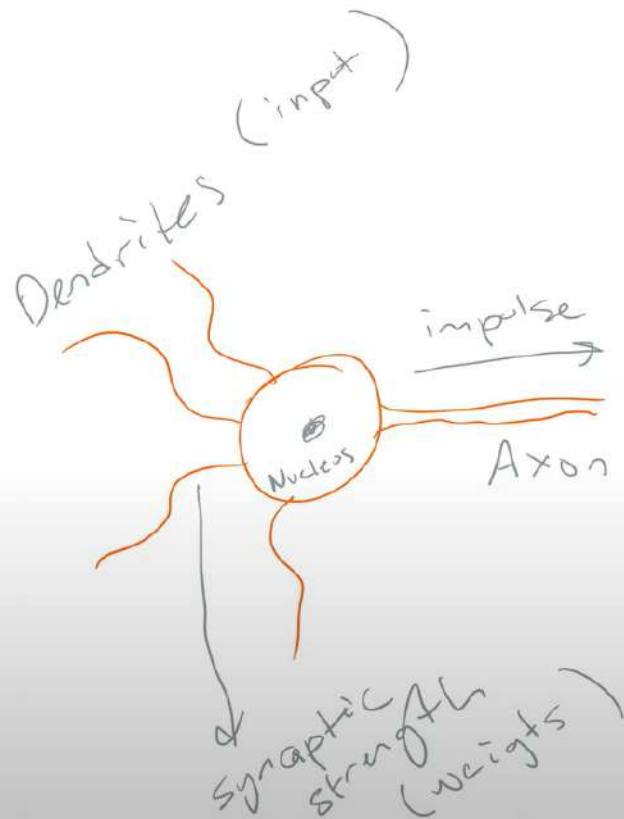


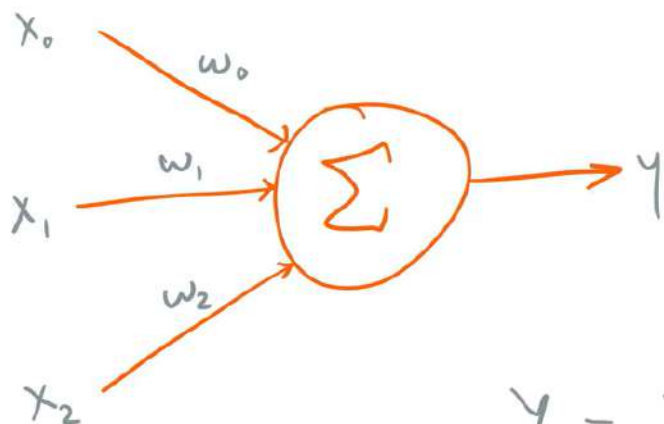
More on this side



# Neurons



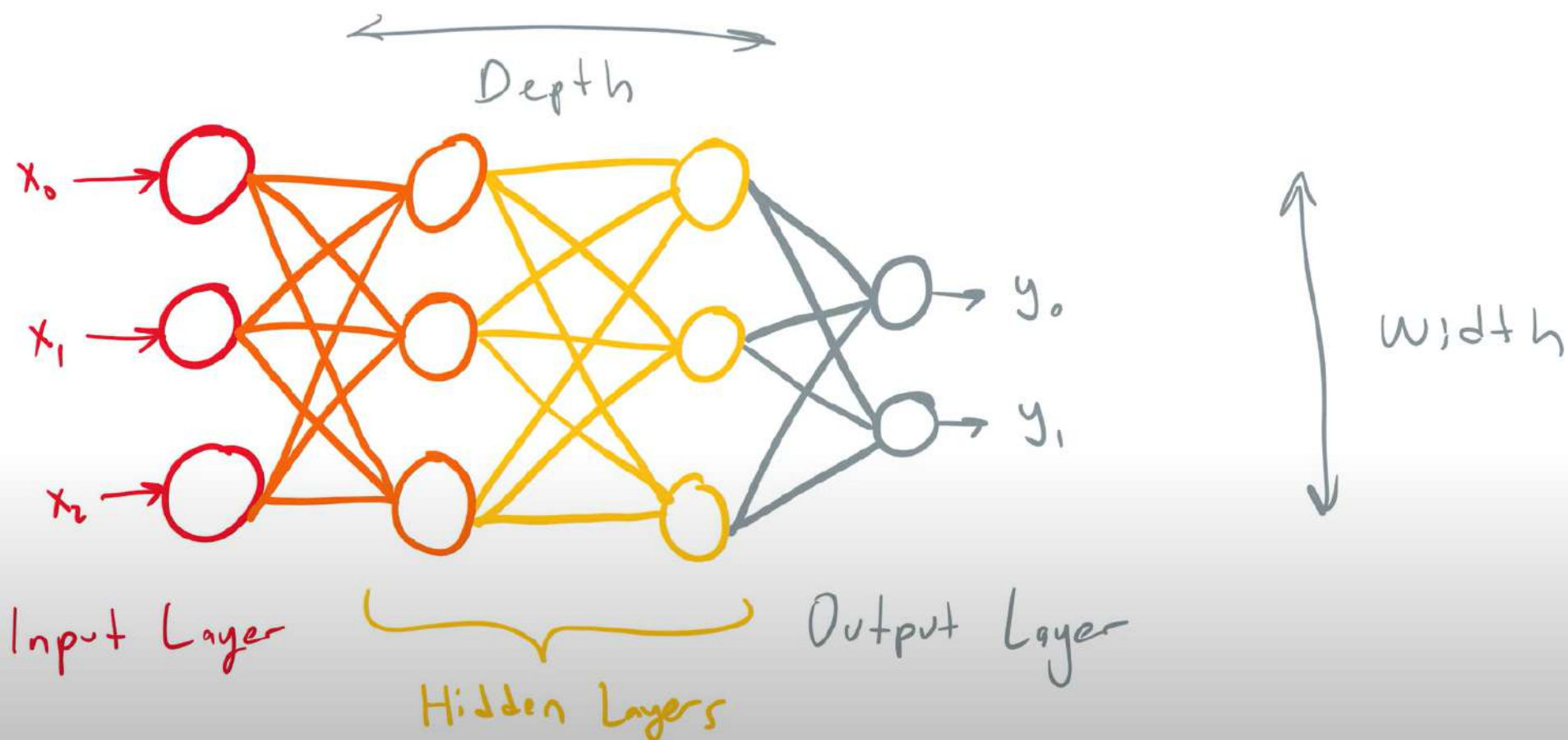
# Artificial Neuron



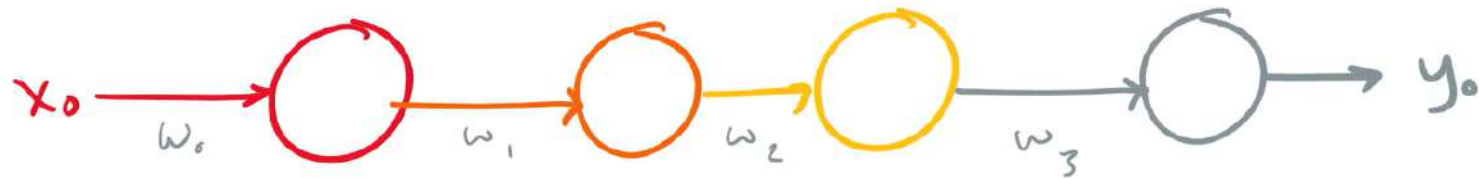
$$y = x_0 w_0 + x_1 w_1 + x_2 w_2 = \sum_i x_i w_i = \mathbf{w}^T \mathbf{x}$$

$$\begin{bmatrix} w_0 & w_1 & w_2 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = x_0 w_0 + x_1 w_1 + x_2 w_2$$

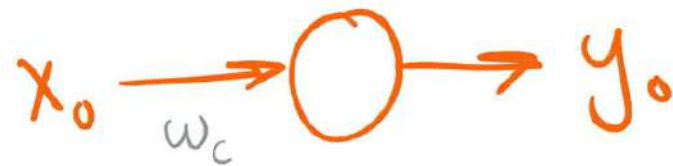
# Artificial Neural Networks



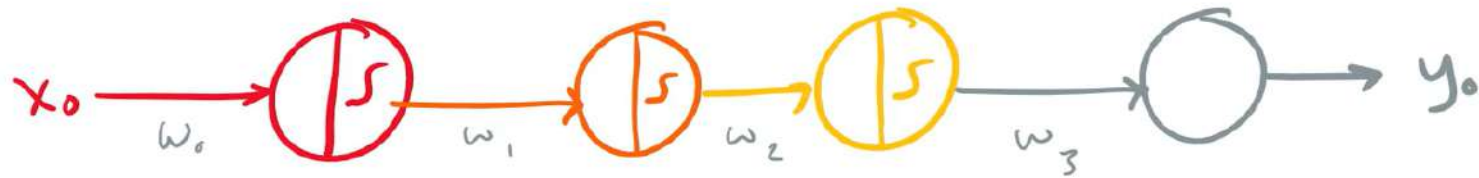
# Activation Function



$$y_0 = x_0 w_0 w_1 w_2 w_3 = x_0 w_c$$



# Activation Function



$$w^T x \rightarrow \sigma(w^T x)$$



# Activation Function

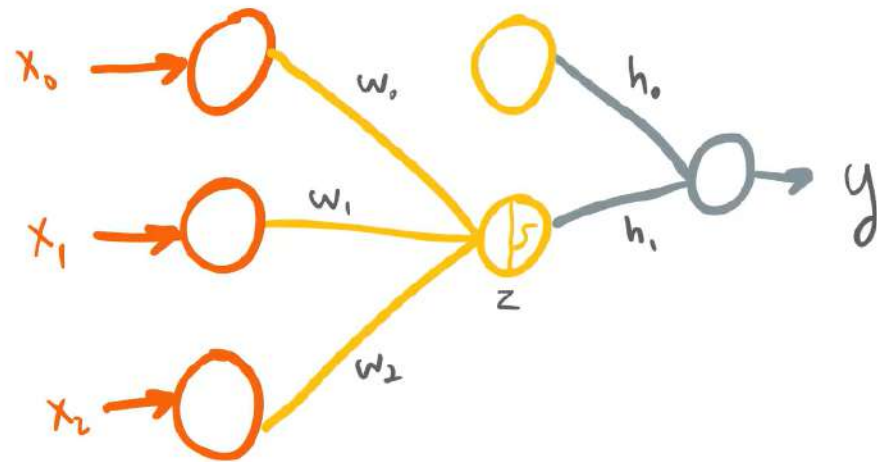


$$w^T x \rightarrow \sigma(w^T x) \quad \sigma(x) = \frac{1}{1 + e^{-x}}$$



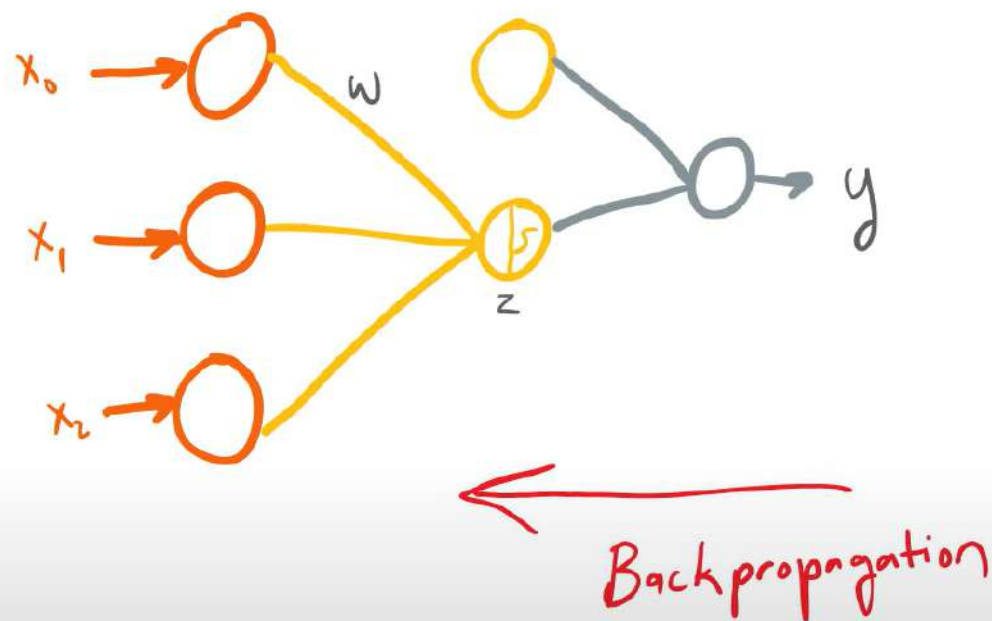
$$y_0 = \sigma(\sigma(\sigma(x_0 w_0) w_1) w_2) w_3$$

# Learning the weights



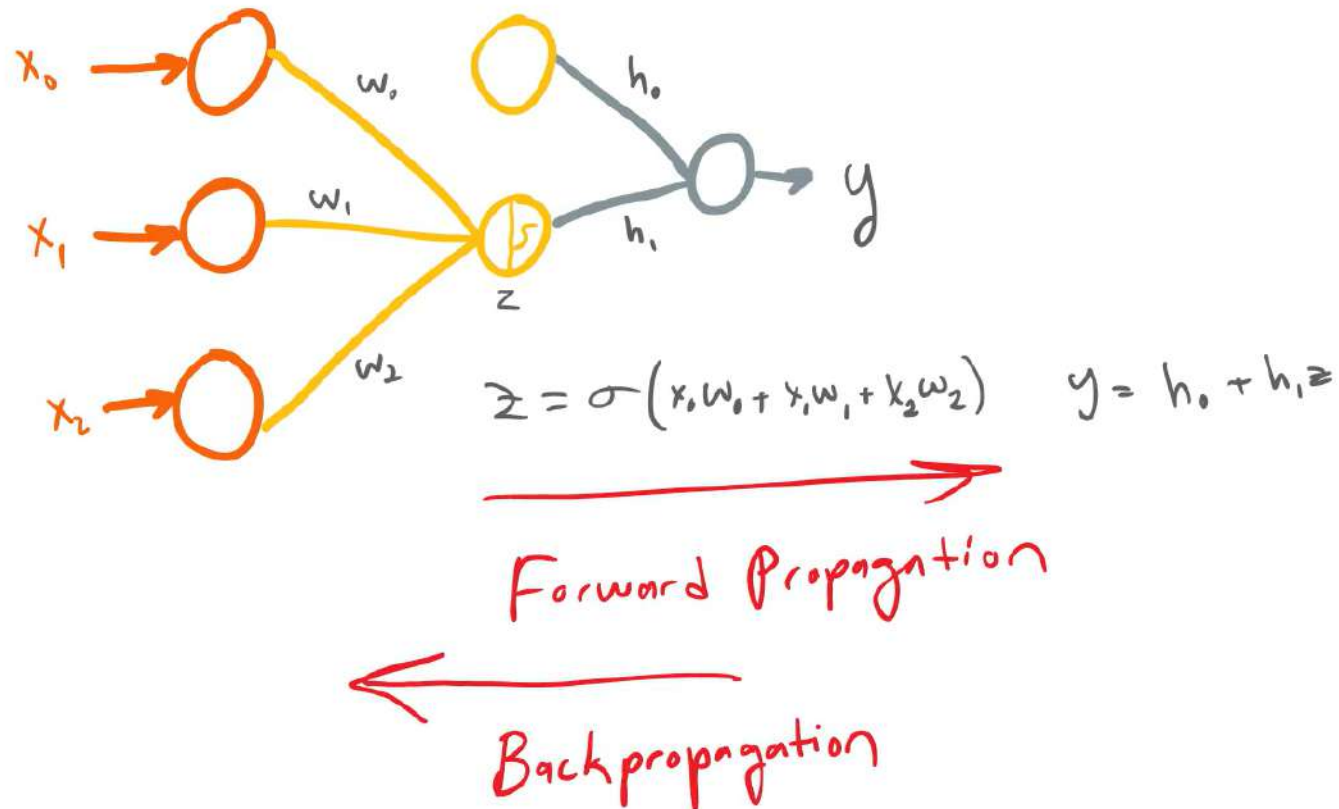


# Learning the weights

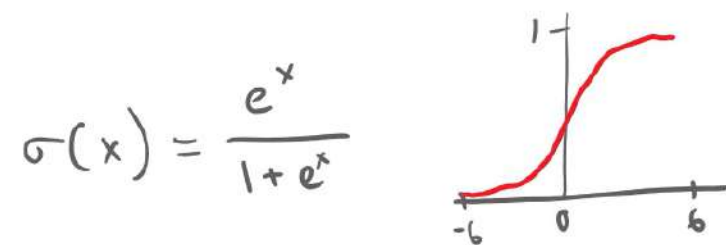


$$\frac{\partial E}{\partial w} = \underbrace{\frac{\partial E}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial w}}_{\text{Chain Rule}}$$

