

Announcement

- Homework 1 has been set out
- There will be total 3 homework assignments this semester
- Final project topic: Google Research Football

Google Research Football (GRF)

- A Novel Reinforcement Learning framework
- Building agents to master the football game
- Compatible with OpenAI Gym API



(a) Kickoff



(b) Yellow card



(c) Corner kick

GRF Modes

- Single-agent mode: 11 vs 11
 - Each team has 11 players
 - At each moment, AI only controls an active player who controls the ball
- Multi-agent mode: 5 vs 5
 - Each team has 5 players, including a rule-based keeper
 - AI controls 4 players
- Football academy
 - Curriculum learning
 - Scenarios: scoring, passing, running with the ball, etc.

Single-Agent Mode: 11 vs 11



Multi-Agent Mode: 5 vs 5



Final Project on GRF

- We will have a tournament
- Tentative mode: Multi-Agent 5 vs 5
- Project grading: reports and tournament results
 - Research ideas are more important than tournament results
- Website:
<https://ai.googleblog.com/2019/06/introducing-google-research-football.html>
- The Framework Paper:
<https://arxiv.org/pdf/1907.11180.pdf>

Deep Reinforcement Learning

Lecture 6: Deep Learning for RL

Instructor: Chongjie Zhang

Tsinghua University

Large-Scale Reinforcement Learning

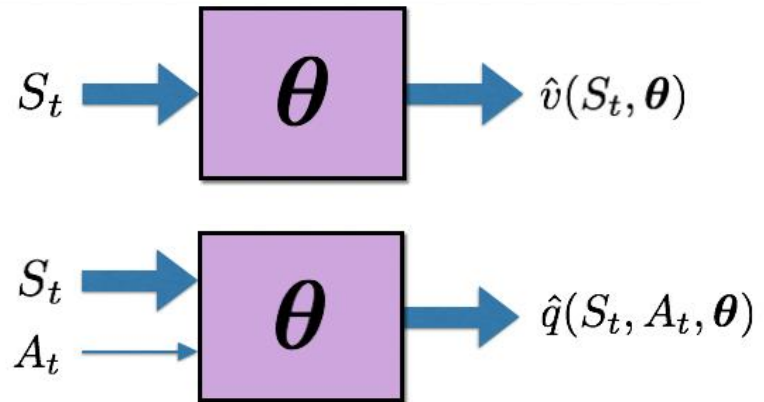
- So far we have represented value functions by a lookup table
 - Every **state** s has an entry $V(s)$, or
 - Every **state-action** pair (s, a) has an entry $Q(s, a)$
- Reinforcement learning should be used to solve large problems, e.g.
 - Backgammon: 10^{20} states
 - Computer Go: 10^{170} states
 - Helicopter, robot, ...: enormous continuous state space
 -
- Tabular methods clearly cannot handle this.. **why?**
 - There are too many states and/or actions to store in memory
 - It is too slow to learn the value of each state individually
 - **You cannot generalize across states!**

Value Function Approximation (VFA)

- Solution for large MDPs:
 - Estimate value function with function approximation
- Value function approximation (VFA) replaces the table with a general parameterized form:

$$\hat{v}(s, \theta) \approx v_{\pi}(s)$$

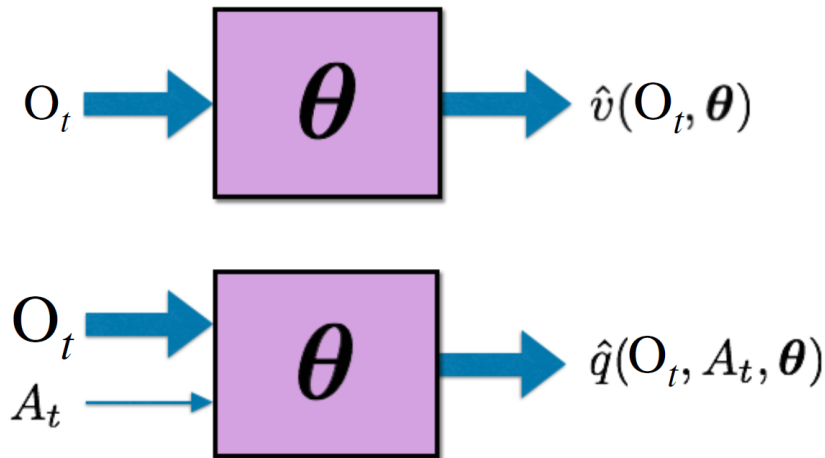
$$\text{or } \hat{q}(s, a, \theta) \approx q_{\pi}(s, a)$$



- Why this is a good idea?
 - Generalization: those functions can be trained to map *similar* states to similar values.

End-to-End RL

- End-to-end RL methods replace the hand-designed state representation with raw observations.



- Good: We get rid of manual design of state representations
- Bad: we need tons of data to train the network since O_t usually WAY more high dimensional than hand-designed S_t

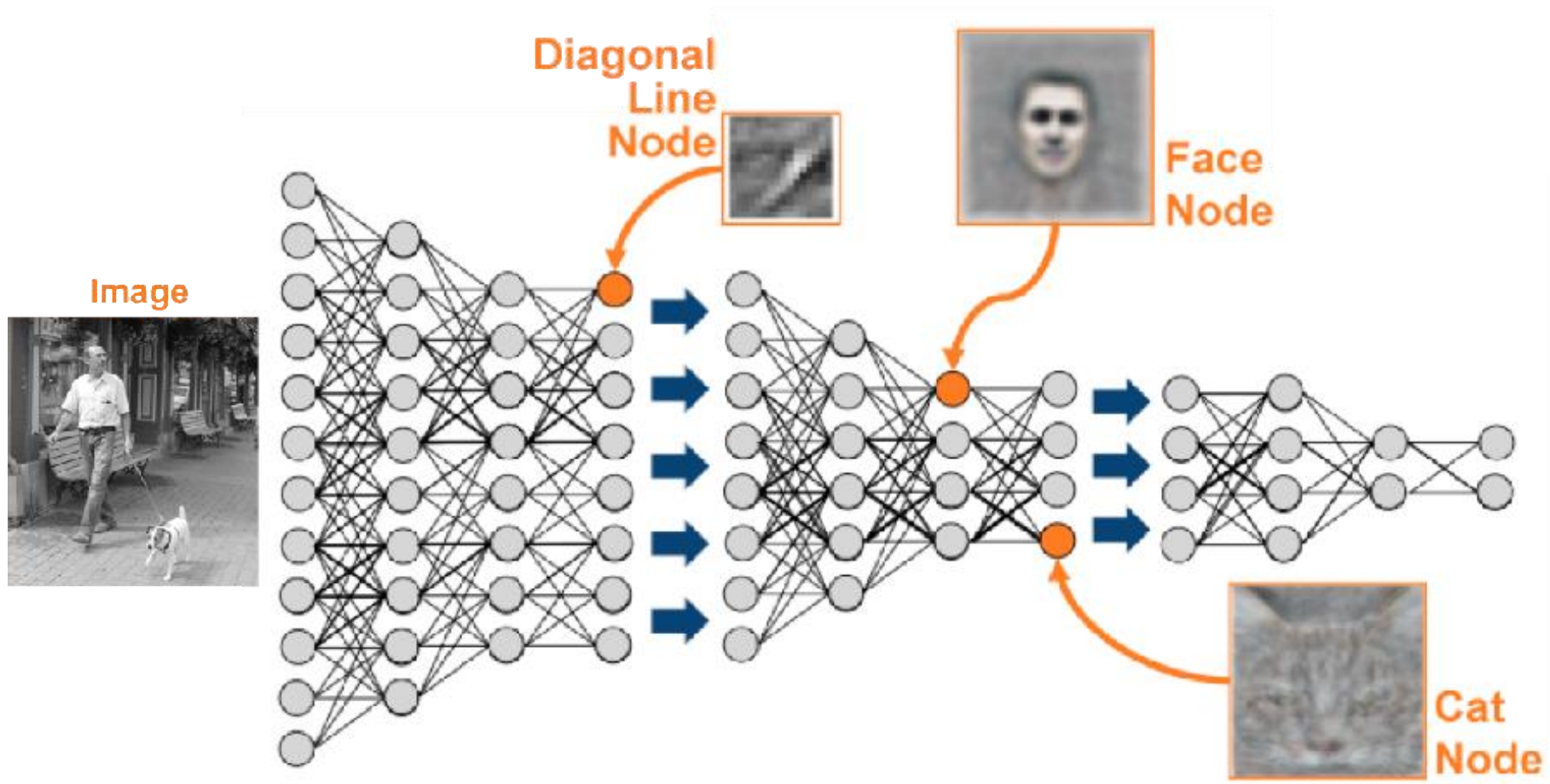
Which Function Approximation?

- There are many function approximators, e.g.
 - Linear combinations of features
 - Neural networks
 - Decision tree
 - Nearest neighbour
 - ...
- In this lecture we will consider:
 - Linear combinations of features
 - **Neural networks**

Today's Lecture

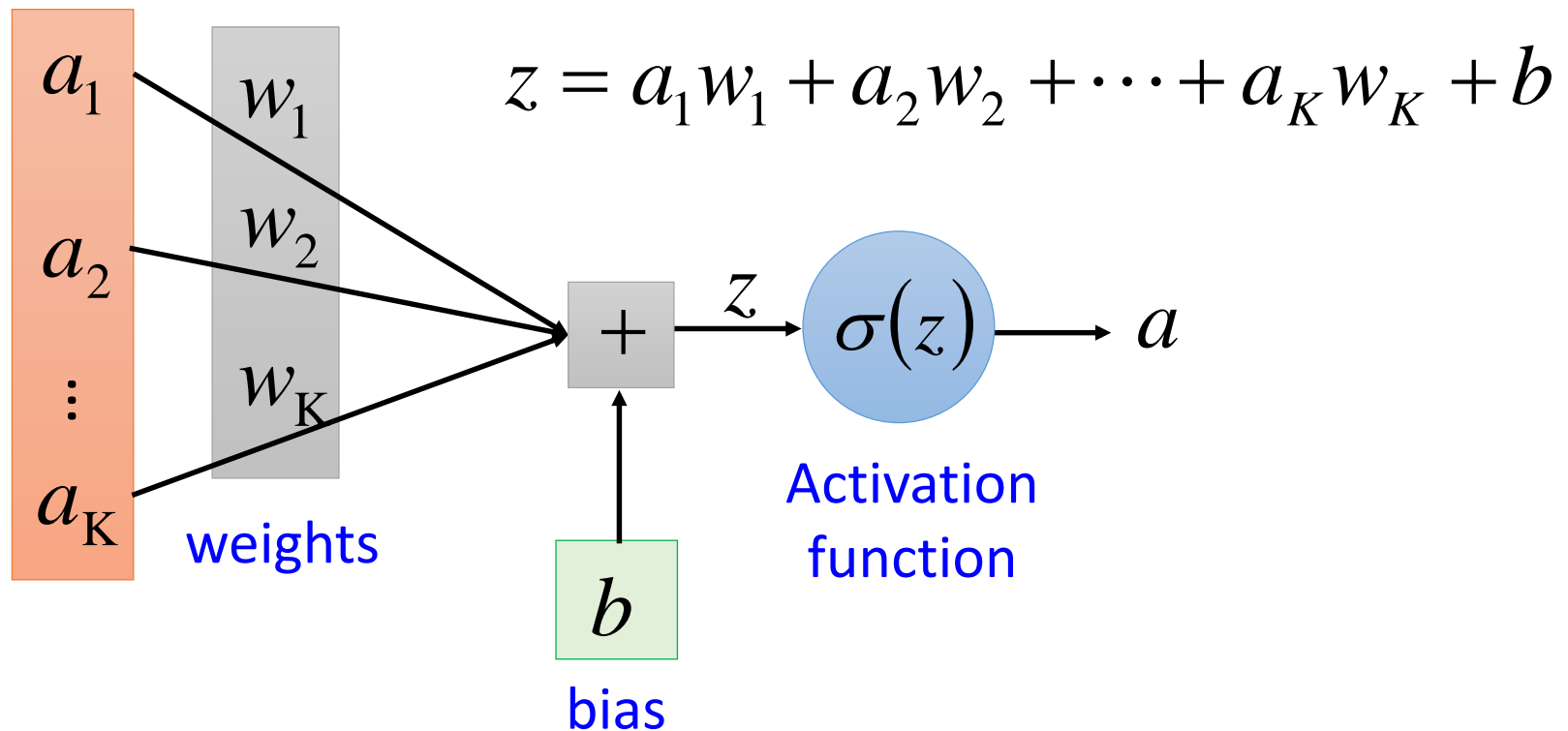
- Neural Networks
- Training Neural Networks
- Convolutional Neural Networks
- Recurrent Neural Networks

Deep Neural Networks



Element of Neural Network

Neuron $f: R^K \rightarrow R$

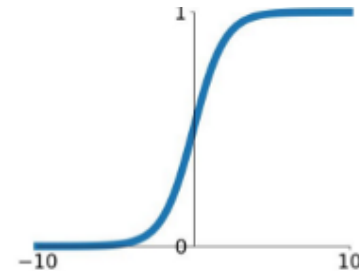


Activation Function - Sigmoid

- Sigmoid activation function:

Sigmoid

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

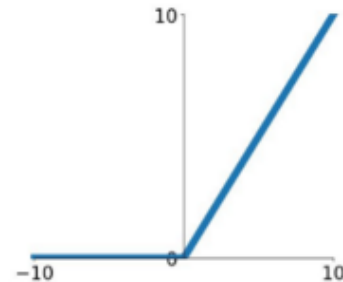


- Squashes the neuron's output between 0 and 1
- Always positive
- Bounded
- Strictly Increasing

Activation Function - ReLU

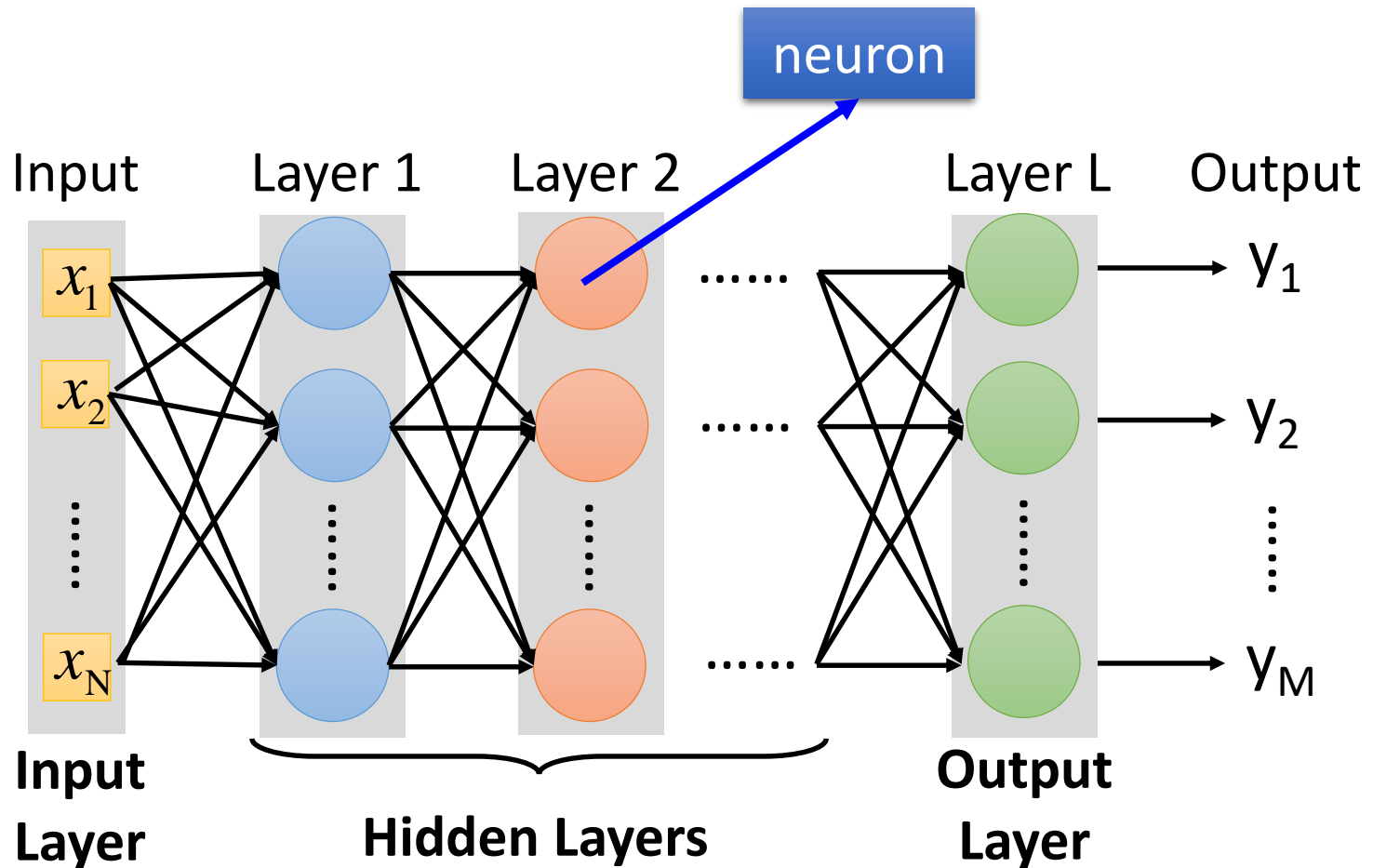
- Rectified linear (ReLU) activation function

$$\text{ReLU} \\ \max(0, x)$$



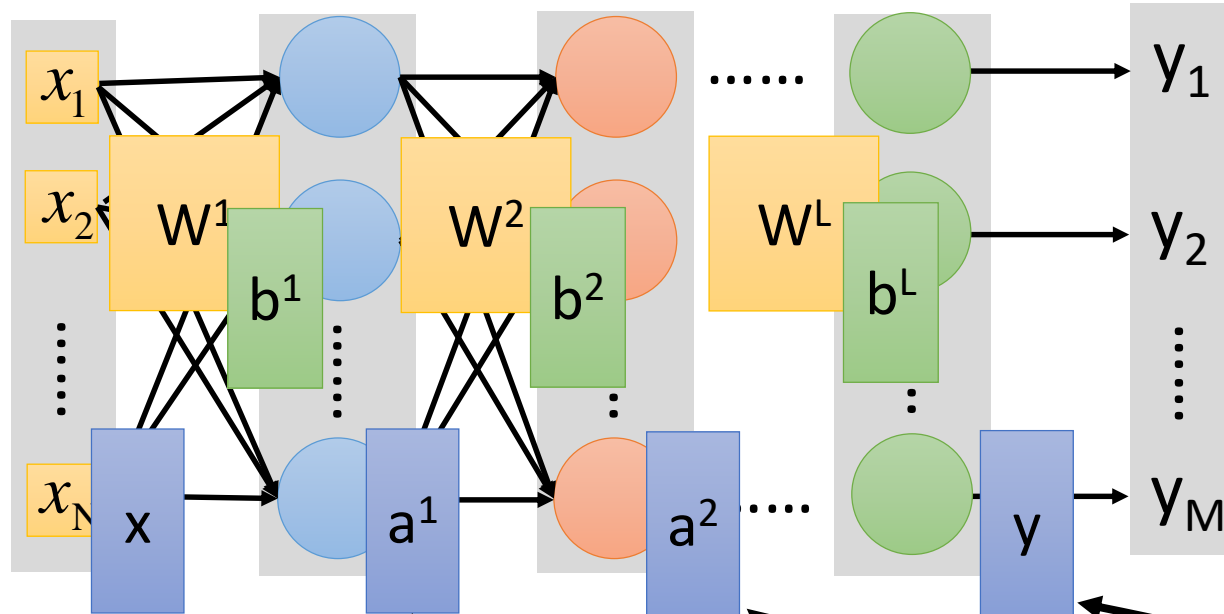
- Bounded below by 0 (always non-negative)
- Tends to produce units with sparse activities
- Not upper bounded
- Increasing

Neural Network



Deep means many hidden layers

Neural Network



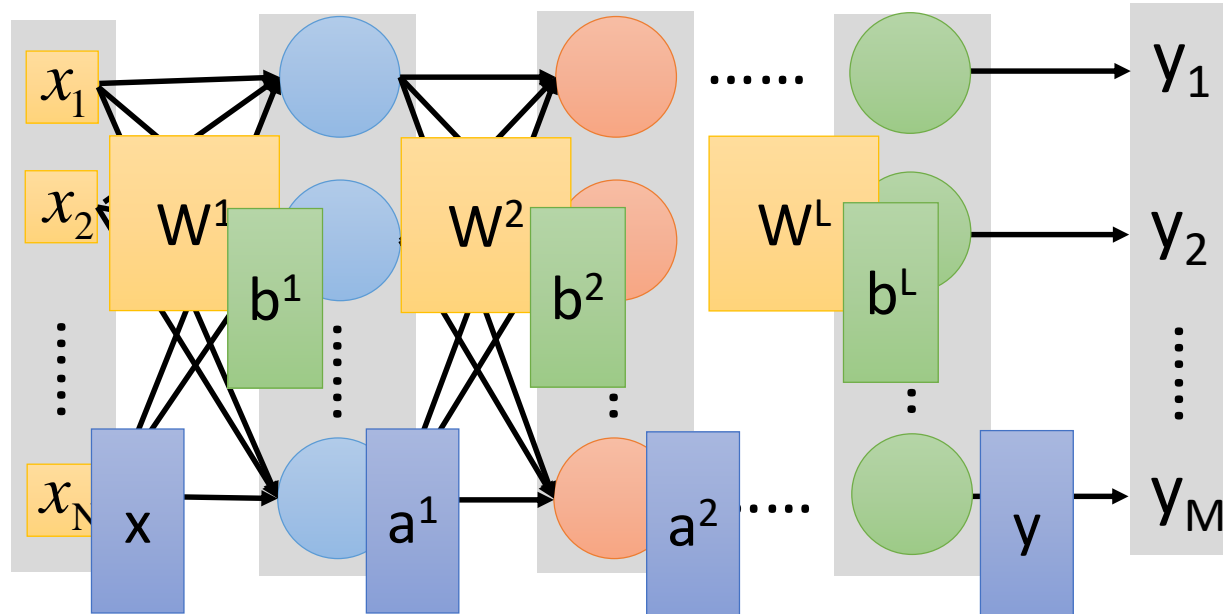
$$\sigma(W^1 x + b^1)$$

$$\sigma(W^2 a^1 + b^2)$$

$$\sigma(W^L a^{L-1} + b^L)$$

Arrows indicate the flow of information from the input layer to the first hidden layer, from the first hidden layer to the second hidden layer, and from the final hidden layer to the output layer.

