

A Gentle Introduction to Deep Learning

Dr. Rodrigue Rizk

Rodrigue.Rizk@usd.edu

Learning Objectives

- ❑ Understand the intuition behind deep neural networks particularly convolutional neural networks (CNN)
- ❑ Learn the machine vision foundations
- ❑ Understand the building blocks of CNN
- ❑ Get hands-on experience with CNN

What Do You See?



How can we help computers see?



What computers ‘see’: Images as Numbers

What you see



Input Image

What you both see

157	153	174	168	160	152	129	151	172	161	165	166
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	94	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	257	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Input Image + values

What the computer “sees”

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	94	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Pixel intensity values
("pix-el"=picture-element)

An image is just a matrix of numbers [0,255]

What computers ‘see’: Images as Numbers

What you see



Input Image

What you both see

157	153	174	168	160	152	129	151	172	161	165	166
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	94	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	257	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Input Image + values

What the computer “sees”

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	83	17	110	210	180	154
180	180	50	14	94	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Pixel intensity values
("pix-el"=picture-element)

An image is just a matrix of numbers [0,255]

Can I just do classification on this long image vector directly?

What computers ‘see’: Images as Numbers

What you see



Input Image

What you both see

157	153	174	168	160	152	129	151	172	161	165	166
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	257	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Input Image + values

What the computer “sees”

157	153	174	168	150	152	129	151	172	161	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	6	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	105	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	36	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	235	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	200	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

Pixel intensity values
("pix-el"=picture-element)

An image is just a matrix of numbers [0,255]

Can I just do classification on this long image vector directly?

No. Instead: exploit image spatial structure. Learn patches. Build them up

High Level Feature Detection

Let's identify key features in each image category



Nose, Eyes, Mouth



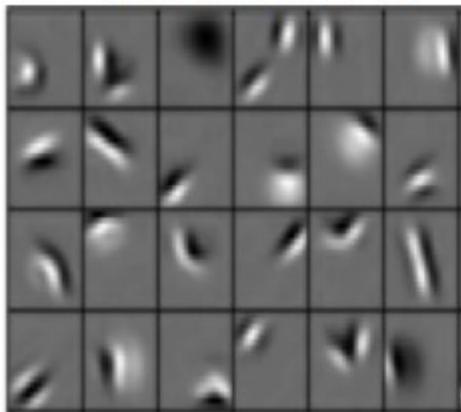
Wheels, License Plate,
Headlights



Door, Windows, Steps

High Level Feature Detection

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

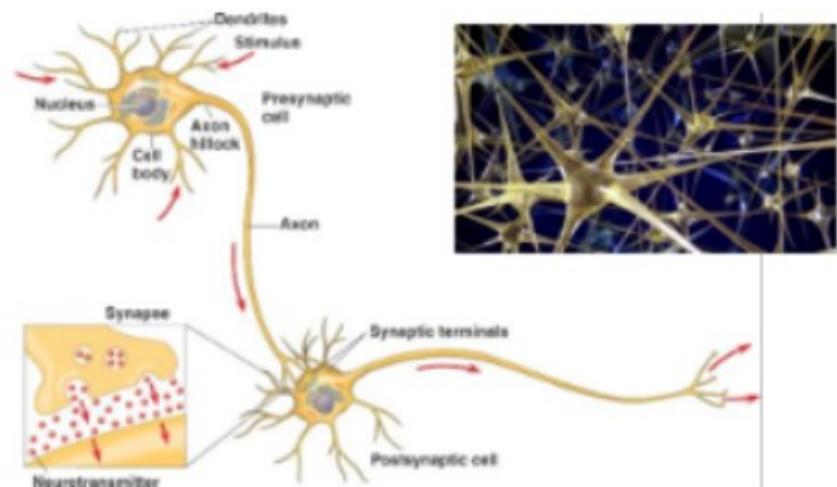
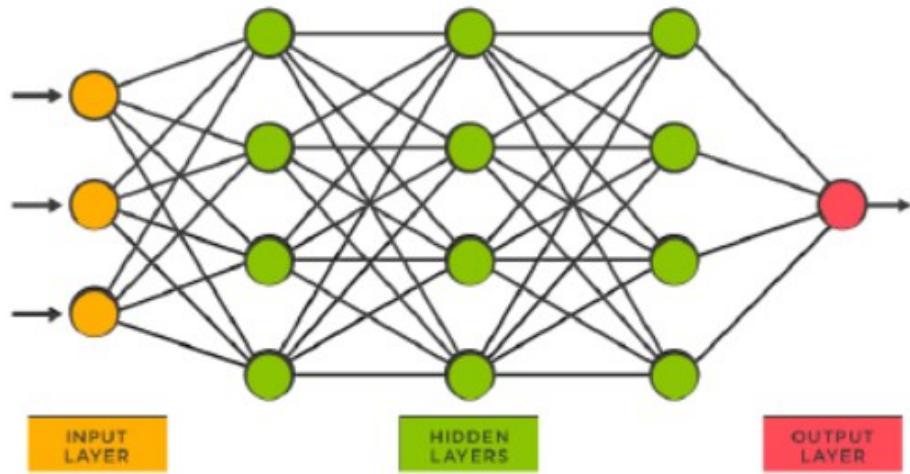
High level features



Facial structure

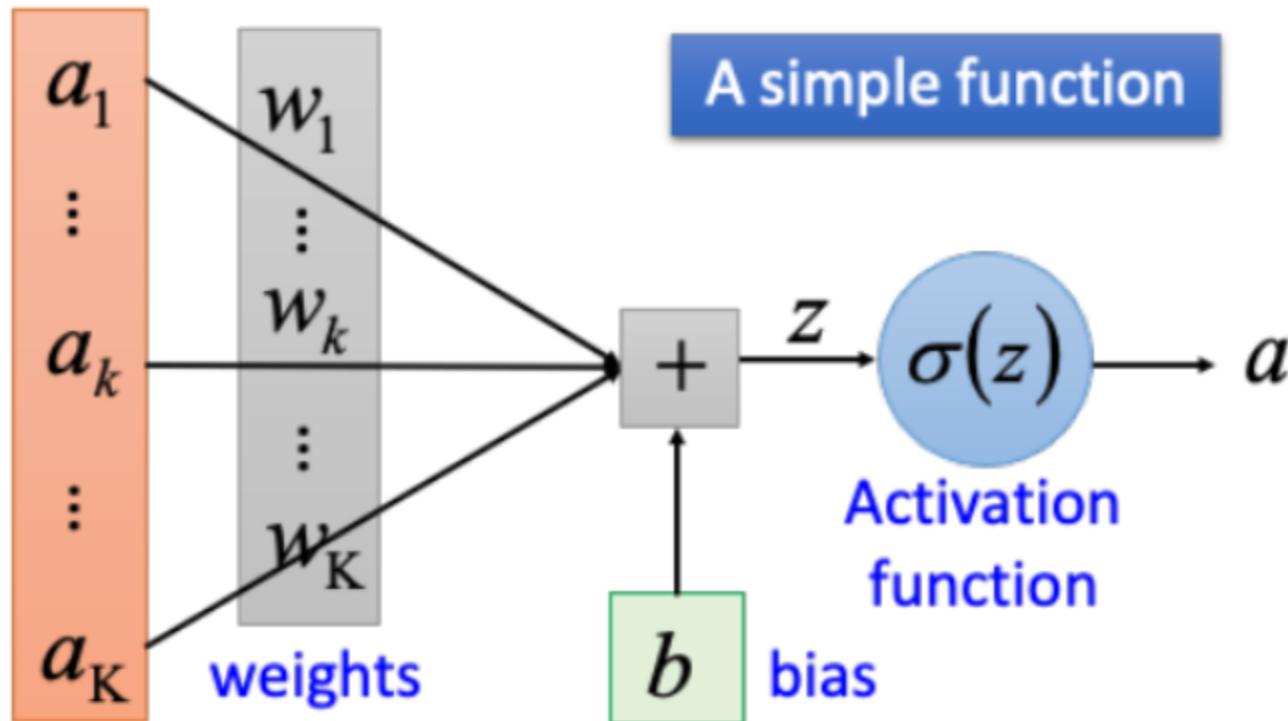
Learn the hierarchy of features directly from the data

Neural Network



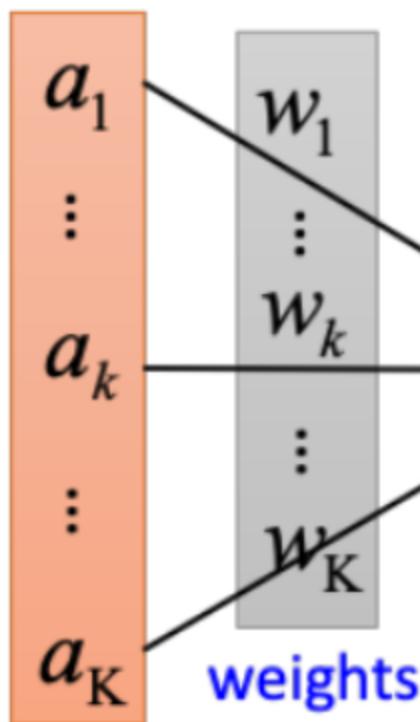
At Each Neuron

$$z = a_1 w_1 + \cdots + a_k w_k + \cdots + a_K w_K + b$$



At Each Neuron

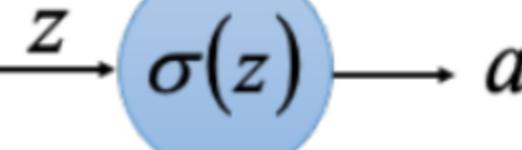
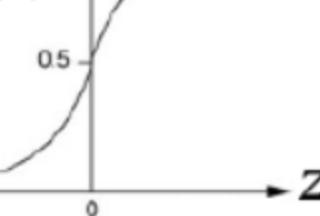
$$z = a_1 w_1 + \cdots + a_k w_k$$