# Learning the weights

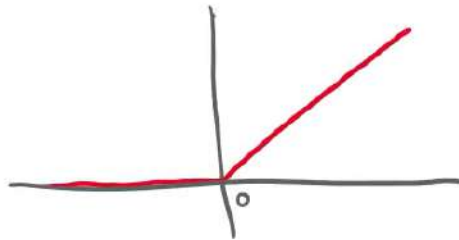


# Activation functions

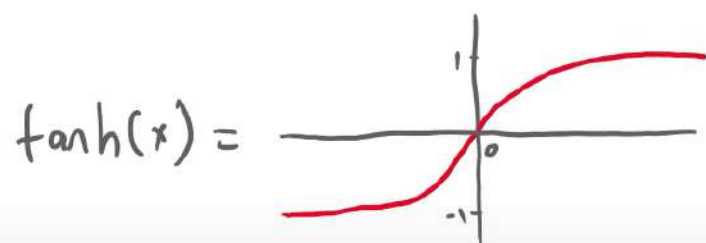
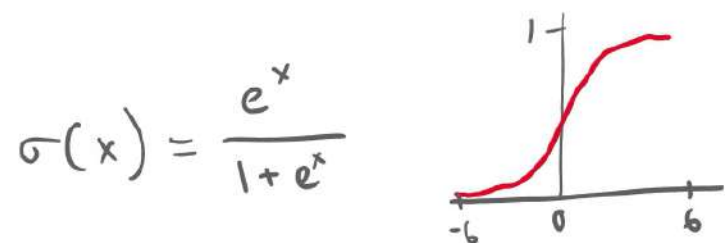


# ReLU: Rectified Linear Unit

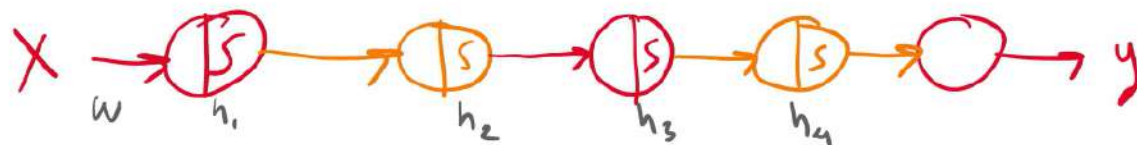
$$y = \max(0, x)$$



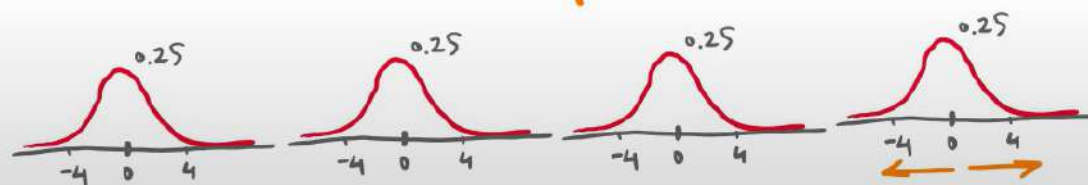
# Activation functions



Vanishing Gradient

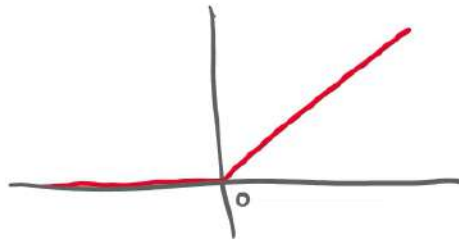


$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w}$$

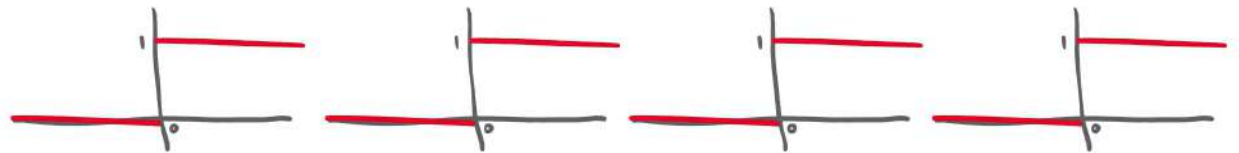


# ReLU: Rectified Linear Unit

$$y = \max(0, x)$$

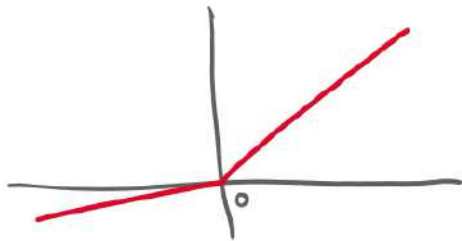
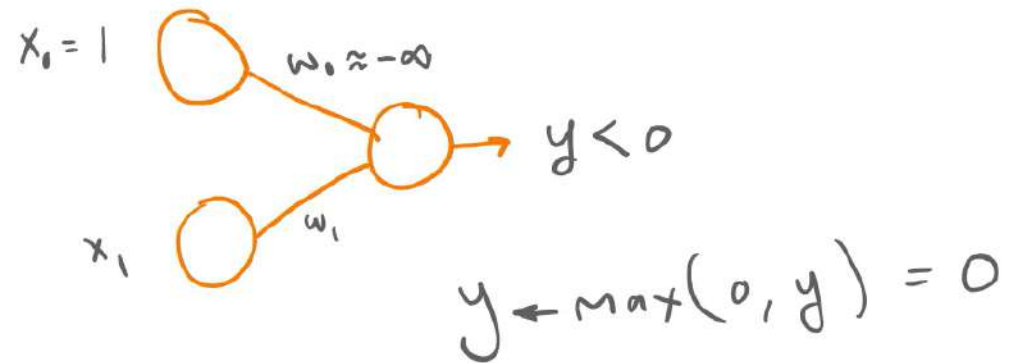
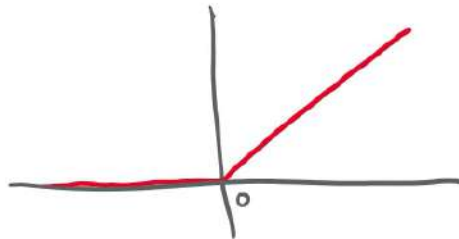


$$\frac{dE}{dw} = \frac{\partial E}{\partial h_4} \frac{\partial h_4}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial w}$$

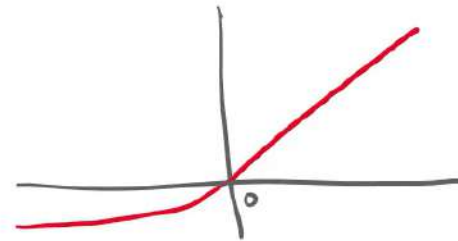


# ReLU variants

$$y = \max(0, x)$$



Leaky ReLU  
Parametric ReLU

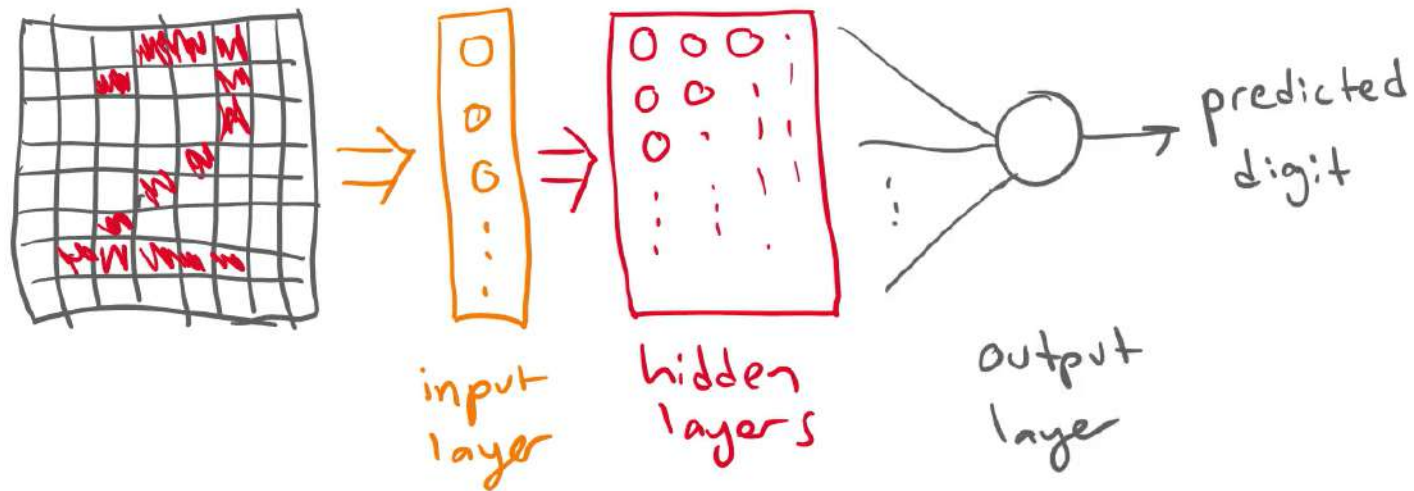


ELU  
(Exponential Linear Unit)

# Loss functions

$$MSE = \frac{1}{n} \sum_i^n (y - \hat{y})^2$$

Classification

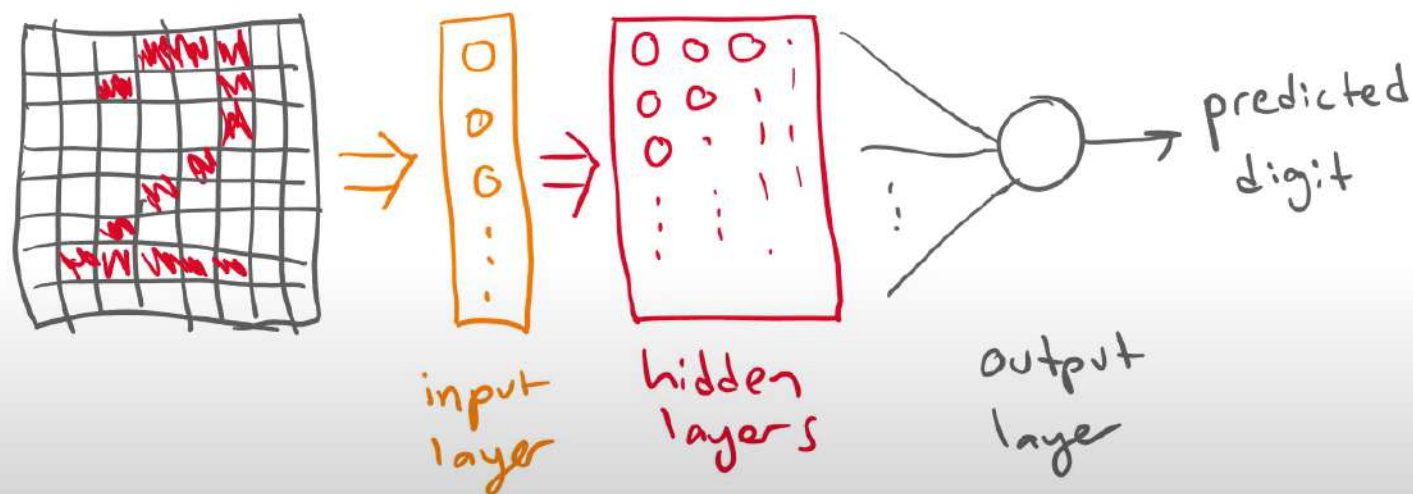




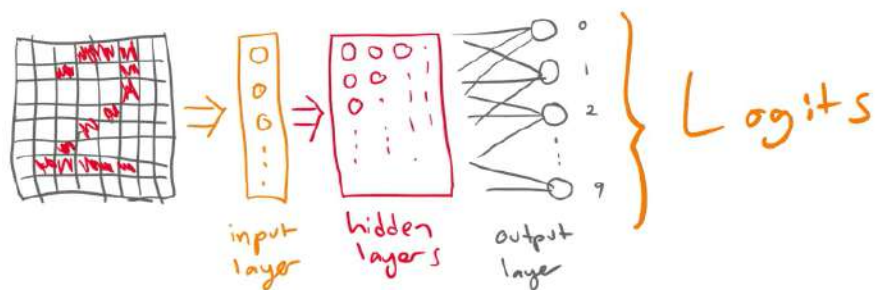
# Loss functions

$$MSE = \frac{1}{n} \sum_i^n (y - \hat{y})^2$$

Classification



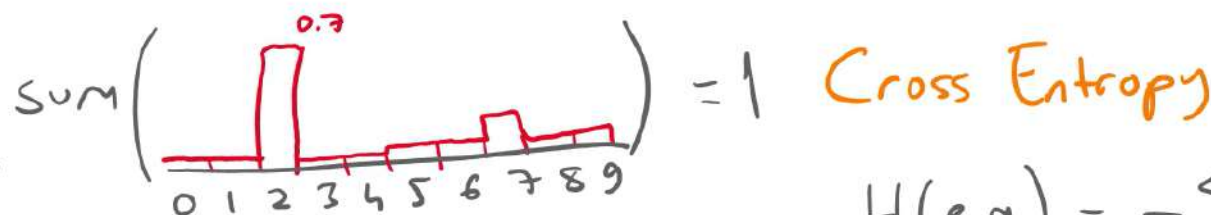
# Softmax function



$$\text{Softmax}(y_i) = \frac{e^{y_i}}{\sum_i e^{y_i}}$$



Ground truth probabilities  
(one-hot encoding)