Neural Network



$$y = f(x)$$

Using parallel computing techniques to speed up matrix operation

$$= \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

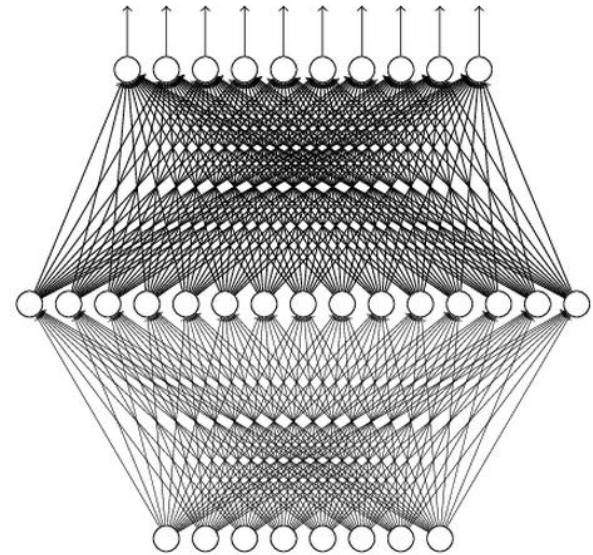
Universal Approximation Theorem

Any continuous function f

$$f : R^N \rightarrow R^M$$

Can be realized by a network
with one hidden layer

(given **enough** hidden
neurons)

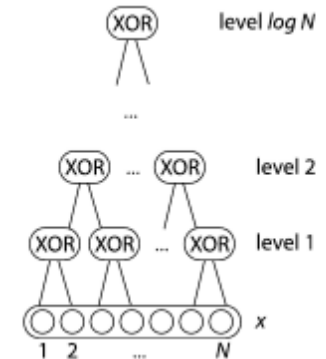
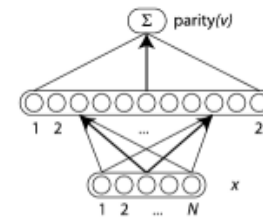
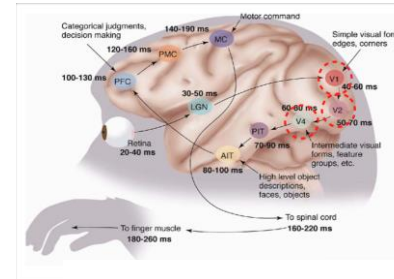


Reference for the reason:
<http://neuralnetworksanddeeplearning.com/chap4.html>

Why “Deep” neural network not “Fat” neural network?

Why Use Deep Networks?

- Motivation from Biology
 - Visual Cortex
- Motivation from Circuit Theory
 - Compact representation
- Modularity
 - More efficiently using data
- In Practice: works better for many domains
 - Hard to argue with results



Today's Lecture

- Neural Networks
- **Training Neural Networks**
- Convolutional Neural Networks
- Recurrent Neural Networks

Training

- Empirical Risk Minimization:

$$\arg \min_{\theta} \frac{1}{T} \sum_t \underbrace{l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})}_{\text{Loss function}} + \underbrace{\lambda \Omega(\theta)}_{\text{Regularizer}}$$

- Learning is cast as optimization.
 - For classification problems, we would like to minimize classification error, e.g., logistic or cross entropy loss.
 - For regression problems, we would like to minimize regression error, e.g., L1 or L2 distance from ground-truth

Stochastic Gradient Descent

- Perform updates after seeing each example:

- Initialize: $\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$

- For $t=1:T$

- for each training example $(\mathbf{x}^{(t)}, y^{(t)})$

$$\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) - \lambda \nabla_{\theta} \Omega(\theta)$$

$$\theta \leftarrow \theta + \alpha \Delta$$

Training epoch
=
Iteration of all examples

- Training a neural network, we need:

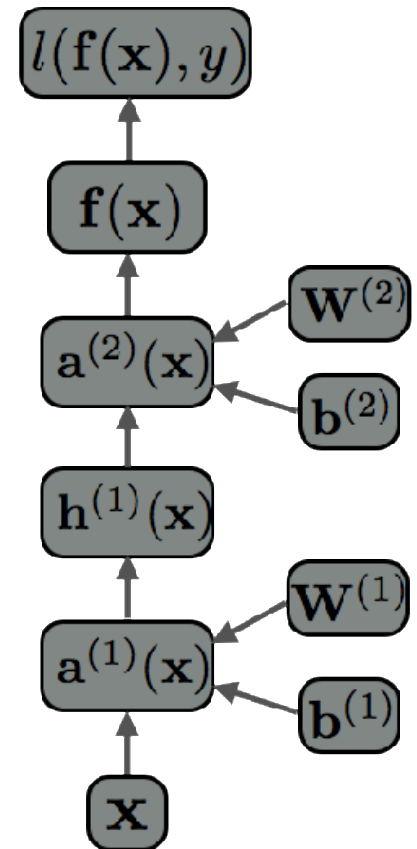
- **Loss function**: $l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$

- A procedure to **compute gradients**: $\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$

- **Regularizer** and its gradient: $\Omega(\theta), \nabla_{\theta} \Omega(\theta)$

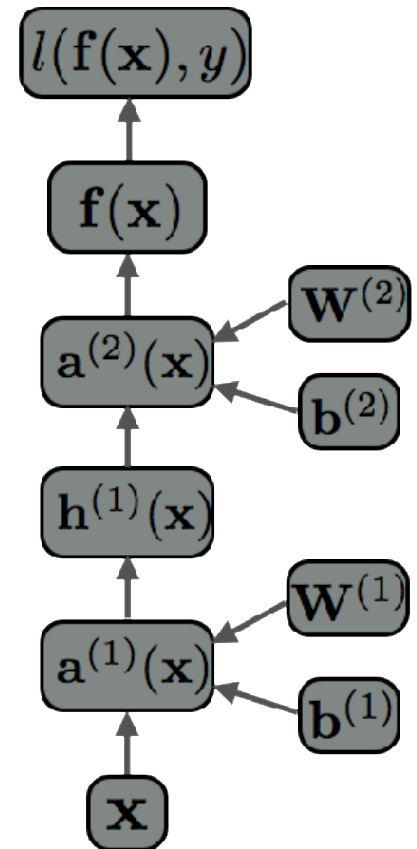
Computational Flow Graph

- Forward propagation can be represented as an acyclic flow graph
- Forward propagation can be implemented in a modular way:
 - Each box can be an object with an **fprop method**, that computes the value of the box given its children
 - Calling the **fprop method** of each box in the right order yields forward propagation



Computational Flow Graph

- Each object also has a **bprop** method
 - it computes the gradient of the loss with respect to each child box.
- By calling **bprop** in the reverse order, we obtain backpropagation



Mini-batch and Momentum

- Make updates based on a mini-batch of examples (instead of a single example)
 - the gradient is the average regularized loss for that mini-batch
 - can give a more accurate estimate of the gradient
- Momentum: Can use an exponential average of previous gradients:

$$\overline{\nabla}_{\theta}^{(t)} = \nabla_{\theta} l(\mathbf{f}(\mathbf{x}^{(t)}), y^{(t)}) + \beta \overline{\nabla}_{\theta}^{(t-1)}$$

- can get pass plateaus more quickly, by “gaining momentum”

Today's Lecture

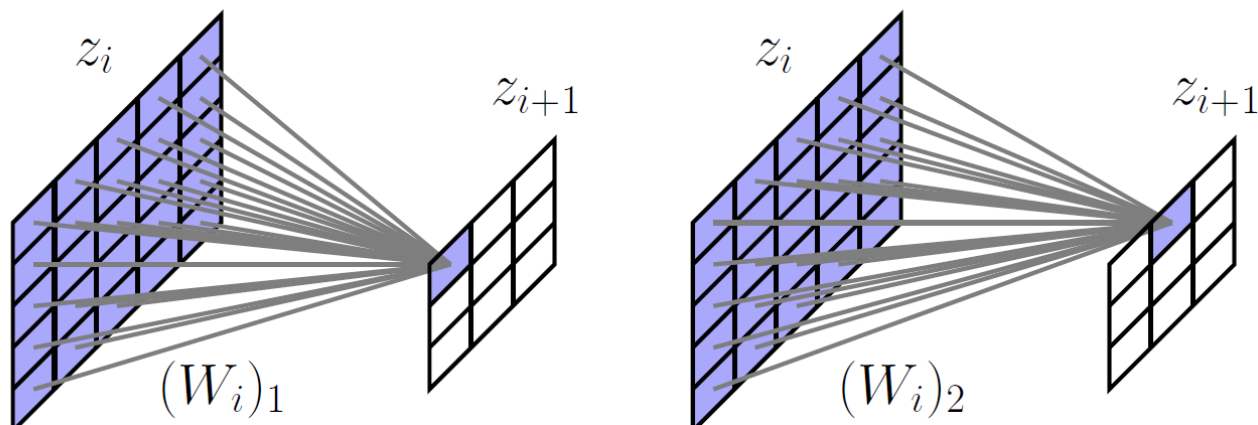
- Neural Networks
- Training Neural Networks
- **Convolutional Neural Networks**
- Recurrent Neural Networks

Problems with Fully Connected Networks

A 256x256 (RGB) image \implies $\sim 200\text{K}$ dimensional input x

A fully connected network would need a very large number of parameters, very likely to overfit the data

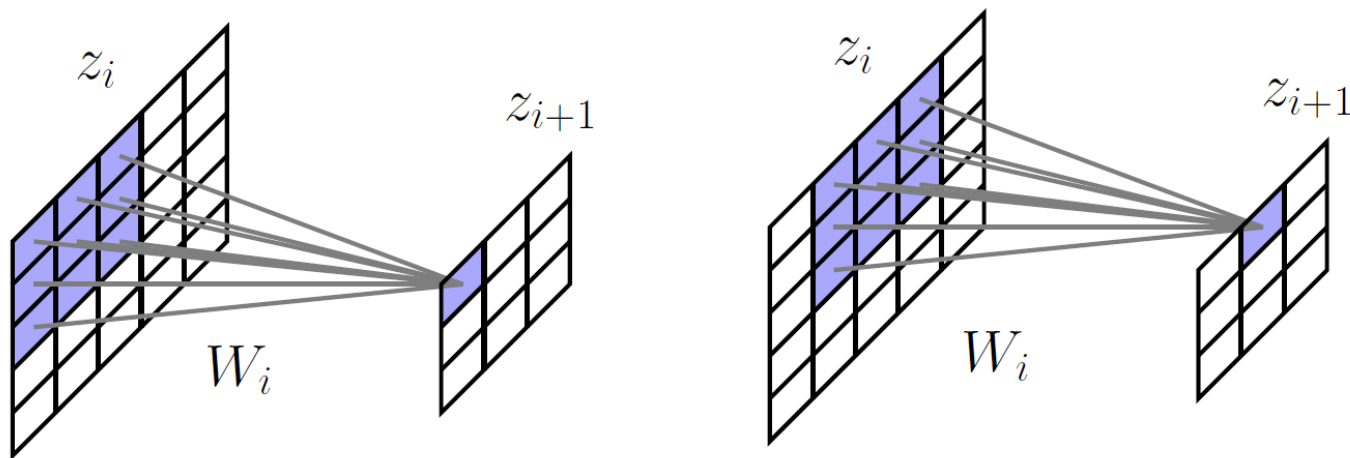
Generic deep network also does not capture the “natural” invariances we expect in images (translation, scale)



Convolutional Neural Networks (CNNs)

To create architectures that can handle large images, restrict the weights in two ways

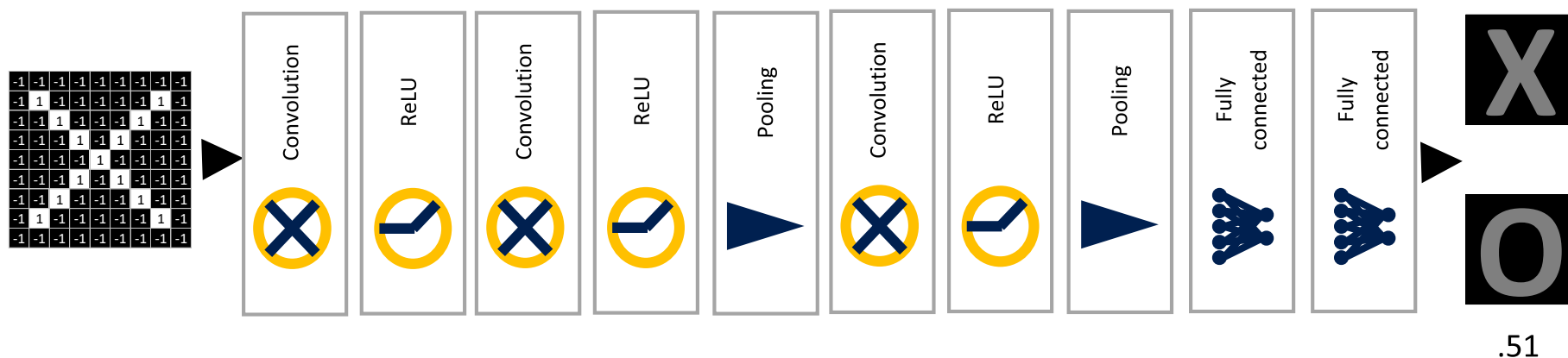
1. Require that activations between layers only occur in “local” manner
2. Require that all activations share the same weights



These lead to an architecture known as a convolutional neural network

Convolutional Neural Networks (CNNs)

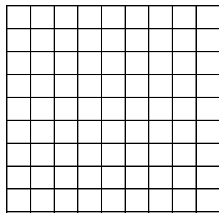
- Containing different types of layers
 - Convolution
 - Non-linearity
 - Pooling (or downsampling)
 - Fully connected layer



A toy ConvNet: X's and O's

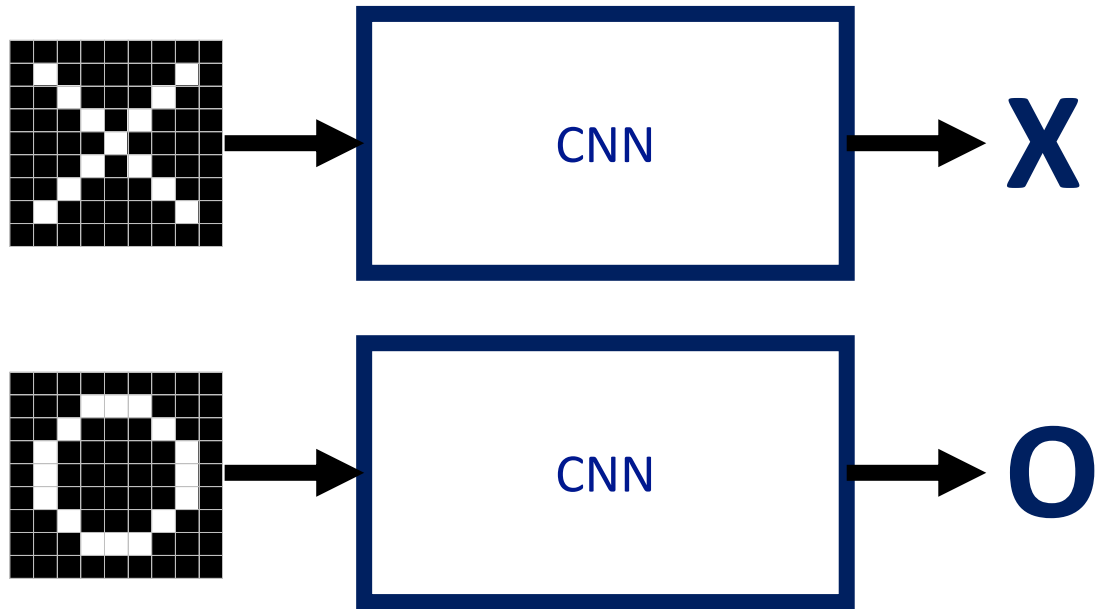
Says whether a picture is of an X or an O

