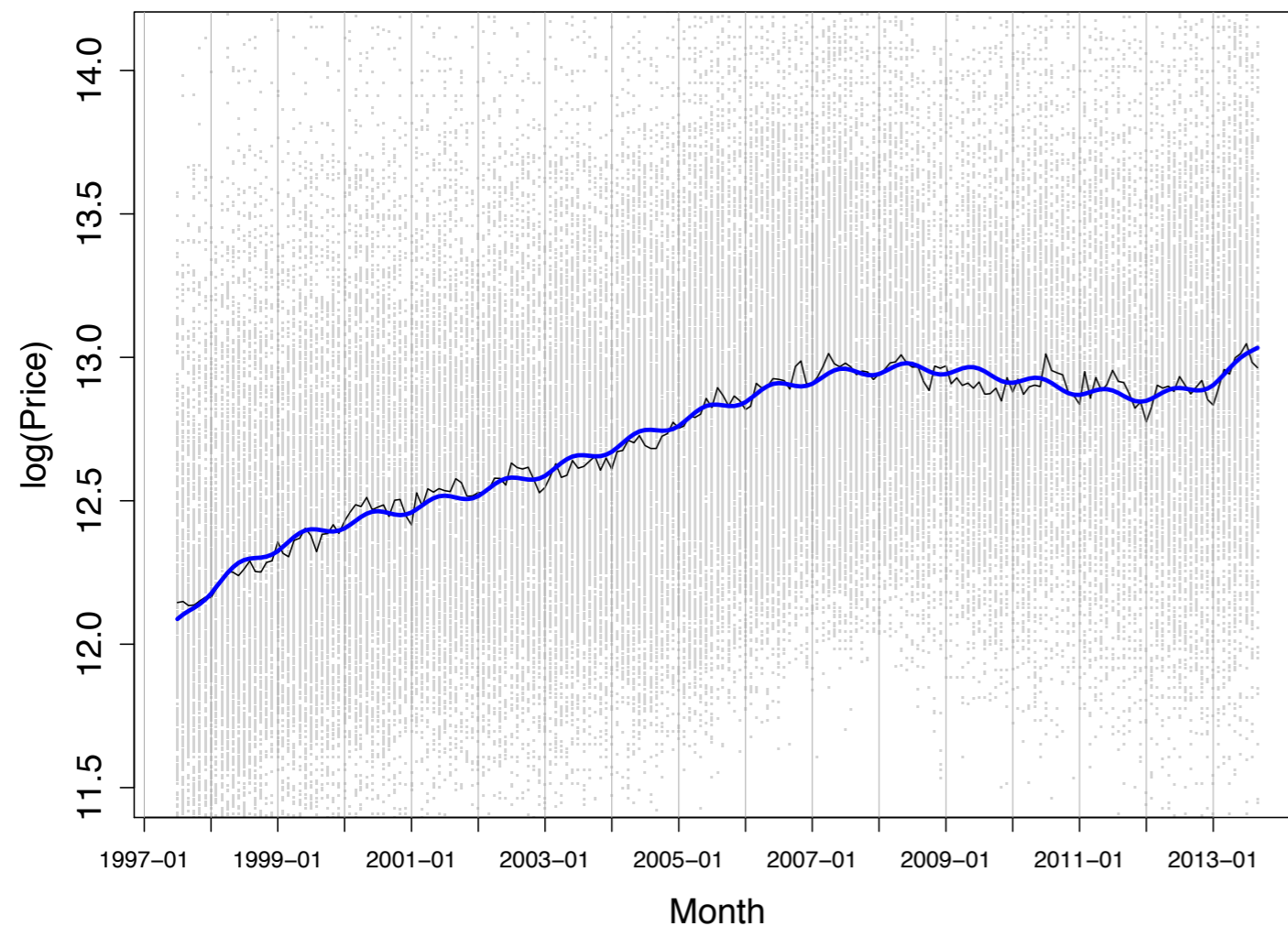# Deep Learning

Sewoong Oh

CSE/STAT 416
University of Washington

- Feature engineering is critical in achieving good performance

- e.g. seasonal trends captured by sinusoids

$$f(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \textcolor{red}{\tilde{w}_4} \sin\left(\frac{2\pi x}{12}\right) + \textcolor{red}{\tilde{w}_5} \cos\left(\frac{2\pi x}{12}\right)$$

# Image classification



Input: **x**
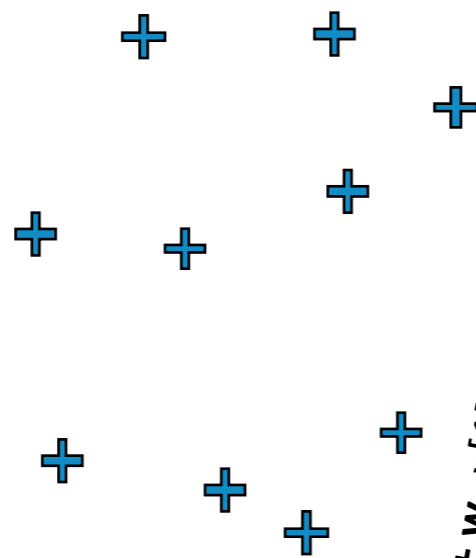**Image pixels**

Output: y
Predicted object

- Feature engineering is extremely challenging
  - For real-data that is high-dimensional and complex
- Neural networks allow us to learn features that are non-linear

# Recall: linear classification

- Input is d-dimensional data
- Output is a partition of the space into two, separated by a hyperplane (line in 2-d)
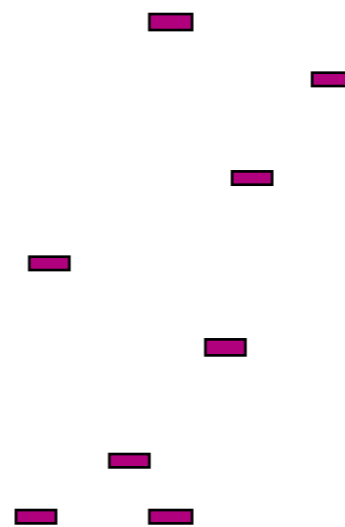- Training searches for the best line

$$\text{Score}(x) = w_0 + \mathbf{w}_1 x[1] + \mathbf{w}_2 x[2] + \dots + \mathbf{w}_d x[d]$$
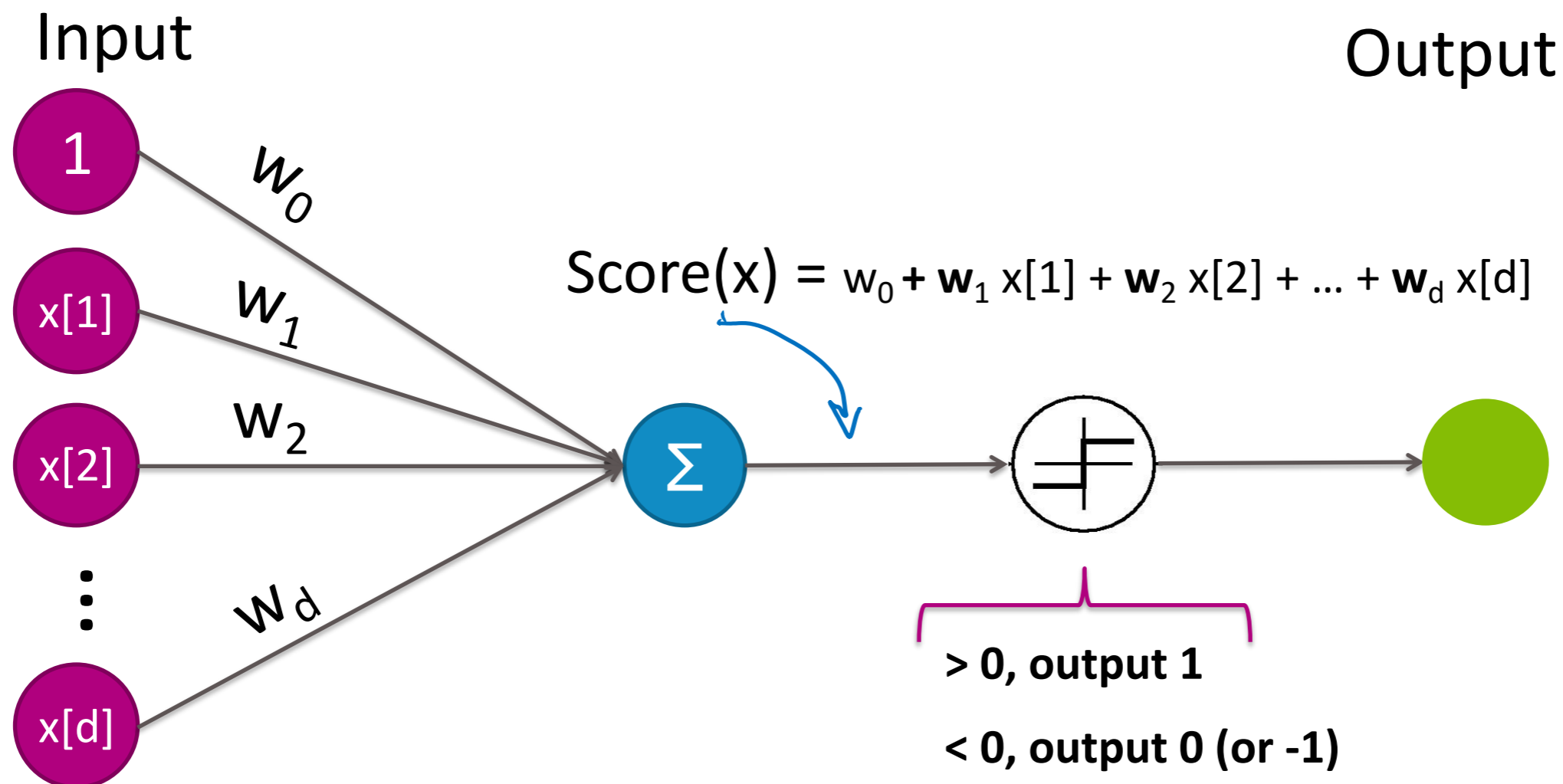
Score(x) **> 0**