Sigmoid Function

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

$$\sigma(z)$$



Activation
function

bias

b

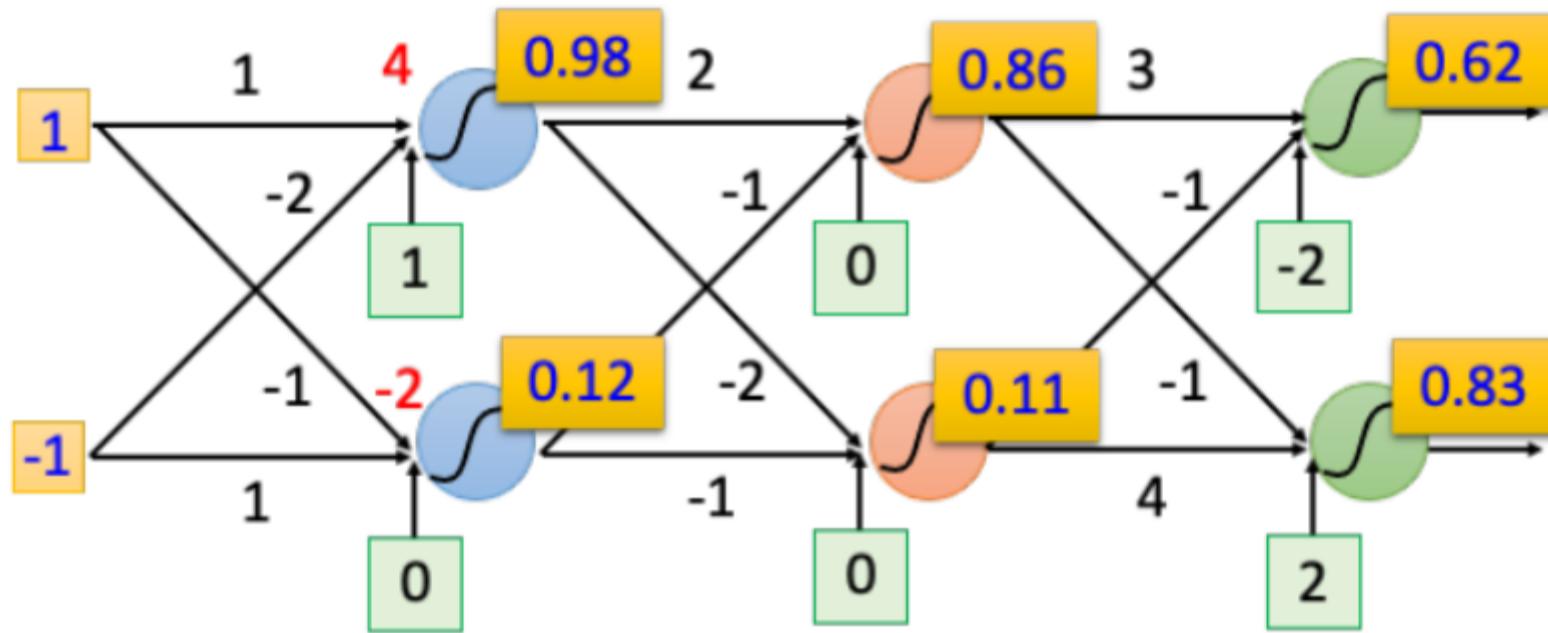
$$a$$

z

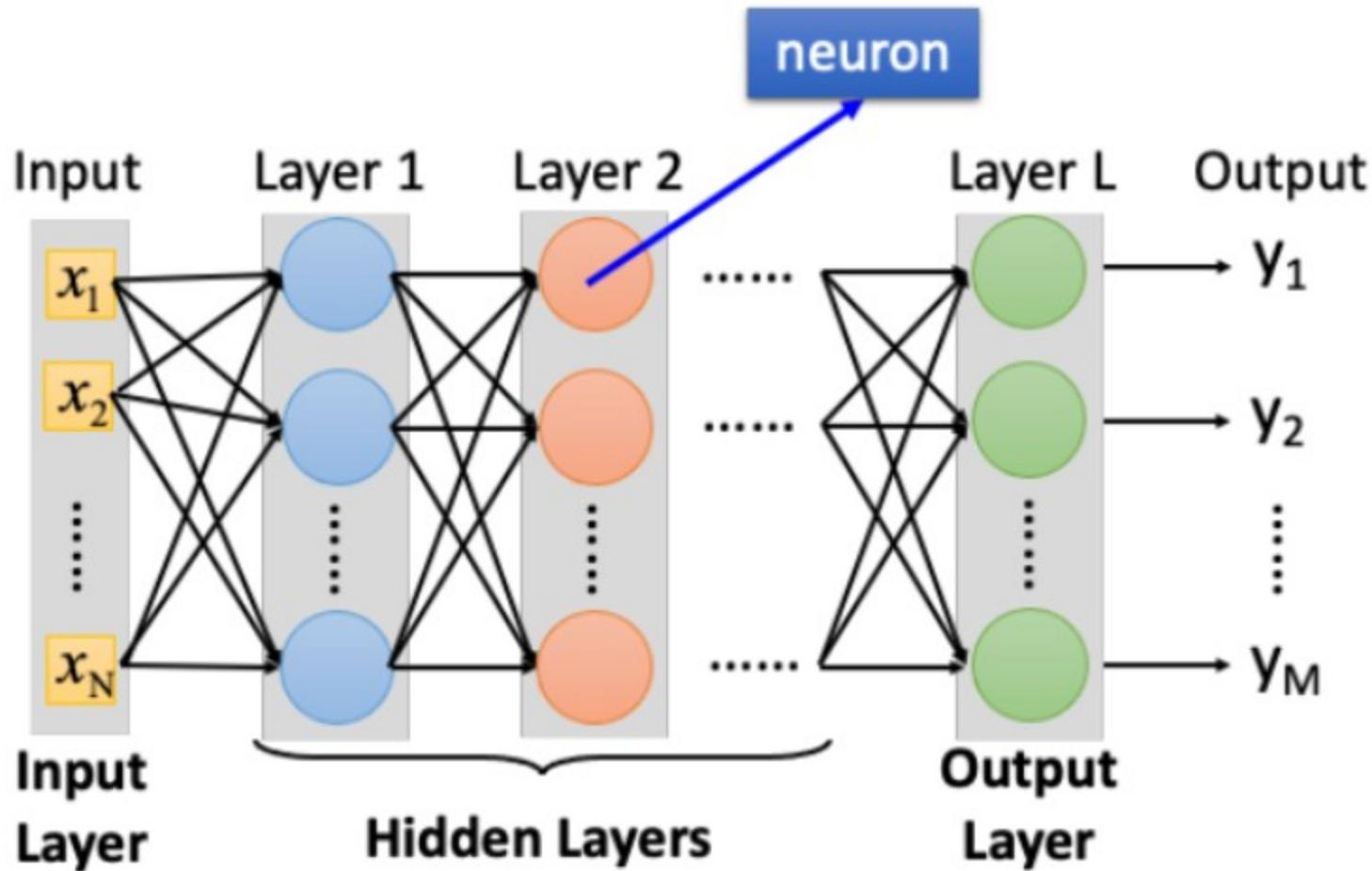
+

z

Fully Connected Neural Network



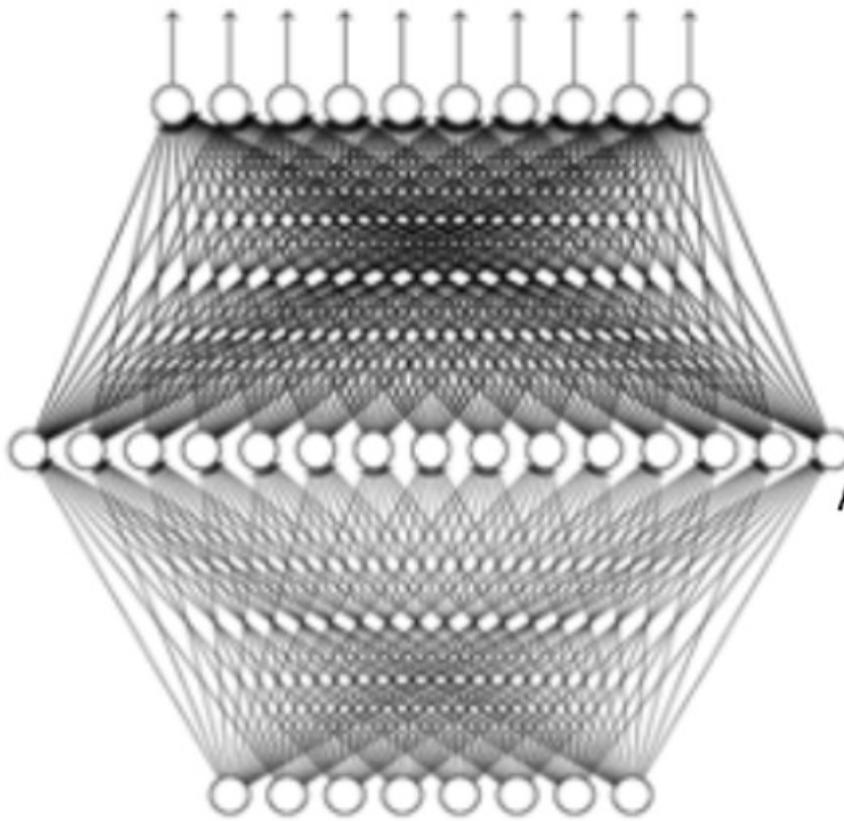
Fully Connected Neural Network



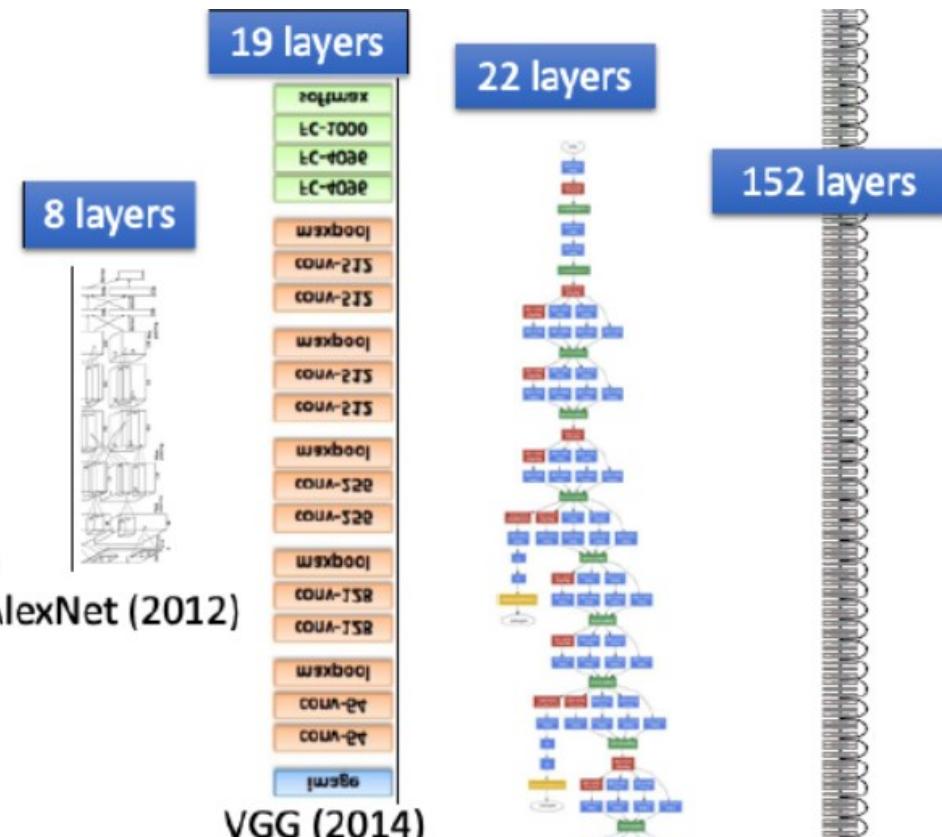
Deep means many hidden layers

Deep Neural Network

Deep = Many hidden layers



AlexNet (2012)



Using multiple layers of neurons to represent some functions

GoogleNet (2015) ResidualNet (2015)