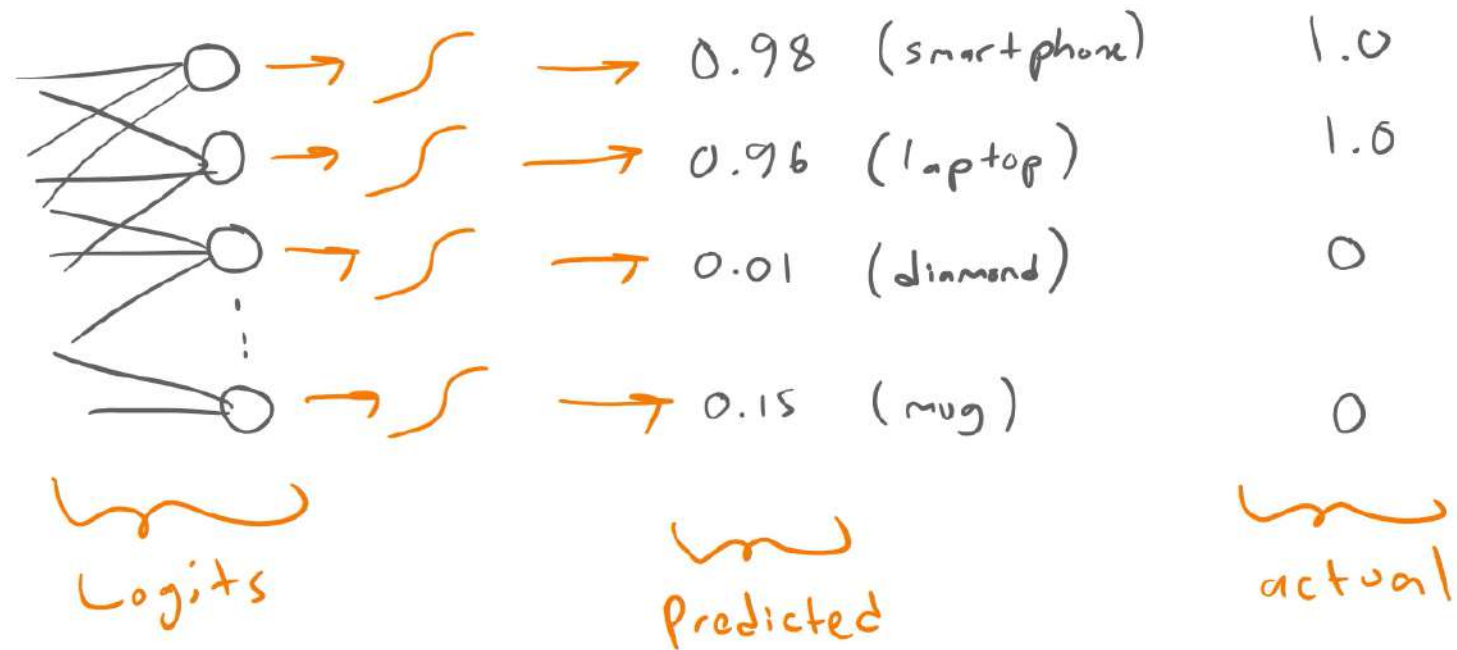


Predicted probabilities

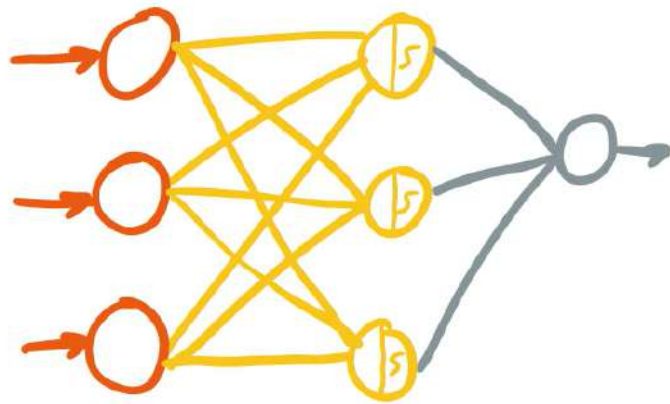
$$H(p, q) = -\sum_i p_i \log q_i$$

Cross Entropy

# Multi-label classification



One more thing: do we really need deep models?

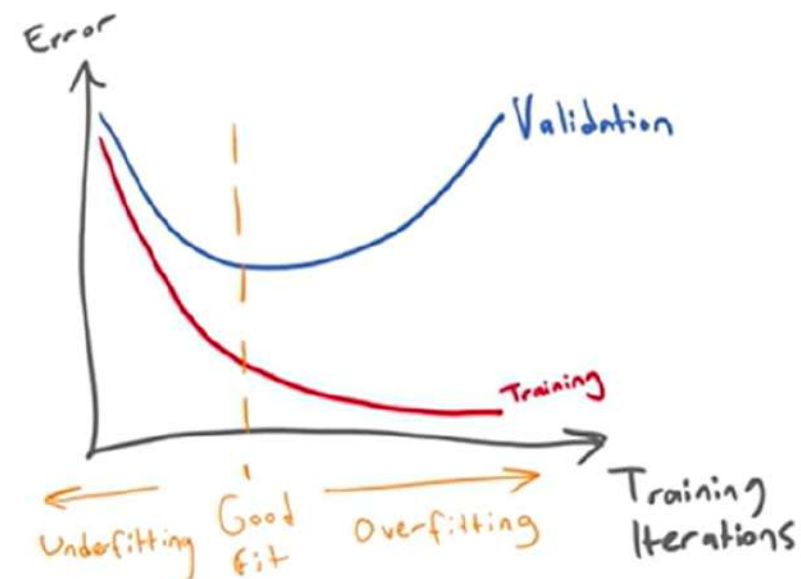


might be  
too wide

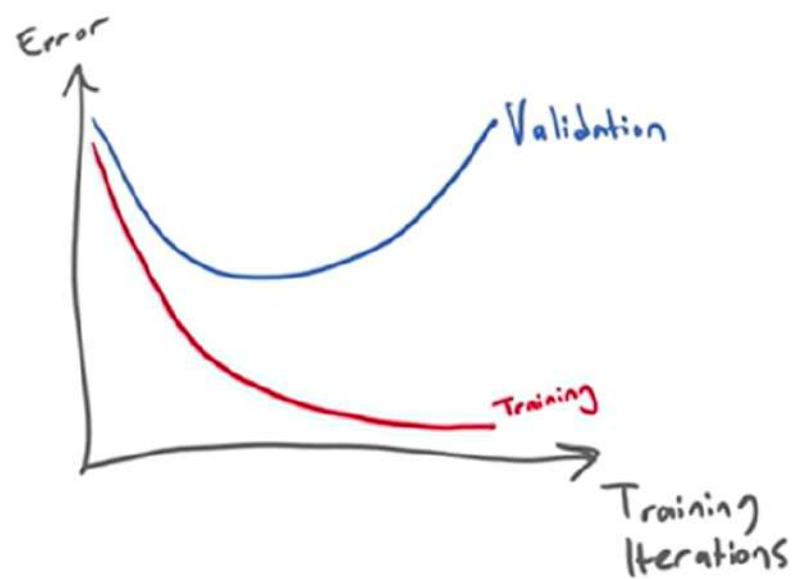
Universal Approximator

can represent  $\neq$  can learn

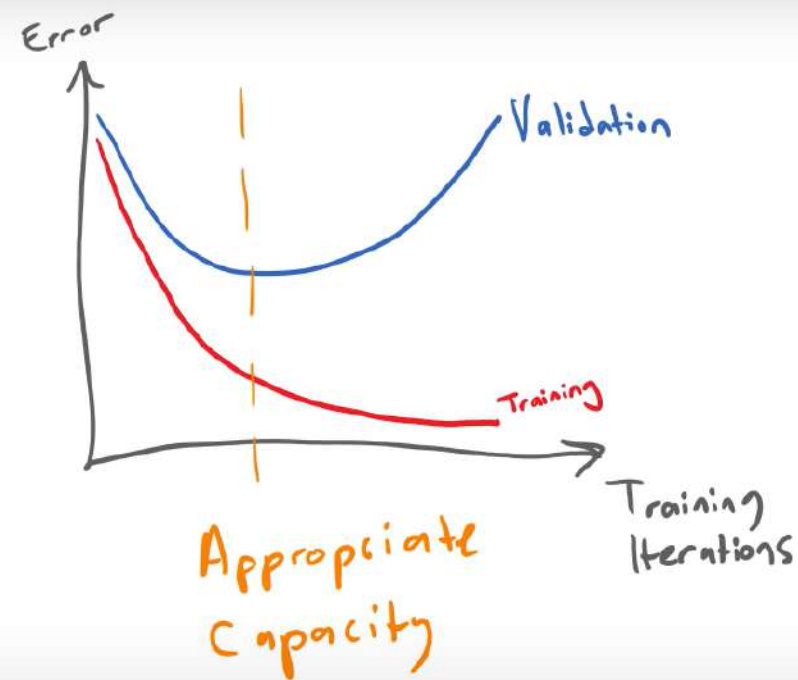
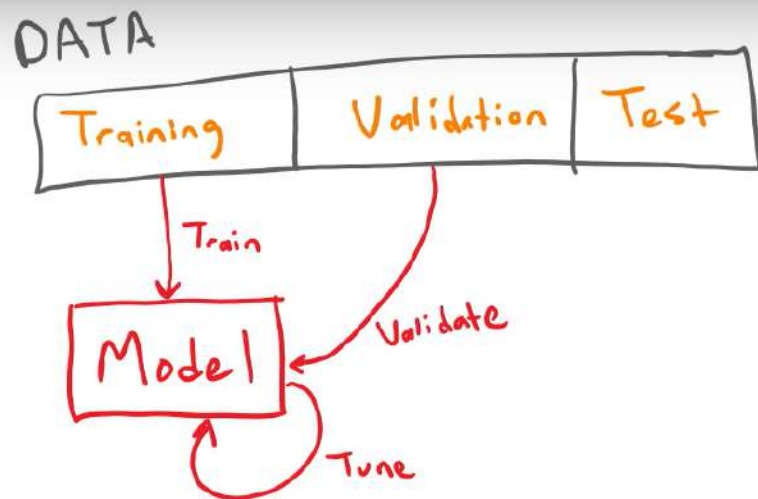
Training loss	Validation loss	
High	High	Underfitting: increase capacity
Low	High	Overfitting: decrease capacity
Low	Low	Good fit: run test
High	Low	Unlikely: debug



# What's next?



## Overfitting, Underfitting, and Model Capacity



### Hyperparameters

- ⚙ # of training iterations
- ⚙ # of layers (depth)
- ⚙ # of hidden units (width)
- ⚙ types of layers
- ⚙ regularization methods
- ⚙ regularization strength
- ⚙ optimization algorithm
- ⚙ learning rate
- ⚙ batch size
- ⚙ ...



DATA



5:06 / 9:06 • Partitioning Data >

Scroll for details



SUBSCRIBE

## ☰ Regularization and Data Augmentation

Fewer parameters

*demand less computational power*

Best performing models tend to be

*large but trained in a way that  
restricts the utilization of their entire potential*

Regularization

- *encourages models to have a preference towards simpler models*
- *reduces the risk of overfitting*



⏮ ⏪ ⏩ ⏭ 🔊 1:05 / 8:08

Scroll for details  
▼

$$\text{loss} \leftarrow \underbrace{\text{loss}}_{\sum_i (y_i - \hat{y}_i)^2} + \underbrace{(\text{weight penalty})}_{\alpha \sum_i |w_i|}$$

$$\text{Loss} \leftarrow \underbrace{\text{Loss}}_{\sum_i (y_i - \hat{y}_i)^2} + \underbrace{(\text{weight penalty})}_{\alpha \sum_i w_i^2}$$

actual      predicted
regularization strength

L2 Regularization

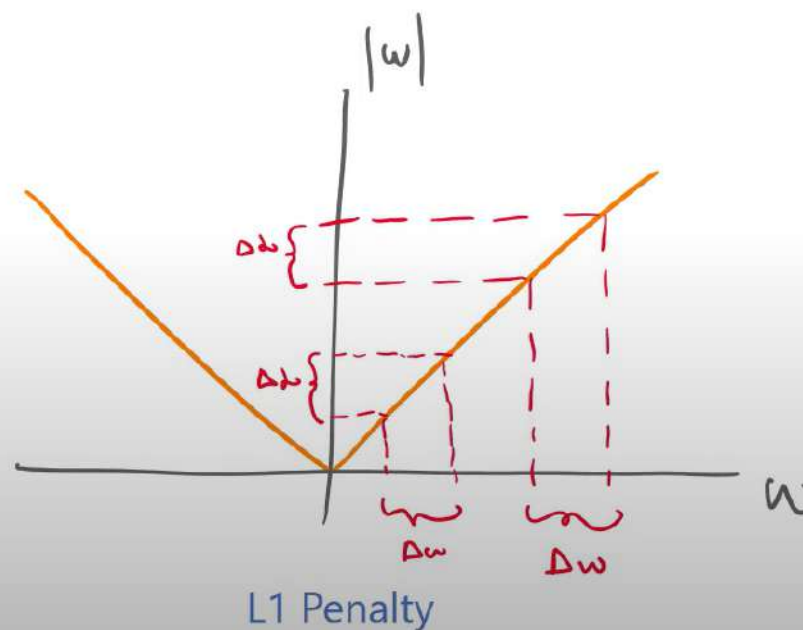
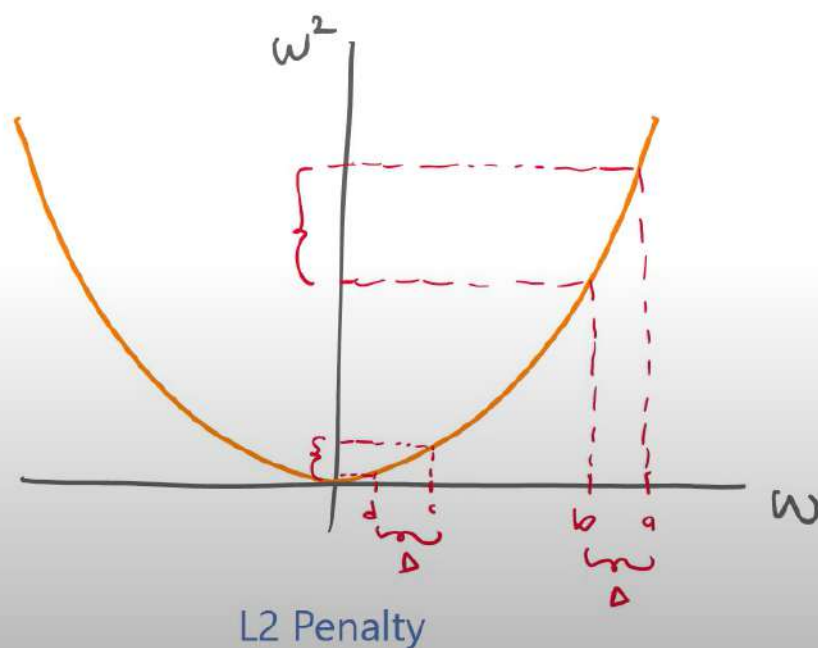
Minimize  $\left( \downarrow \sum_i (y_i - \hat{y}_i)^2 + \alpha \downarrow \sum_i w_i^2 \right)$

Ridge Regression

$$\text{Loss} \leftarrow \underbrace{\text{Loss}}_{\sum_i (y_i - \hat{y}_i)^2} + \underbrace{(\text{weight penalty})}_{\alpha \sum_i |w_i|}$$

L1 Regularization LASSO

Leads to sparser results



## Regularization and Data Augmentation

### Data Augmentation



Original Image



Scaled Image



Translated Image



Rotated Image



Mirrored Image



Brightness / Contrast  
Shifted Image



Noise Added  
Image



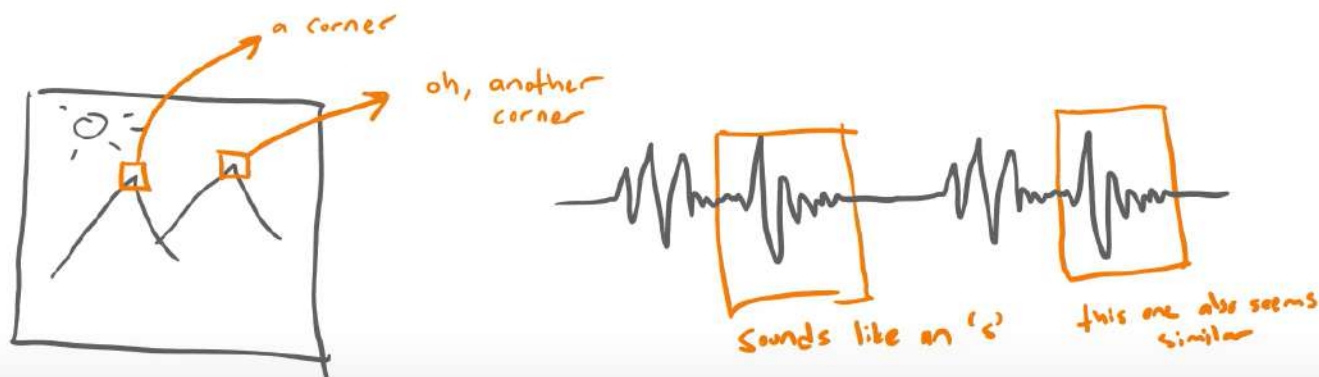
Regularization methods introduce additional information such as