A two-dimensional
array of pixels

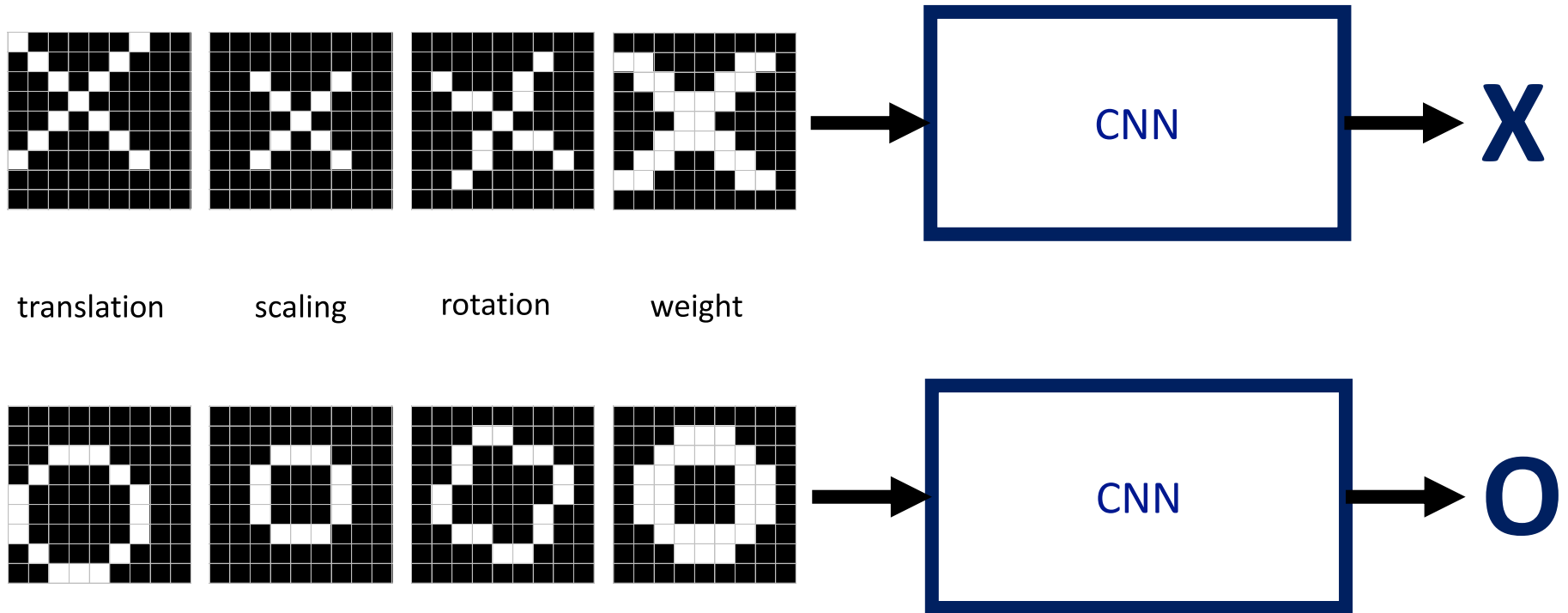


X or **O**

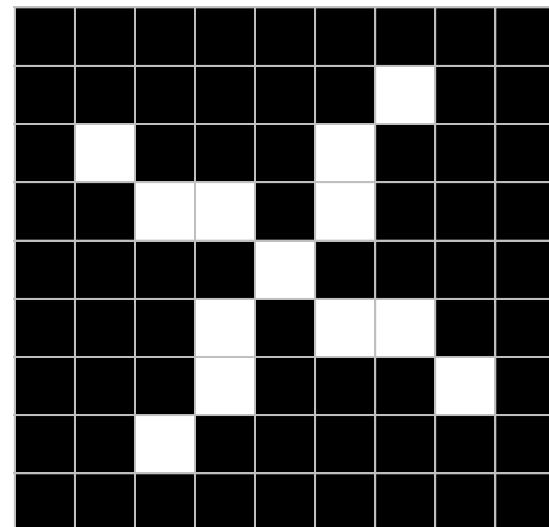
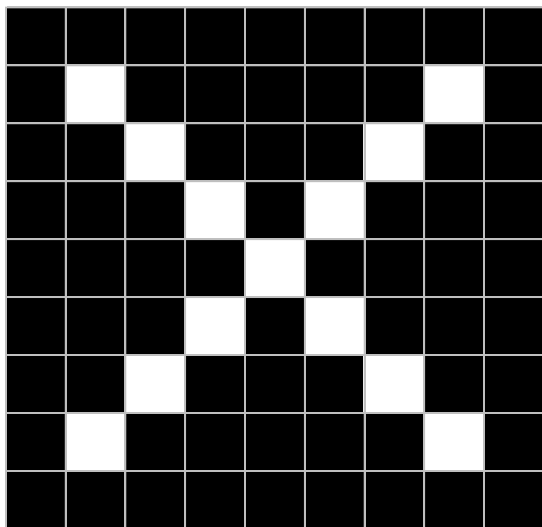
For example



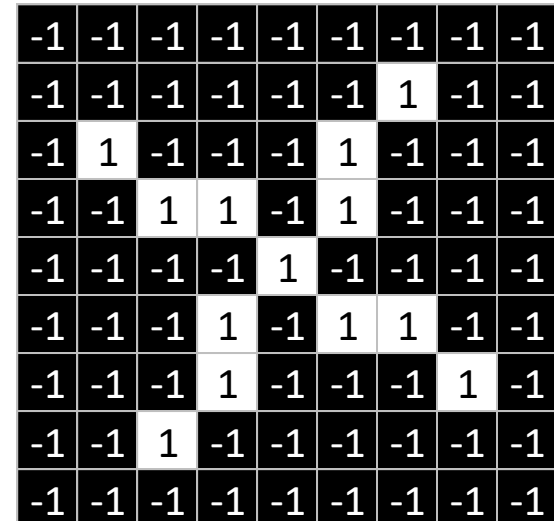
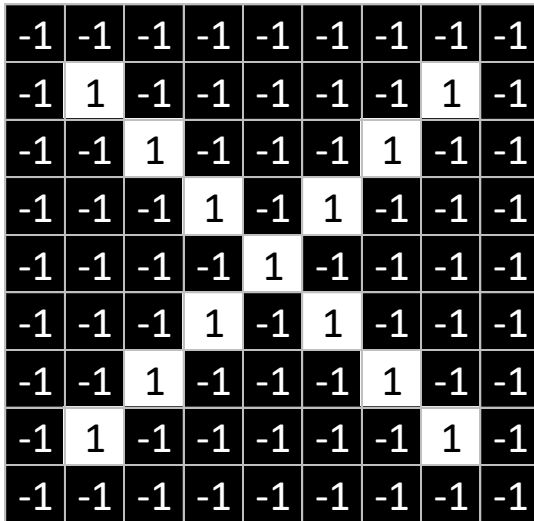
Trickier cases



Deciding is hard



What computers see

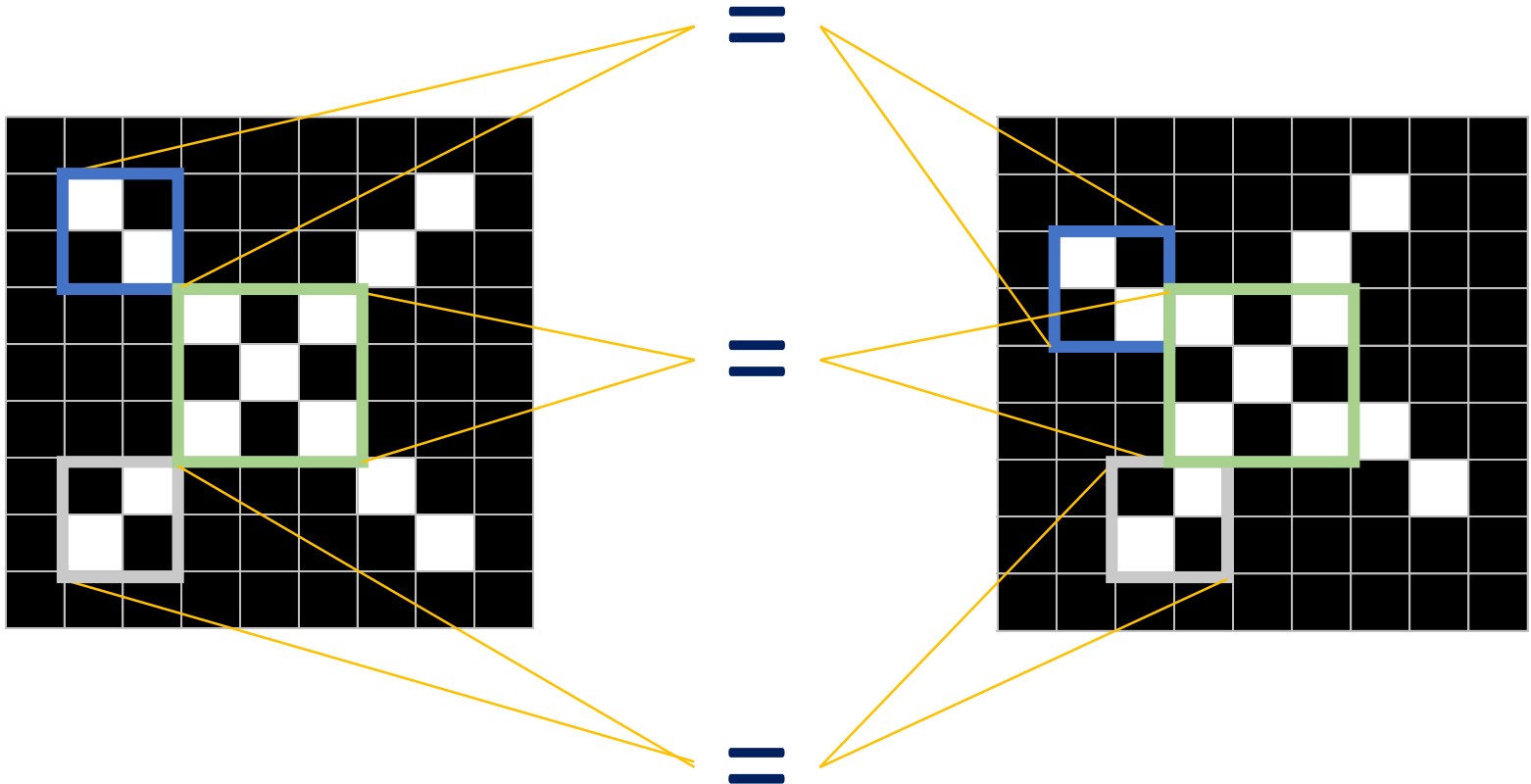


Computers are literal

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1

[illegible]

ConvNets match pieces of the image



Features match pieces of the image

1	-1	-1
-1	1	-1
-1	-1	1

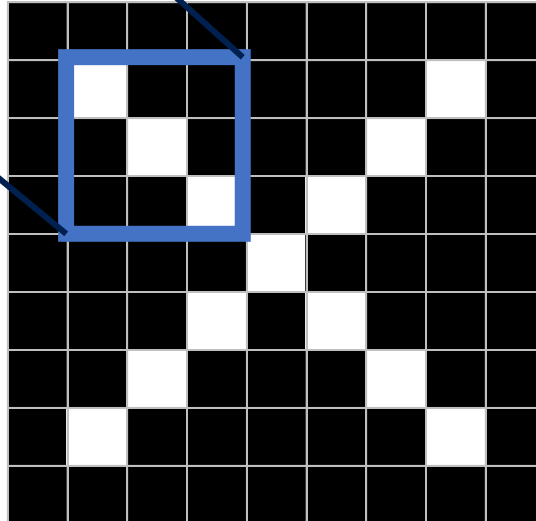
1	-1	1
-1	1	-1
1	-1	1

-1	-1	1
-1	1	-1
1	-1	-1

1	-1	-1
-1	1	-1
-1	-1	1

1	-1	1
-1	1	-1
1	-1	1

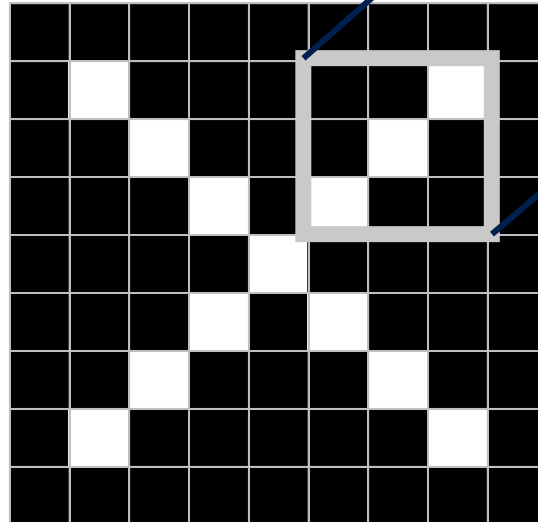
-1	-1	1
-1	1	-1
1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

1	-1	1
-1	1	-1
1	-1	1

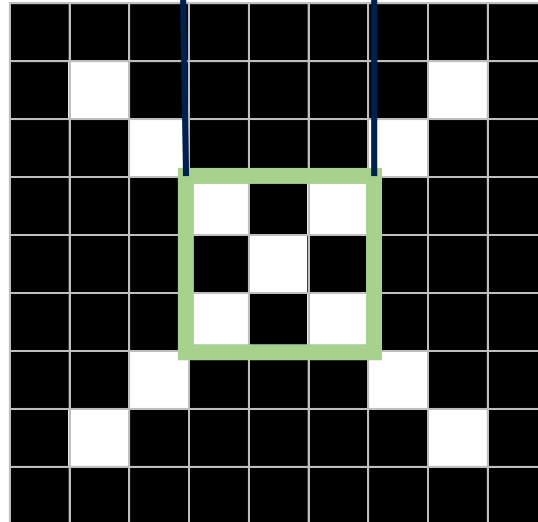
-1	-1	1
-1	1	-1
1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

1	-1	1
-1	1	-1
1	-1	1

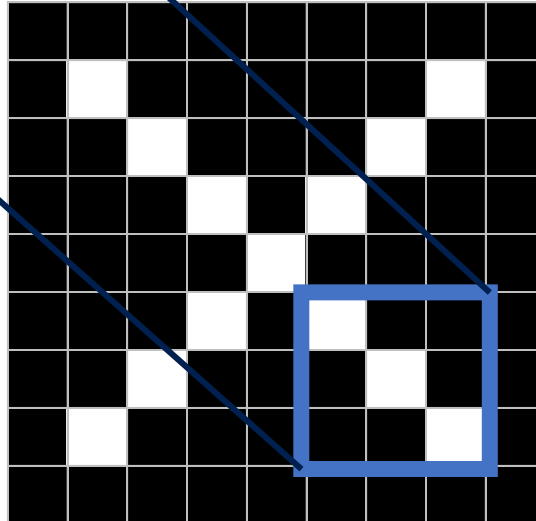
-1	-1	1
-1	1	-1
1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

1	-1	1
-1	1	-1
1	-1	1

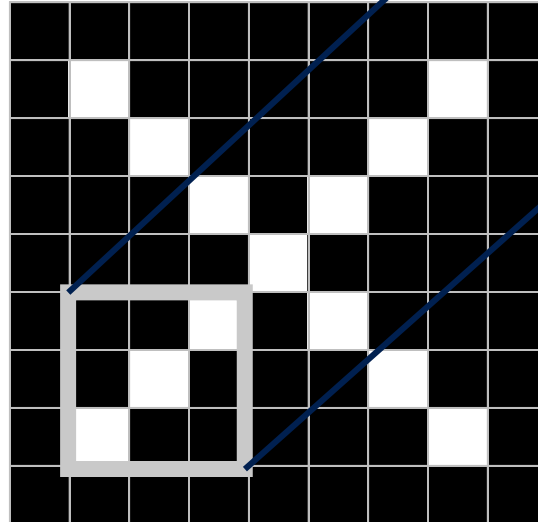
-1	-1	1
-1	1	-1
1	-1	-1



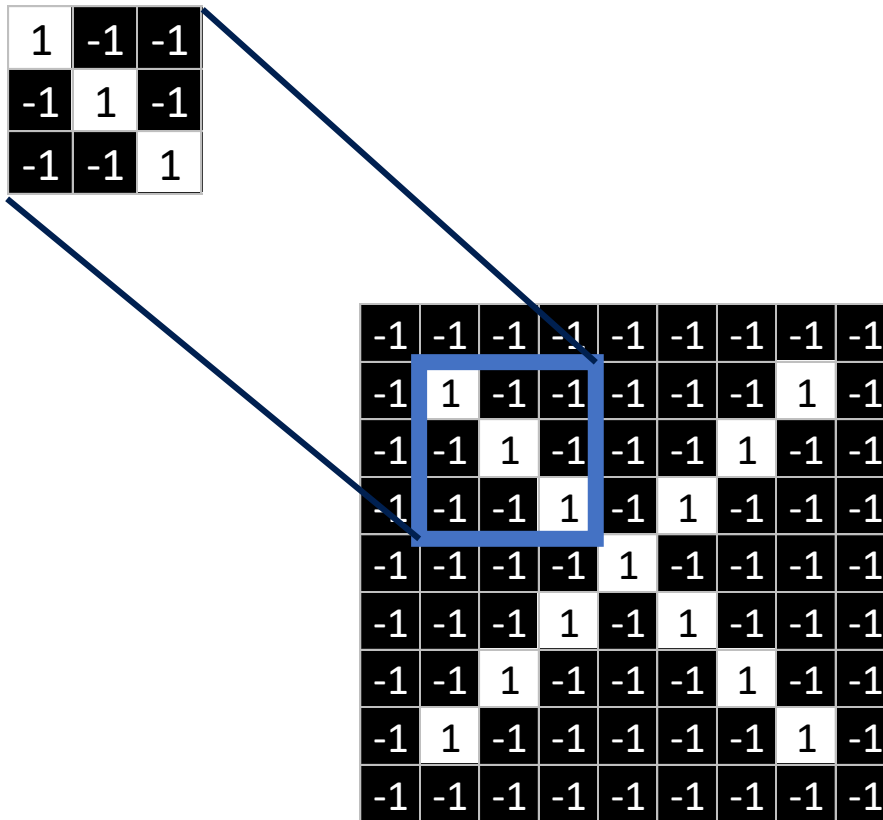
1	-1	-1
-1	1	-1
-1	-1	1

1	-1	1
-1	1	-1
1	-1	1

-1	-1	1
-1	1	-1
1	-1	-1



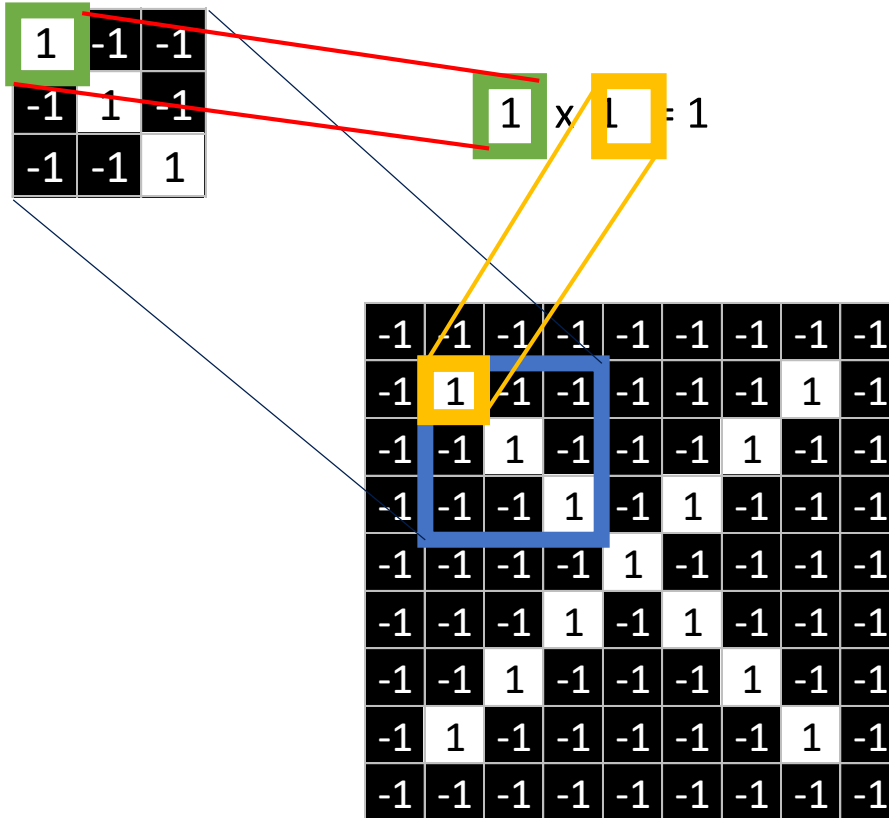
Filtering: The math behind the match



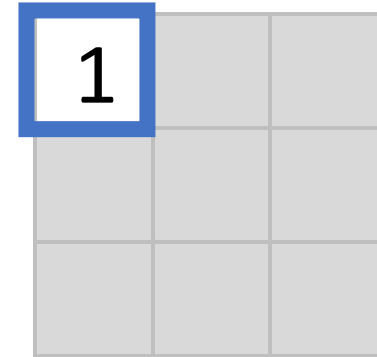
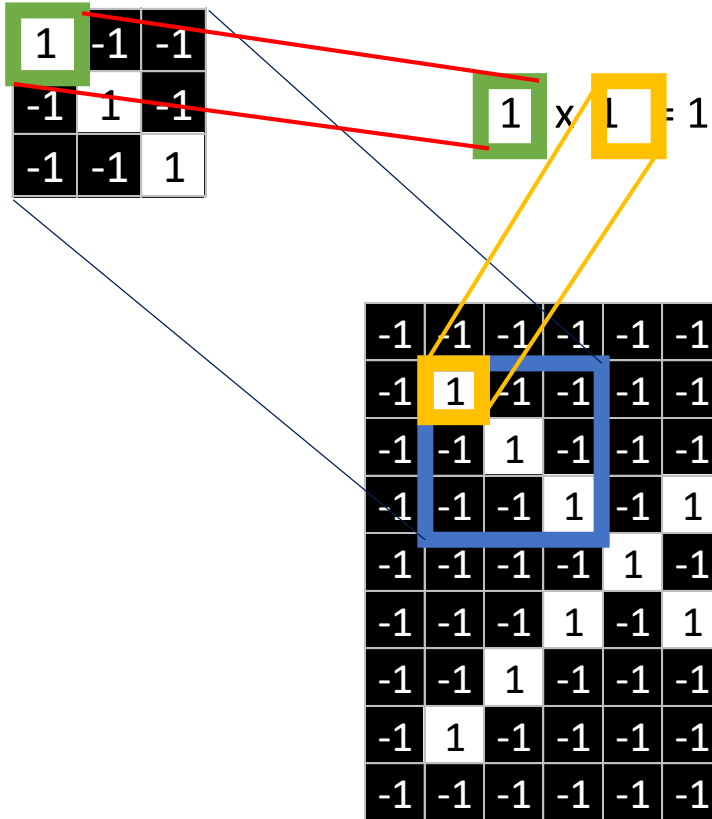
Filtering: The math behind the match

1. Line up the feature and the image patch.
2. Multiply each image pixel by the corresponding feature pixel.
3. Add them up.
4. Divide by the total number of pixels in the feature.