Score(x) **< 0**

$w_0 + \mathbf{w}_1 x[1] + \mathbf{w}_2 x[2] + \dots + \mathbf{w}_d x[d] = 0$

# Graph representation of classifier: useful for defining neural networks
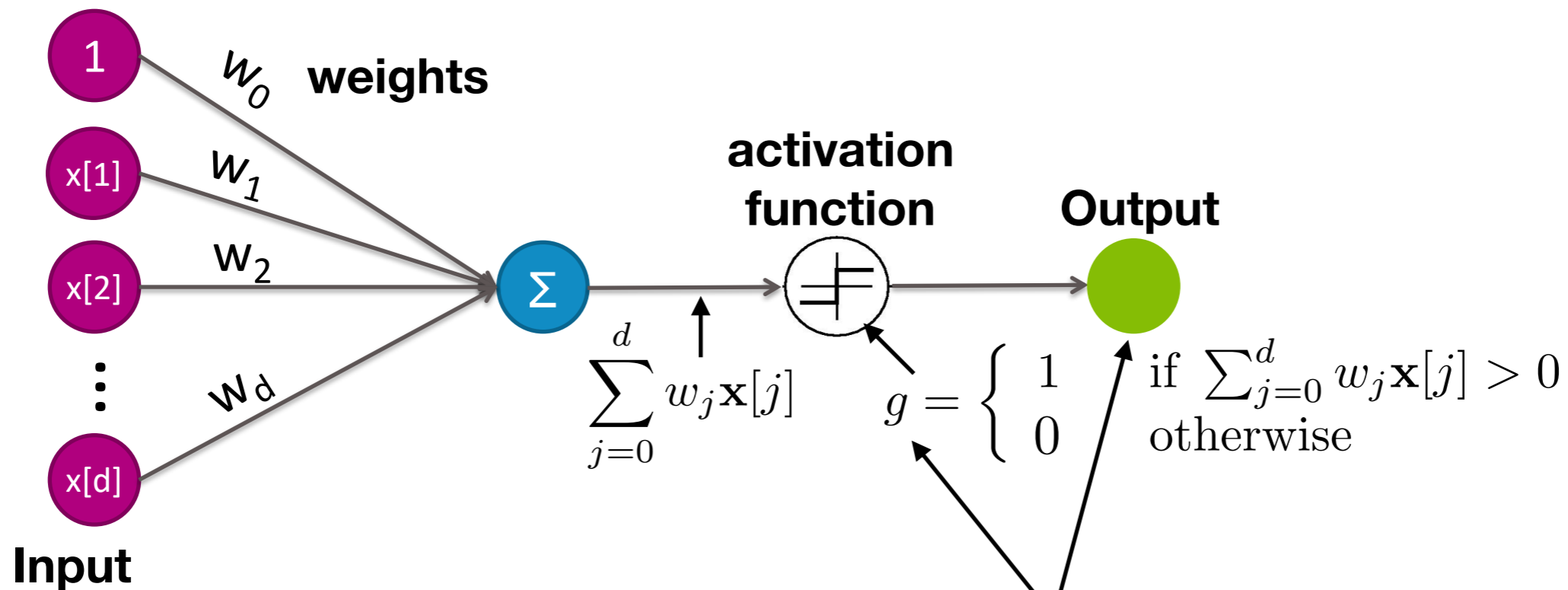
- We study an alternative representation of a linear classifier
- This graphical representation paves the way for designing deep neural networks
- This allow one to compactly represent a function (as a composition of many simple operations)

Input

Output

1

$w_0$

x[1]

$w_1$

Score(x) = $w_0$ + $w_1$ x[1] + $w_2$ x[2] + … + $w_d$ x[d]

$w_2$

x[2]

$\Sigma$

⋮

$w_d$

x[d]

> 0, output 1

< 0, output 0 (or -1)

# Single-layer neural network

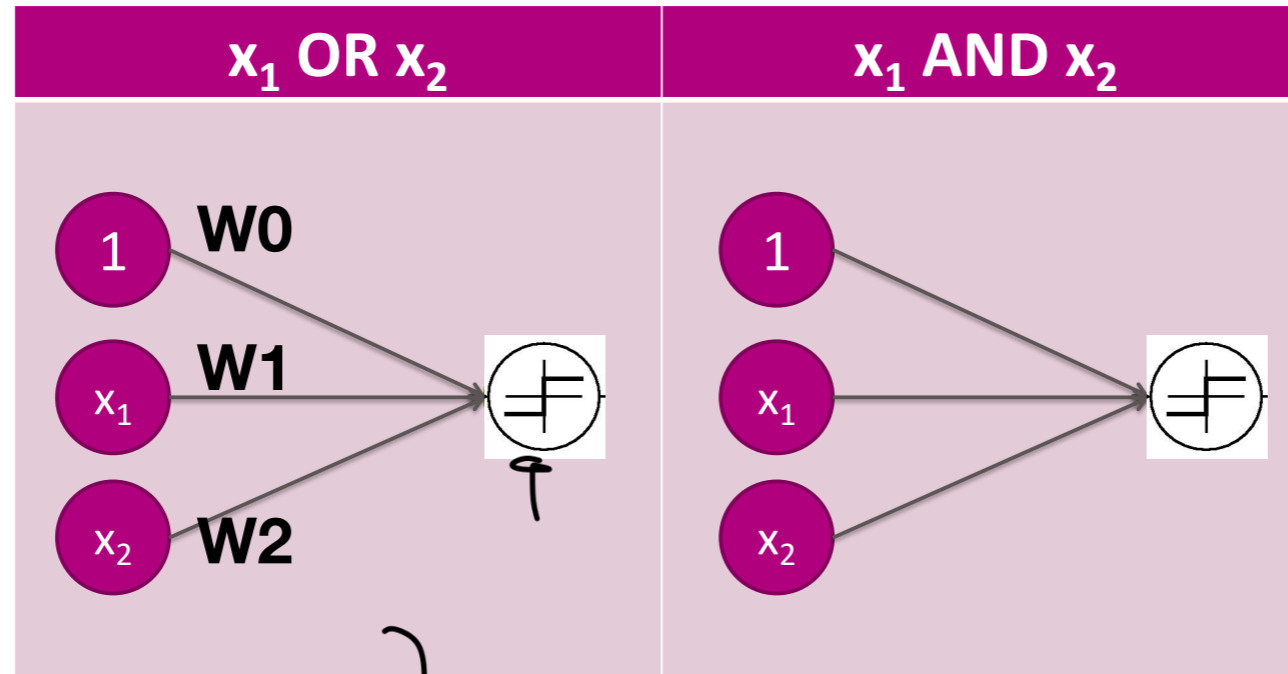This is a **single**-layer and **one-neuron** neural network