- *smaller weights are better*
- *parameter sharing is useful*
- *distorting the input in some ways shouldn't change the results*

Prevent overfitting

- *limit model capacity*
- *get more data*

Data augmentation



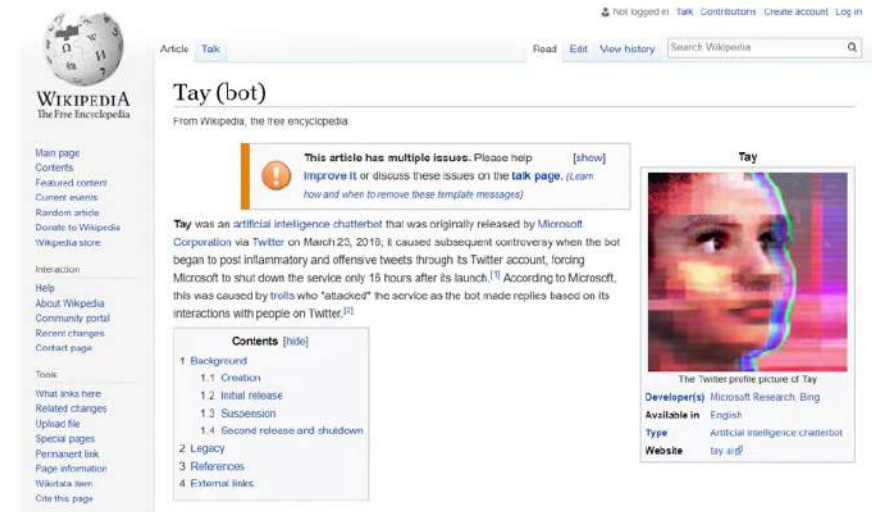
## Dataset Bias

*data can be **biased** no matter how **big** it is*

## Models used for

***medical, financial, and legal** purposes  
affect people's lives*

Biased models can further reinforce biases





## Subjective Studies

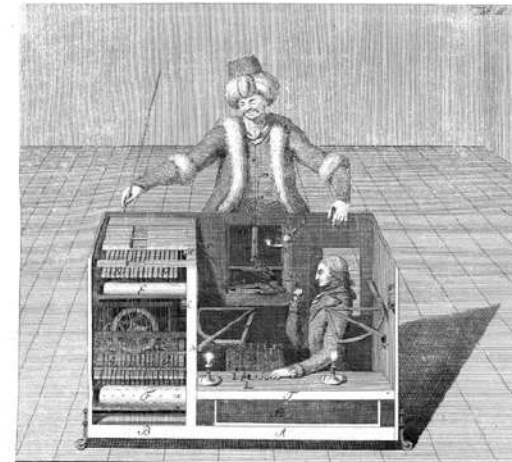
*Surveys*

*Games*

*Crowdsourcing*

## Data collection guidelines

- Ask *unbiased* questions
- Design *an easy to use* interface
- Make it *fun* for the users
- Ensure that the entire process is *ethical*



## Data Imputation

**Caveat:** *data might not be missing at random*

Gender	Likes
M	Cats
M	Dogs
F	Missing
M	Cats
F	Dogs
F	Dogs
F	Cats

Classes
Cats
Dogs
Missing

## Feature scaling

Variable	Range of values
Age	0 – 100+
Annual income	0 – 1,000,000+
Years of experience	0 – 40+

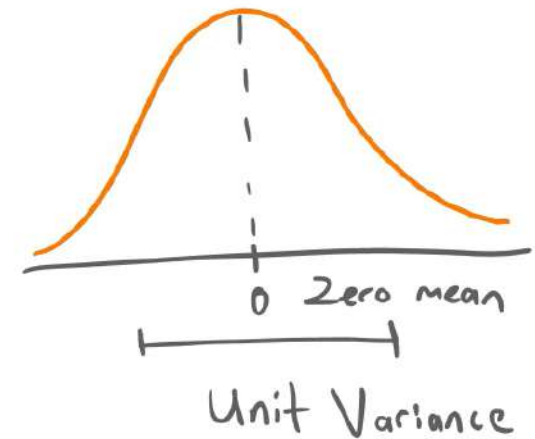


Range of values
0 – 1
0 – 1
0 – 1

$$\hat{x} = \frac{x - \min(x)}{\max(x) - \min(x)}$$

$$\hat{age} = \frac{age - 18}{99 - 18}$$

## Standardization



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

## Data Imbalance

Classes	Number of samples
Cats	5000
Dogs	5000
Tigers	150
Caracals	50
Axolotls	25



### Undersampling

Classes	# samples
Cats	25
Dogs	25
Tigers	25
Caracals	25
Axolotls	25

### Oversampling

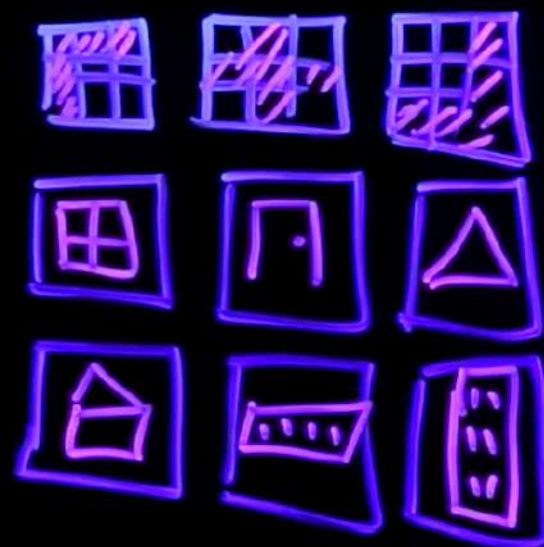
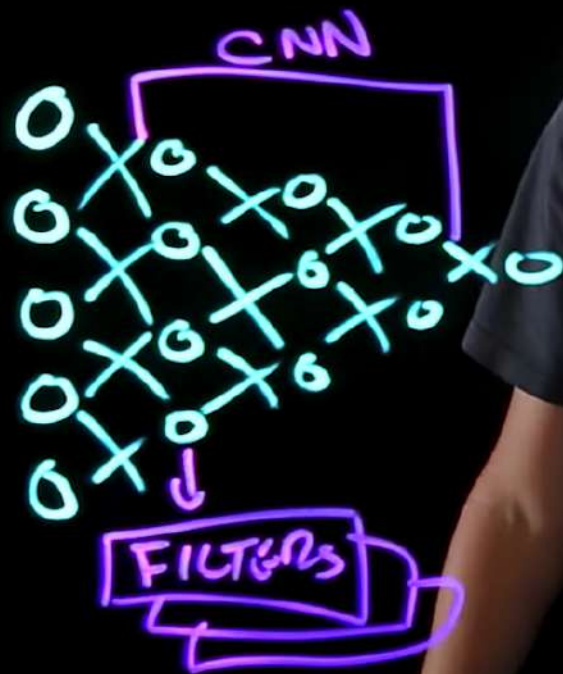
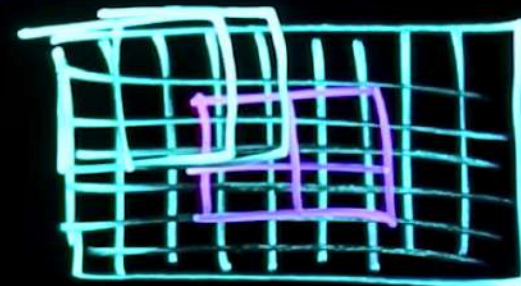
Classes	# samples
Cats	5000
Dogs	5000
Tigers	5000
Caracals	5000
Axolotls	5000

### Class-weighted loss function

$$\mathcal{L}_{\text{loss}} = \sum w q \log p$$

↓  
Class weights

# CONVOLUTIONAL NEURAL NETWORK (CNN)



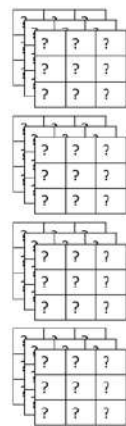


# Convolutional Neural Networks Explained

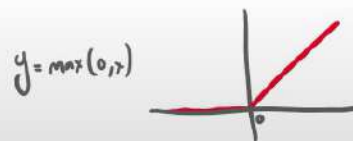


Convolutional Layer

Filters



Activation Function



5:11 / 14:30

Scroll for details

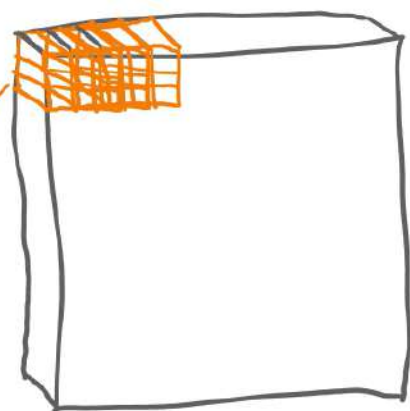


Settings

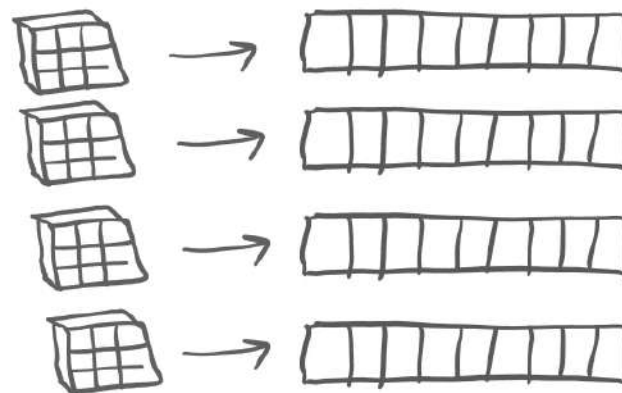
SCRIBES

# Convolutional Neural Networks Explained

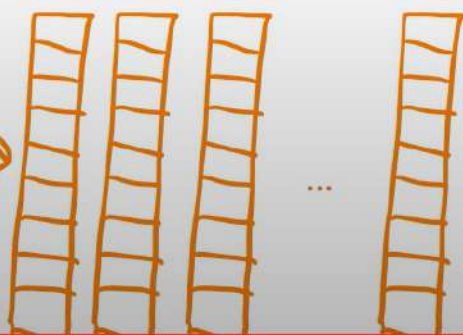
## Implementation as a matrix multiplication



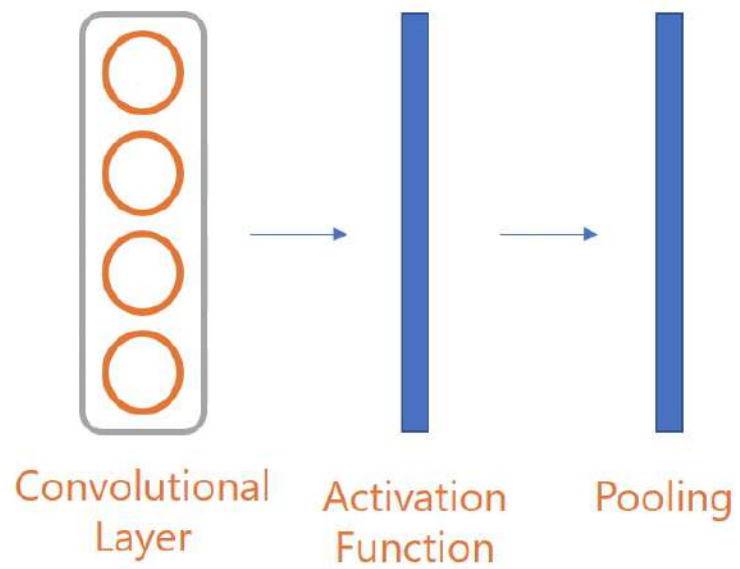
Input



Filters



$$\begin{matrix} 4 \\ \begin{matrix} \text{---} \\ | \\ \text{---} \end{matrix} \\ 27 \end{matrix} \begin{matrix} 27 \\ \text{---} \\ | \\ \text{---} \end{matrix} = \begin{matrix} N \\ \text{---} \\ | \\ \text{---} \end{matrix}$$

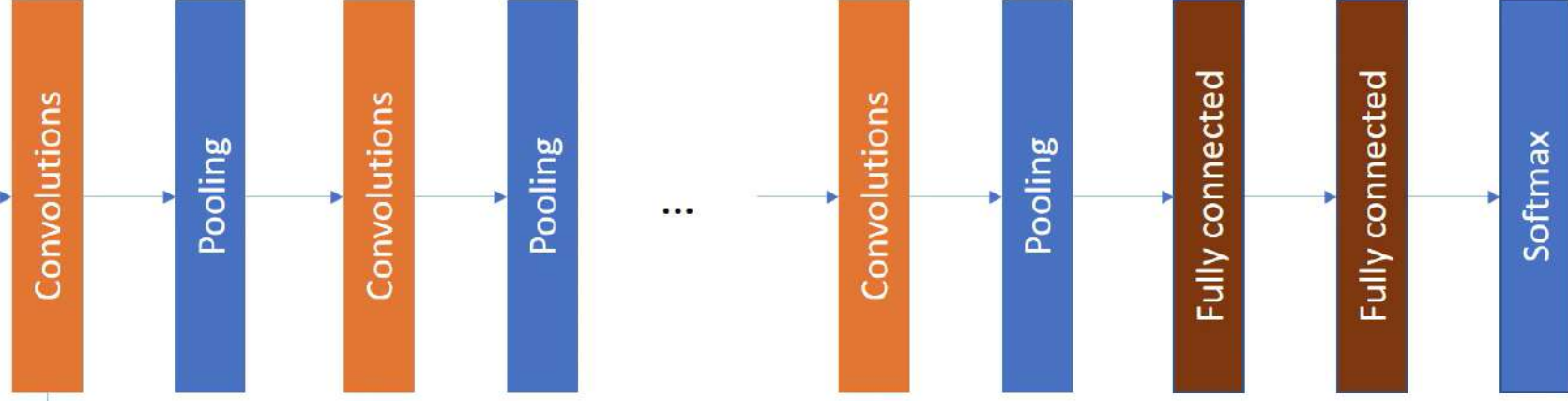
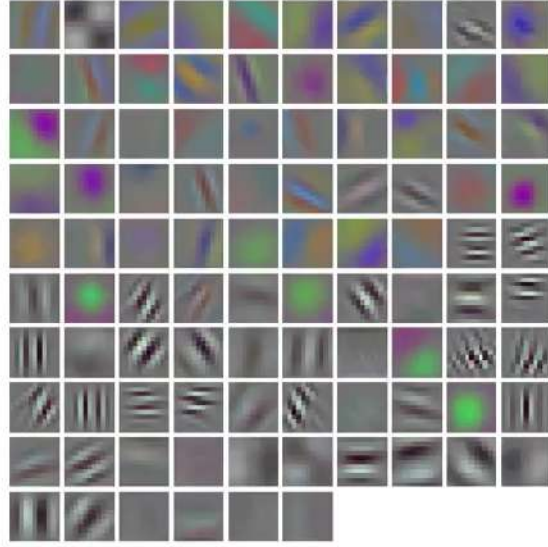