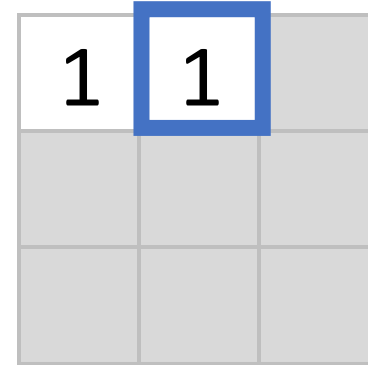
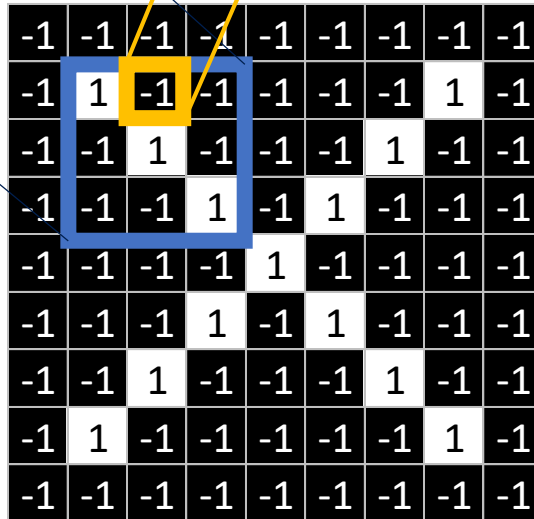
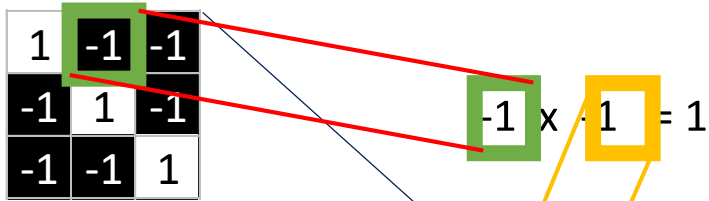
Filtering: The math behind the match



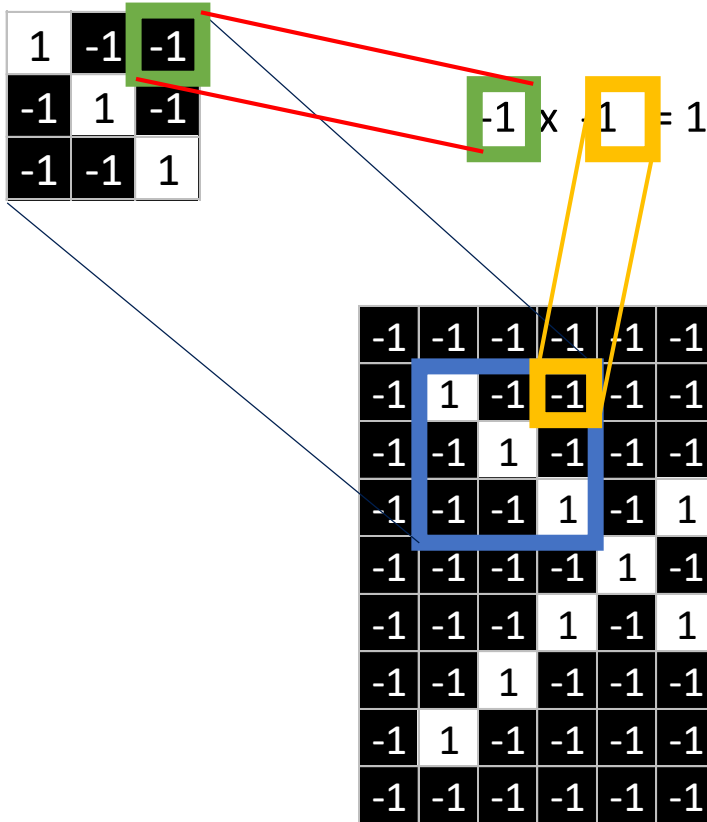
Filtering: The math behind the match



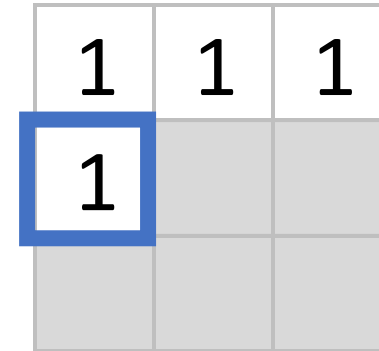
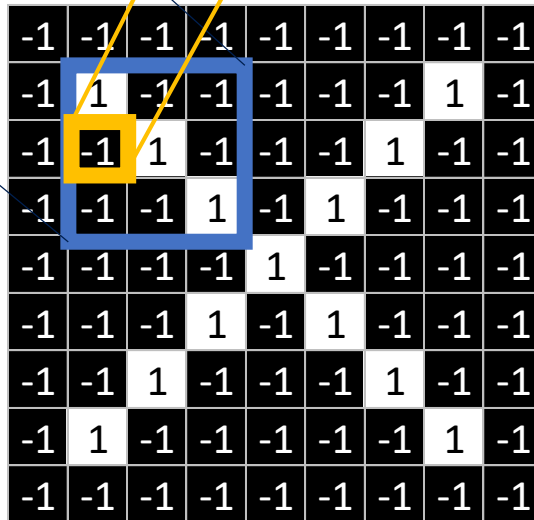
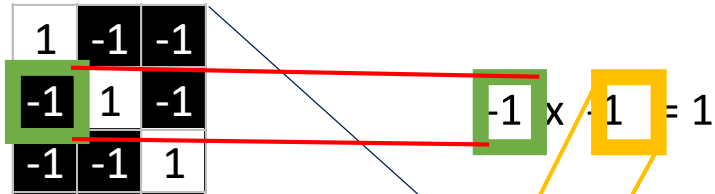
Filtering: The math behind the match



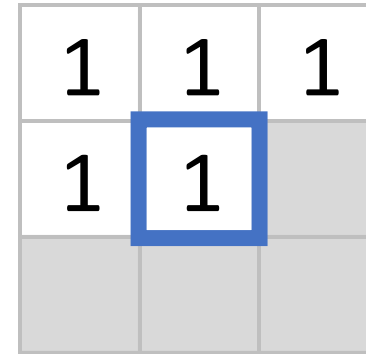
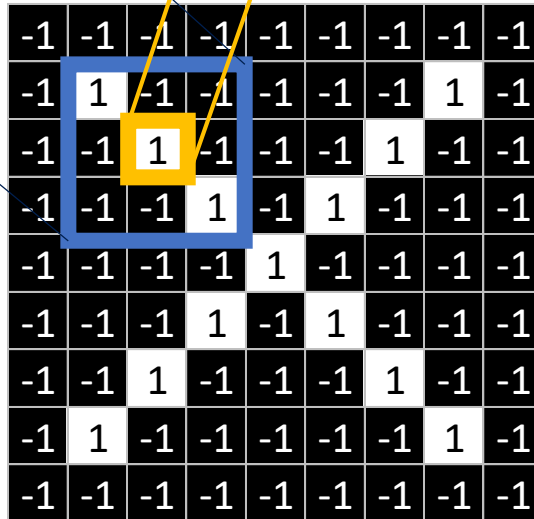
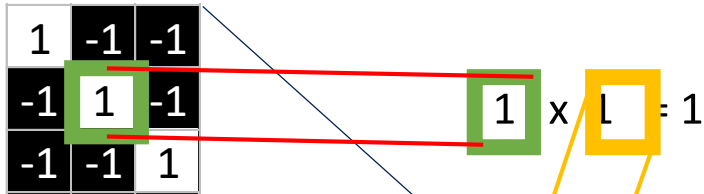
Filtering: The math behind the match



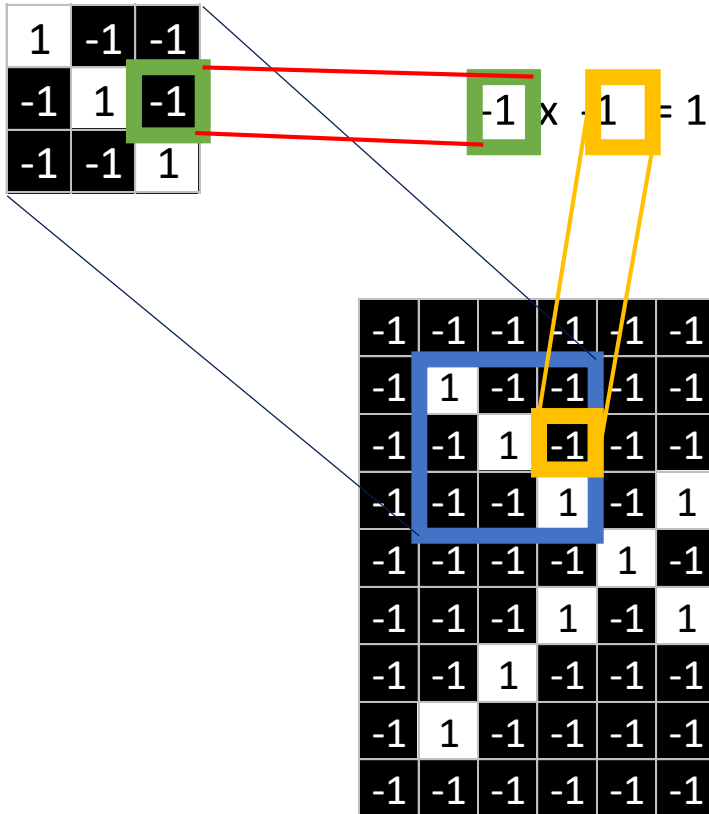
Filtering: The math behind the match



Filtering: The math behind the match

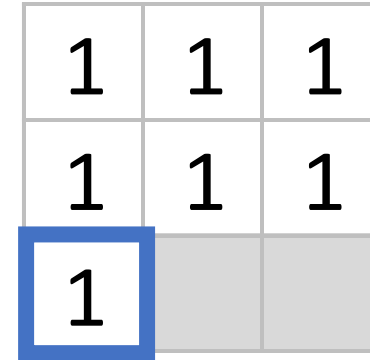
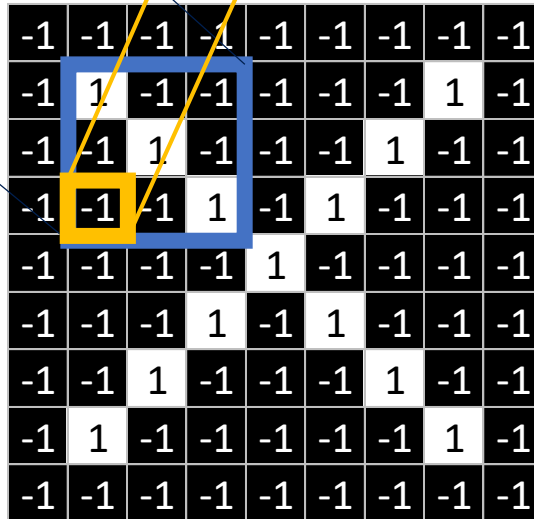
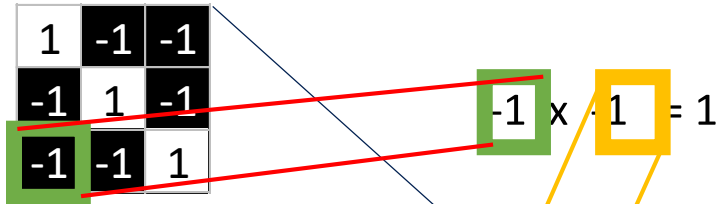