Output Layer

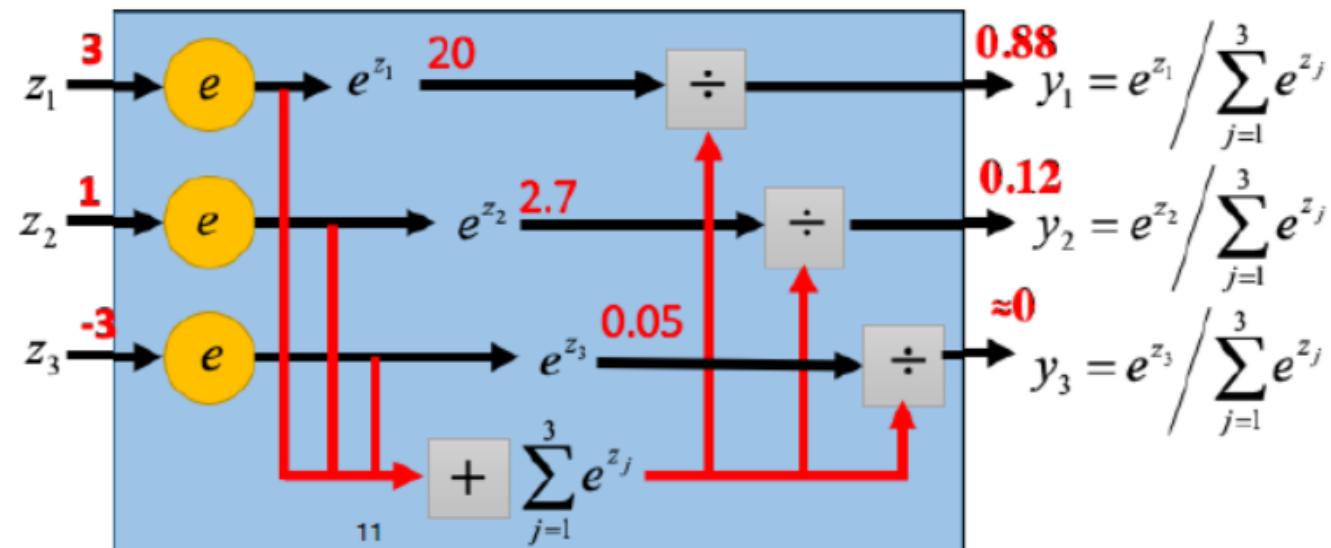
Ordinary Layer

$$z_1 \rightarrow \sigma \rightarrow y_1 = \sigma(z_1)$$

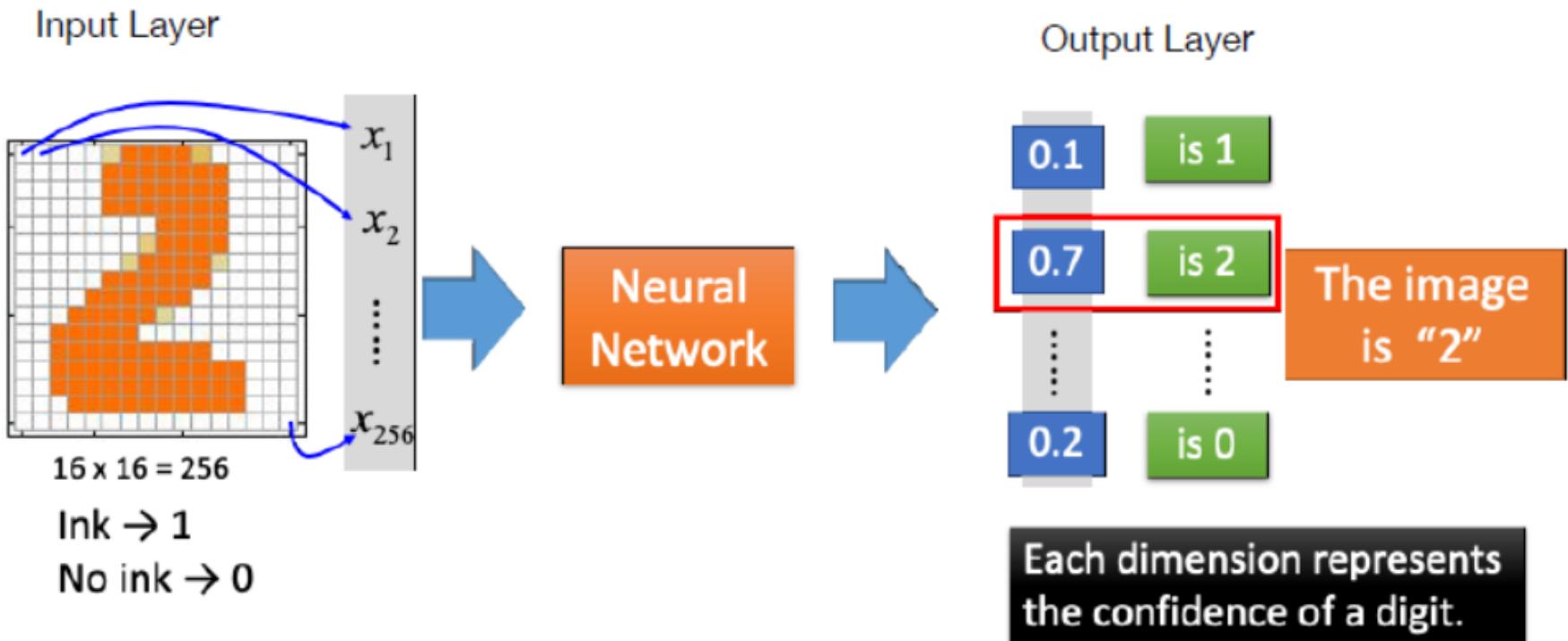
$$z_2 \rightarrow \sigma \rightarrow y_2 = \sigma(z_2)$$

$$z_3 \rightarrow \sigma \rightarrow y_3 = \sigma(z_3)$$

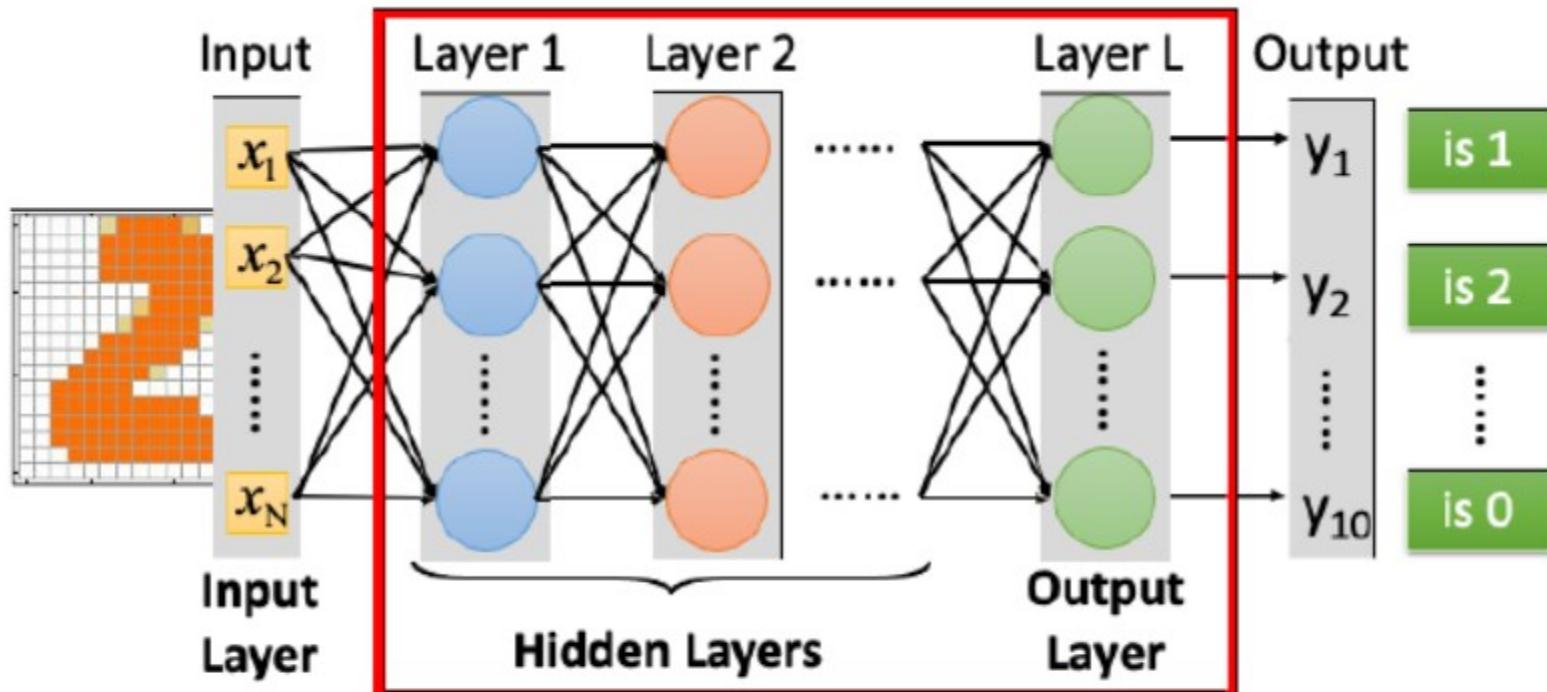
Softmax Layer



An Example



An Example

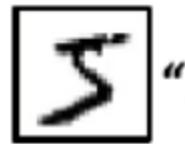


Needs to determine the network structure

How many layers?

How many neurons for each layer?

Learning Target



"5"



"0"



"4"



"1"



"9"



"2"



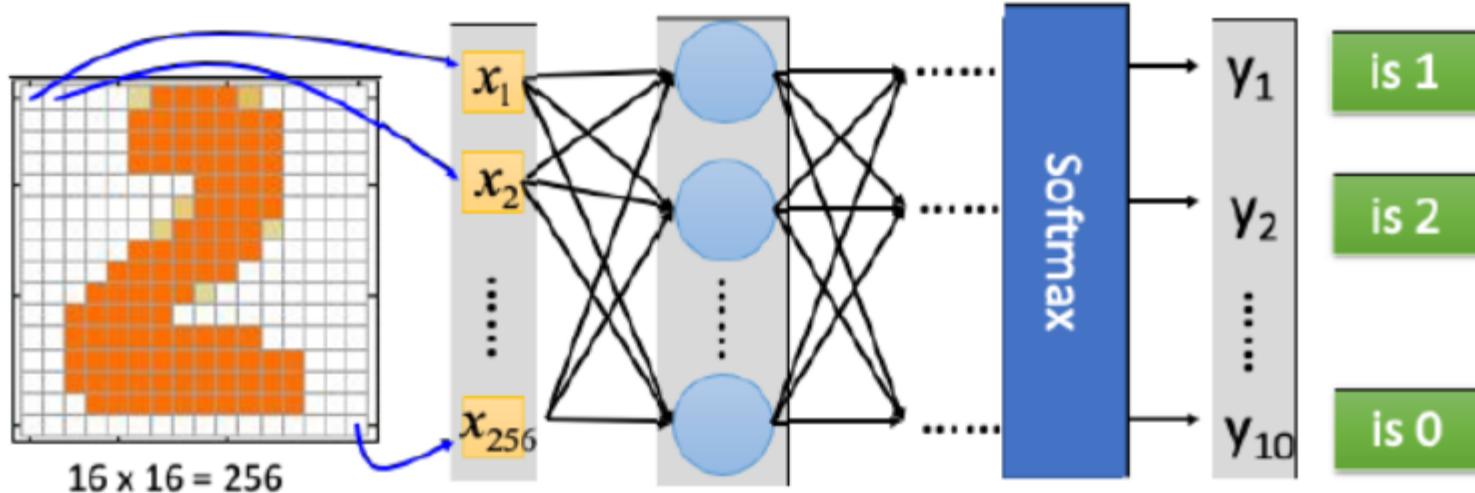
"1"



"3"

The learning target is defined on
the training data.