0	0	5	1
0	2	4	0
1	0	1	0
0	0	0	3

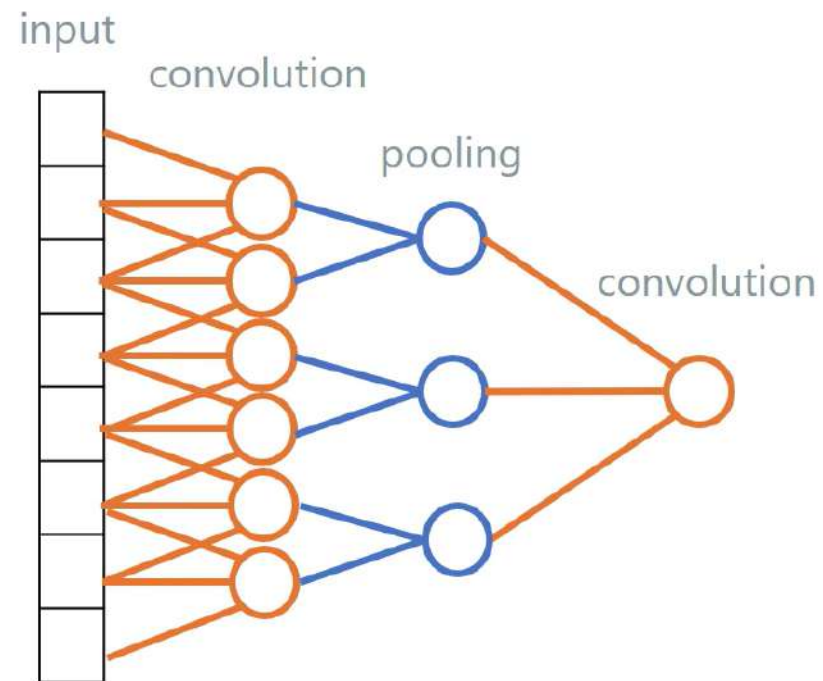
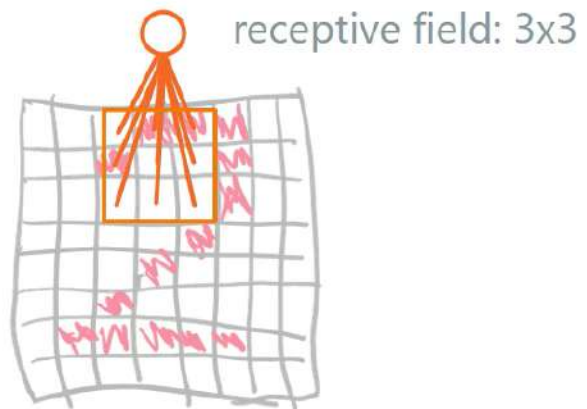
Max-pooling

2	5
1	3



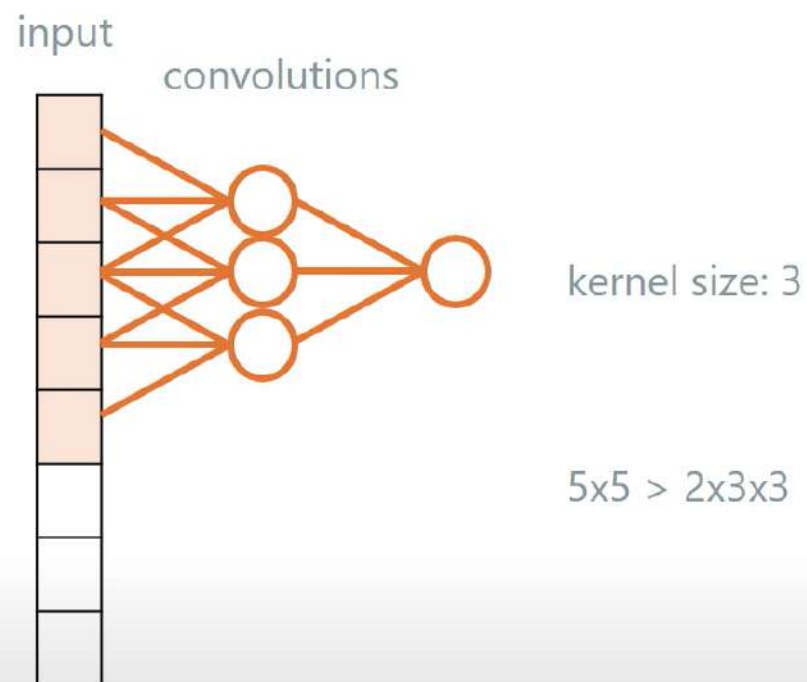
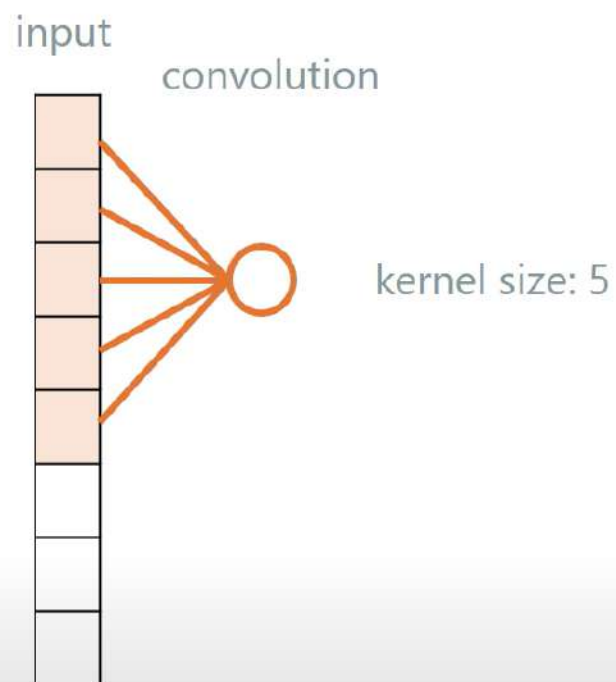


## Receptive field

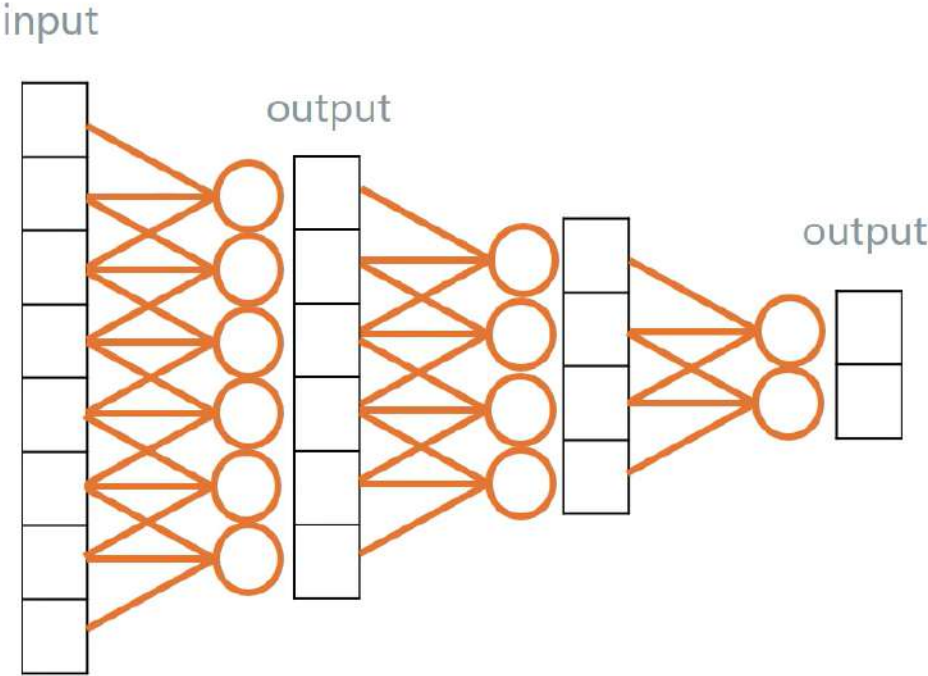


## Convolutional Neural Networks Explained

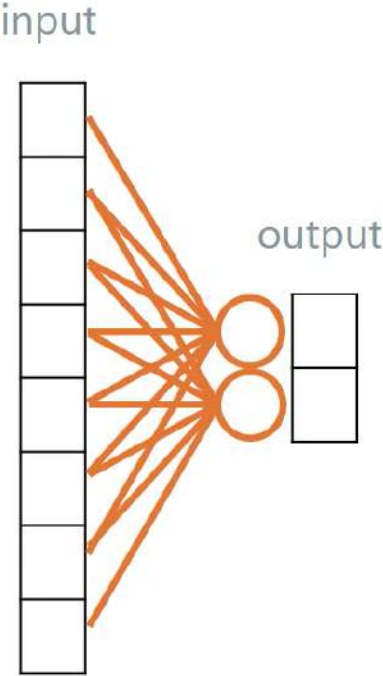
### Kernel size



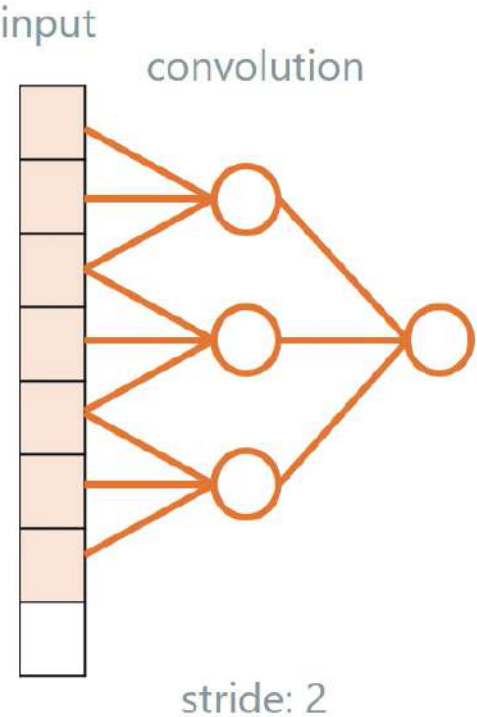
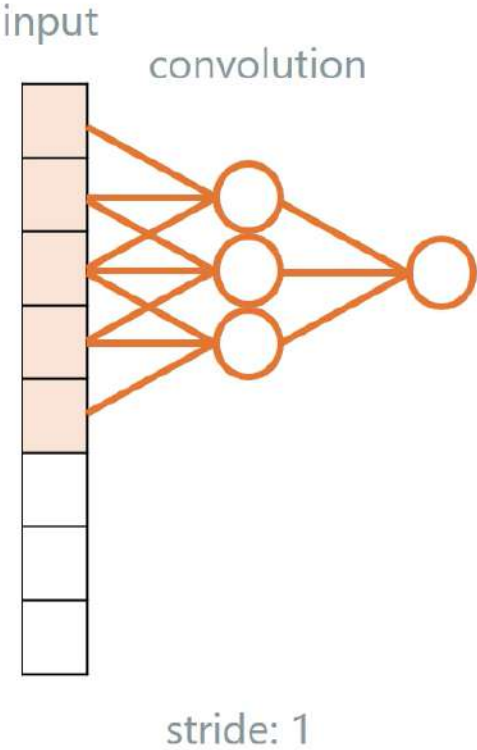
# Padding



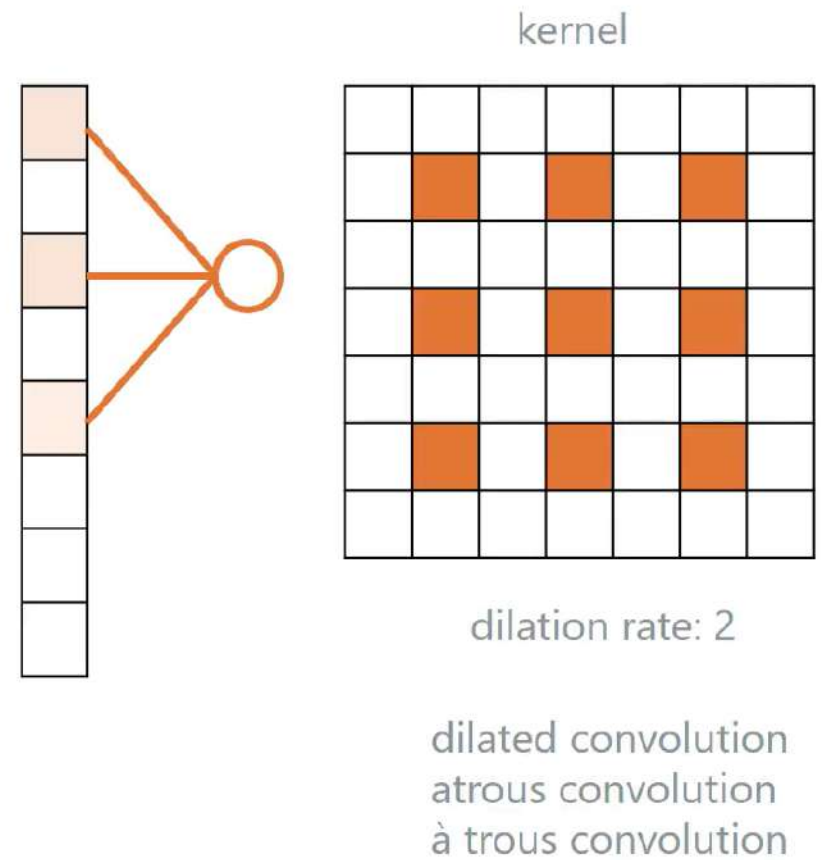
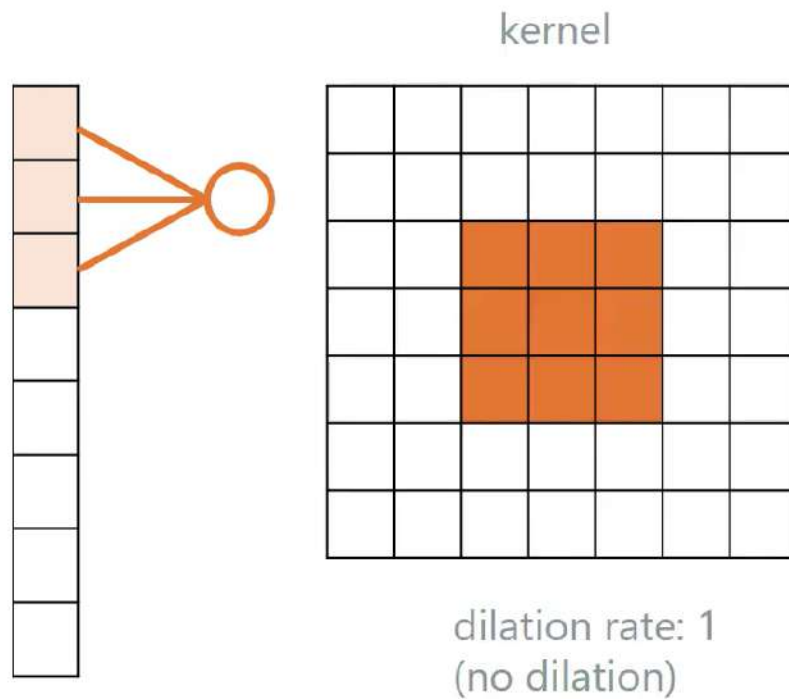
Padding='VALID'



# Stride



## Dilation rate



## Building blocks of CNNs

- Convolution
- Pooling

## Hyperparameters

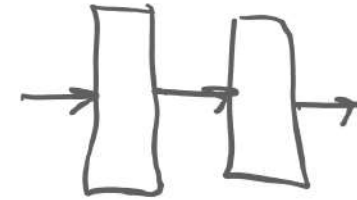
- Kernel size, stride, dilation rate
- Number of filters, number of layers
- Pooling size & type
- Padding type
- The way we arrange the layers

## What's next?

- How to design a CNN
- Popular CNN architectures

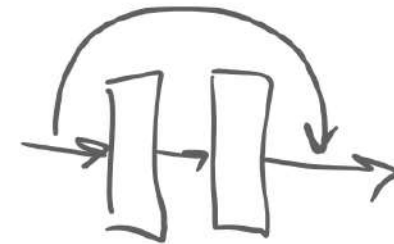
## No human intervention in deep learning?

- humans are still in the loop!



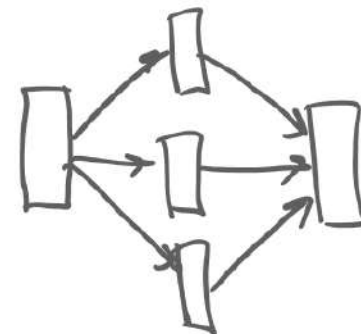
## Grid search on all hyperparameters?