Filtering: The math behind the match

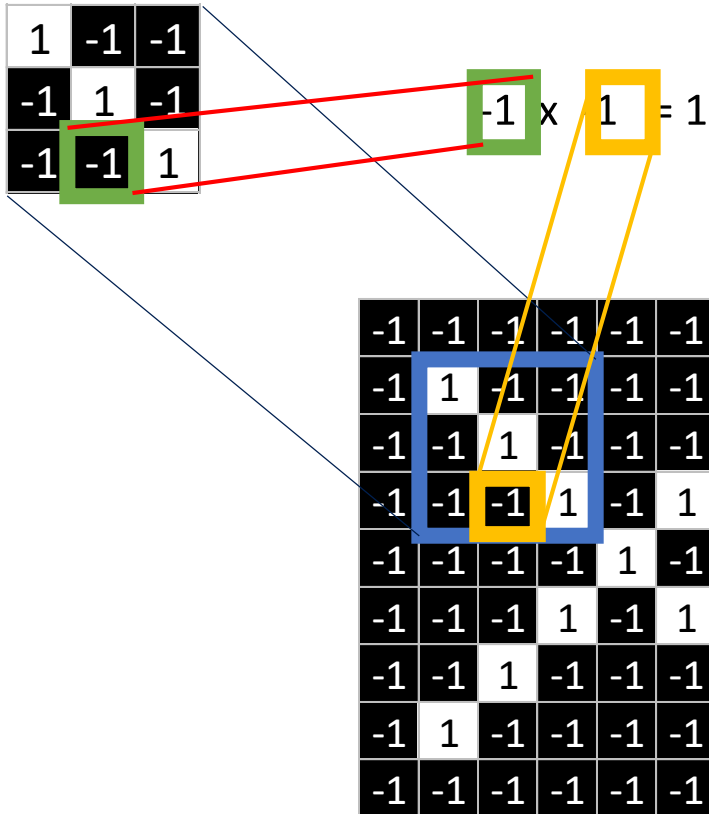


1	1	1
1	1	1

Filtering: The math behind the match

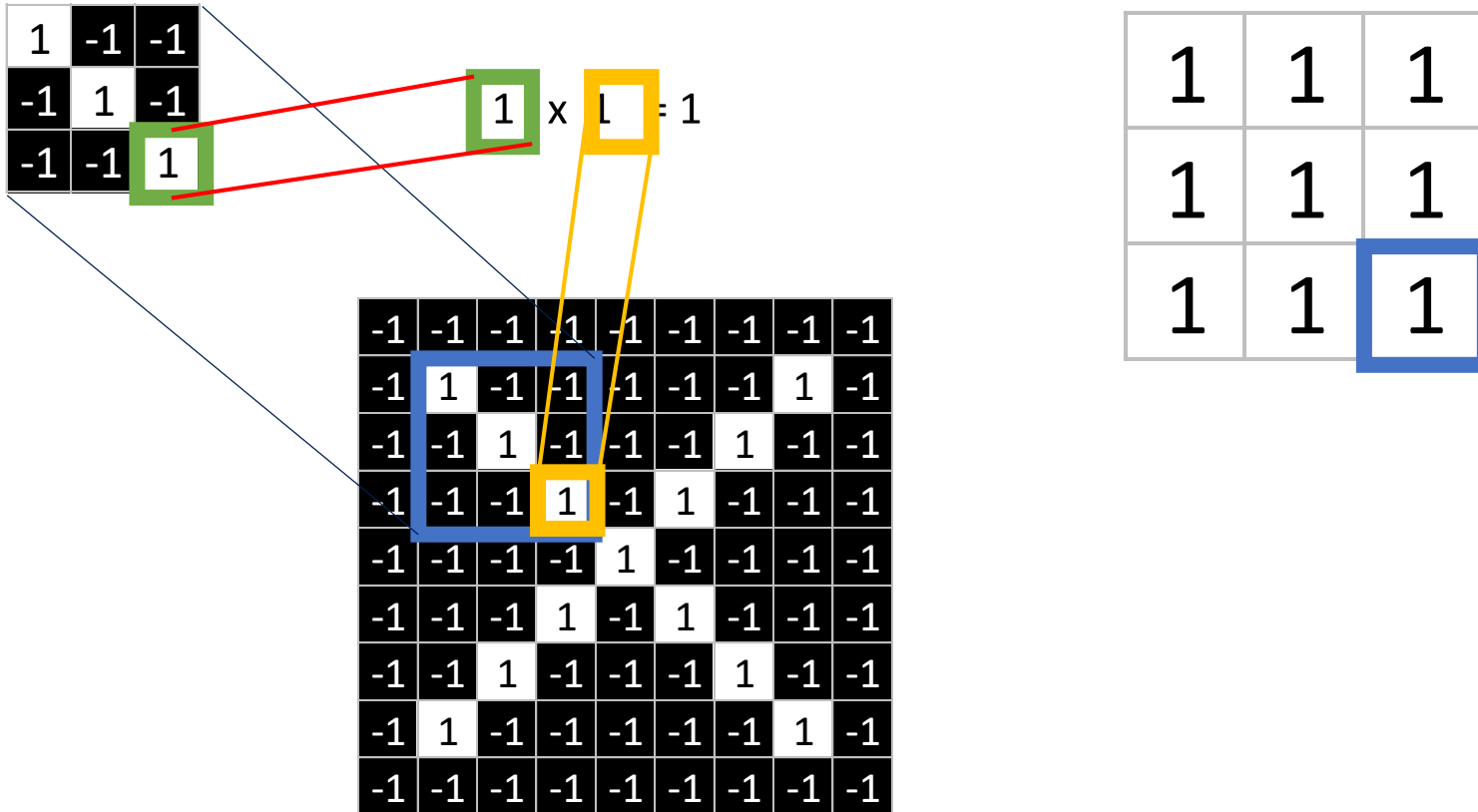


Filtering: The math behind the match

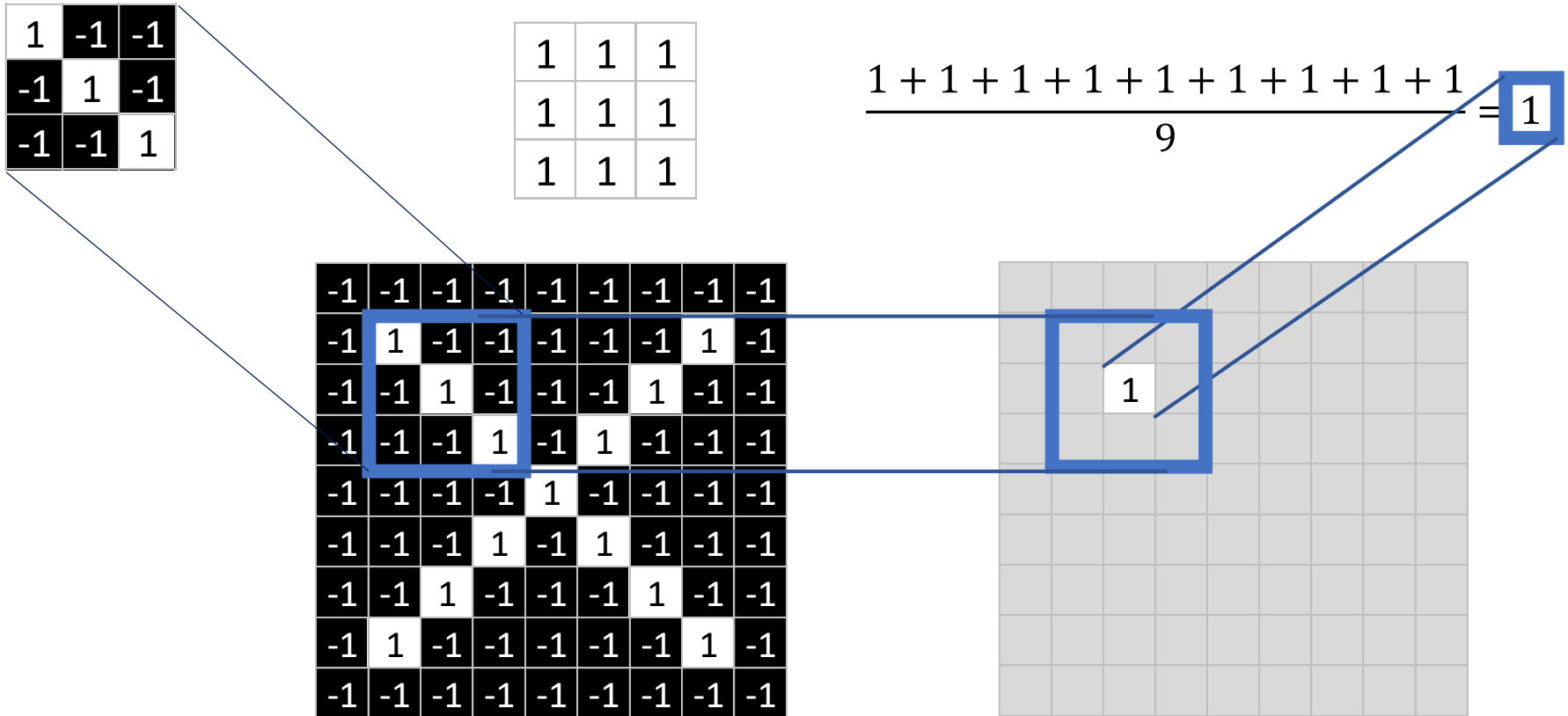


1	1	1
1	1	1
1	1	

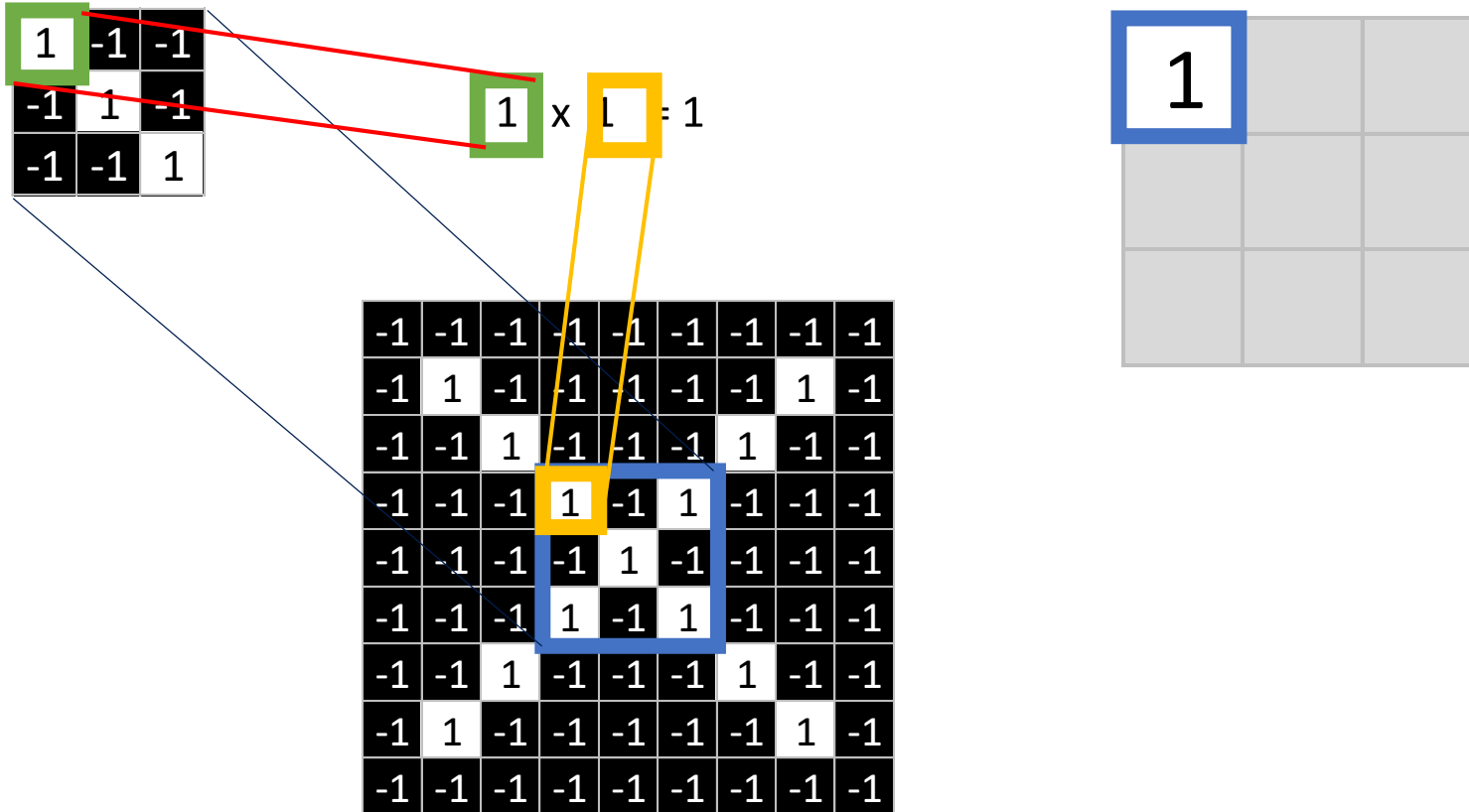
Filtering: The math behind the match



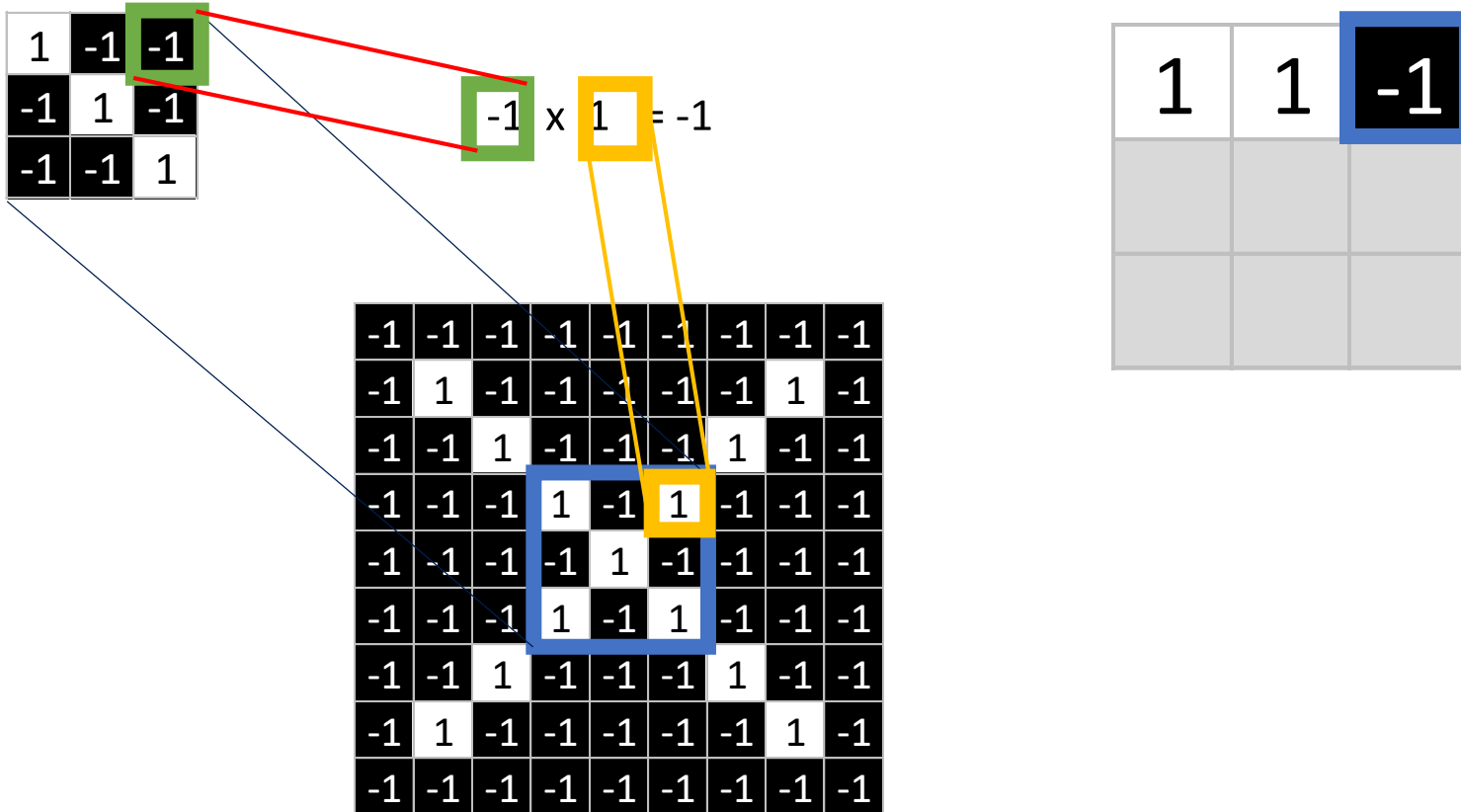
Filtering: The math behind the match



Filtering: The math behind the match



Filtering: The math behind the match



Filtering: The math behind the match

1	-1	-1
-1	1	-1
-1	-1	1

1	1	-1
1	1	1
-1	1	1

$$\frac{1+1-1+1+1+1-1+1+1}{9} = .55$$

[illegible]

Convolution: Trying every possible match

1	-1	-1
-1	1	-1
-1	-1	1

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1



0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

Convolution: Trying every possible match

-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1



1	-1	-1
-1	1	-1
-1	-1	1

=

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1



1	-1	1
-1	1	-1
1	-1	1

=

0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	-0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33

-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	1	-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1	-1	-1	-1
-1	-1	-1	1	-1	1	-1	-1	-1	-1
-1	-1	1	-1	-1	-1	1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	-1	-1	-1	-1	-1



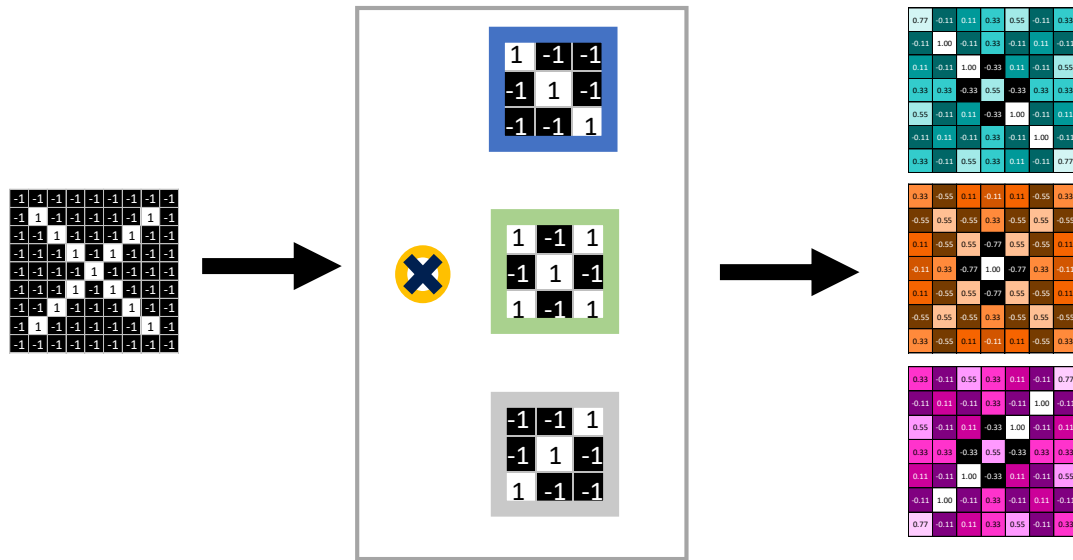
-1	-1	1
-1	1	-1
1	-1	-1

=

0.33	-0.11	0.55	0.33	0.11	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.77	-0.11	0.11	0.33	0.55	-0.11	0.33

Convolution layer

One image becomes a stack of filtered images



Convolution layer

One image becomes a stack of filtered images

-1	-1	-1	-1	-1	-1	-1
-1	1	-1	-1	-1	-1	-1
-1	-1	1	-1	-1	-1	-1
-1	-1	-1	1	-1	-1	-1
-1	-1	-1	-1	1	-1	-1
-1	-1	1	-1	-1	1	-1
-1	-1	-1	-1	-1	1	-1



0.77	-0.11	0.33	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	0.33	0.11	-0.11	0.55
0.33	0.33	-0.11	0.55	-0.11	0.33	0.33
0.55	-0.11	0.11	0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

0.33	0.55	0.11	-0.11	0.11	-0.55	0.33
0.55	0.33	0.55	0.33	-0.11	0.55	0.33
0.11	0.55	0.55	-0.77	0.55	0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	0.55	0.33	-0.55	0.55	0.55
0.33	-0.55	0.11	0.11	0.11	-0.55	0.33

0.33	-0.11	0.55	0.33	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00
0.55	-0.11	0.11	-0.22	1.00	-0.11
0.33	0.33	-0.33	0.55	-0.33	0.33
0.11	-0.11	1.00	0.33	0.11	-0.11
-0.11	1.00	-0.11	0.33	-0.11	0.11
0.77	-0.11	0.11	0.33	0.55	0.11

Convolution in Image Processing

Convolutions (typically with *prespecified* filters) are a common operation in many computer vision applications



Original image z



Gaussian blur



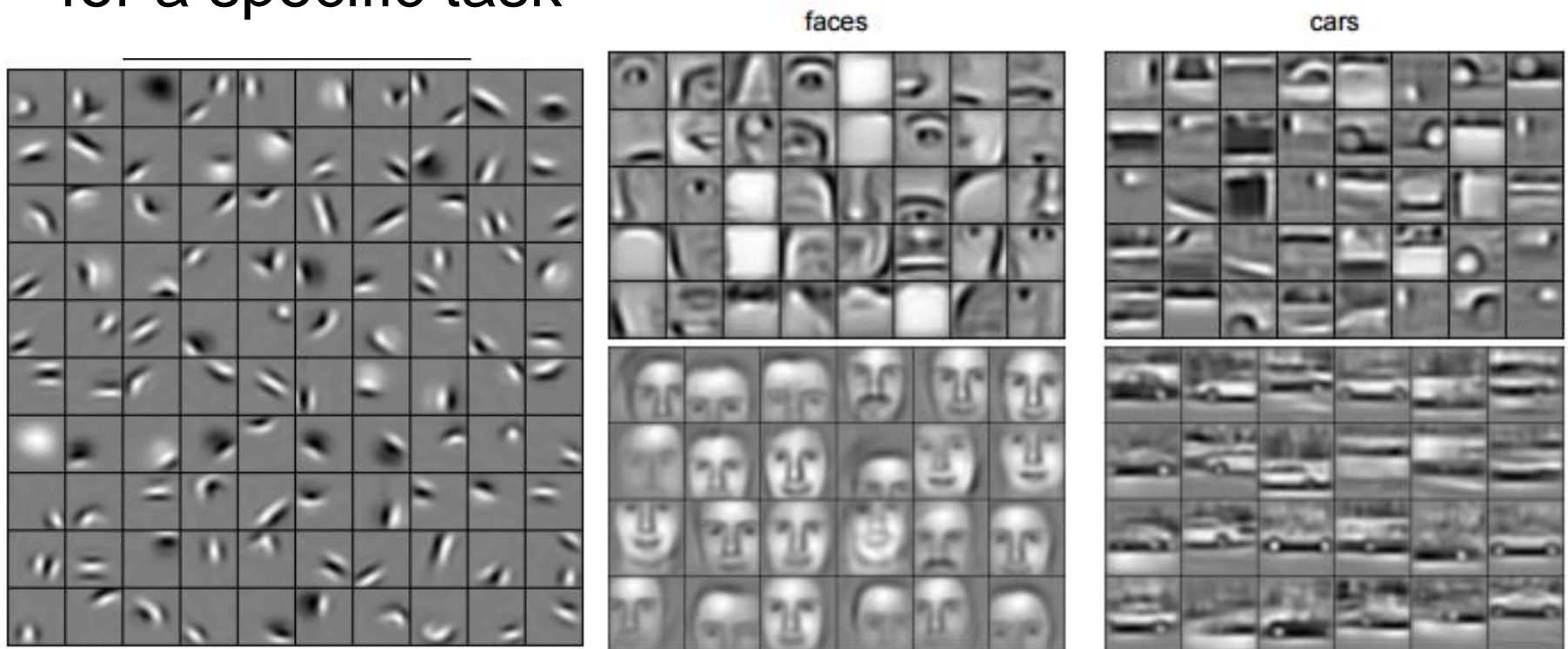
Image gradient

$$z * \begin{bmatrix} 1 & 4 & 7 & 4 & 1 \\ 4 & 16 & 26 & 16 & 4 \\ 7 & 26 & 41 & 26 & 7 \\ 4 & 16 & 26 & 16 & 4 \\ 1 & 4 & 4 & 4 & 1 \end{bmatrix} / 273$$

$$\left(\left(z * \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \right)^2 + \left(z * \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} \right)^2 \right)^{\frac{1}{2}}$$

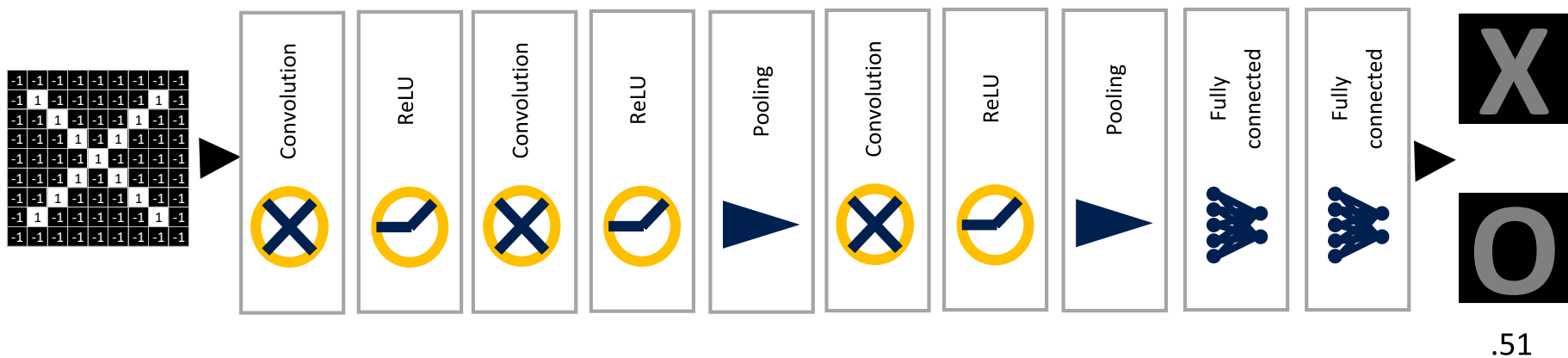
Learning CNNs

- Idea of a convolutional neural network, in some sense, is to let the network “learn” the right filters for a specific task



Convolutional Neural Networks (CNNs)

- Containing different types of layers
 - Convolution
 - **Non-linearity**
 - Pooling (or downsampling)
 - Fully connected layer

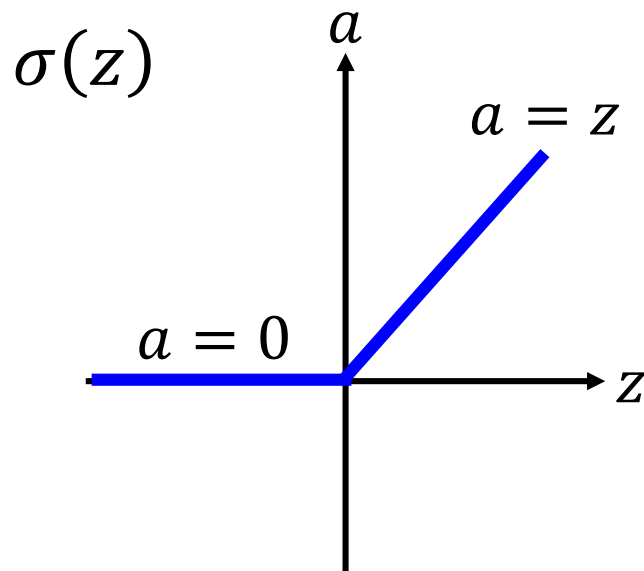


Non-linearity Layer

- Convolution is a linear operation
- Non-linearity layer creates an activation map from the feature map generated by the convolutional layer
- Consisting an activation function (an element-wise operation)
- Rectified linear units (ReLus) is advantageous over the traditional sigmoid or tanh activation functions

A Common Activation Function in CNNs

- Rectified Linear Unit (ReLU)