$$f(x) = \text{sign}(w_0 + w_1 x[1] + \cdots w_d x[d])$$



- This is a generic formula for one **neuron**:
  - input: 1,x[1],..,x[d]
  - We take weighted sum
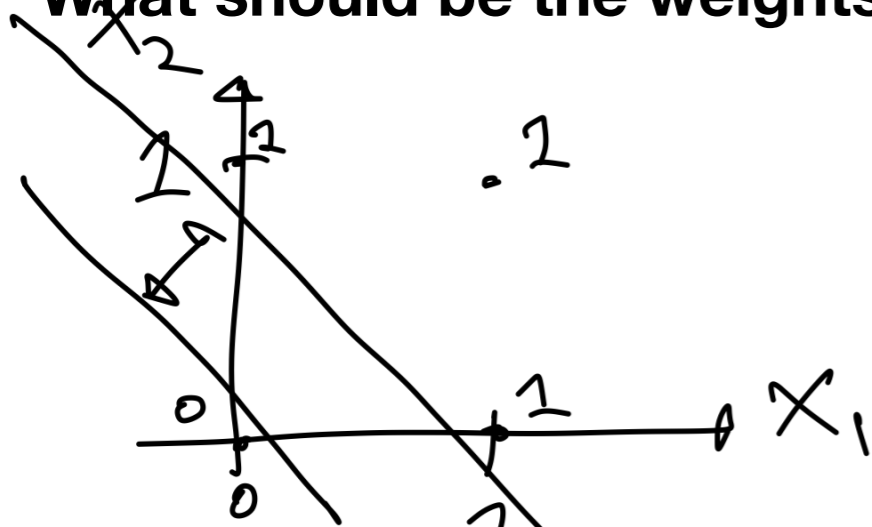  - And pass it through an **activation function** *g()*

# What can be represented by a linear classifier?

- x1 x2 y
- 0 0 0
- 0 1 1
- 1 0 1
- 1 1 1

| $x_1$ OR $x_2$ | $x_1$ AND $x_2$ |
|---|---|



- x1 x2 y
- 0 0 0
- 0 1 0
- 1 0 0
- 1 1 1

**What should be the weights?**

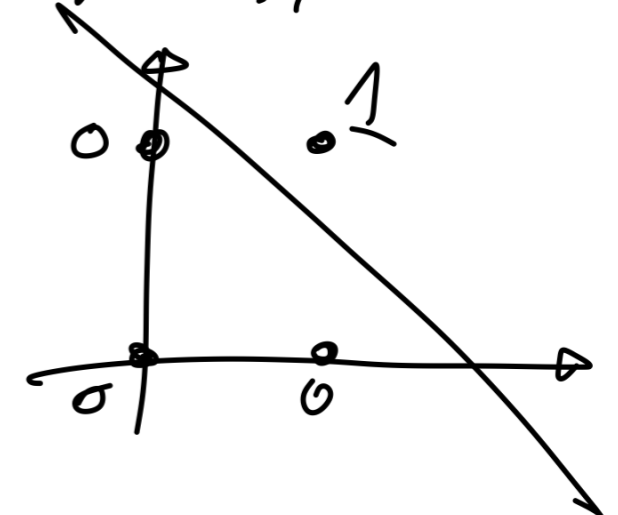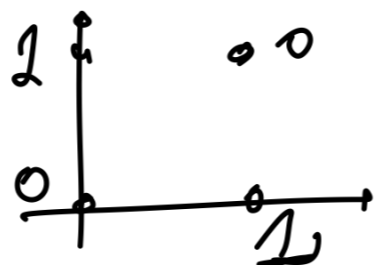$$\text{sign}(W_0 + W_1 X_1 + W_2 X_2) = Y$$

$-0.5 + 0.5X_1 + 0.5X_2 \rightarrow OR$

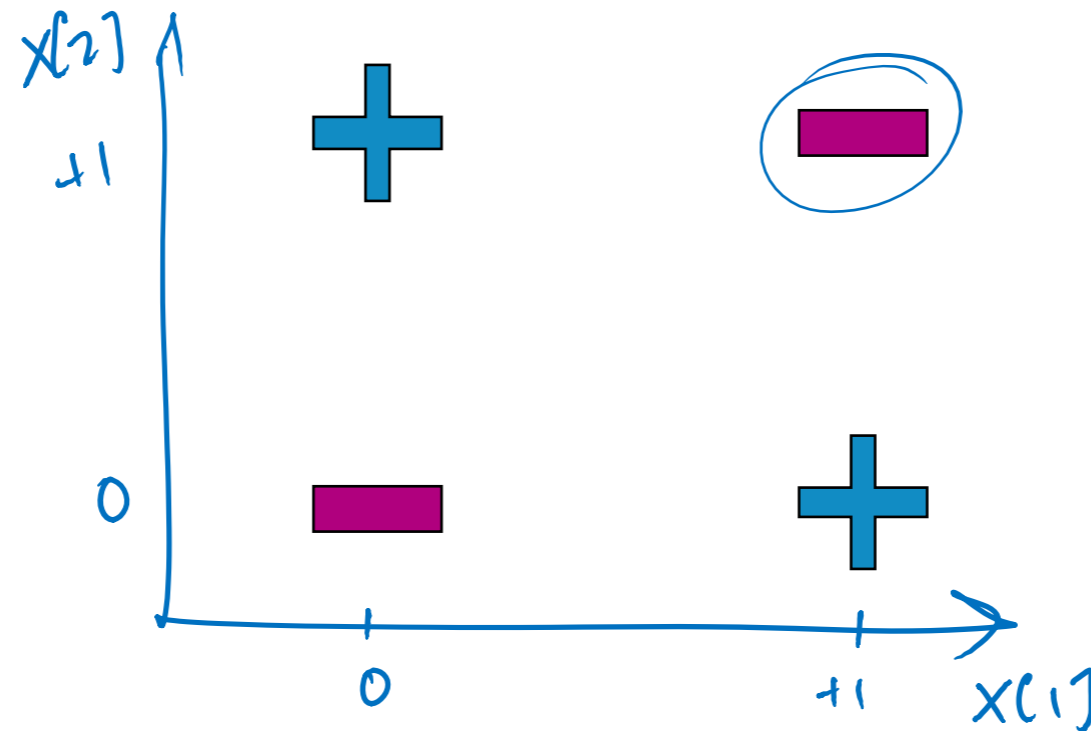$+0.5 + 0.5X_1 + 0.5X_2 \rightarrow AND$

**Note that there is a one-to-one correspondence between a linear classifier and a neural network of the above form**
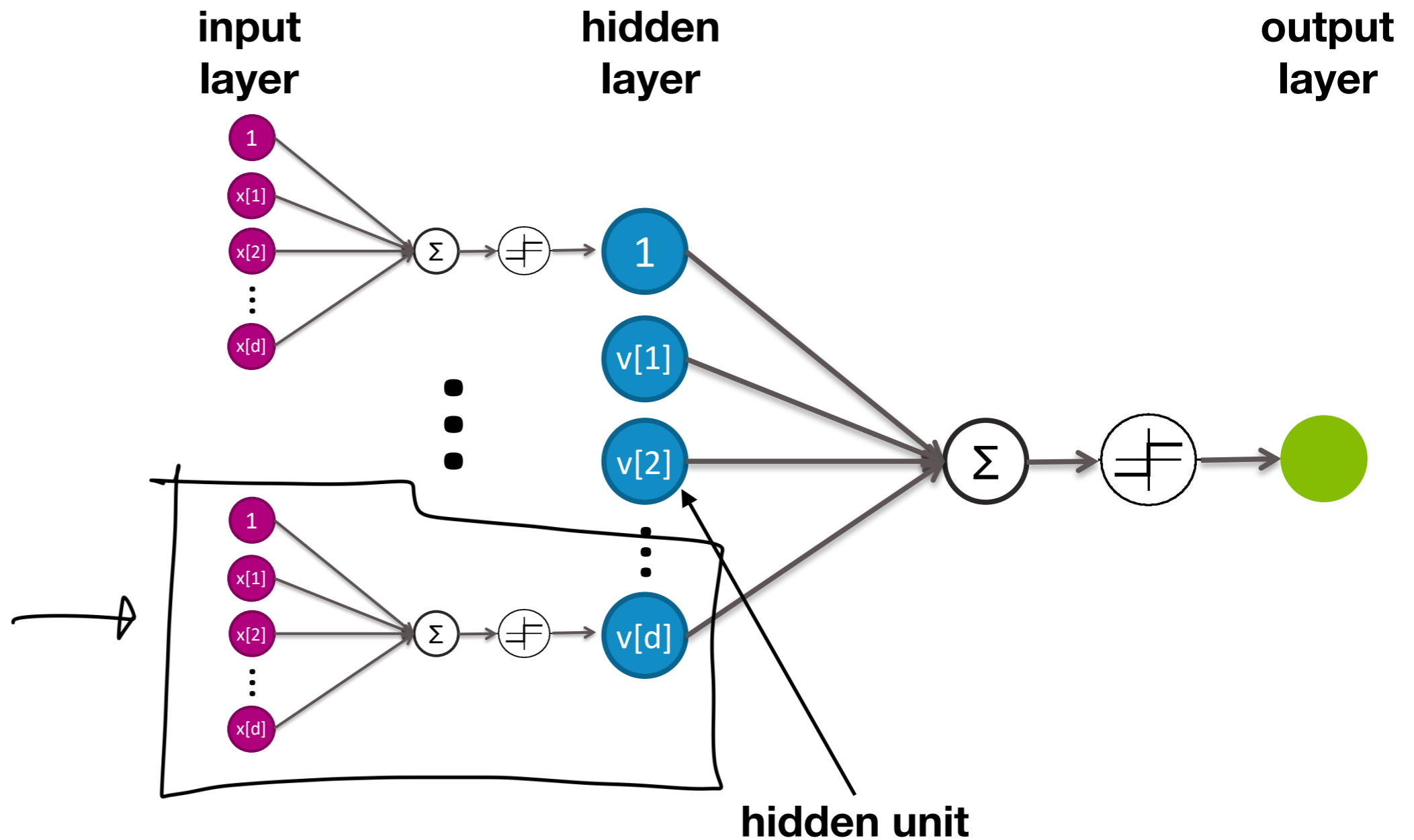
**What cannot be learned?**