- too many hyperparameters
- infinitely many ways to design a network



## Designing a deep neural network involves

- human expertise
- trial and error





## How to design a ConvNet?

### TL;DR

- you don't design
- pick something that works and use it

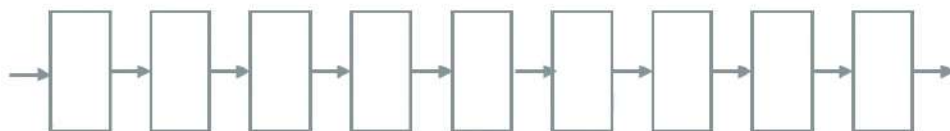
### Working on novel problems?

- borrow ideas from successful models
- design your own model

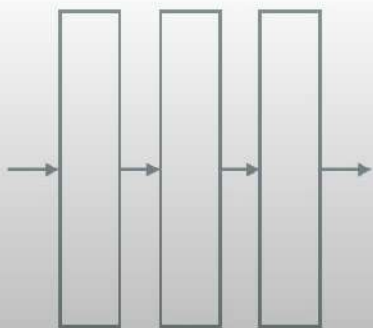
## How to Design a Convolutional Neural Network

### How do choose the number of layers and units?

- start small
- gradually increase model size
- smallest model: linear regression
- deeper:



- wider:

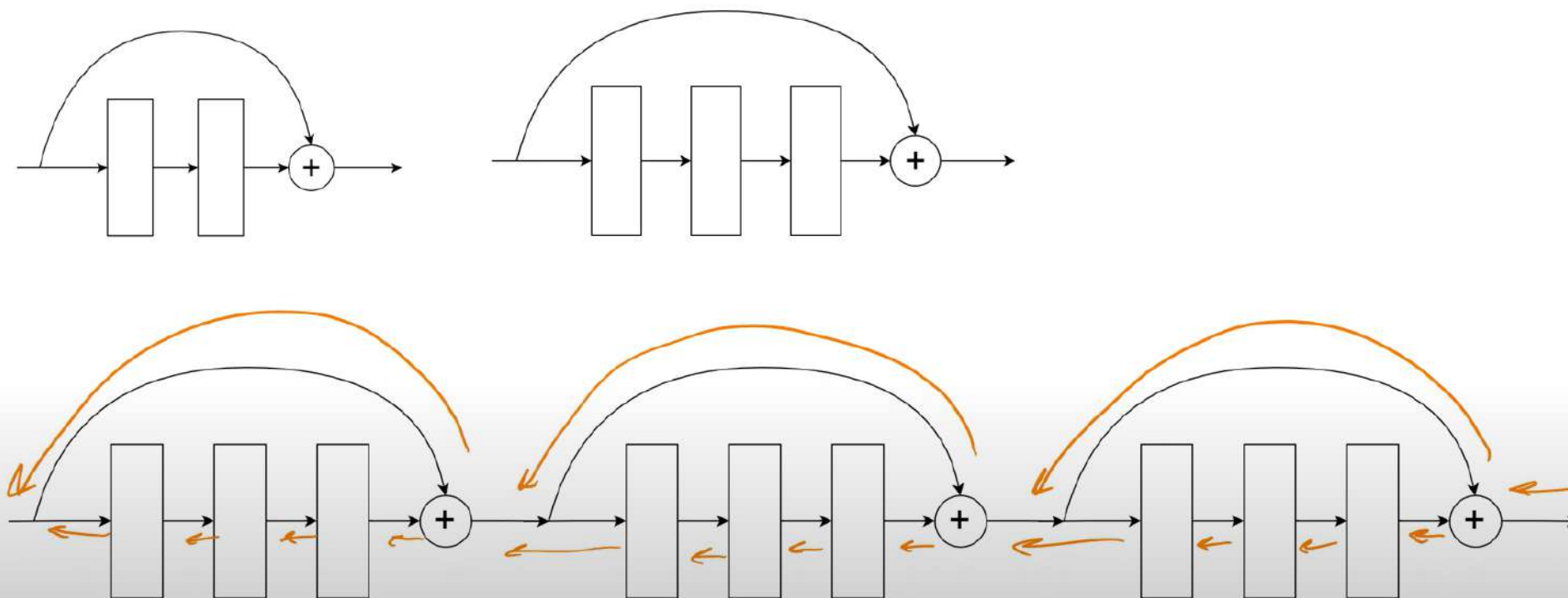


## Skip connections

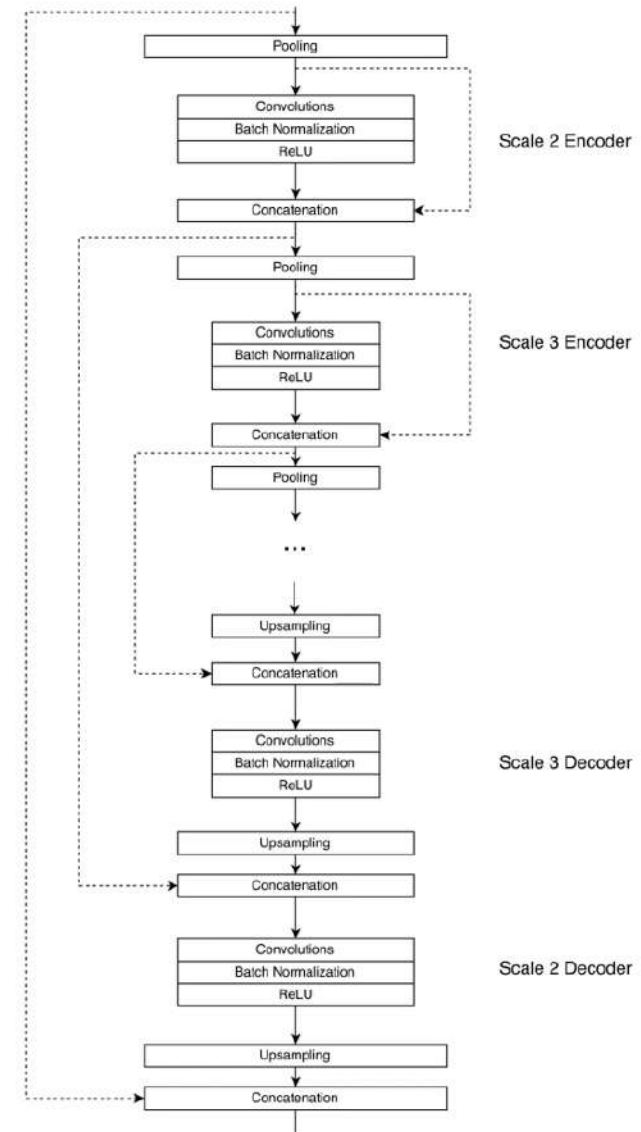
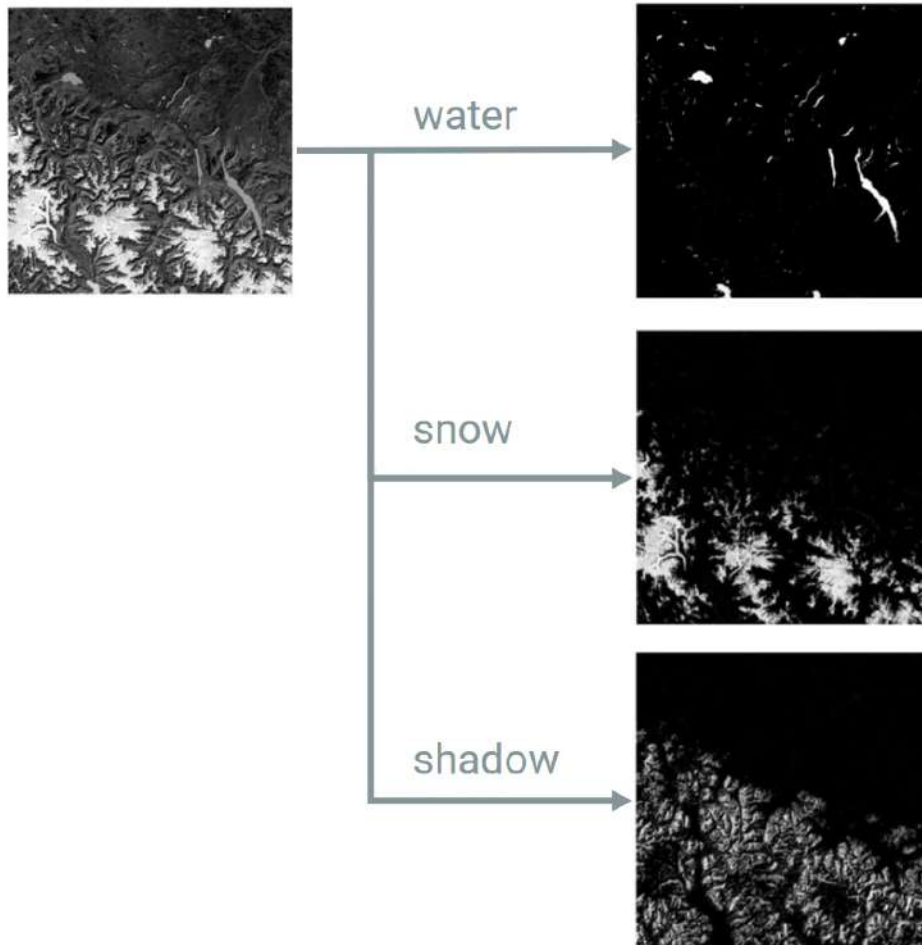
- Fully Convolutional Networks

## Skip connections

- building blocks of the ResNet architecture



# Skip connections



## Skip connections

- Fully Convolutional Networks

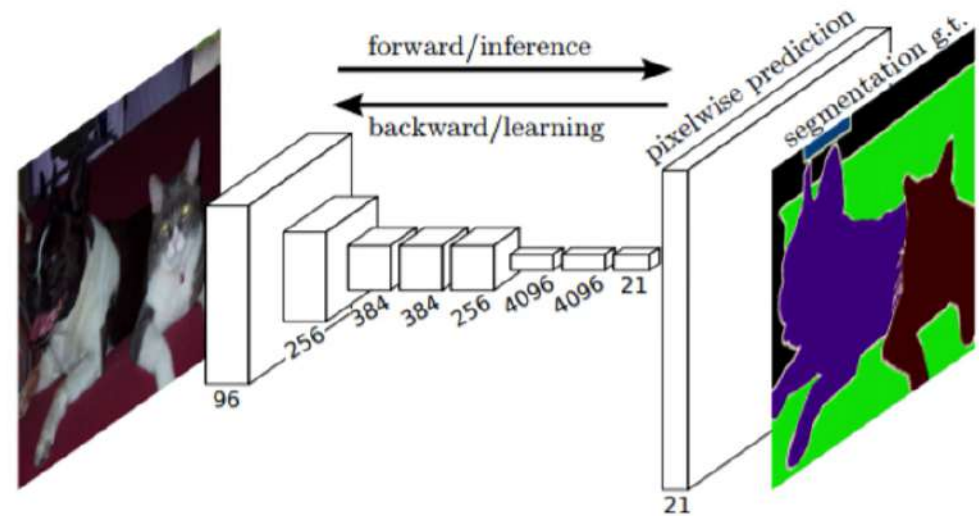
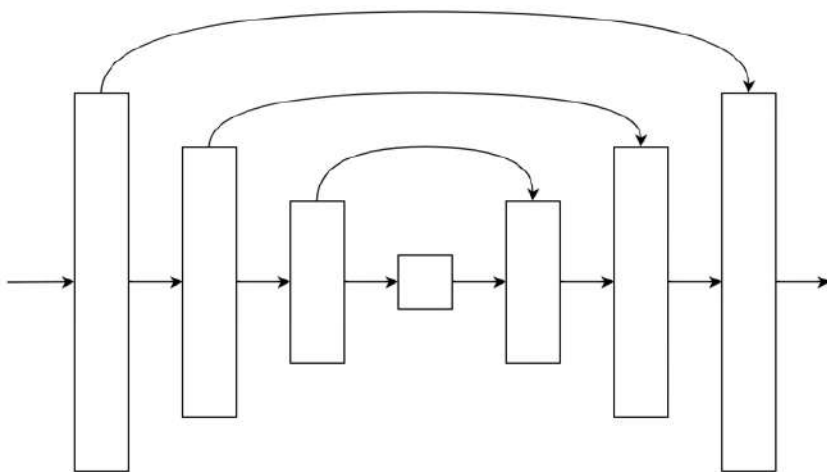


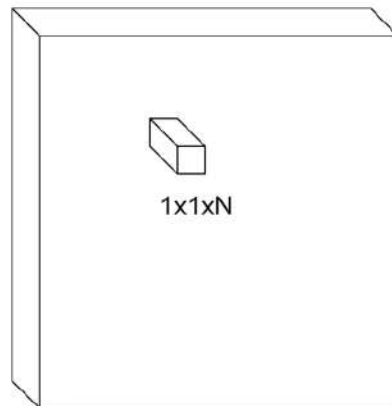
Figure credit: Long and Shelhamer (2015)

## How to choose kernel size?

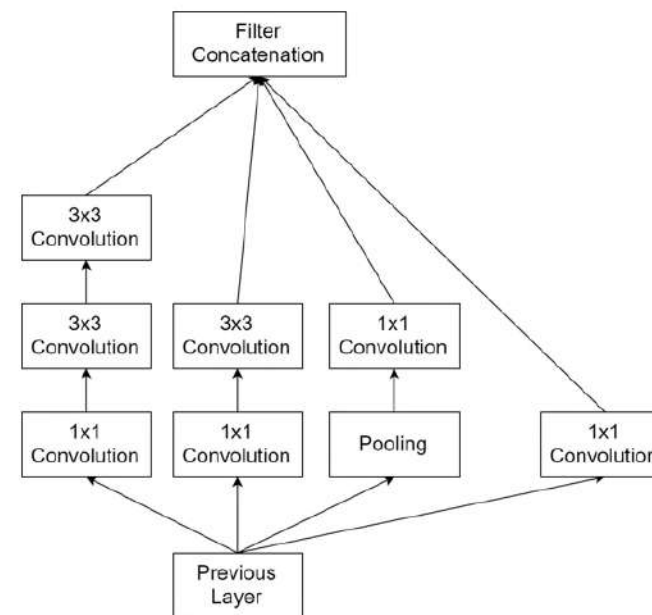
- 3x3 and 1x1 kernels usually work the best