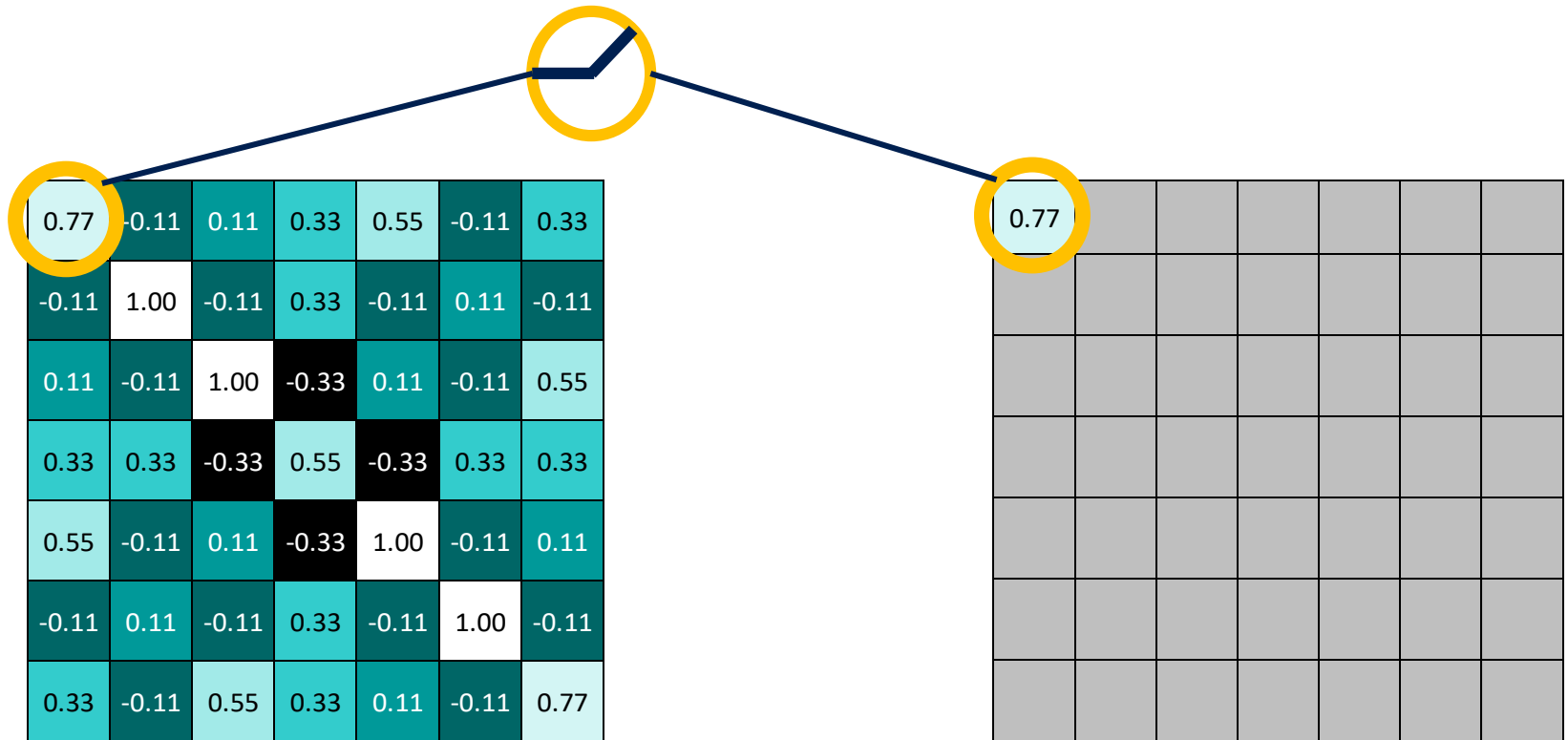


[Xavier Glorot, AISTATS'11]
[Andrew L. Maas, ICML'13]
[Kaiming He, arXiv'15]

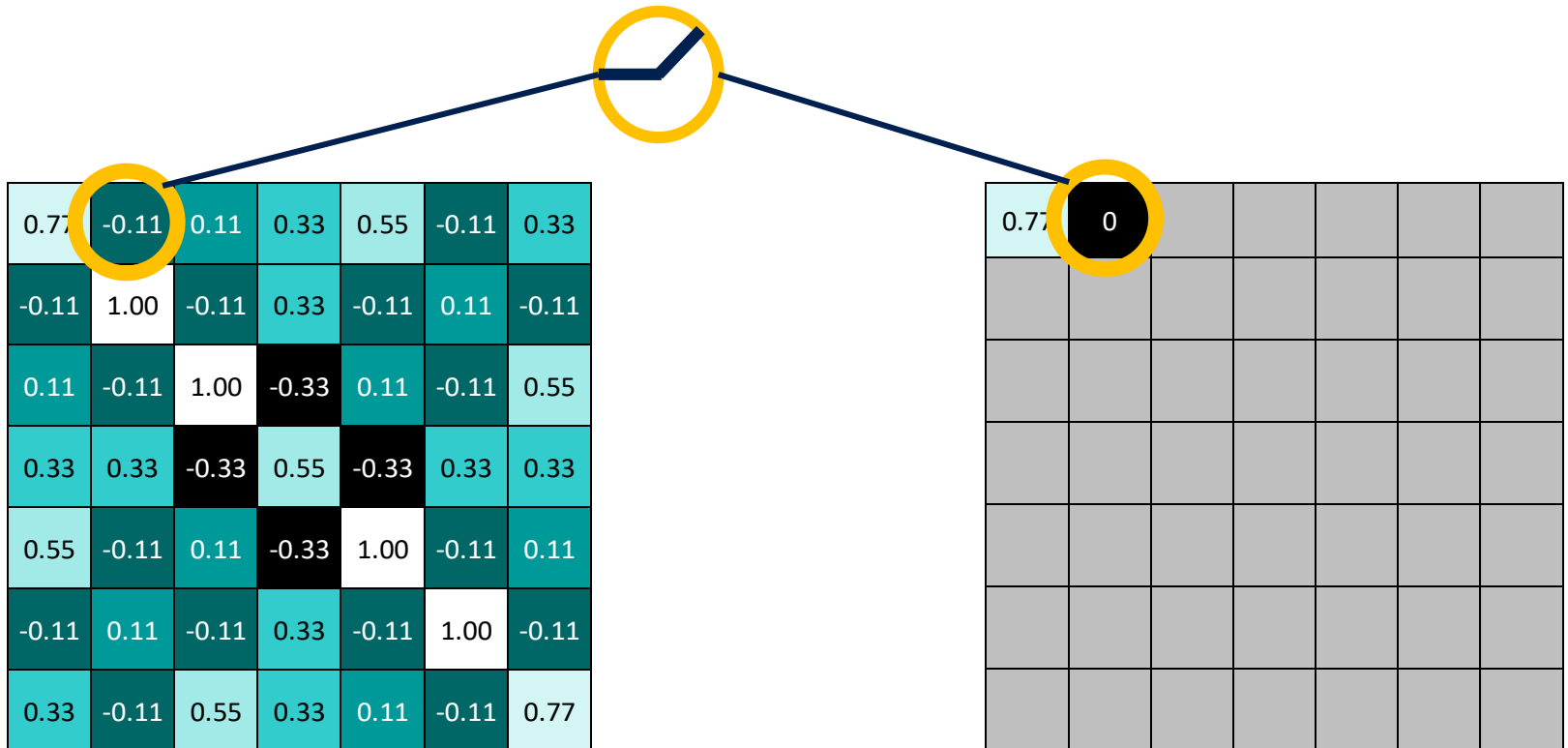
Reason:

1. Fast to compute
2. Cancellation problem
3. More sparse activation volume
4. Vanishing gradient problem

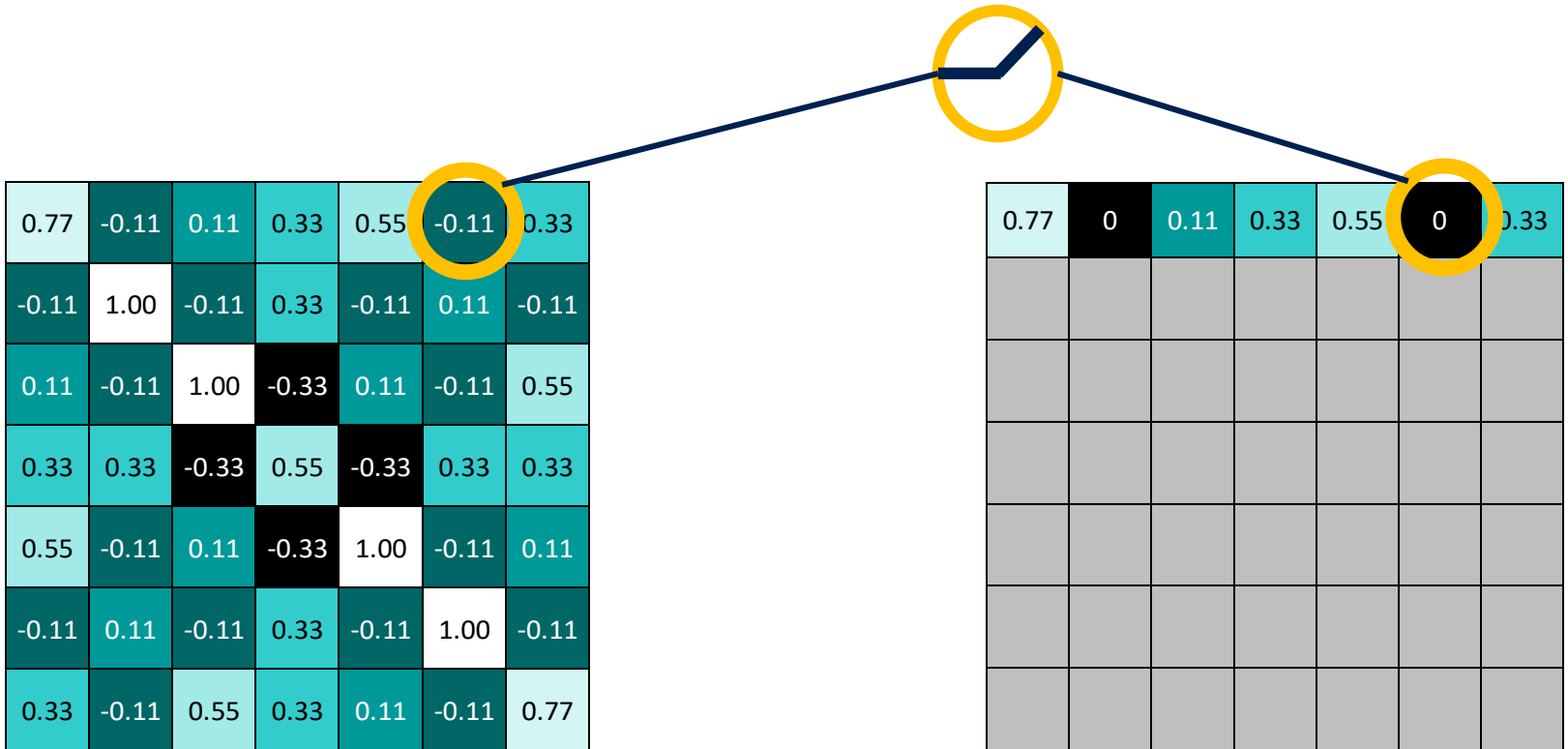
Rectified Linear Units (ReLUs)



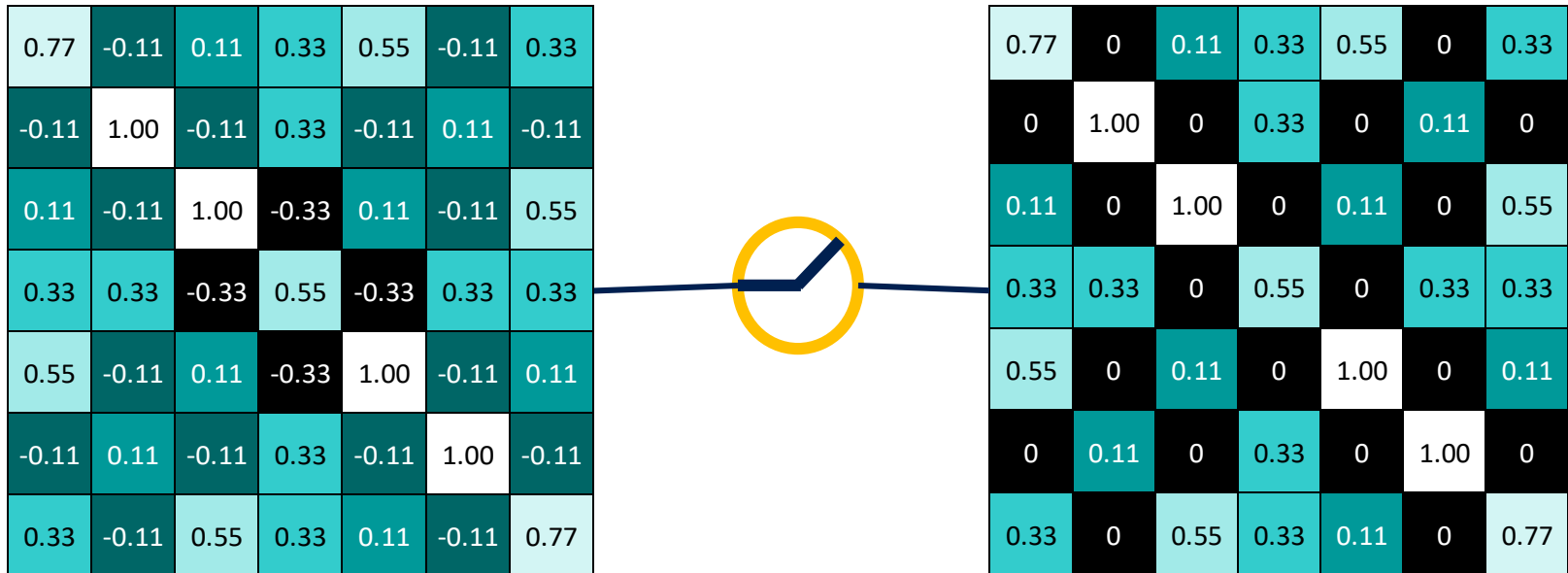
Rectified Linear Units (ReLUs)



Rectified Linear Units (ReLUs)



Rectified Linear Units (ReLUs)



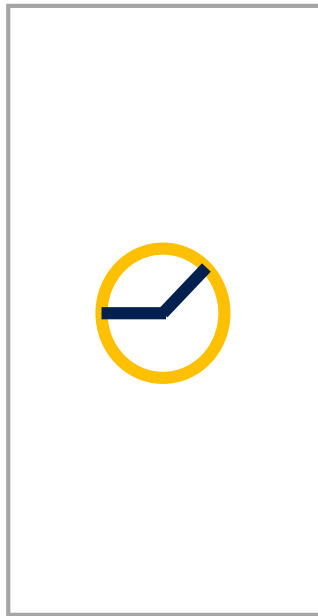
ReLU layer

A stack of images becomes a stack of images with no negative values.

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	-0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33

0.33	-0.11	0.55	0.33	0.11	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.77	-0.11	0.11	0.33	0.55	-0.11	0.33



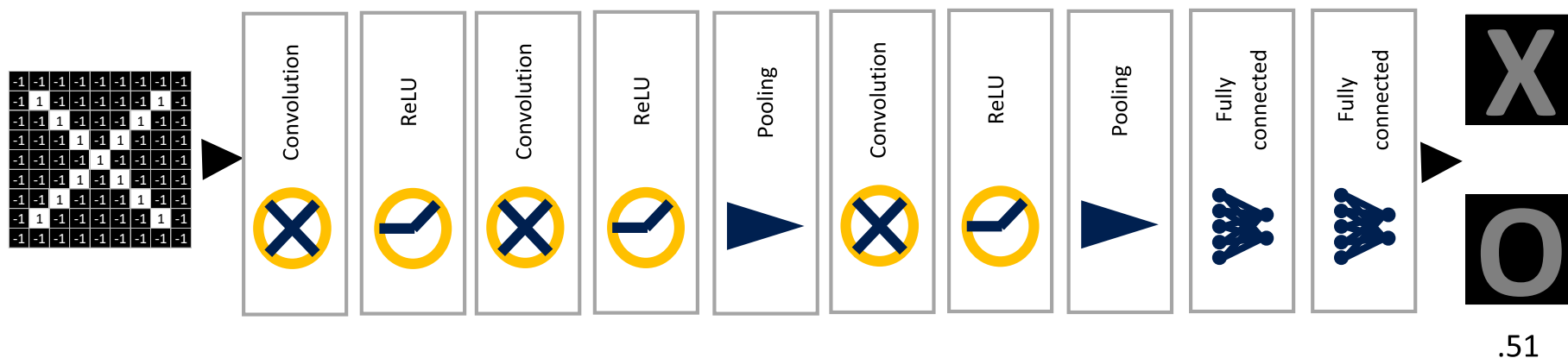
0.77	0	0.11	0.33	0.55	0	0.33
0	1.00	0	0.33	0	-0.11	0
0.11	0	1.00	0	0.11	0	0.55
0.33	0.33	0	0.55	0	0.33	0.33
0.55	0	0.11	0	1.00	0	0.11
0	0.11	0	0.33	0	1.00	0
0.33	0	0.55	0.33	0.11	0	0.77

0.33	0	0.11	0	0.11	0	0.33
0	0.55	0	0.33	0	0.55	0
0.11	0	0.55	0	0.55	0	0.11
0	0.33	0	1.00	0	0.33	0
0.11	0	0.55	0	0.55	0	0.11
0	0.55	0	0.33	0	0.55	0
0.33	0	0.11	0	0.11	0	0.33

0.33	0	0.55	0.33	0.11	0	0.77
0	0.11	0	0.33	0	1.00	0
0.55	0	0.11	0	1.00	0	0.11
0.33	0.33	0	0.55	0	0.33	0.33
0.11	0	1.00	0	0.11	0	0.55
0	1.00	0	0.33	0	0.11	0
0.77	0	0.11	0.33	0.55	0	0.33

Convolutional Neural Networks (CNNs)

- Containing different types of layers
 - Convolution
 - Non-linearity
 - **Pooling (or downsampling)**
 - Fully connected layer



Pooling: Shrinking the Image Stack

- Motivation: the activation maps can be large
- Reducing the spacial size of the activation maps
 - Often after multiple stages of other layers (i.e., convolutional and non-linear layers)
- Steps:
 1. Pick a window size (usually 2 or 3).
 2. Pick a stride (usually 2).
 3. Walk your window across your filtered images.
 4. From each window, take the maximum value.