# How can we get higher representation power?



- How can we build upon the single-layer, one-neuron function, to get a class of functions that can represent more complex functions?

# Hidden layer

- We compose neurons to create a network of neurons
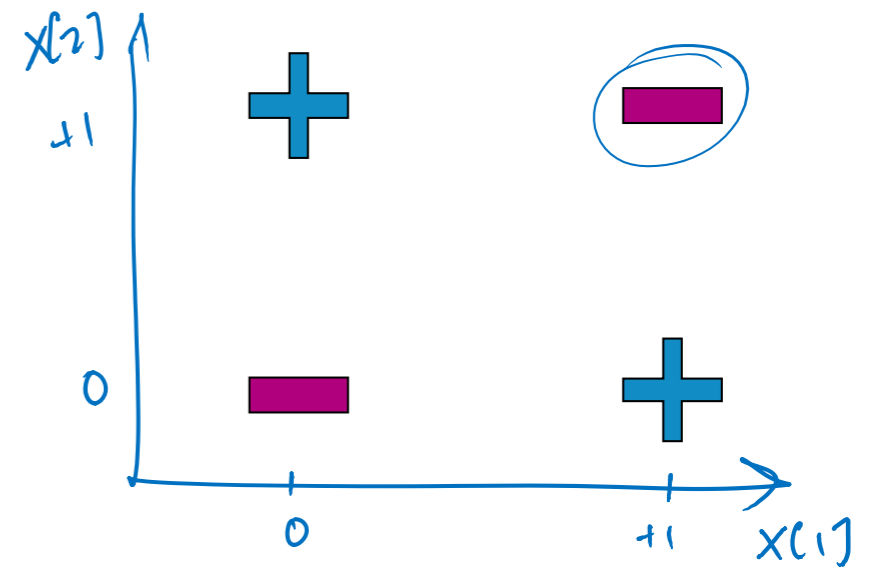  -> neural network



$$h_2(x) = \text{sign}(\underbrace{w_{20} + w_{21}x[1] + \cdots w_{2d}x[d]}_{w_2^T x})$$

$$f(x) = \text{sign}\Big( (w^{(2)})^T \underbrace{\text{sign}\big((W^{(1)})^T x\big)}_{=h(x)} \Big)$$

# Example: XOR function

XOR = x[1] AND NOT x[2]  OR  NOT x[1] AND x[2]

v[1]                v[2]

# XOR as a 2-layer neural network

$y = x[1]$ XOR $x[2]$ $= (x[1]$ AND $\neg x[2])$ OR $(x[2]$ AND $\neg x[1])$

$v[1] = (x[1]$ AND $\neg x[2])$

$\quad = g(-0.5 + x[1] - x[2])$

$v[2] = (x[2]$ AND $\neg x[1])$

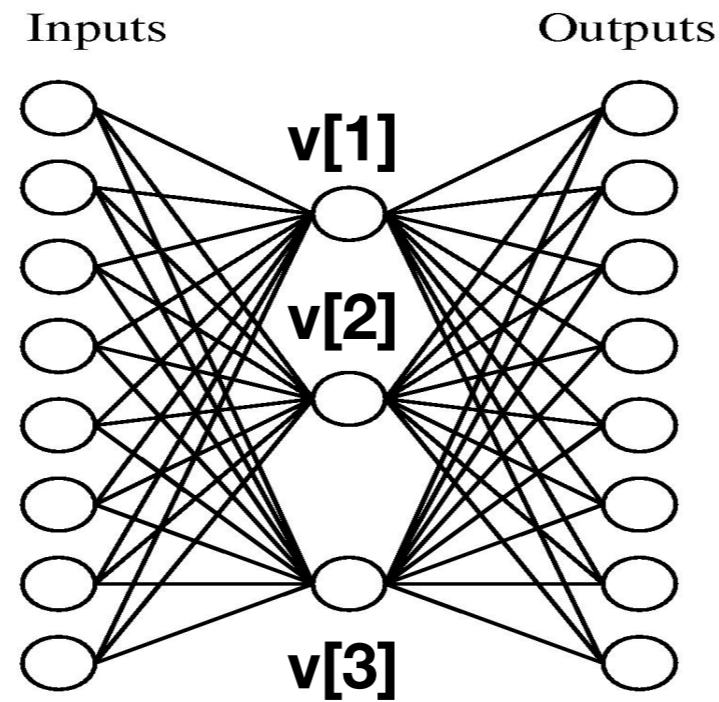$\quad = g(-0.5 + x[2] - x[1])$

$y = v[1]$ OR $v[2]$

$\quad = g(-0.5 + v[1] + v[2])$

# Two-layer neural network ( = one-hidden layer neural network)

Inputs       Outputs

**v[1]**

**v[2]**

**v[3]**

Single unit:

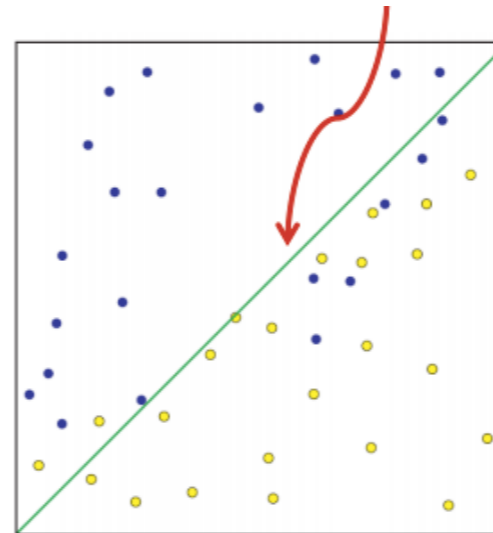$$out(\mathbf{x}) = g(w_0 + \sum_j w_j \mathbf{x}[j])$$