## Pointwise (1x1) filters

- channel-wise dense layers
- learn cross-channel features



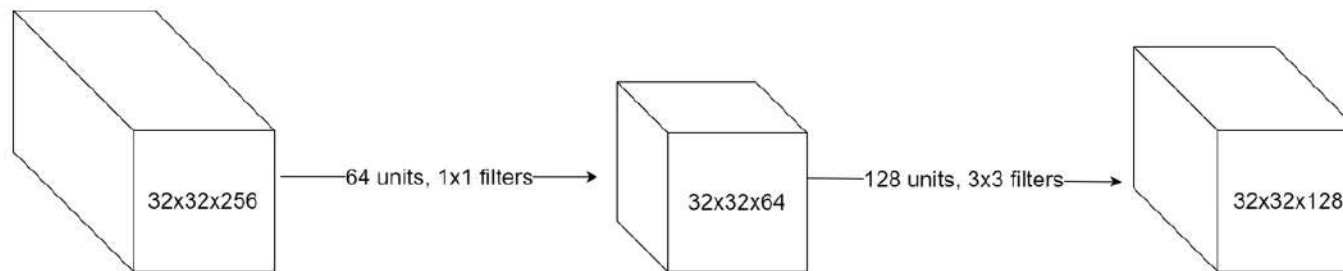
## Inception module



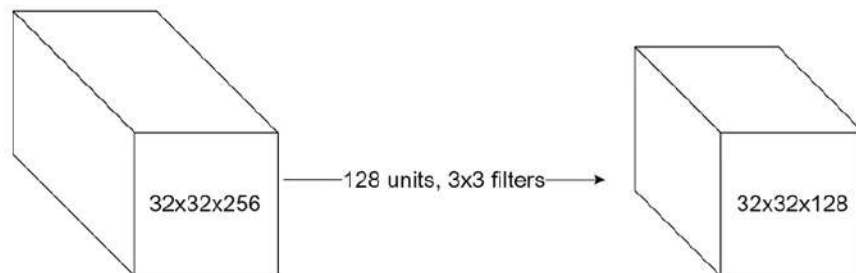


## Pointwise (1x1) filters

- Dimensionality reduction



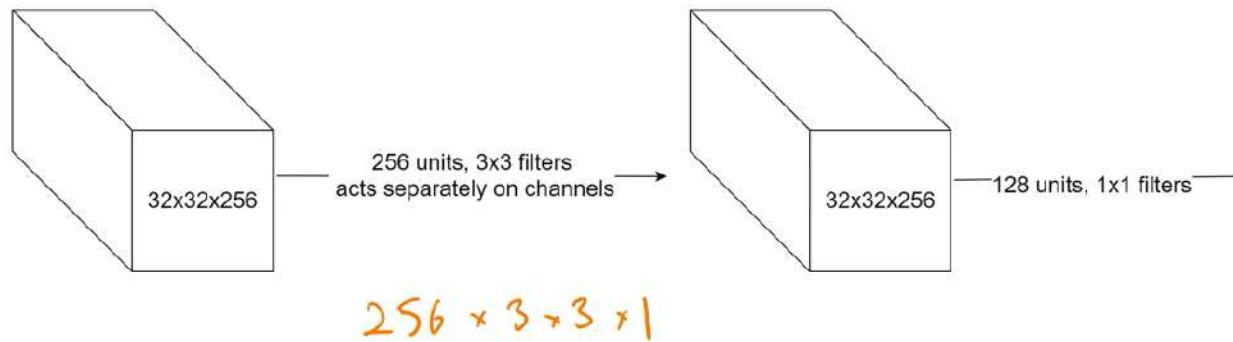
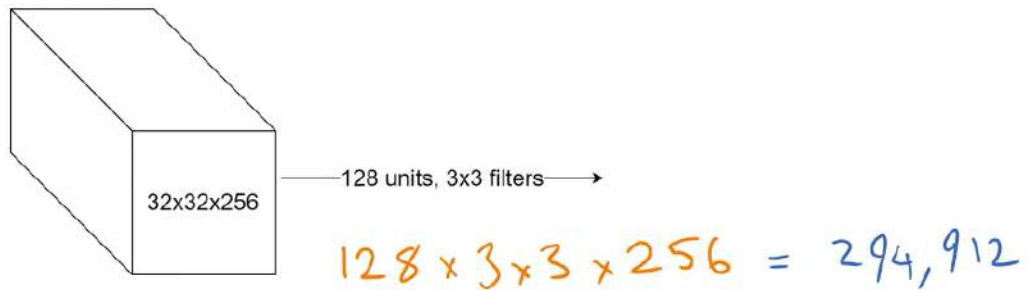
$$32 \times 32 \times 64 \times 1 \times 1 \times 256 + 32 \times 32 \times 128 \times 3 \times 3 \times 64 \approx 92 \text{ M}$$



$$32 \times 32 \times 128 \times 3 \times 3 \times 256 \approx 300 \text{ M}$$

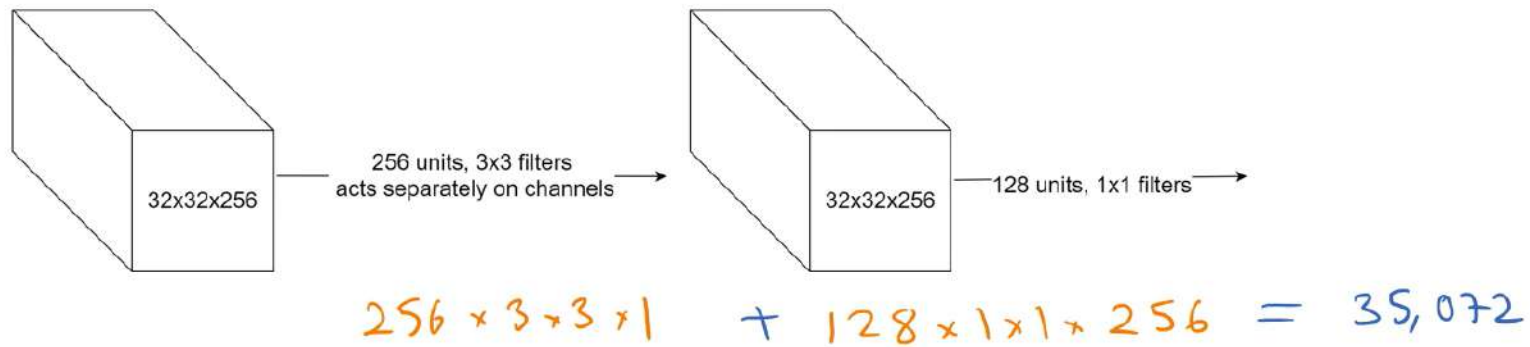
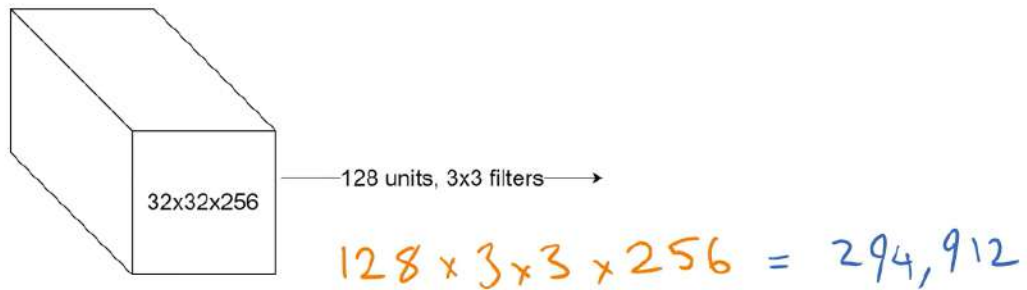
## Pointwise (1x1) filters

- Depthwise separable convolution

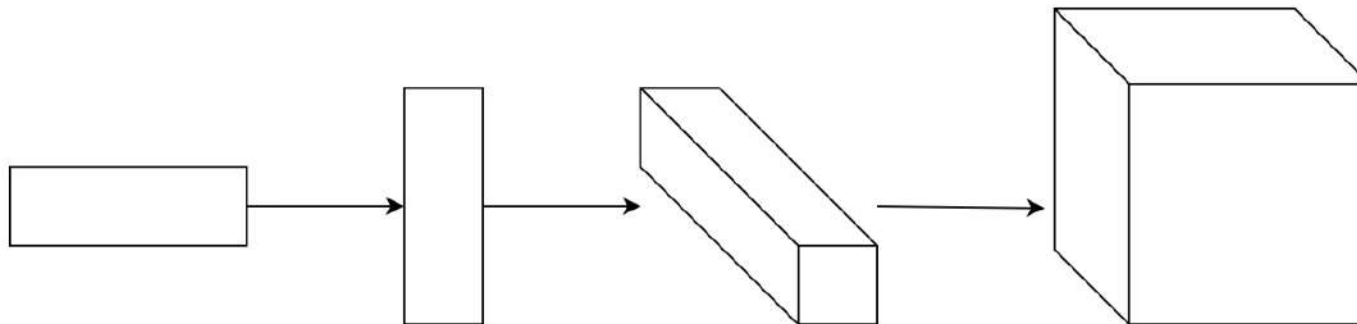
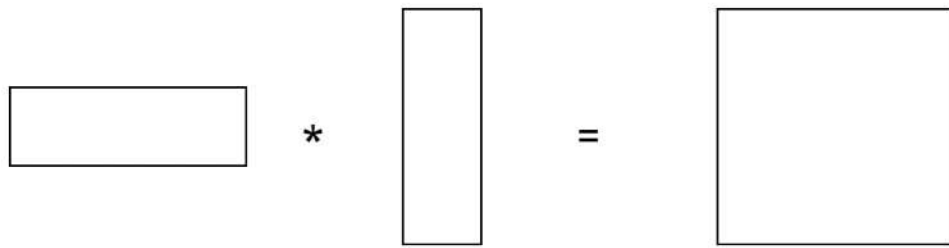


## Pointwise (1x1) filters

- Depthwise separable convolution

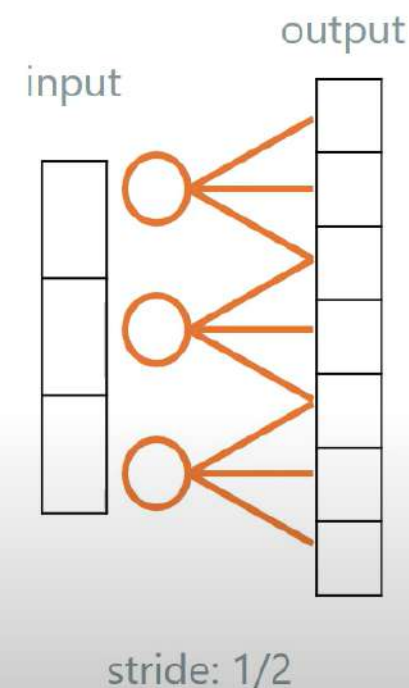
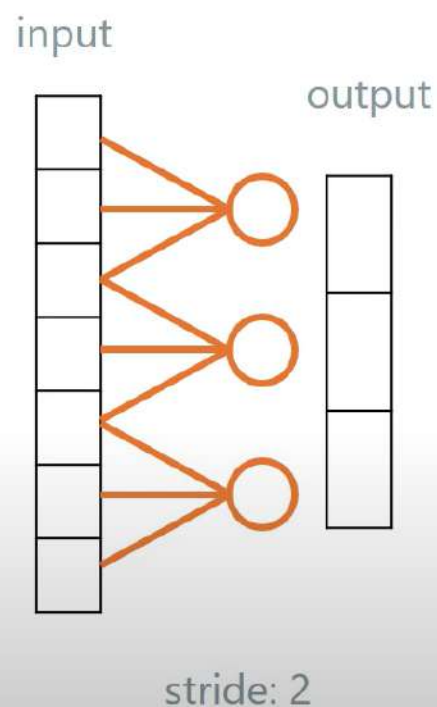


## Separable convolution



## How to choose stride?

- **Stride 1:** preserve spatial resolution
- **Stride 2:** downsample
- **Fractional stride (1/2):** upsample
  - transposed convolution
  - deconvolution



## How to choose pooling parameters?

- A very common setting
  - max pooling
  - pooling size: 2x2
  - same padding
- Global average pooling / pooling to fixed size
  - handles variable-sized inputs

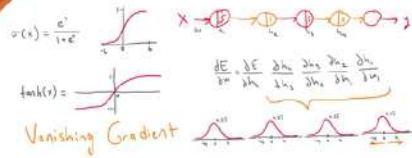
## How to choose activation functions?

- Short answer: **choose ReLU**
- For more information:

Deep Learning Crash Course

### ARTIFICIAL NEURAL NETWORKS

#### PART II Going Deeper





## What type of regularization to use?

- Short answer: use L2 weight decay and dropout
- For more information:



## How to choose the batch size?

- Image recognition tasks
  - batch size: 32
- Noisy gradient?
  - use larger batches
- Stuck in local minima? Out of memory?
  - use smaller batches


# Optimization

What we want to achieve:

Maximize **accuracy** = (#correct samples) / (#all samples)  
on a test set

What we actually do:

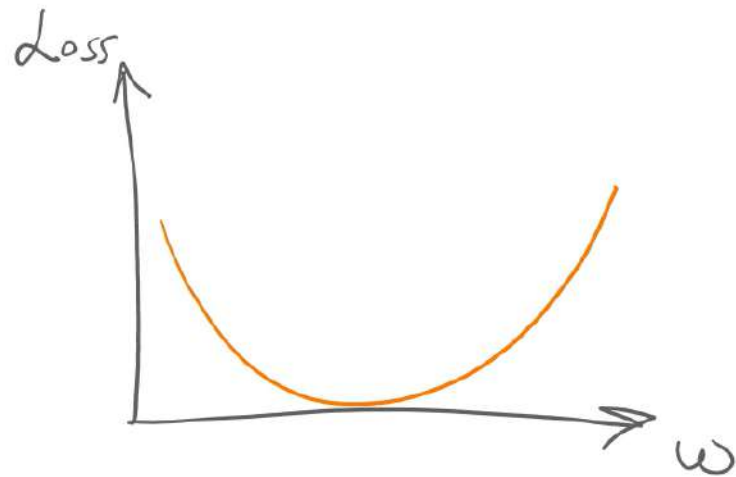
Minimize **cross entropy** =  $\sum p \log(q)$



True class probabilities

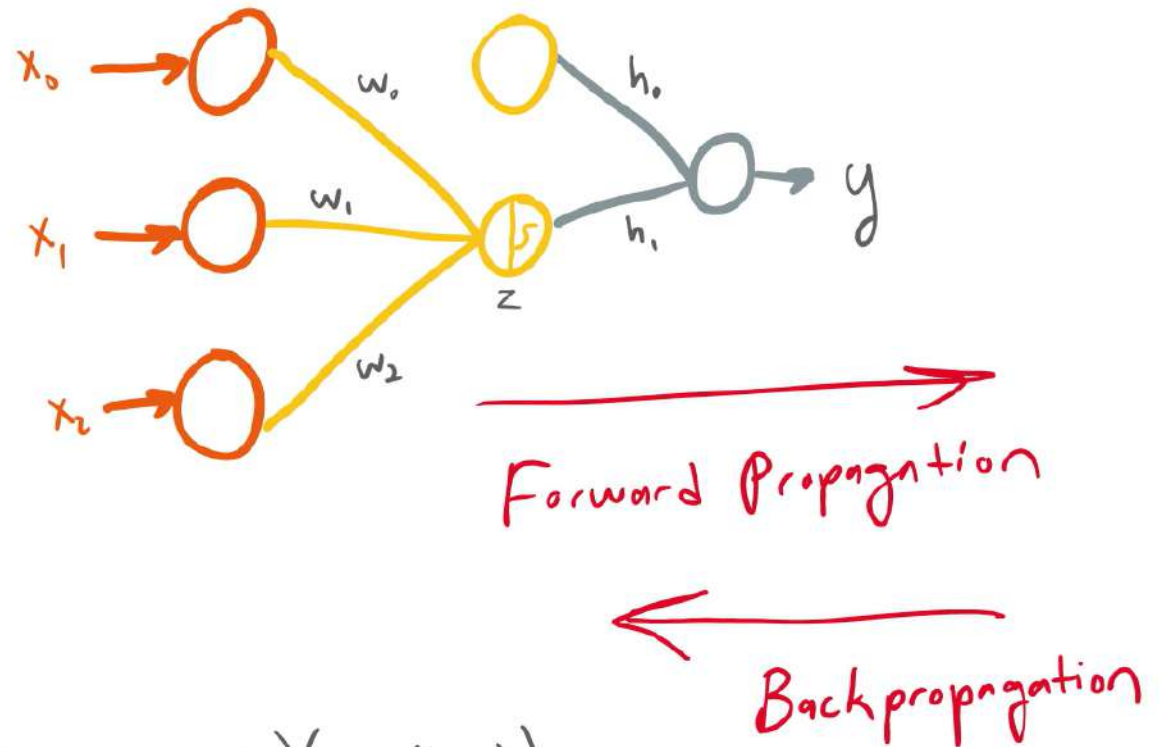
Predicted class probabilities

# Vanilla Stochastic Gradient Descent



$$w \leftarrow w - \alpha \frac{\partial E}{\partial w}$$

new weights  $\leftarrow$  (old weights) - (learning rate)(gradient)



# Momentum

(new weights) <- (old weights) – (learning rate) (gradient)