Pooling

maximum

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

1.00			

Pooling

maximum

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

1.00	0.33		

Pooling

maximum

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

1.00	0.33	0.55	

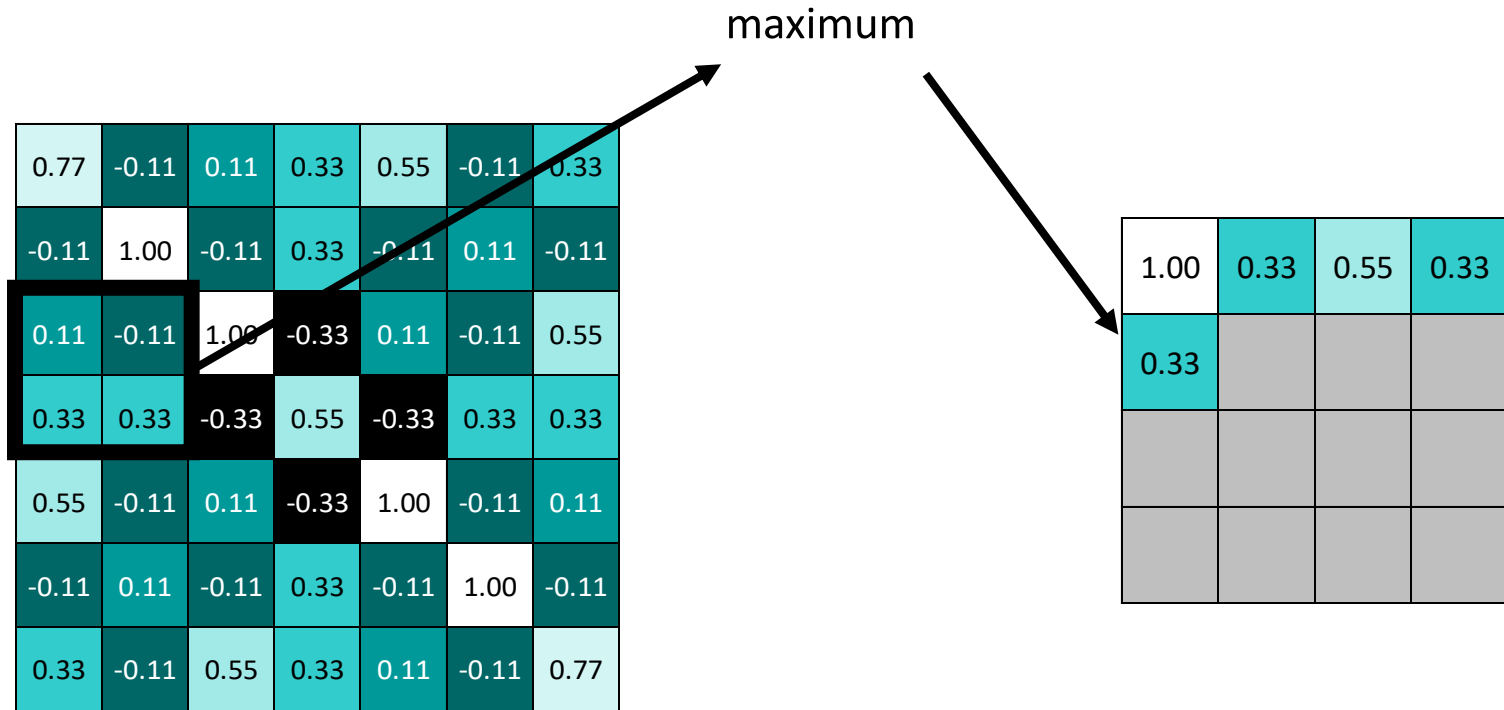
Pooling

maximum

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

1.00	0.33	0.55	0.33

Pooling



Pooling

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77

max pooling



1.00	0.33	0.55	0.33
0.33	1.00	0.33	0.55
0.55	0.33	1.00	0.11
0.33	0.55	0.11	0.77

0.77	-0.11	0.11	0.33	0.55	-0.11	0.33
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77



1.00	0.33	0.55	0.33
0.33	1.00	0.33	0.55
0.55	0.33	1.00	0.11
0.33	0.55	0.11	0.77

0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.11	0.33	-0.77	1.00	-0.77	0.33	-0.11
0.11	-0.55	0.55	-0.77	0.55	-0.55	0.11
-0.55	0.55	-0.55	0.33	-0.55	0.55	-0.55
0.33	-0.55	0.11	-0.11	0.11	-0.55	0.33



0.55	0.33	0.55	0.33
0.33	1.00	0.55	0.11
0.55	0.55	0.55	0.11
0.33	0.11	0.11	0.33

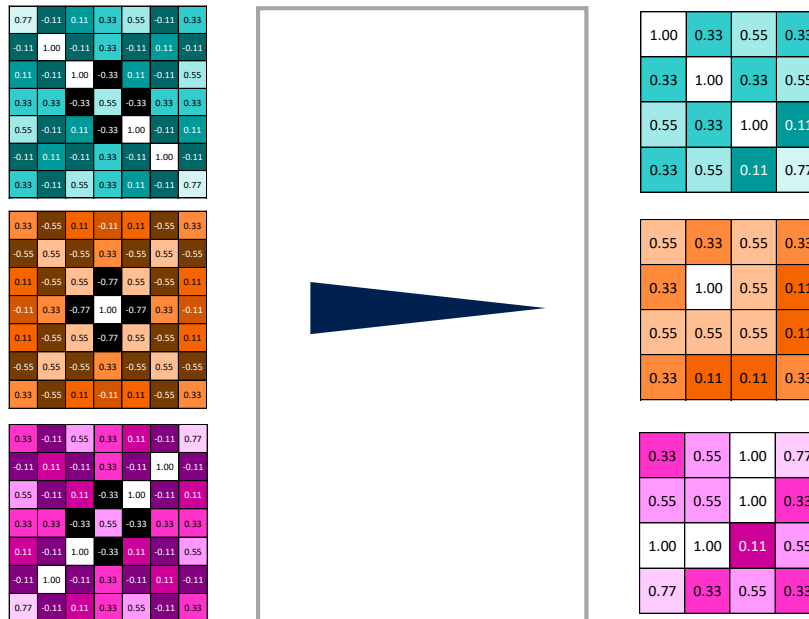
0.33	-0.11	0.55	0.33	0.11	-0.11	0.77
-0.11	0.11	-0.11	0.33	-0.11	1.00	-0.11
0.55	-0.11	0.11	-0.33	1.00	-0.11	0.11
0.33	0.33	-0.33	0.55	-0.33	0.33	0.33
0.11	-0.11	1.00	-0.33	0.11	-0.11	0.55
-0.11	1.00	-0.11	0.33	-0.11	0.11	-0.11
0.77	-0.11	0.11	0.33	0.55	-0.11	0.33



0.33	0.55	1.00	0.77
0.55	0.55	1.00	0.33
1.00	1.00	0.11	0.55
0.77	0.33	0.55	0.33

Pooling layer

A stack of images becomes a stack of smaller images.

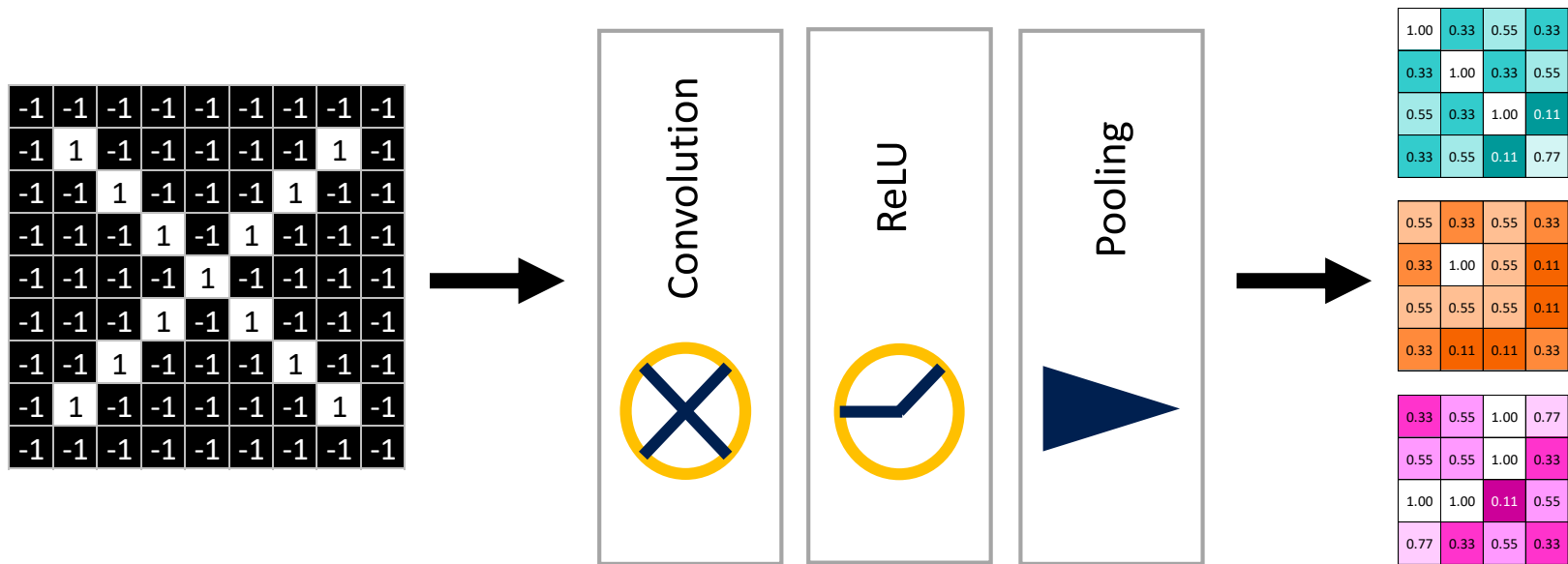


Pros and Cons of Pooling Layer

- Pros:
 - Reducing the computational requirements
 - Minimizing the likelihood of overfitting
- Cons:
 - Aggressive reduction can limit the depth of a network and ultimately limit the performance

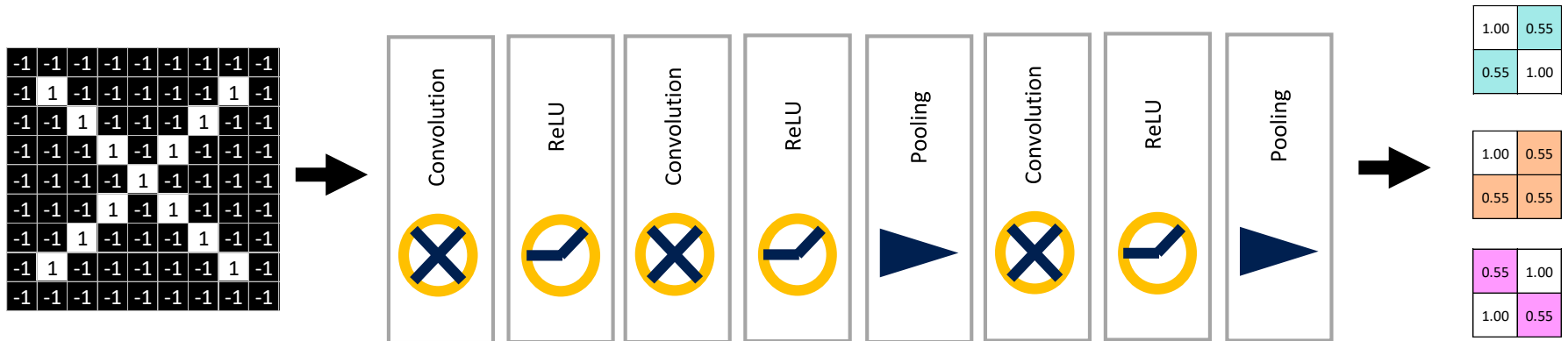
Layers get stacked

The output of one becomes the input of the next.



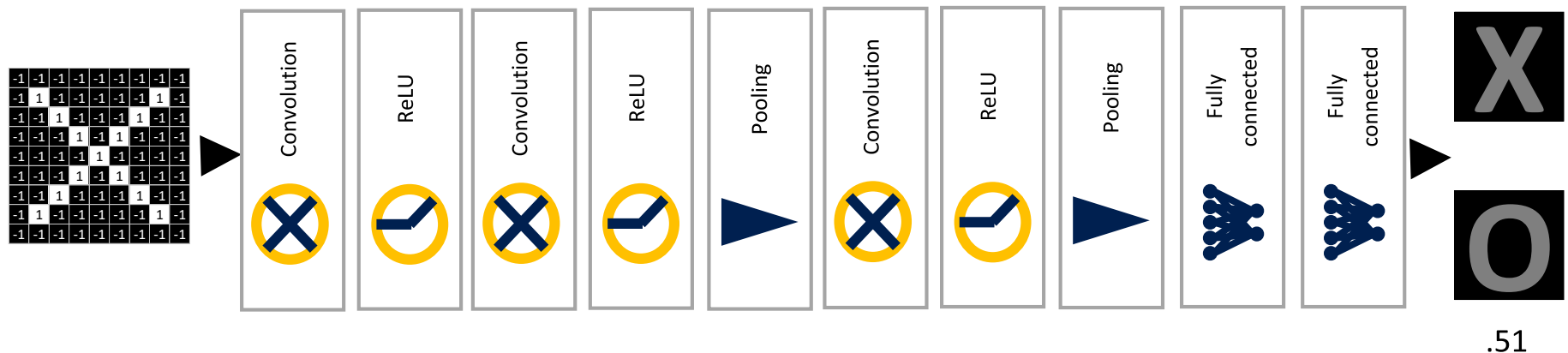
Deep stacking

Layers can be repeated several (or many) times.



Convolutional Neural Networks (CNNs)

- Containing different types of layers
 - Convolution
 - Non-linearity
 - Pooling (or downsampling)
 - **Fully connected layer**

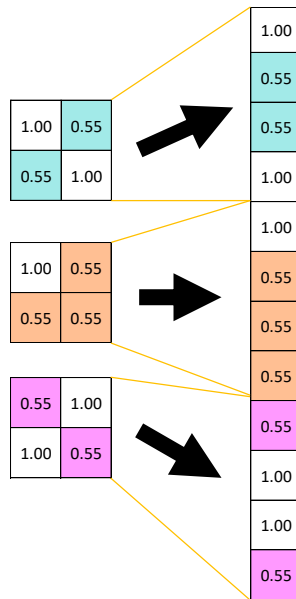


Fully Connected Layer

- Multilayer perceptron
- Mapping the activation volume from previous layers into a class probability distribution
- Non-linearity is built in the neurons, instead of a separate layer
- Viewed as 1×1 convolution kernels