Learning Target



Ink $\rightarrow 1$

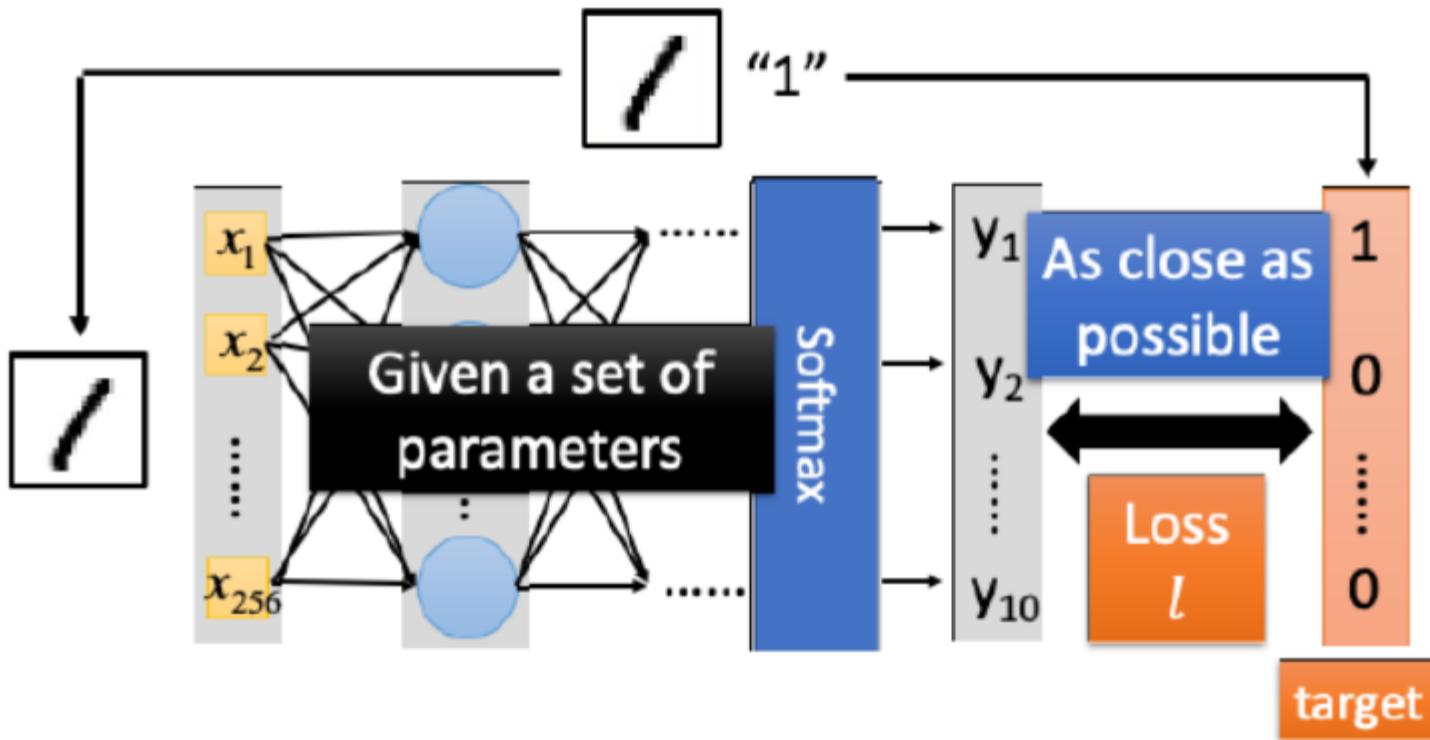
No ink $\rightarrow 0$

The learning target is

Input:  $\rightarrow y_1$ has the maximum value

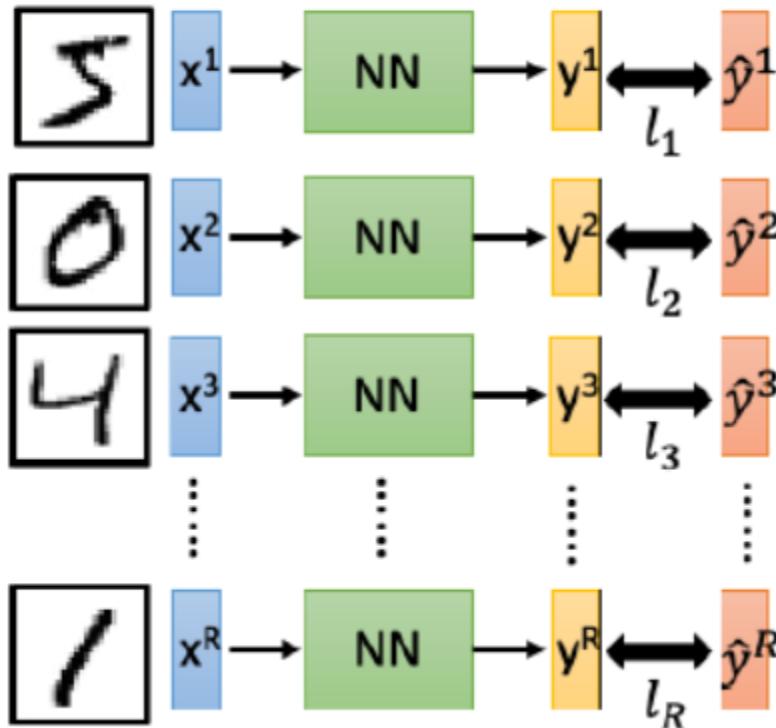
Input:  $\rightarrow y_2$ has the maximum value

Good Function = Loss as Small as Possible



Loss can be square error or cross entropy between the output and the target

Total Loss



$$L = \sum_{r=1}^R l_r$$

As small as possible

Find a function in
function set that
minimizes total loss L

Find the network
parameters θ^* that
minimize total loss L

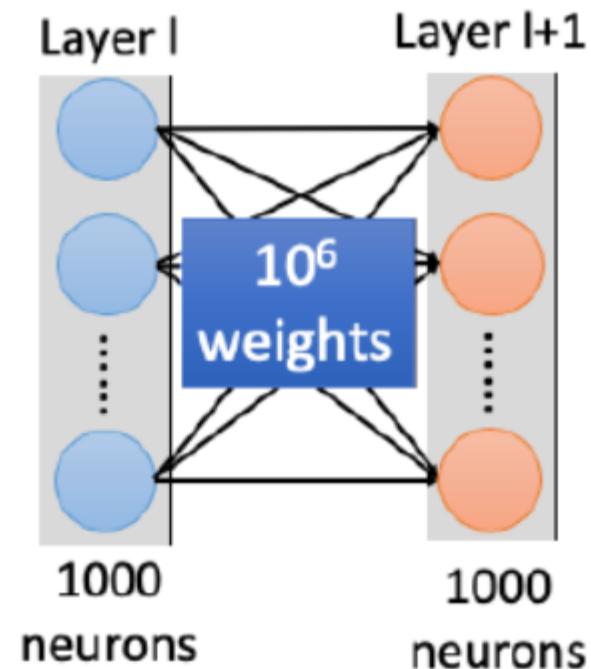
Total Loss

Find **network parameters θ^*** that minimize total loss L

Enumerate all possible values

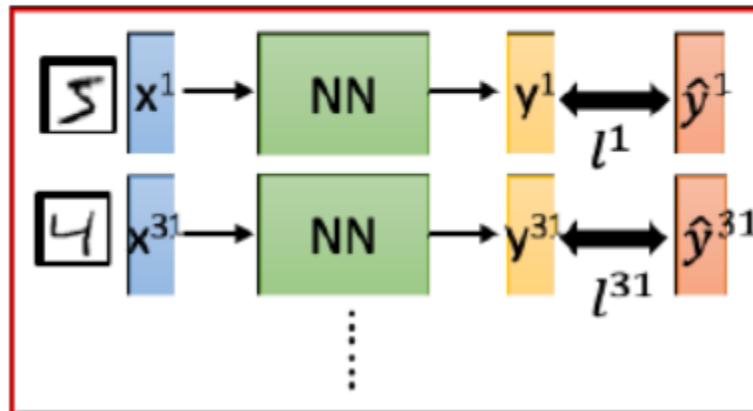
Network parameters $\theta =$
 $\{w_1, w_2, w_3, \dots, b_1, b_2, b_3, \dots\}$

Millions of parameters

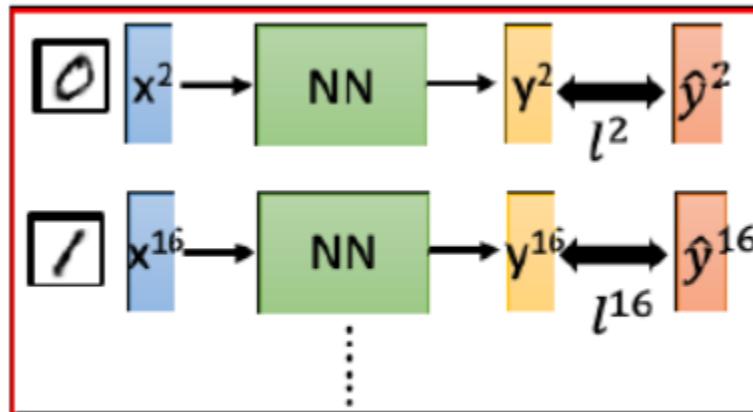


Mini-batch

Mini-batch



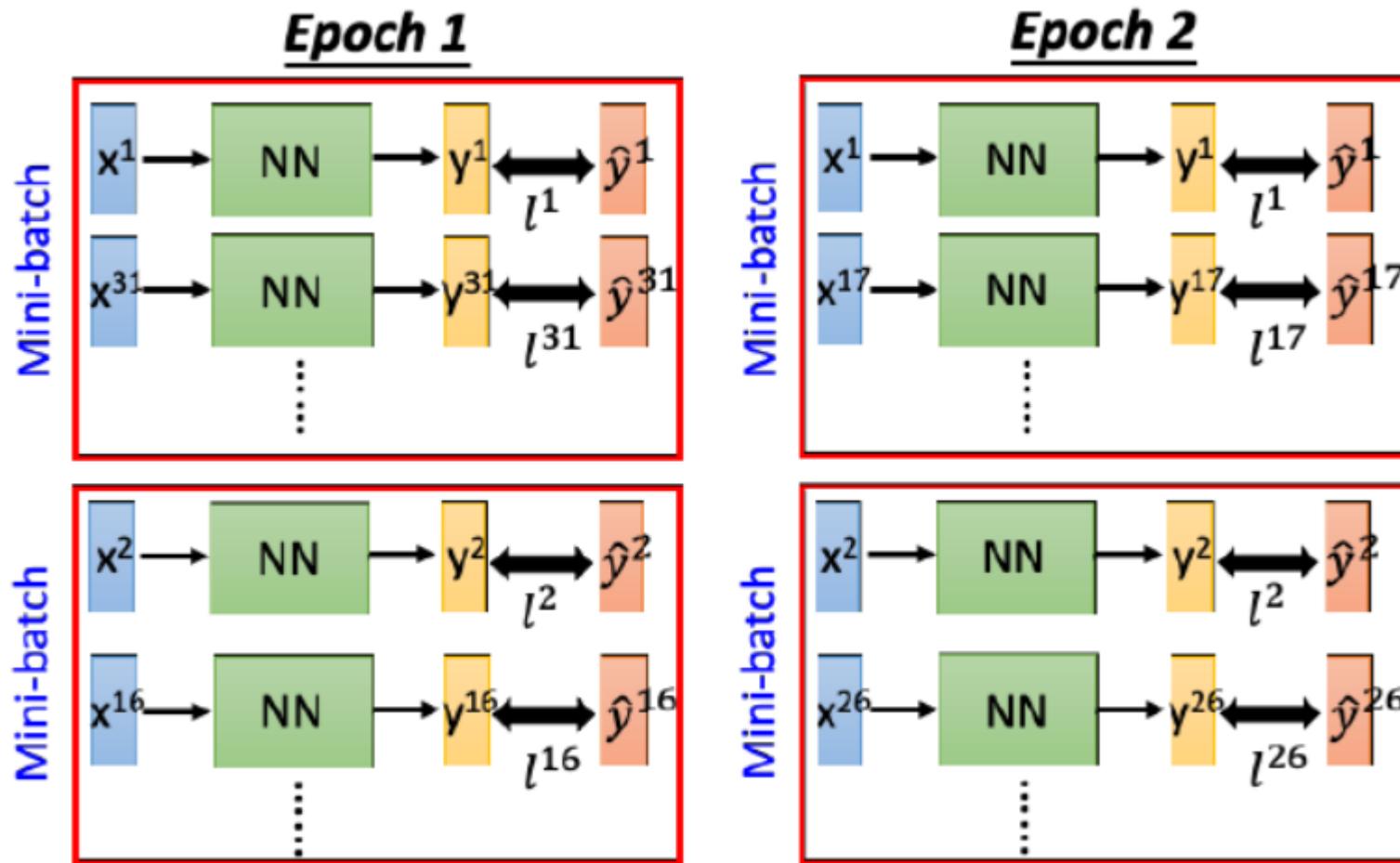
Mini-batch



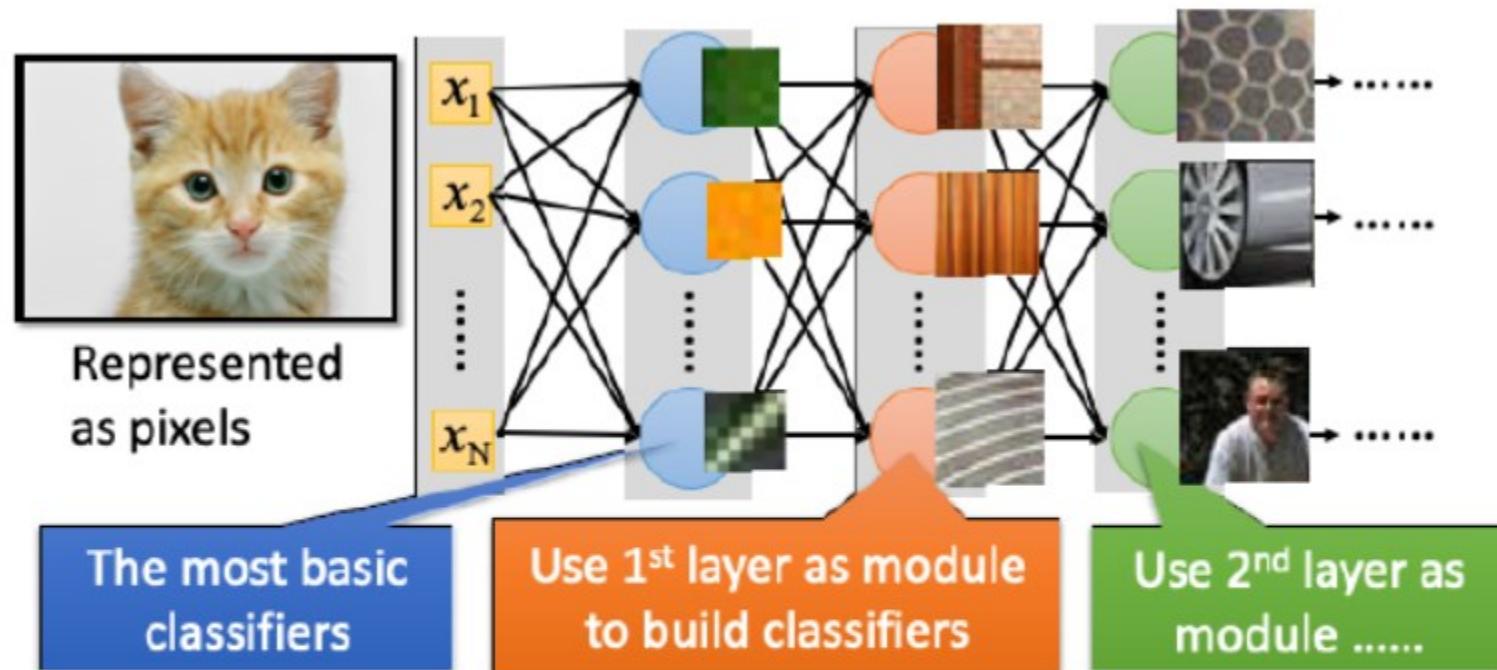
- Randomly initialize network parameters
 - Pick the 1st batch
 $L' = l^1 + l^{31} + \dots$
 Update parameters once
 - Pick the 2nd batch
 $L'' = l^2 + l^{16} + \dots$
 Update parameters once
 :
 - Until all mini-batches have been picked