1-hidden layer:

$$out(\mathbf{x}) = g(w_0 + \sum_k w_k \boxed{g(w_0^k + \sum_j w_j^k \mathbf{x}[j])})$$

$$\mathbf{v}[k]$$

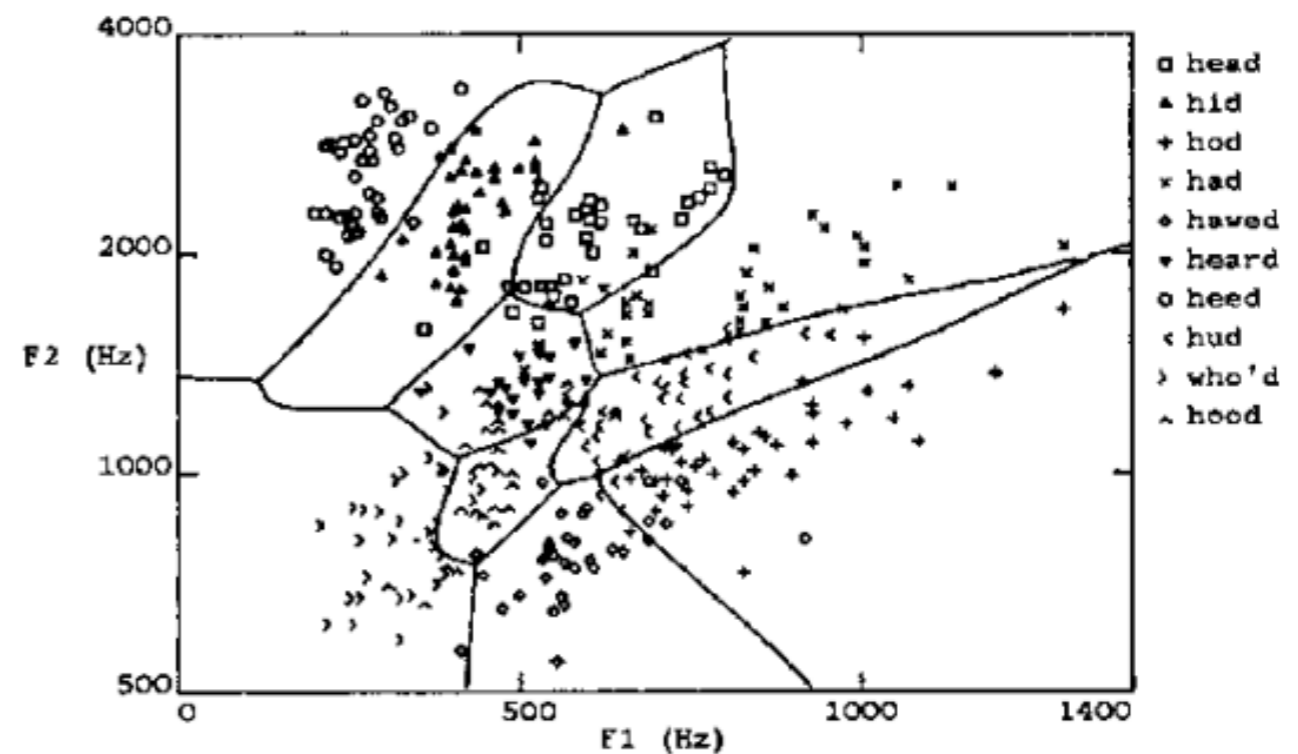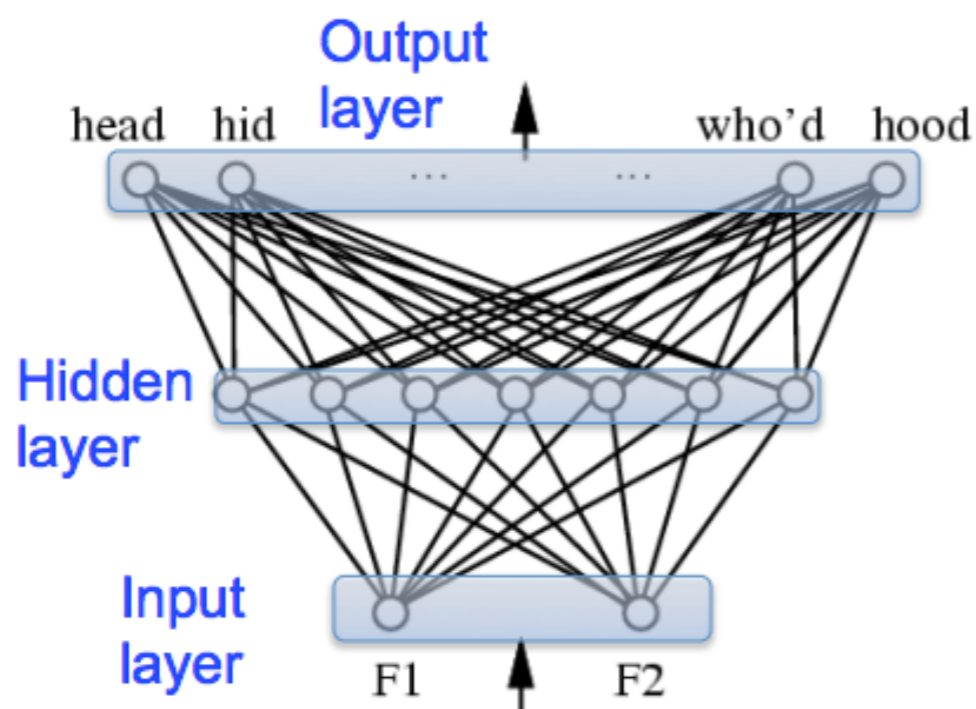# Example of 2-layer neural network in action

**Linear decision boundary**



1-layer neural networks
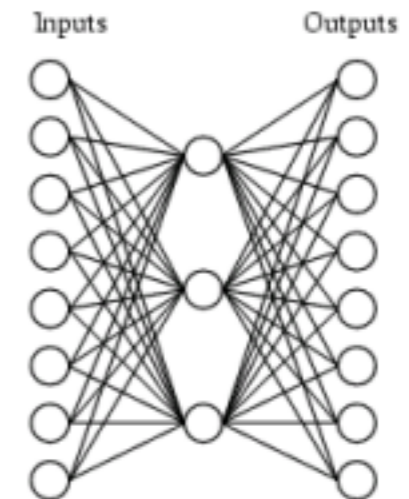only represents linear classifiers

Example: 2-layer neural network trained to distinguish vowel sounds using 2 formants (features)

a highly non-linear decision boundary can be learned from 2-layer neural networks

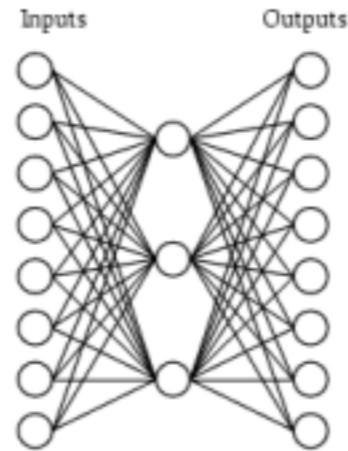# Representation power of a 2-layer neural network

- Can such function be learned?
- If we are manually designing functions, then 3 hidden layer is enough.
- The reason is that there is some simplicity or pattern in the data that we want to represent: it only has basis vectors!

Inputs          Outputs

A target function:

| Input | | Output |
|---|---|---|
| 10000000 | → | 10000000 |
| 01000000 | → | 01000000 |
| 00100000 | → | 00100000 |
| 00010000 | → | 00010000 |
| 00001000 | → | 00001000 |
| 00000100 | → | 00000100 |
| 00000010 | → | 00000010 |
| 00000001 | → | 00000001 |

A network:



Learned hidden layer representation:

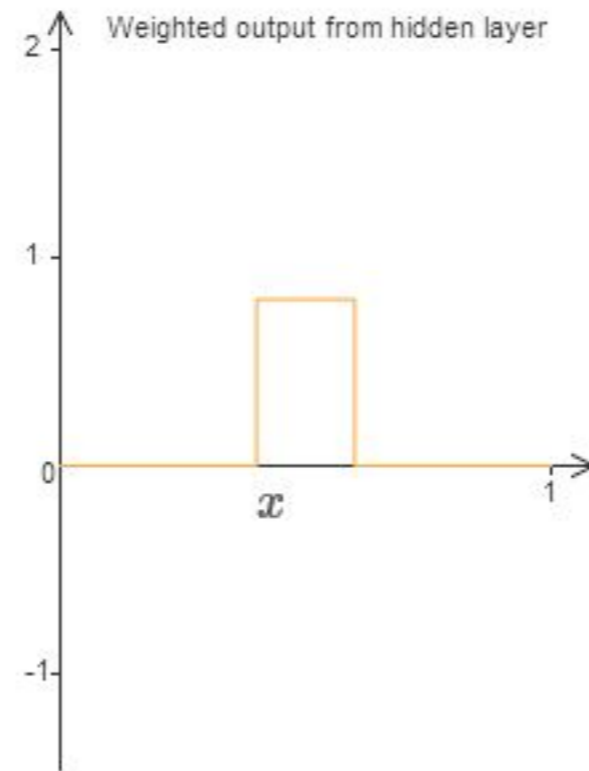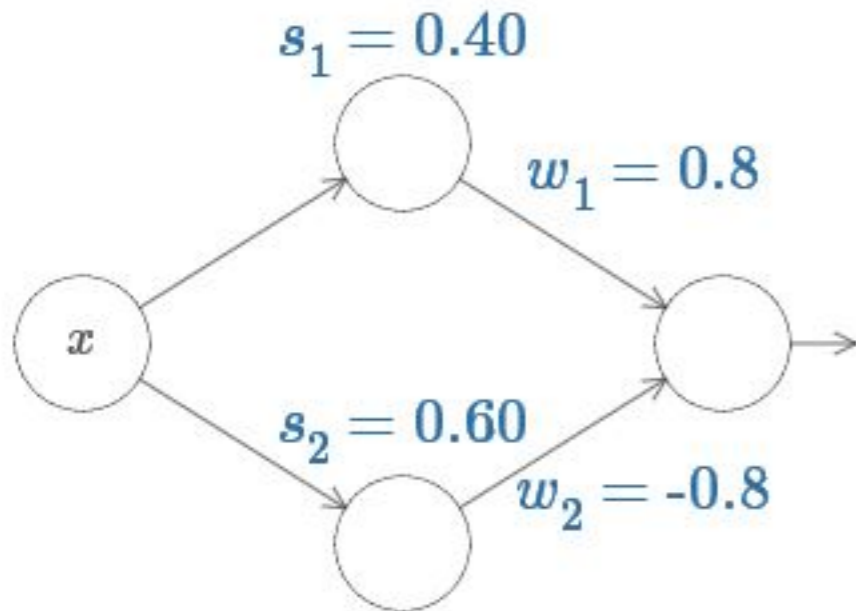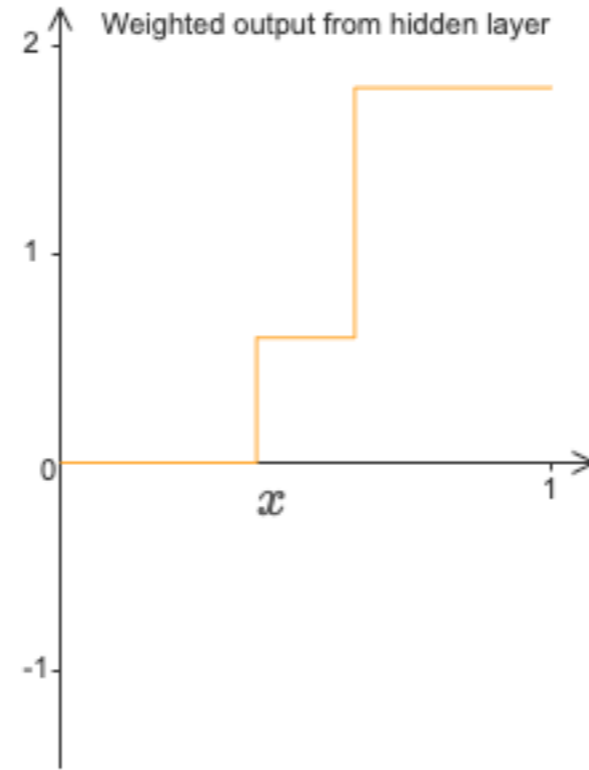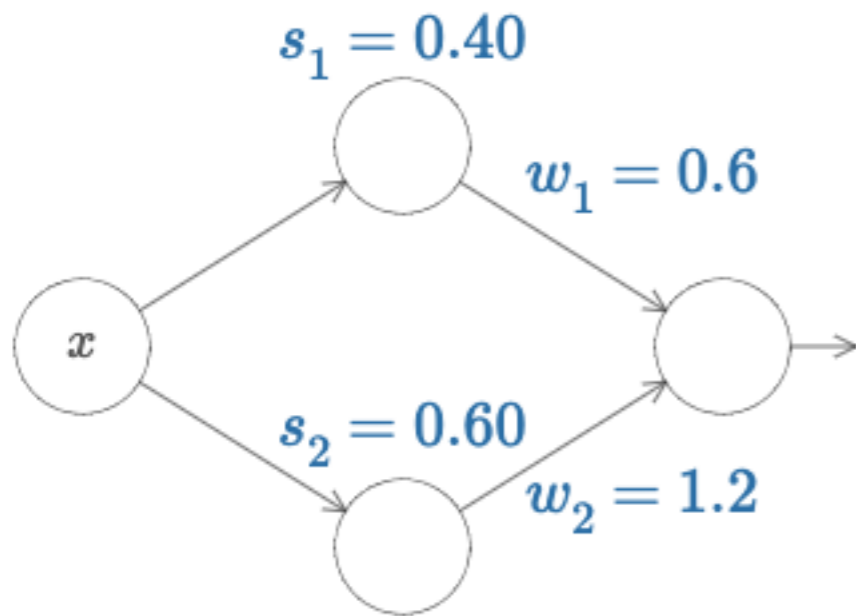| Input | Hidden Values | | | Output |
|---|---|---|---|---|
| 10000000 → | .89 | .04 | .08 | → 10000000 |
| 01000000 → | .01 | .11 | .88 | → 01000000 |
| 00100000 → | .01 | .97 | .27 | → 00100000 |
| 00010000 → | .99 | .97 | .71 | → 00010000 |
| 00001000 → | .03 | .05 | .02 | → 00001000 |
| 00000100 → | .22 | .99 | .99 | → 00000100 |
| 00000010 → | .80 | .01 | .98 | → 00000010 |
| 00000001 → | .60 | .94 | .01 | → 00000001 |