# Momentum

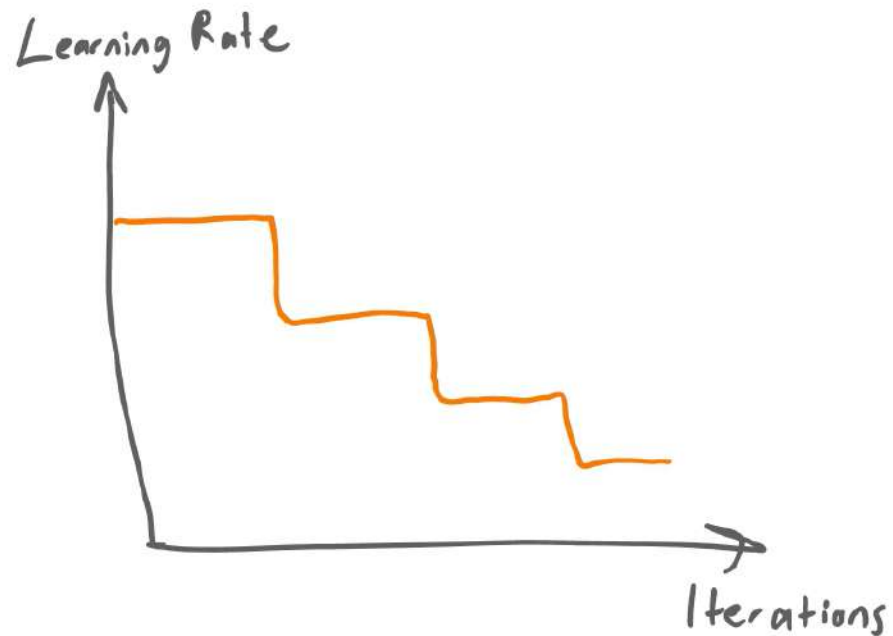
(new weights) <- (old weights) – (learning rate) (gradient)  
+ past gradients

(accumulator) <- (old accumulator) (momentum) + (gradient)

(new weights) <- (old weights) – (learning rate) (accumulator)

# Learning rate schedules

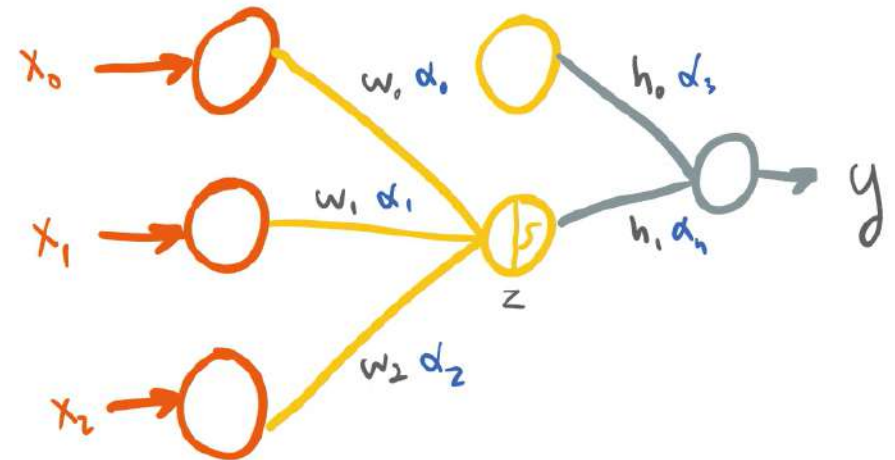
(new weights)  $\leftarrow$  (old weights)  $-$  (learning rate) (gradient)



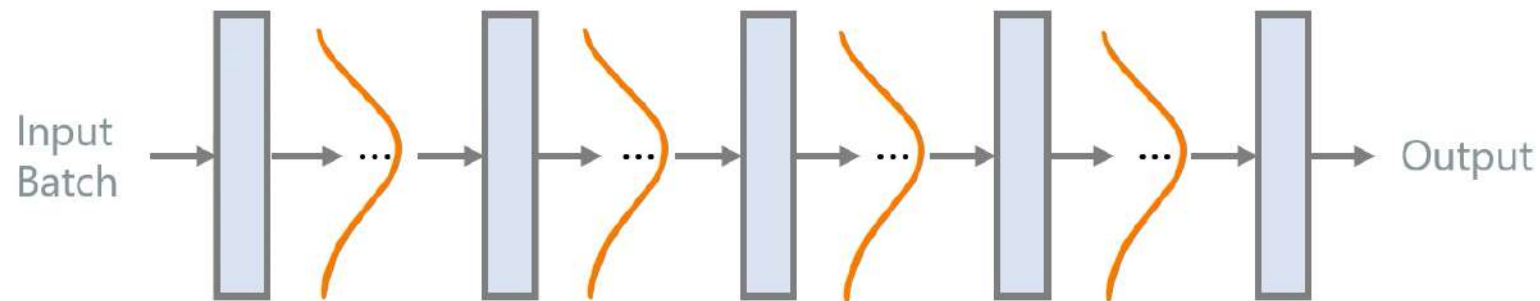


# Adaptive Methods

- AdaGrad
  - Large gradient: decrease  $\alpha$  faster
  - Small gradient: decrease  $\alpha$  slower
- RMSProp
  - Moving average of gradients
- Adam
  - Adaptive Moment Estimation
  - RMSProp + Momentum



# Batch Normalization



Subtract batch mean  
Divide by standard deviation

$$h \leftarrow \frac{h - \mu}{\sigma}$$

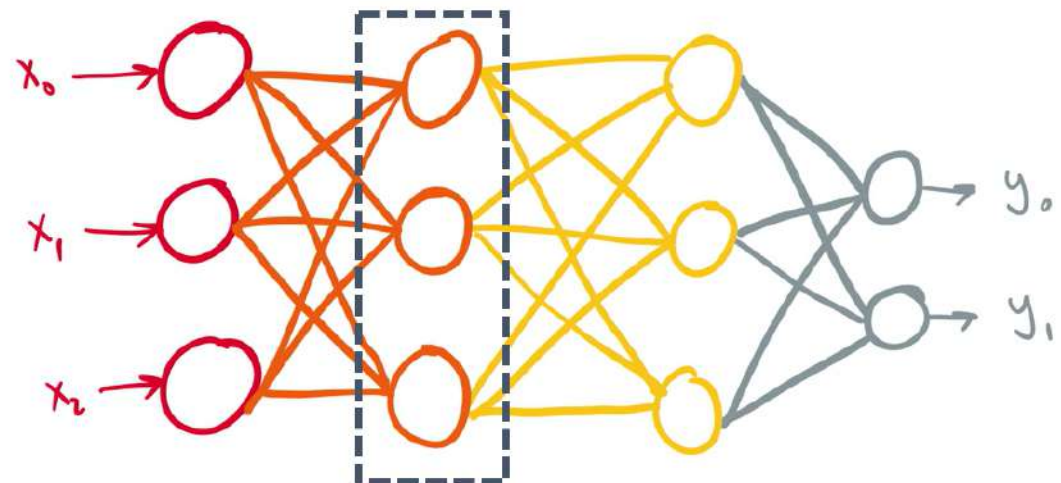
Scale and shift

$$h \leftarrow \gamma h + \beta$$

# Initialization

- Biases
  - Initialize to zero (or another small constant)

- Weights
  - Initialize to zero? No!
  - Break the symmetry
  - Random values



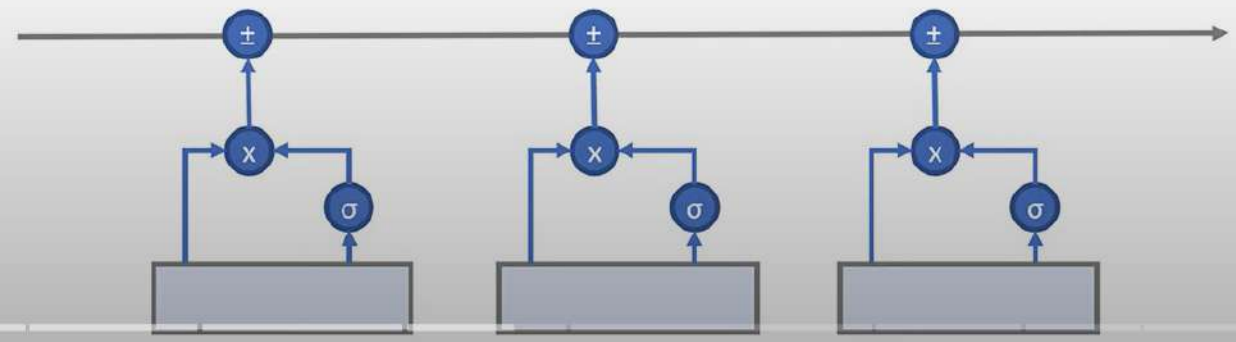
All would get  
the same updates

# Initialization

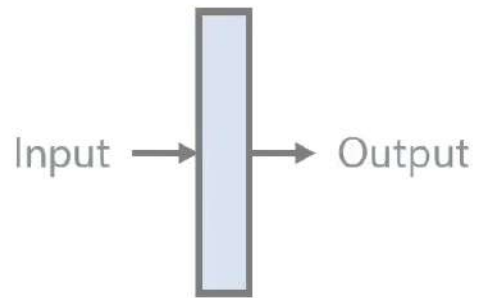
- Glorot (Xavier) initializer

$$W \sim \text{Uniform}(-r, r) \quad r = \sqrt{\frac{6}{n_{in} + n_{out}}}$$

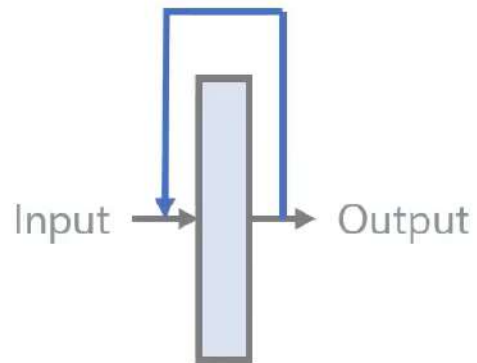
# RECURRENT NEURAL NETWORKS



# Feedforward vs Recurrent



Feedforward Neural Network



Recurrent Neural Network (RNN)

# Sequential Data

- Audio
- Video
- Text
- Biomedical data





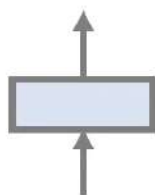
# Types of Inputs and Outputs

Task	Input	Output
Machine translation	Text	Text
Speech recognition	Audio	Text
Speech synthesis	Text	Audio
Sentiment analysis	Text	Categorical variable
Text generator	Random number	Text

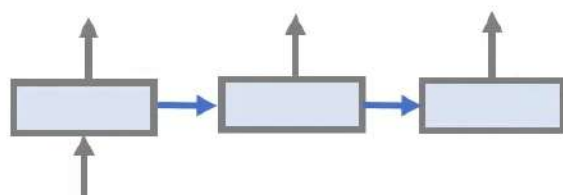


# Types of RNNs

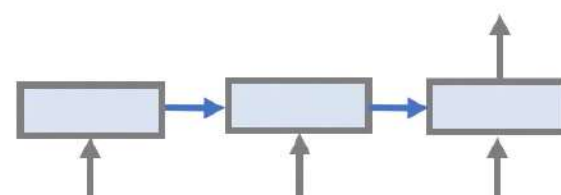
One to one



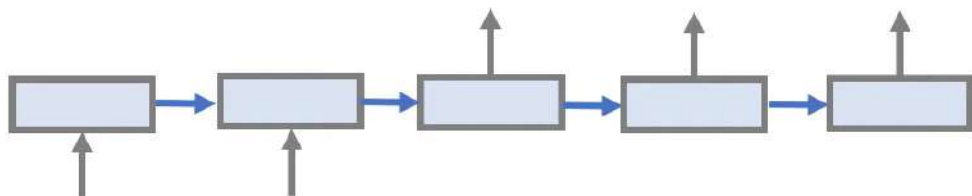
One to many



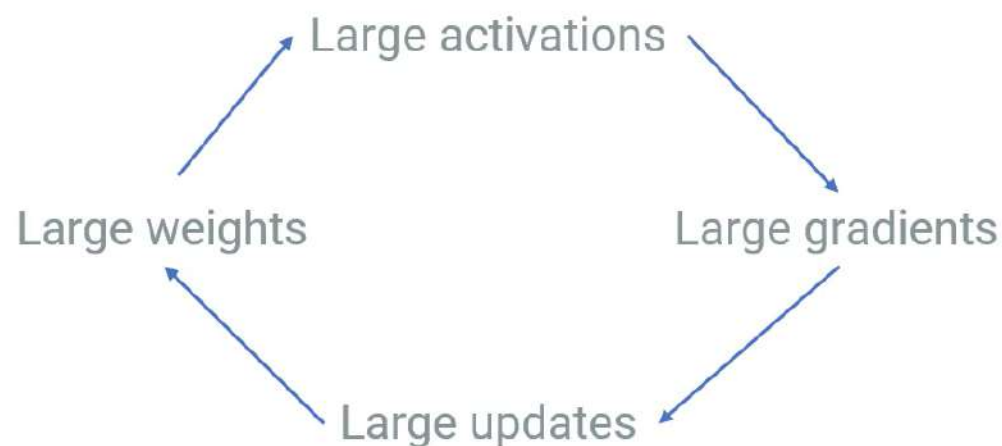
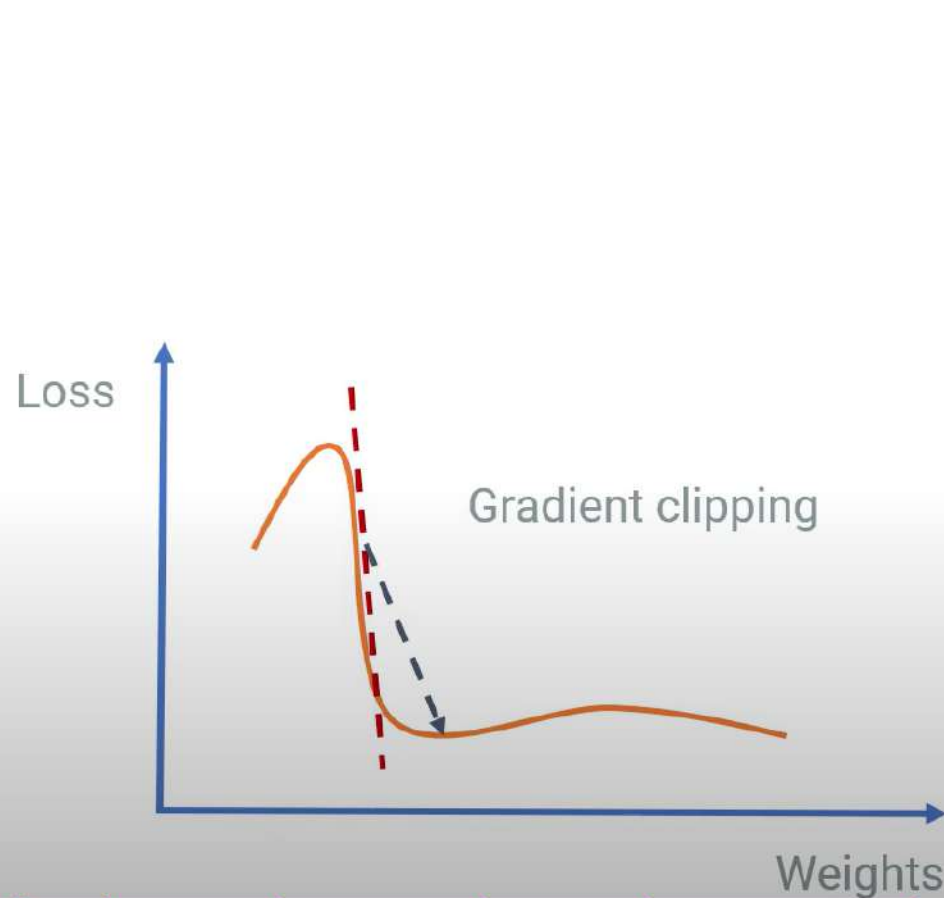
Many to one



Many to many



# Exploding Gradients



# Gated Modules

- LSTM (Long Short Term Memory)
- GRU (Gated Recurrent Unit)