one epoch
- Repeat the above process

Mini-batch



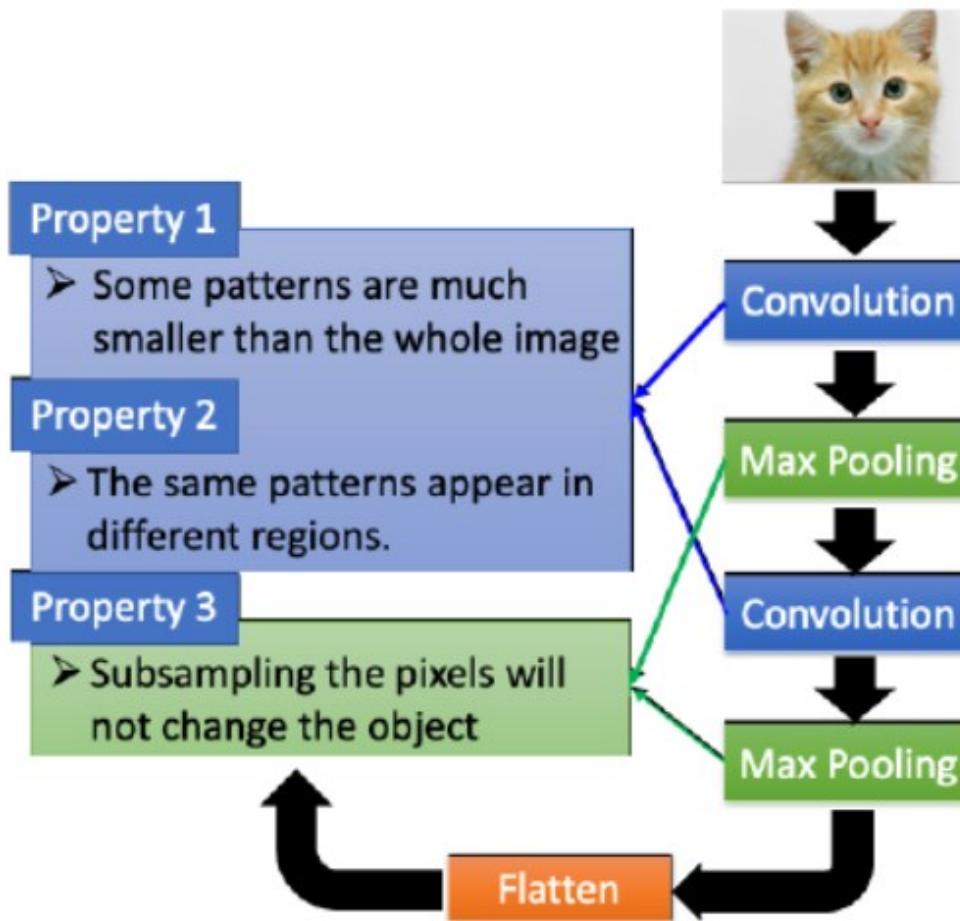
Convolutional Neural Network (CNN)



Each pixel as one input

Can we simplify the network by considering the properties of Images?

Convolutional Neural Network (CNN)



We can subsample the pixels to make image smaller

Less parameters for the network to process the image

CNN - Convolution

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

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

Filter1
Matrix

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

Filter 1
Matrix

...

...

...

...

Each filter detects a small region (3x3) **Property 1**

CNN - Convolution

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

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

Filter 1

 3

CNN - Convolution

Stride = 1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

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

Filter 1

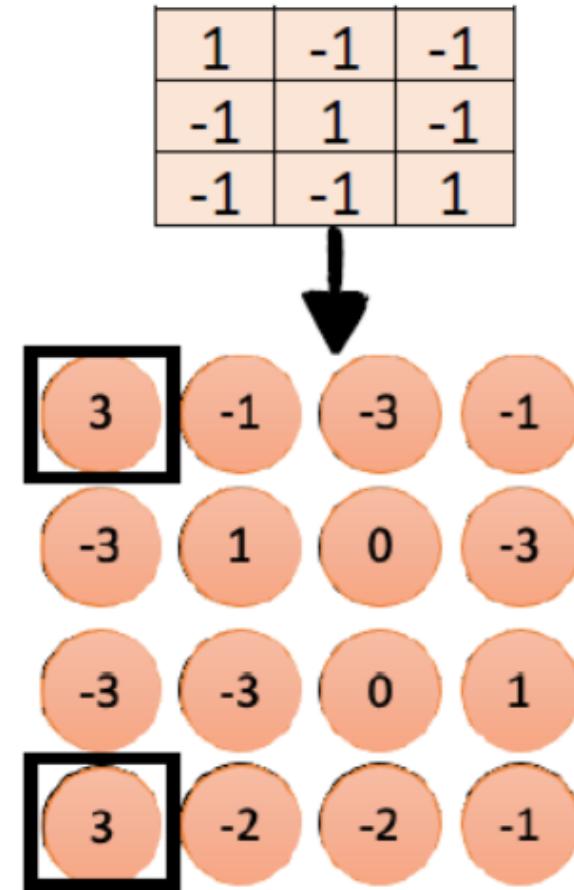
$$3 \quad -1$$

$$-3$$

1	0	0	0	0	0	1
0	1	0	0	1	0	0
0	0	1	1	0	0	0
1	0	0	0	1	0	0
0	1	0	0	1	0	0
0	0	1	0	1	0	0

CNN - Convolution

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0



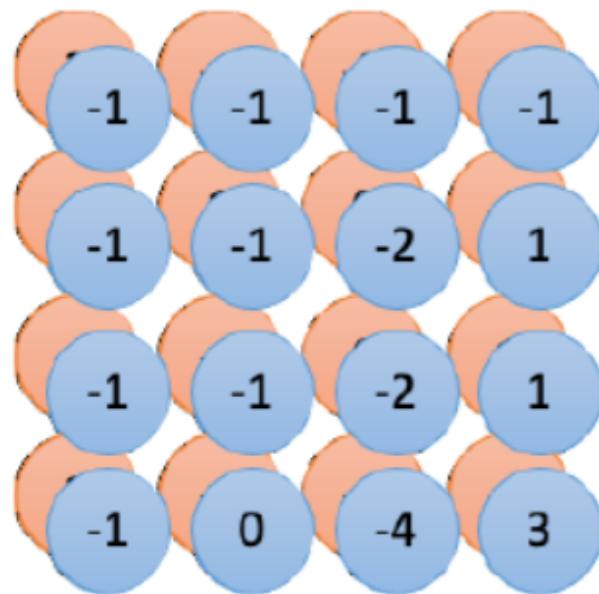
Property 2: Same patterns appears at different region!

CNN - Convolution

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

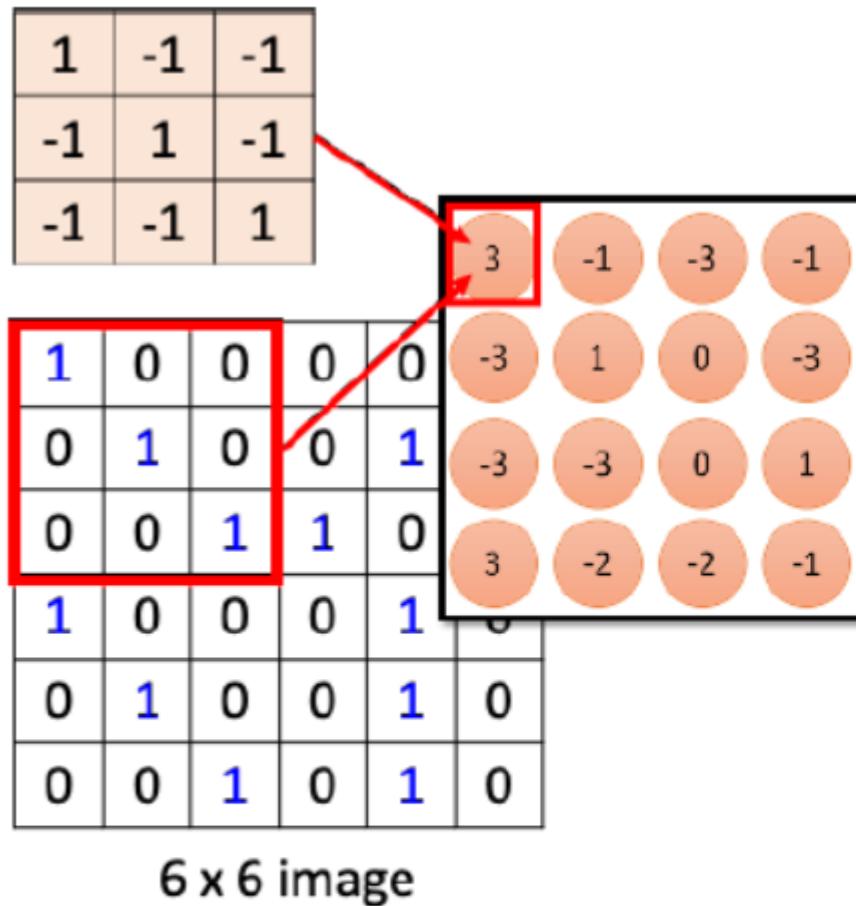
-1	1	-1
-1	1	-1
-1	1	-1

Filter2

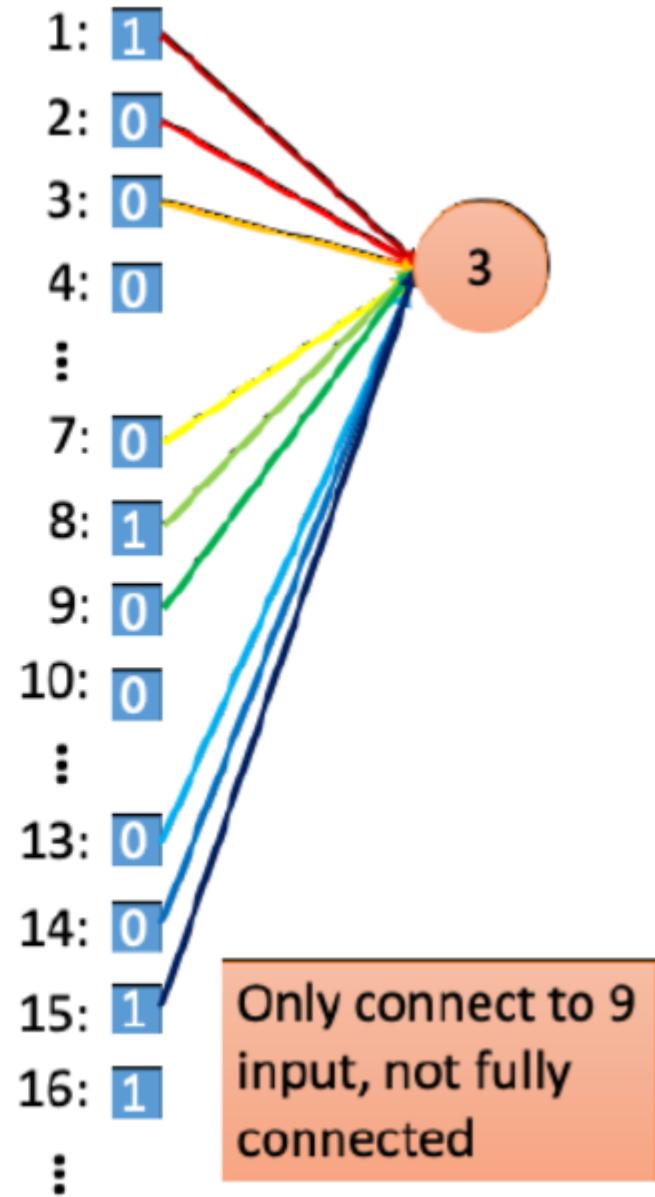


Property 2: Same patterns appears at different region!