# A 2-layer neural network can represent any function, if we allow enough units in the hidden layer

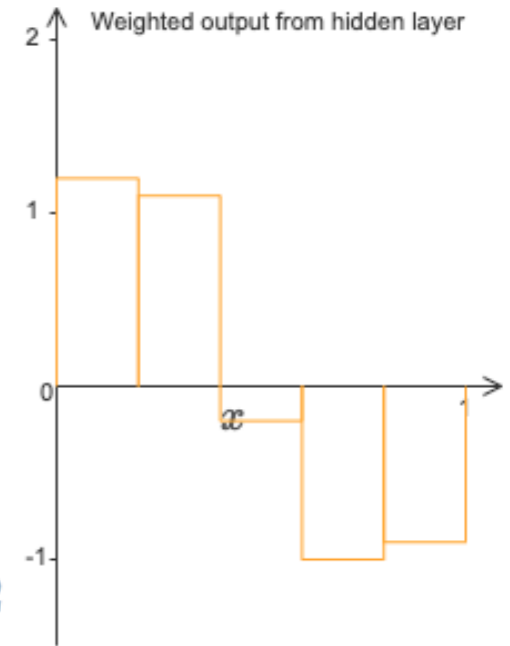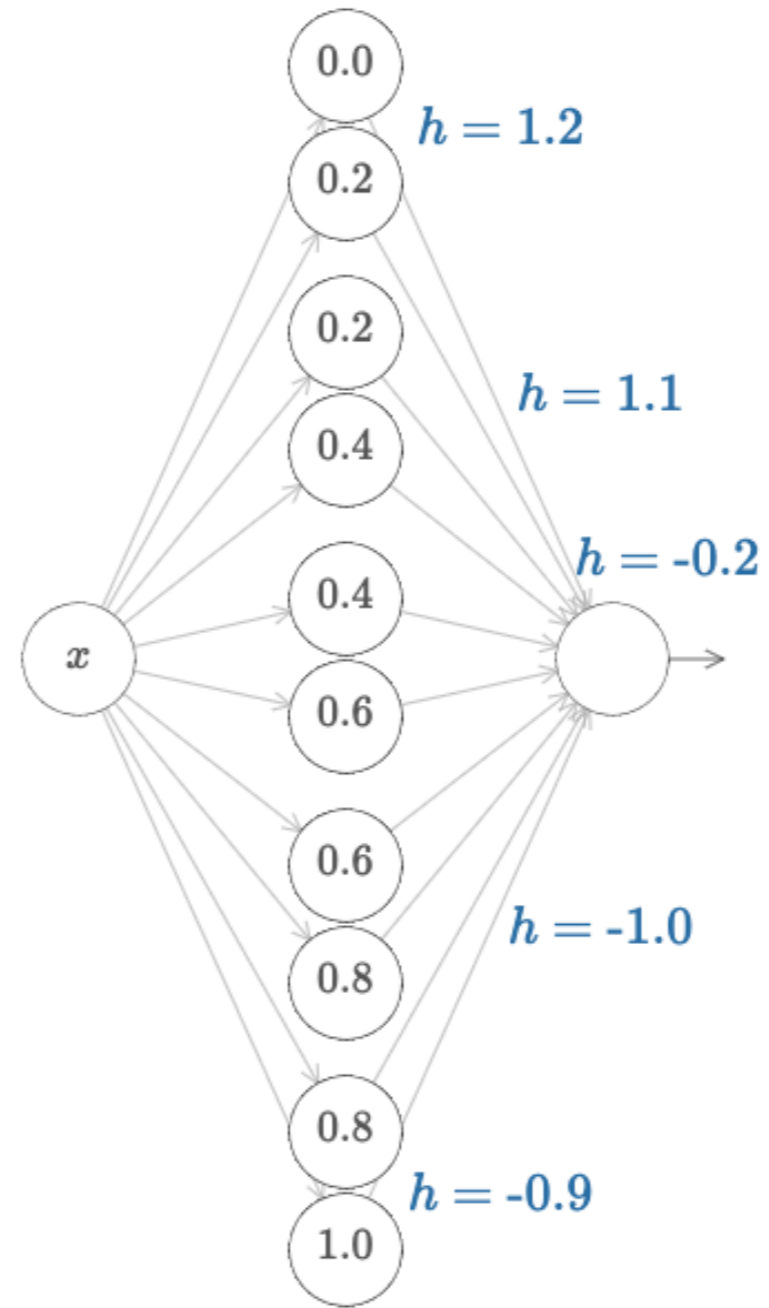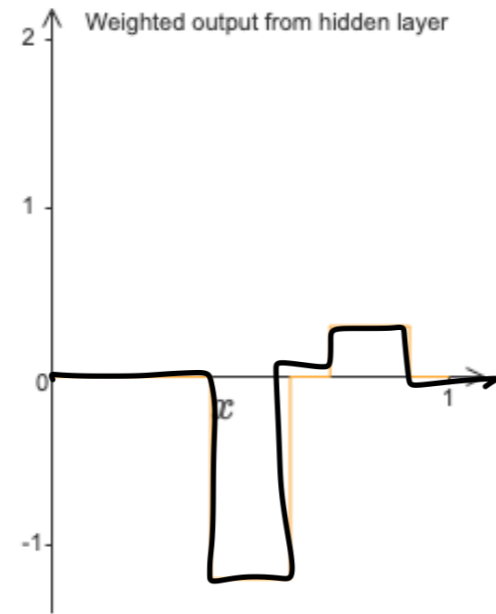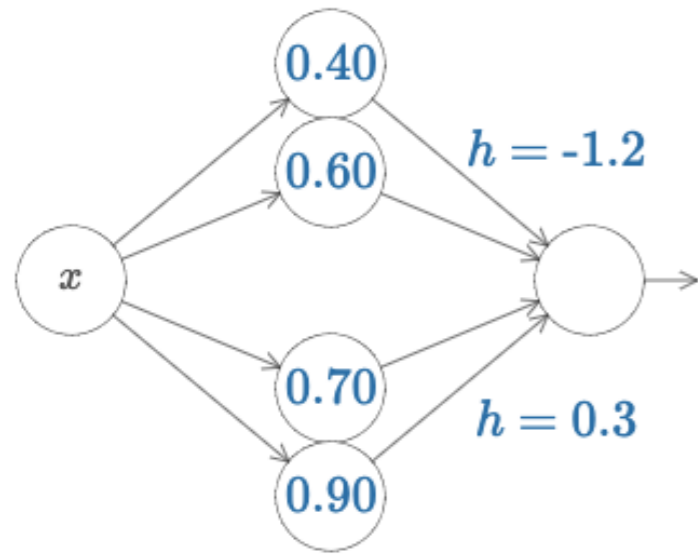**One-dimensional input/output example for illustration**



- We can compose step functions
  to approximate piece constant functions
  and use them to approximate any function

- More pieces (more hidden units) give better approximation

- demo: http://neuralnetworksanddeeplearning.com/chap4.html
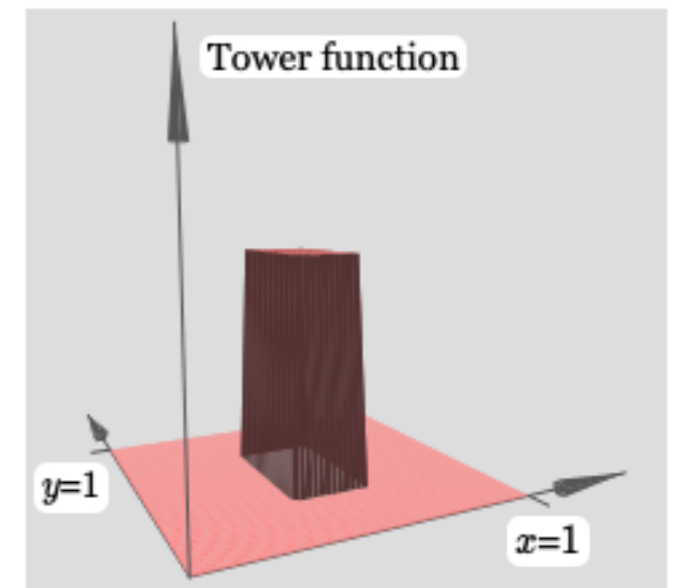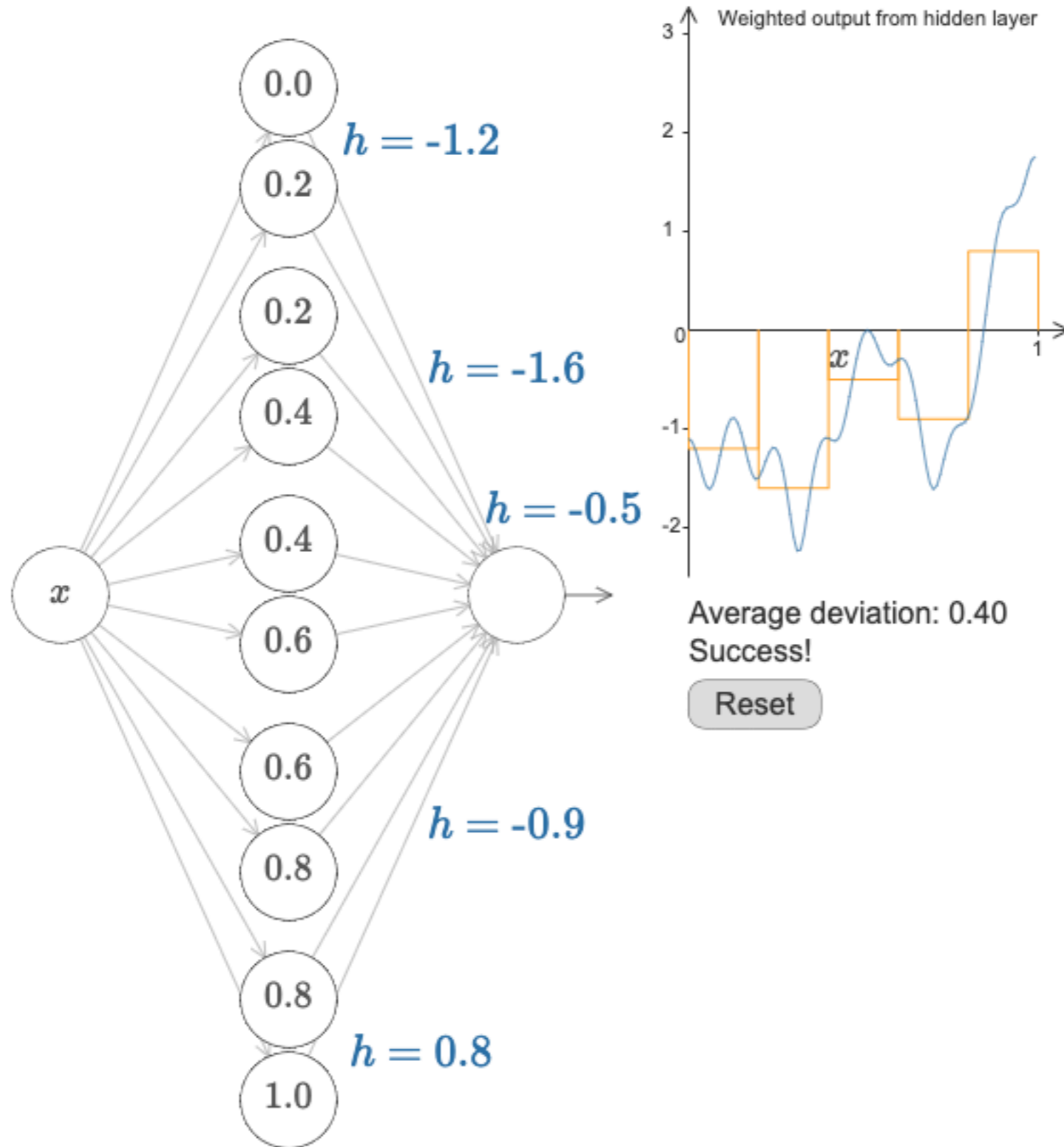
# Example



$s_1 = 0.40$

$w_1 = 0.6$

$x$

$s_2 = 0.60$

$w_2 = 1.2$

Weighted output from hidden layer

$s_1 = 0.40$

$w_1 = 0.8$

$x$

$s_2 = 0.60$

$w_2 = -0.8$

Weighted output from hidden layer

.8

.4    .6

-.8

# Example

# Example



$h = -1.2$

$h = -1.6$

$h = -0.5$

$h = -0.9$

$h = 0.8$

0.0
0.2
0.2
0.4
0.4
0.6
0.6
0.8
0.8
1.0

$x$

Weighted output from hidden layer

$x$

Average deviation: 0.40
Success!