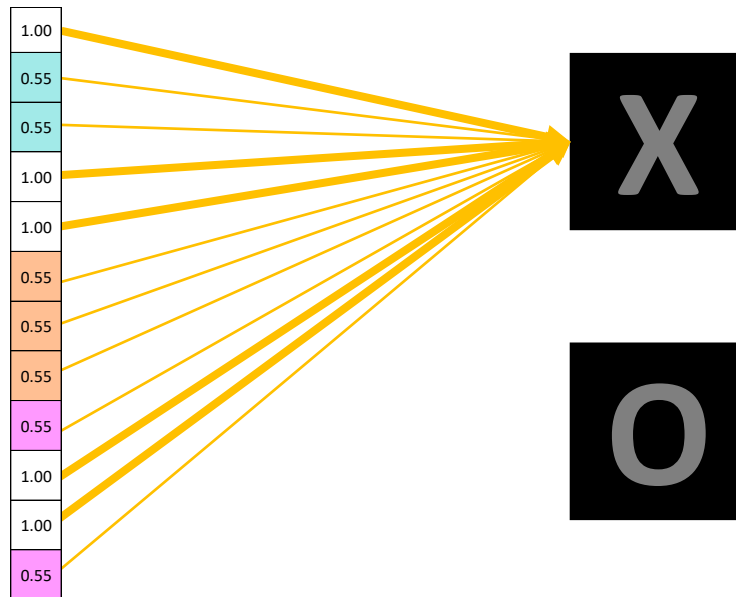
Fully connected layer

Every value gets a vote



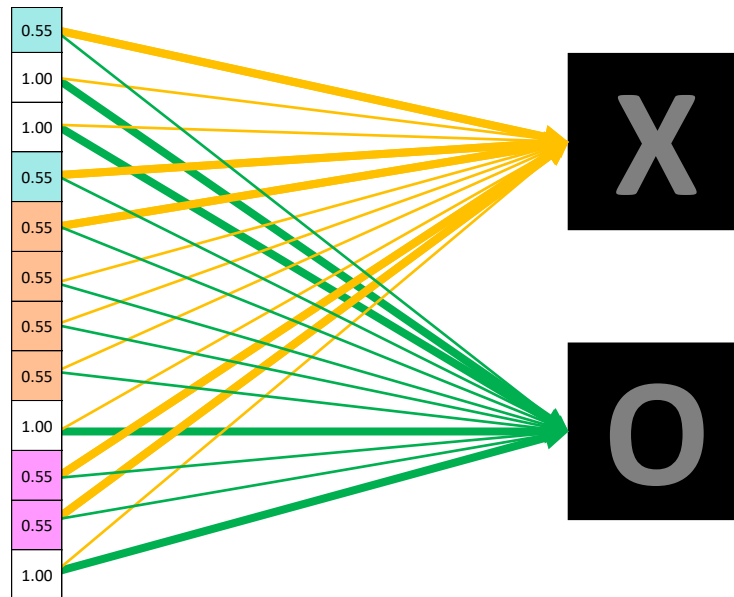
Fully connected layer

Vote depends on how strongly a value predicts X or O



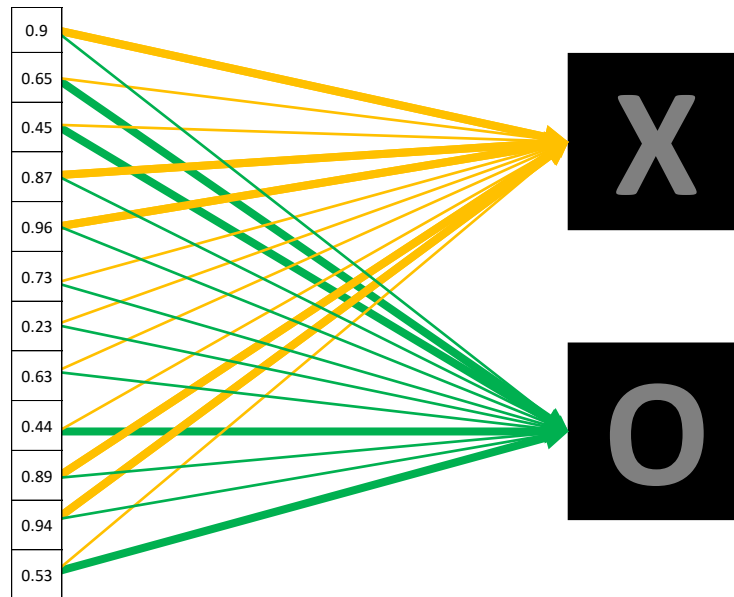
Fully connected layer

Vote depends on how strongly a value predicts X or O



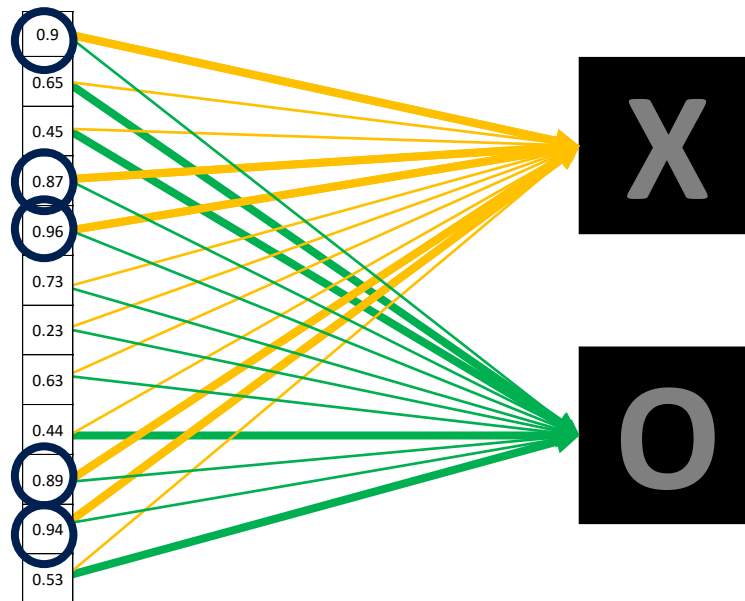
Fully connected layer

Future values vote on X or O



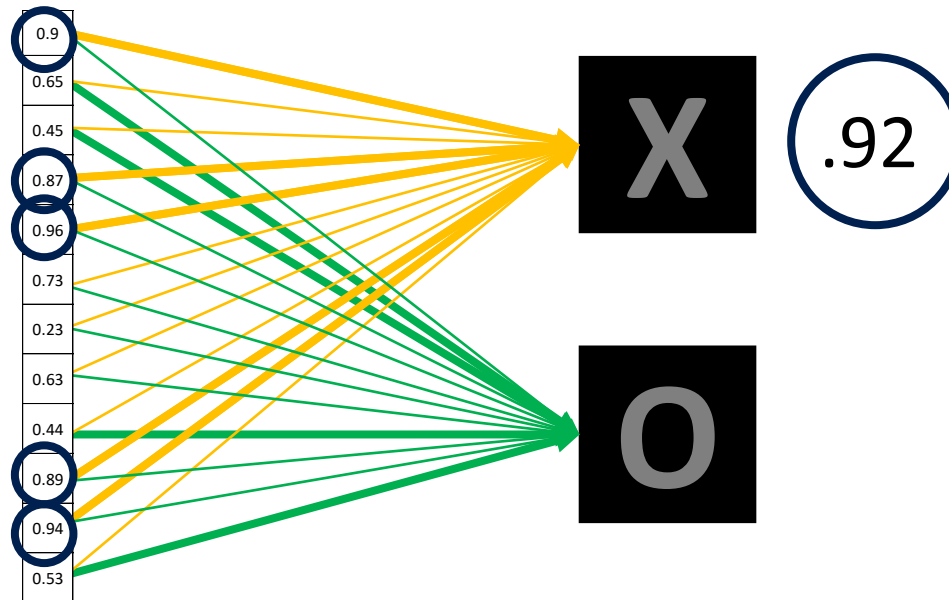
Fully connected layer

Future values vote on X or O



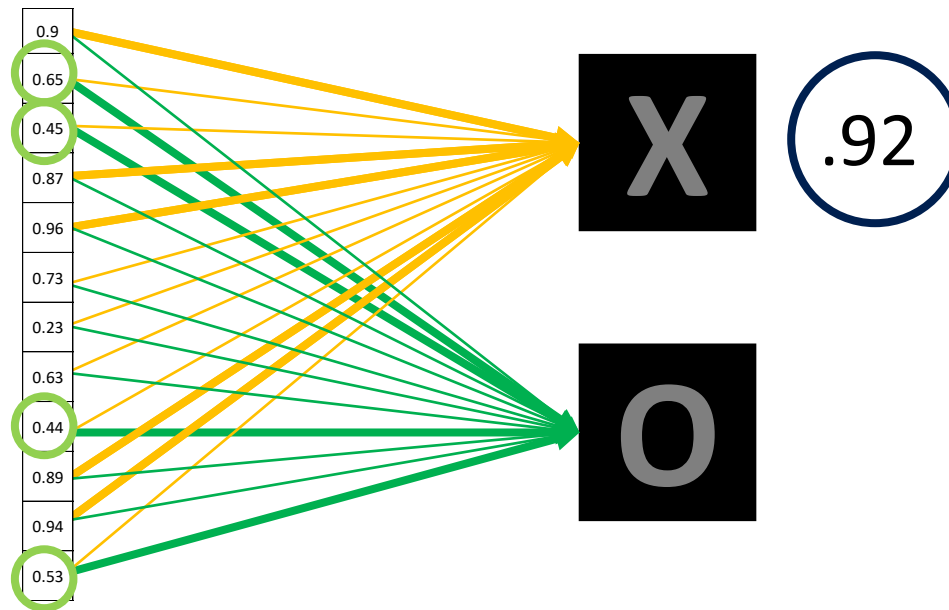
Fully connected layer

Future values vote on X or O



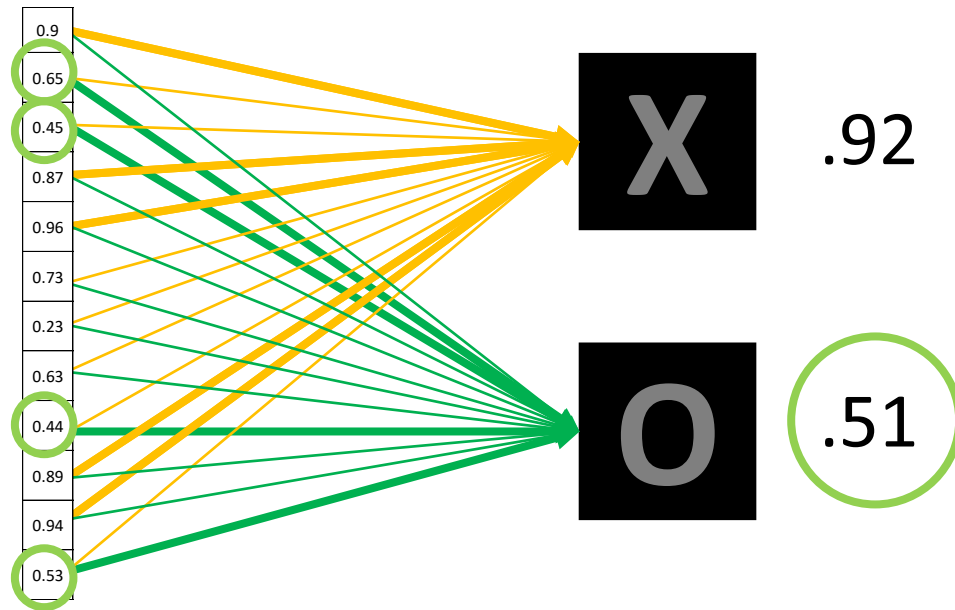
Fully connected layer

Future values vote on X or O



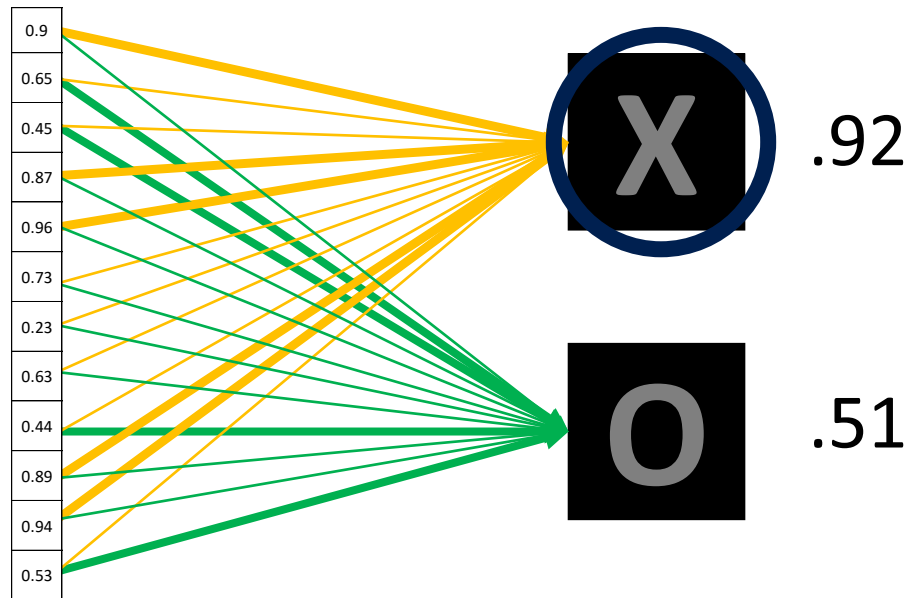
Fully connected layer

Future values vote on X or O



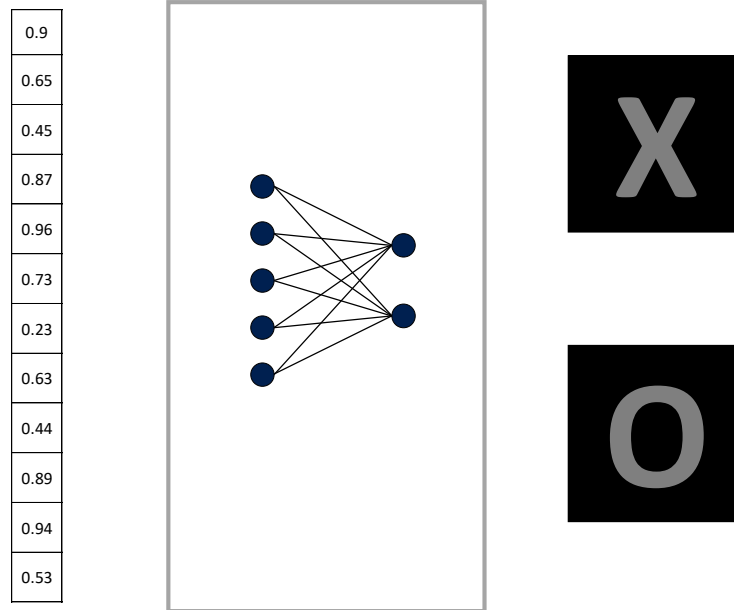
Fully connected layer

Future values vote on X or O



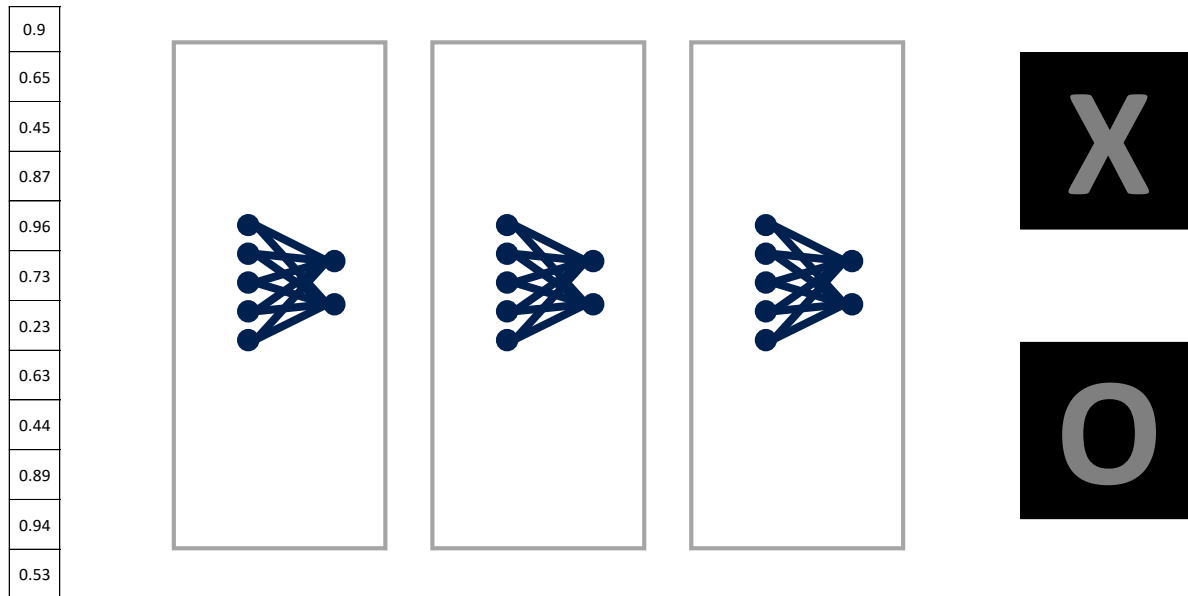
Fully connected layer

A list of feature values becomes a list of votes.



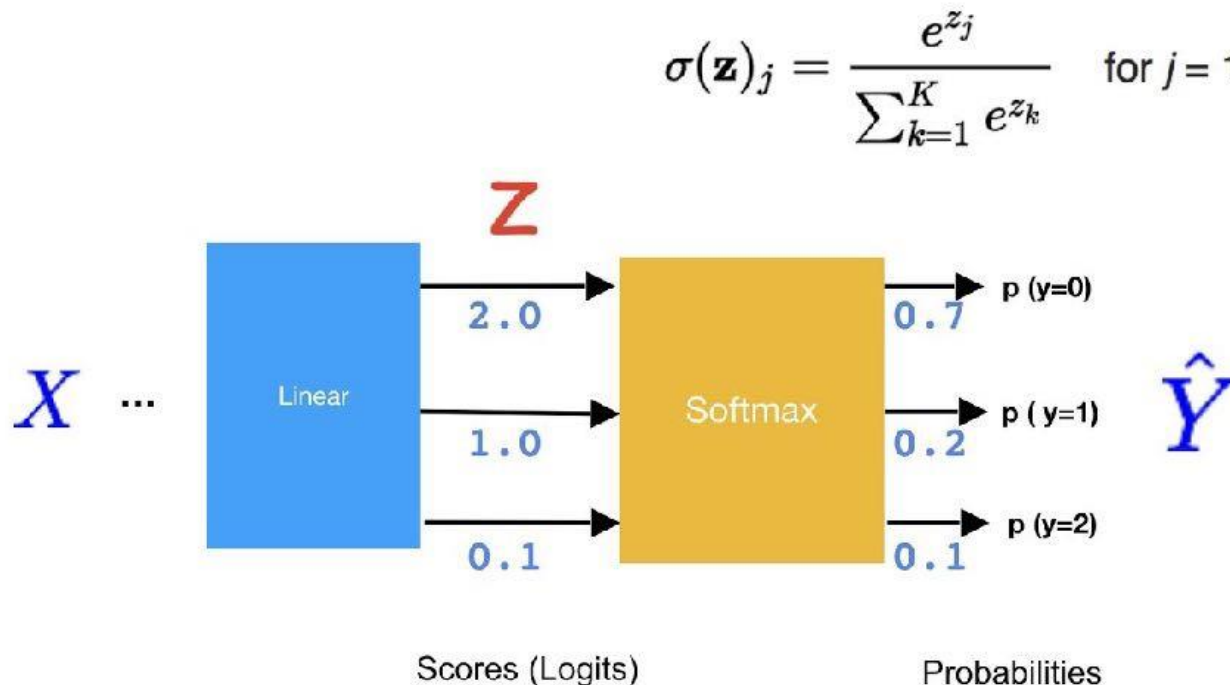
Fully connected layer

These can also be stacked.



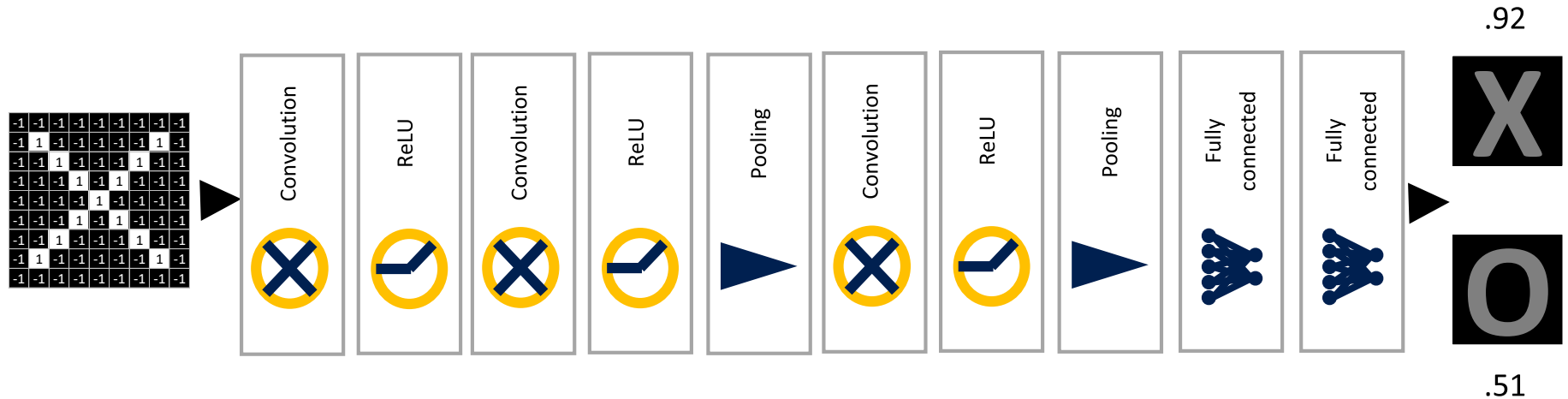
Softmax

- For classification: Output layer is a regular, fully connected layer with softmax non-linearity
 - Output provides an estimate of the conditional probability of each class



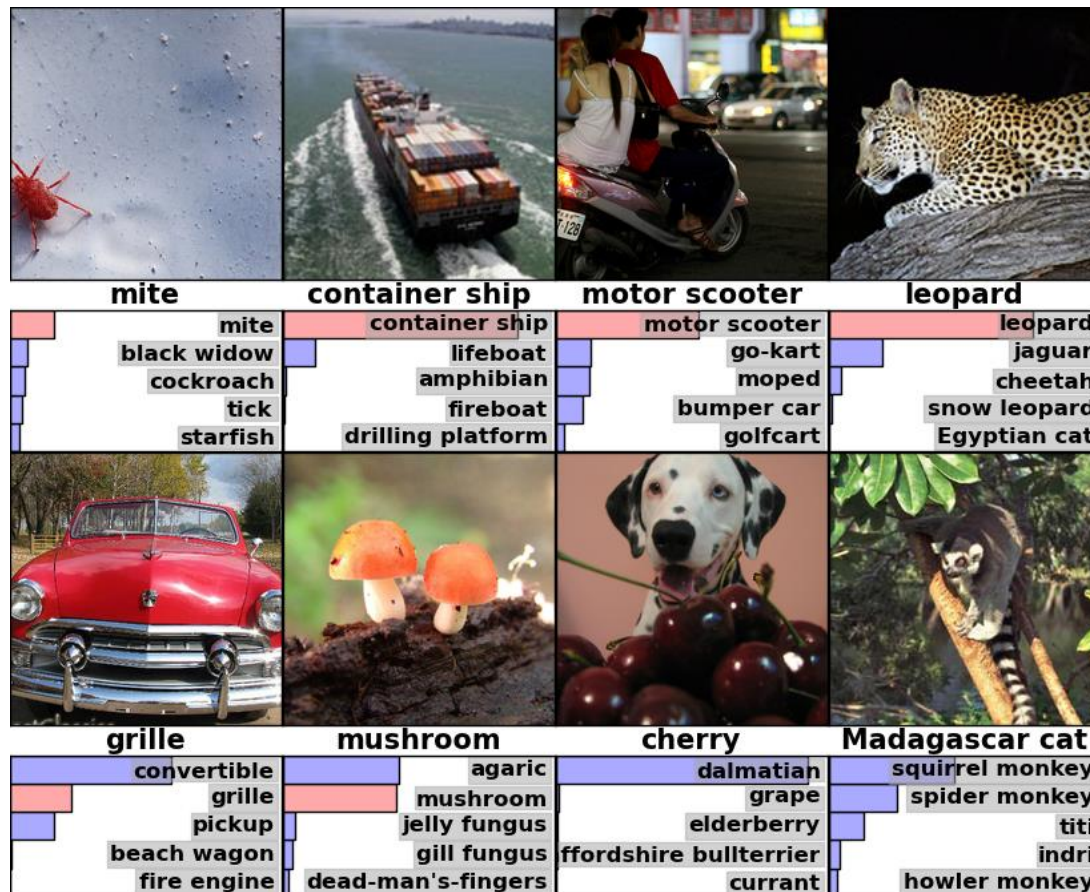
Putting it all together

A set of pixels becomes a set of votes.



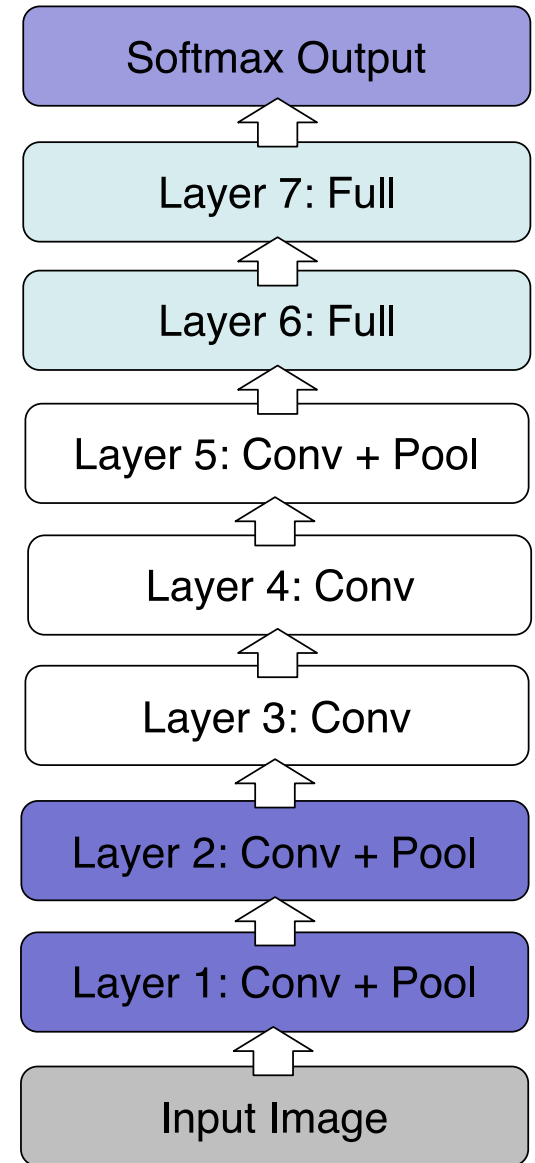
Breakthrough in Computer Vision

“AlexNet” (Krizhevsky et al., 2012), winning entry of ImageNet 2012 competition with a Top-5 error rate of 15.3% (next best system with highly engineered features based got 26.1% error)



AlexNet

- 8 layers total
- Trained on Imagenet dataset [Deng et al. CVPR'09]
- 18.2% top-5 error



[From Rob Fergus' CIFAR 2016 tutorial]

Breakthrough in Computer Vision