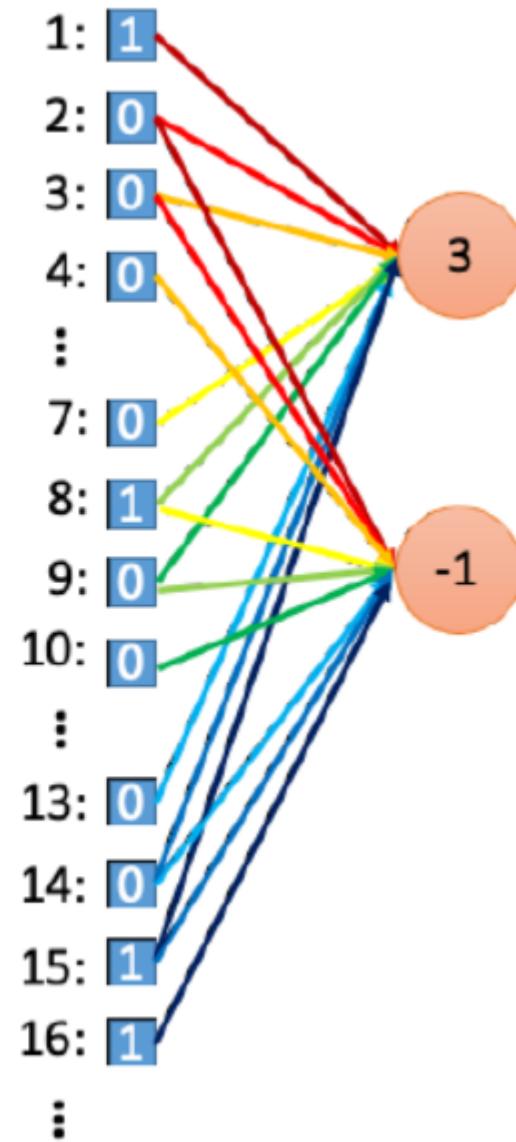
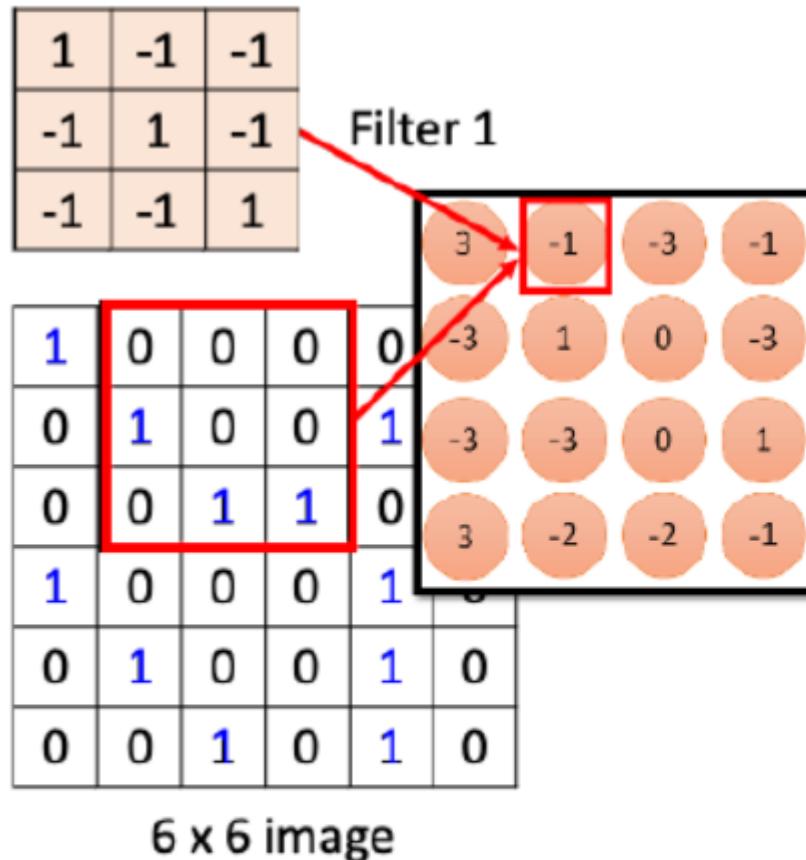
CNN - Convolution



Less parameters!



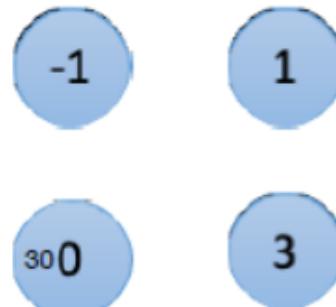
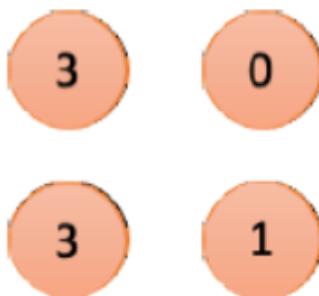
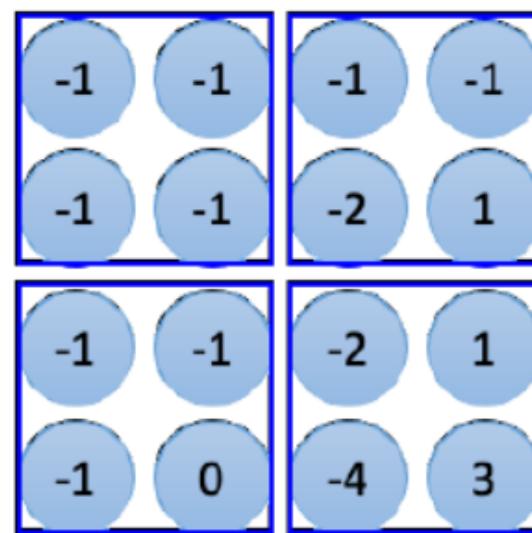
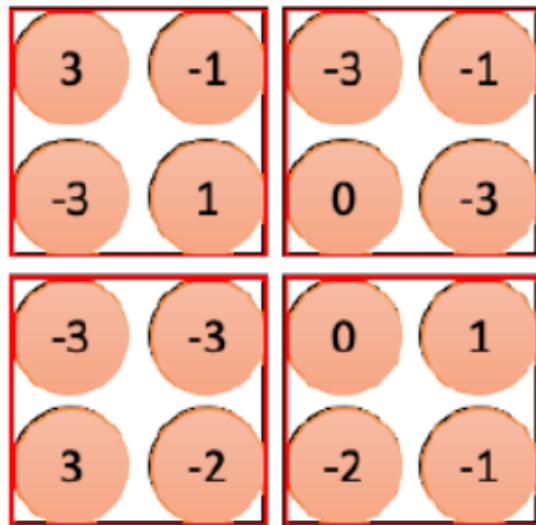
CNN - Convolution



CNN - Max Pooling

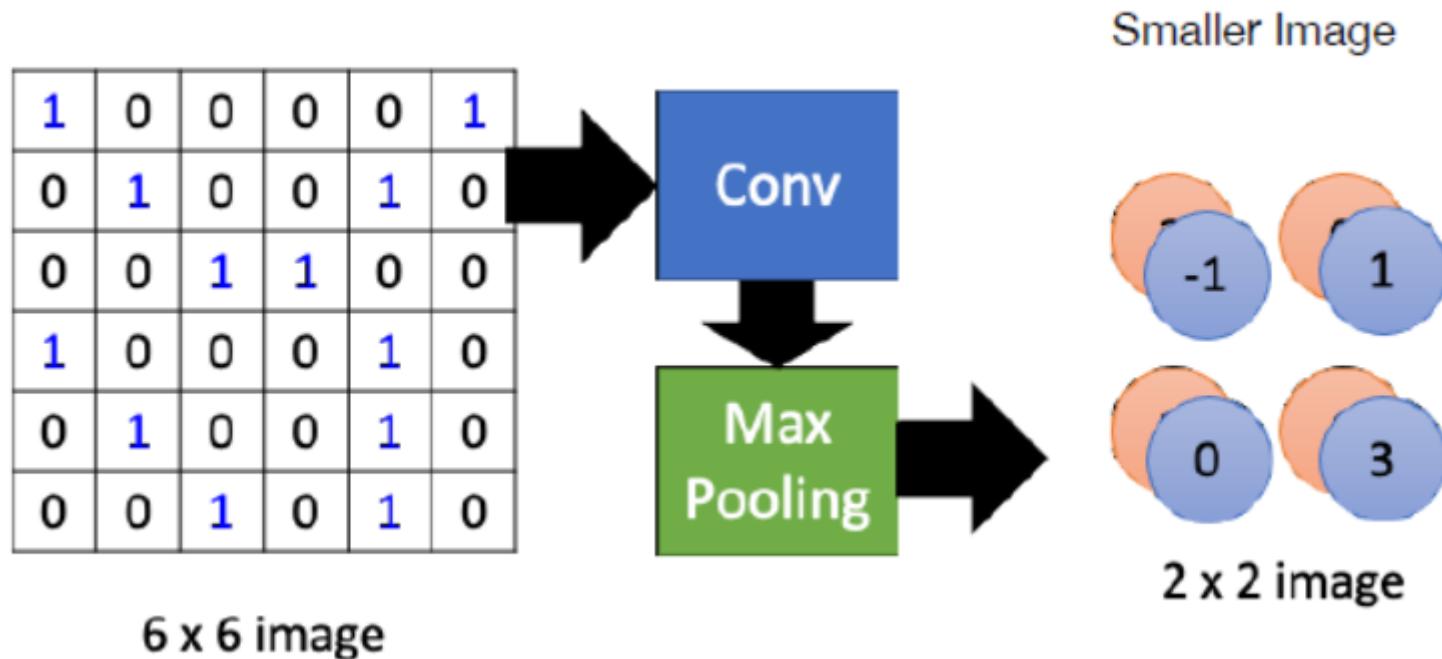
1	-1	-1
-1	1	-1
-1	-1	1

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

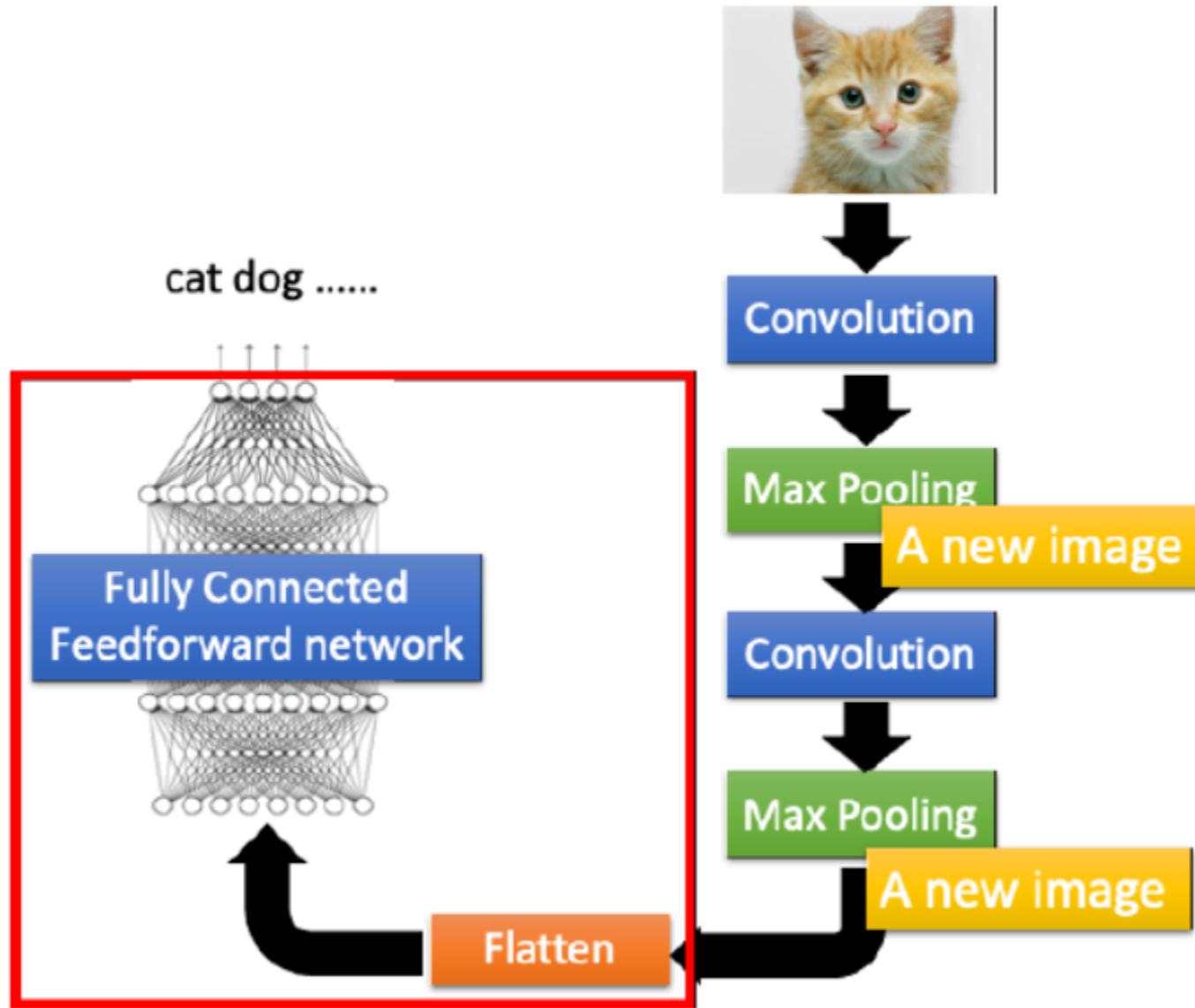


Take the dominant Features Feature Dimensionality Reduction

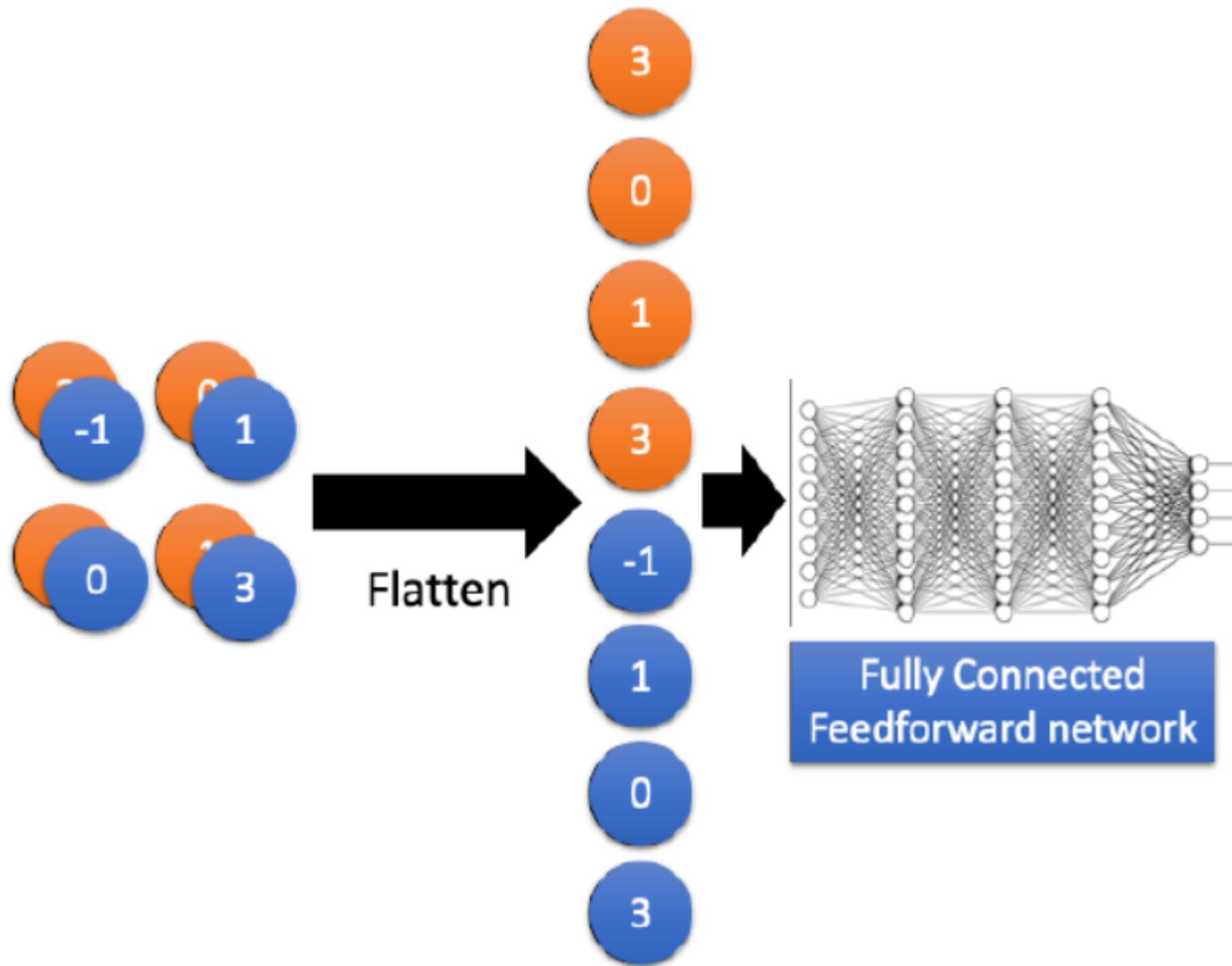
CNN - Convolution



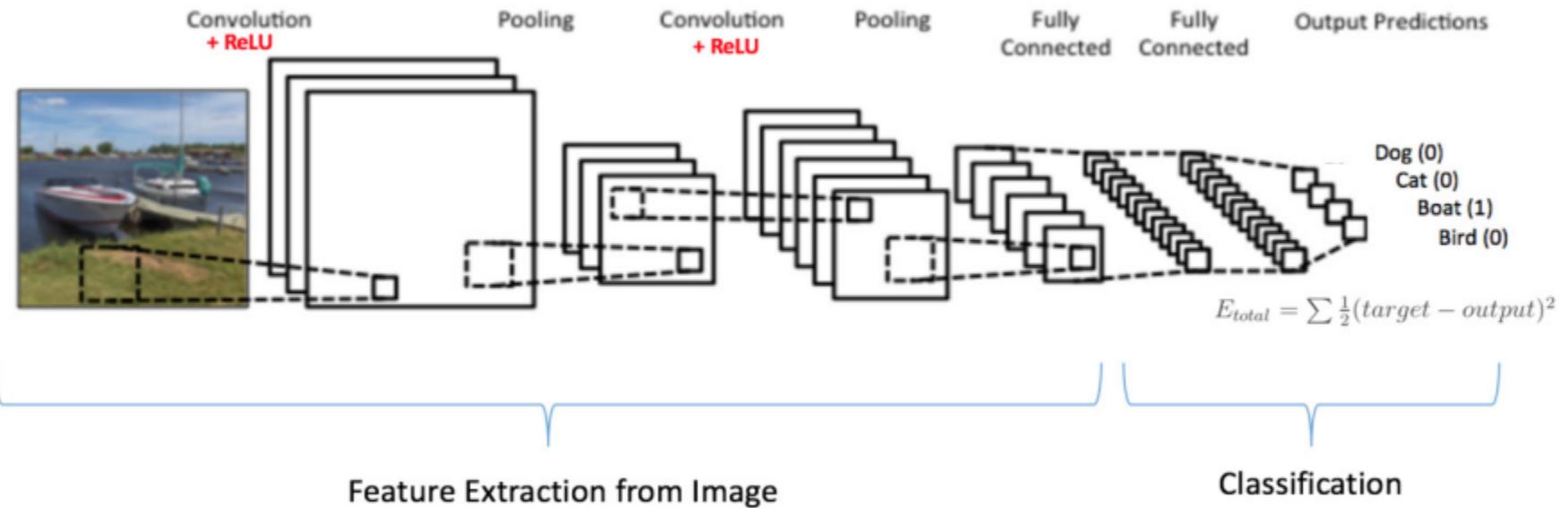
The Whole CNN



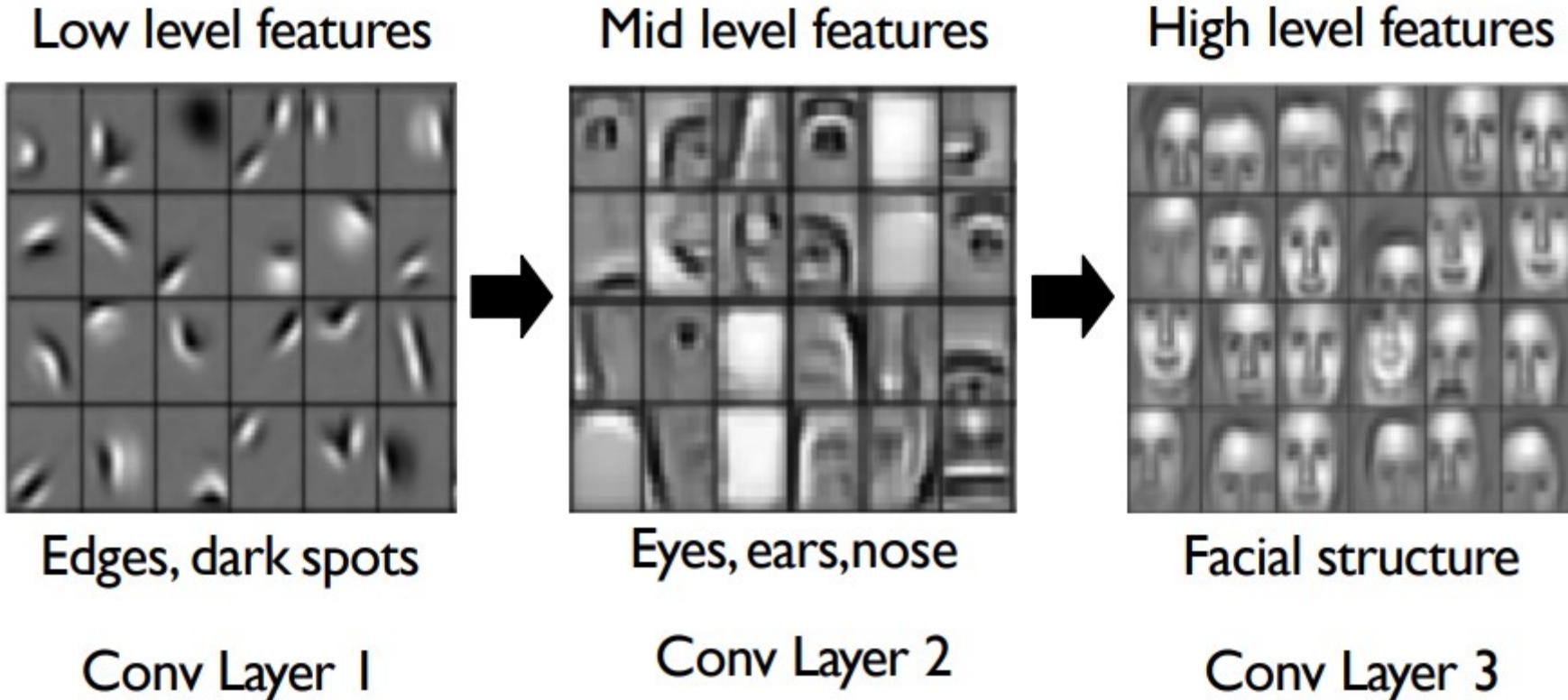
Flatten



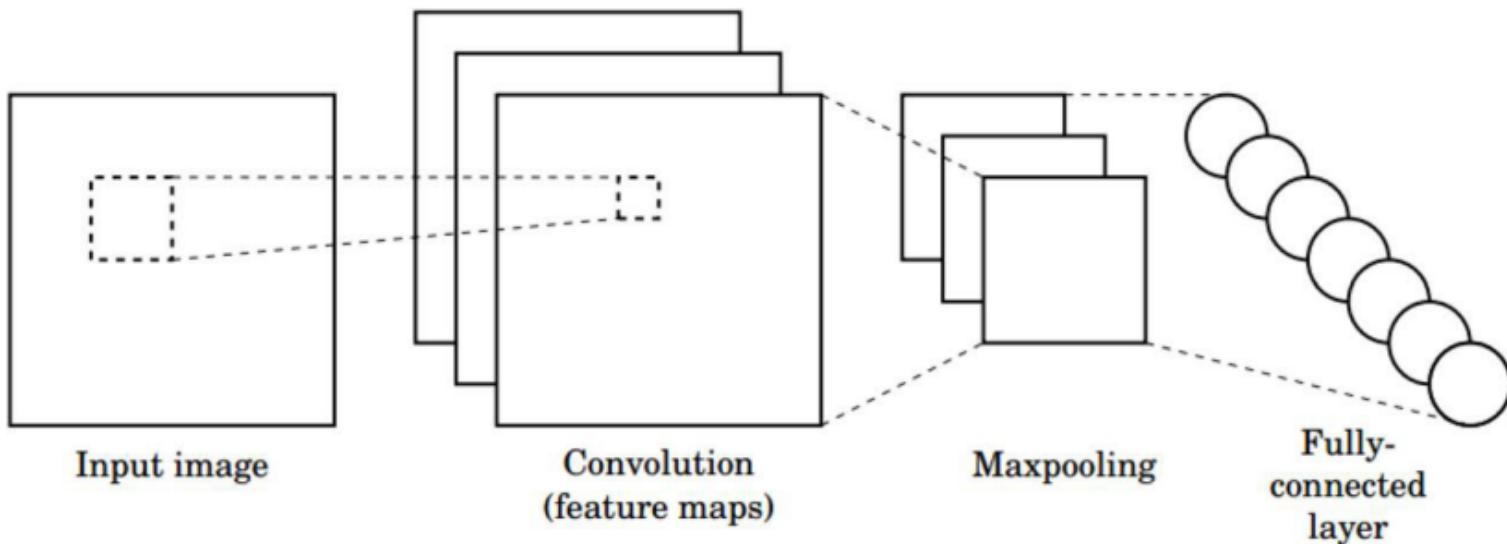
An image classification CNN



Representation Learning in Deep CNNs



CNNs for Classification

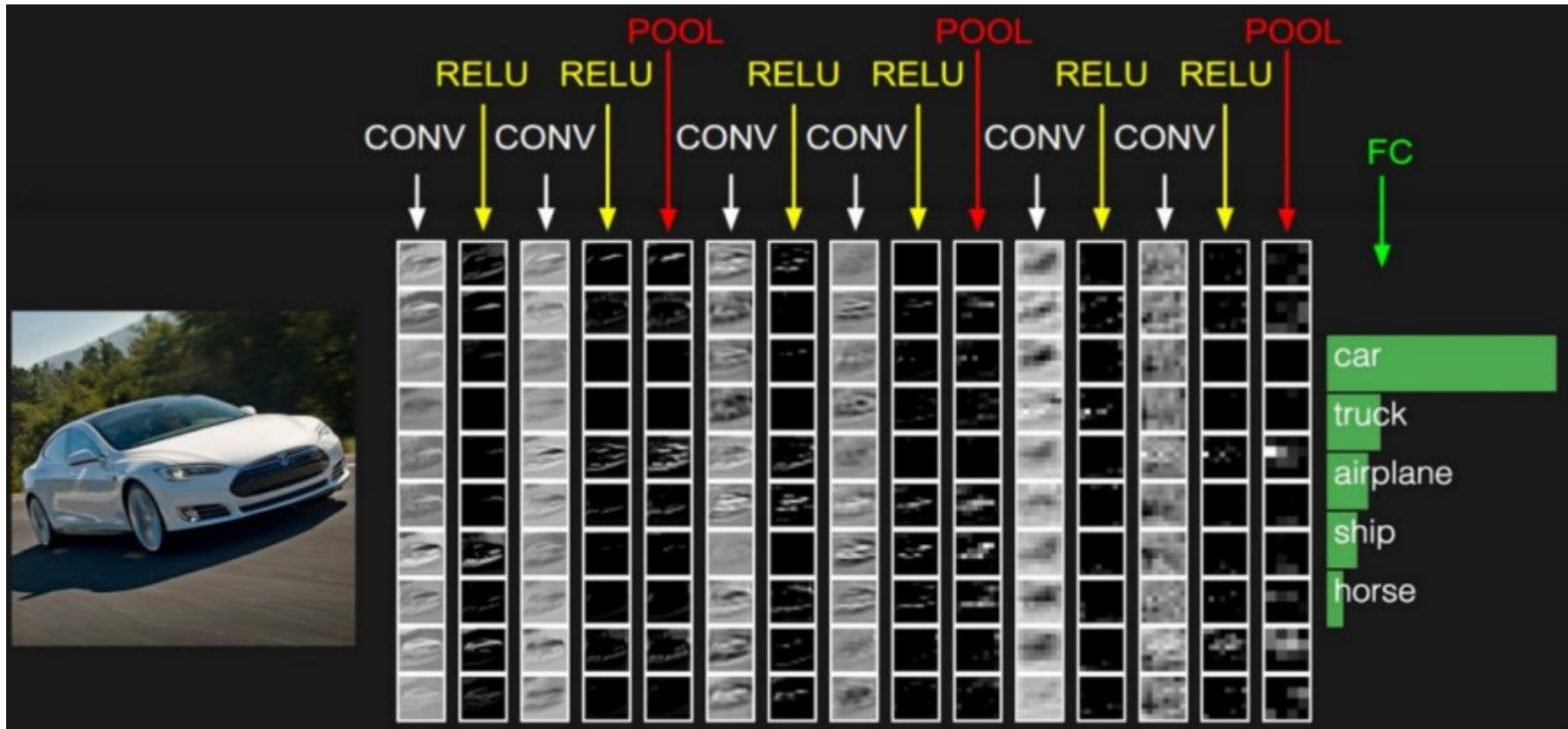


- 1. Convolution:** Apply filters to generate feature maps.
- 2. Non-linearity:** Often ReLU.
- 3. Pooling:** Downsampling operation on each feature map.

Train model with image data.

Learn weights of filters in convolutional layers.

Example – Six convolutional layers



Hands-On Example

- ❑ Python provides a set of libraries including different deep learning packets
- ❑ Standard libraries provide the ready-to-use implementation of algorithm
- ❑ **Tensorflow** is the one we will use in this lecture



Hands-On Example

```

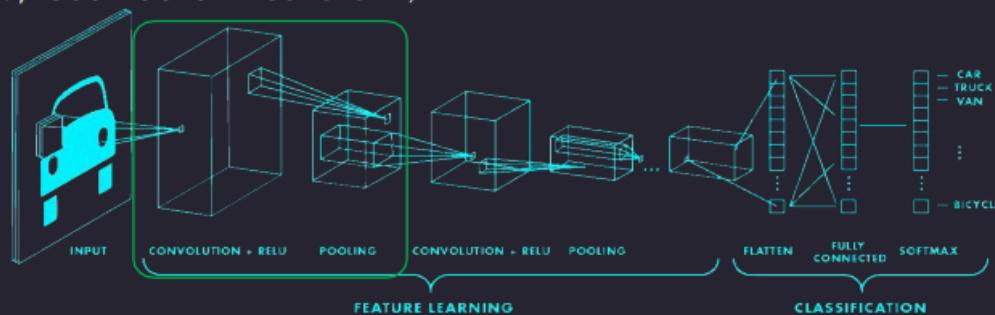
import tensorflow as tf

def generate_model():
    model = tf.keras.Sequential([
        # first convolutional layer
        tf.keras.layers.Conv2D(32, filter_size=3, activation='relu'),
        tf.keras.layers.MaxPool2D(pool_size=2, strides=2),
    ])

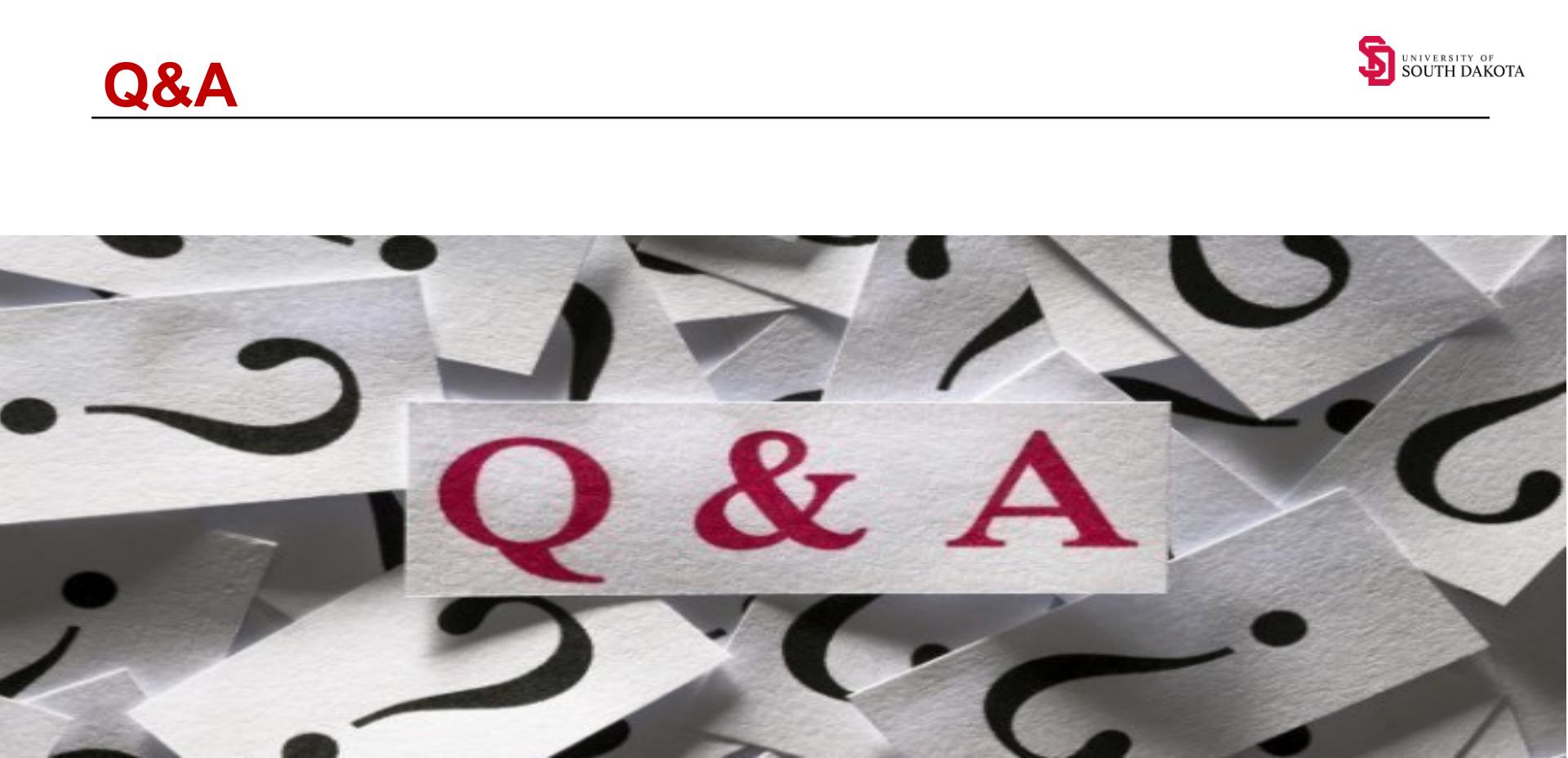
    # second convolutional layer
    tf.keras.layers.Conv2D(64, filter_size=3, activation='relu'),
    tf.keras.layers.MaxPool2D(pool_size=2, strides=2),

    # fully connected classifier
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(1024, activation='relu'),
    tf.keras.layers.Dense(10, activation='softmax')
    # 10 outputs
])

return model
  
```



Q&A



Q & A

The only stupid question is the one you were afraid to ask but never did.

-Rich Sutton