Reset

Tower function

$y=1$

$x=1$
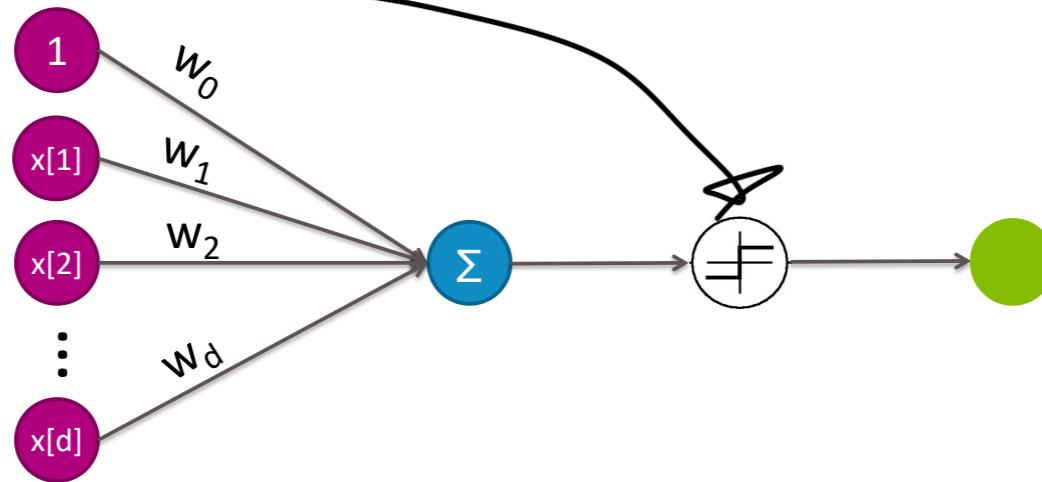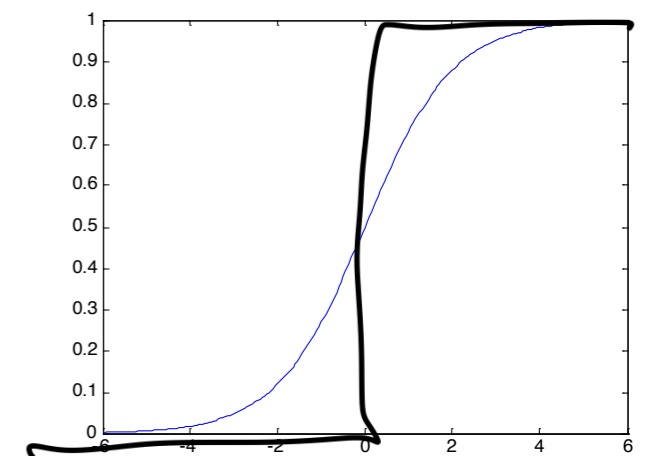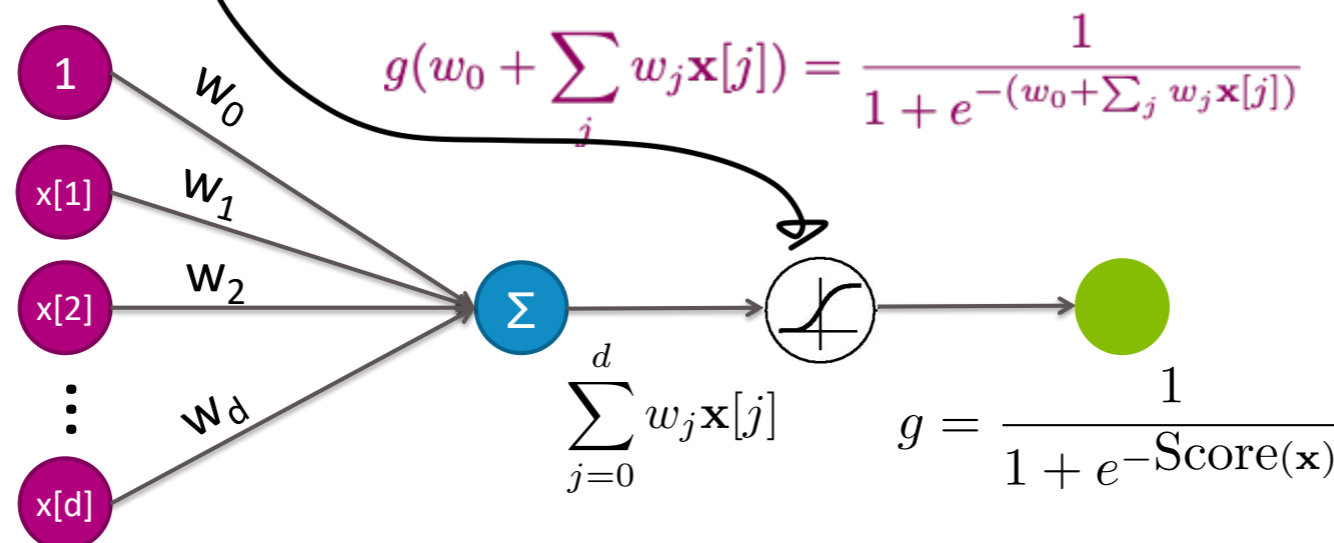
# General neural networks

- **Sign** activation function is never used in practice because the gradient is zero almost everywhere



- instead, **sigmoids** can be used because it is differentiable, and can approximate the sign function

$$g\left(w_0 + \sum_j w_j \mathbf{x}[j]\right) = \frac{1}{1 + e^{-(w_0 + \sum_j w_j \mathbf{x}[j])}}$$

$$\sum_{j=0}^{d} w_j \mathbf{x}[j] \qquad g = \frac{1}{1 + e^{-\mathrm{Score}(\mathbf{x})}}$$

20

# Activation functions

•Sigmoid
-Historically popular, but (mostly) fallen out of favor
•Neuron's activation saturates
(weights get very large -> gradients get small)
•Not zero-centered -> other issues in the gradient steps
-When put on the output layer, called "softmax" because
interpreted as class probability (soft assignment)

•Hyperbolic tangent  $g(x) = \tanh(x)$
-Saturates like sigmoid unit, but zero-centered

•Rectified linear unit (ReLU)  $g(x) = x^+ = \max(0,x)$
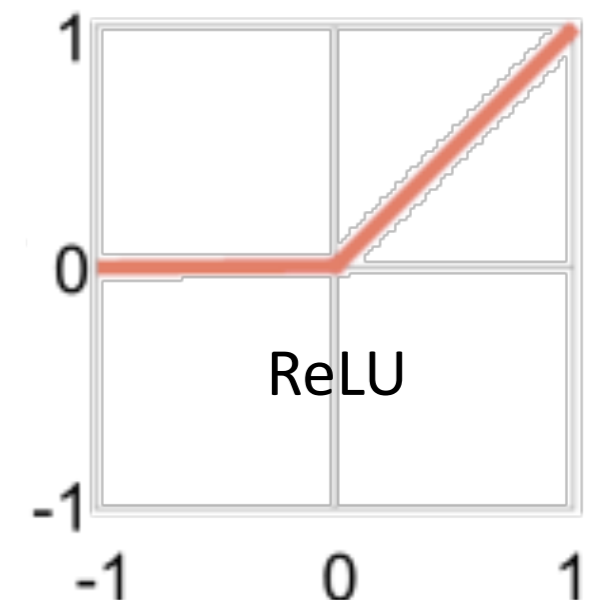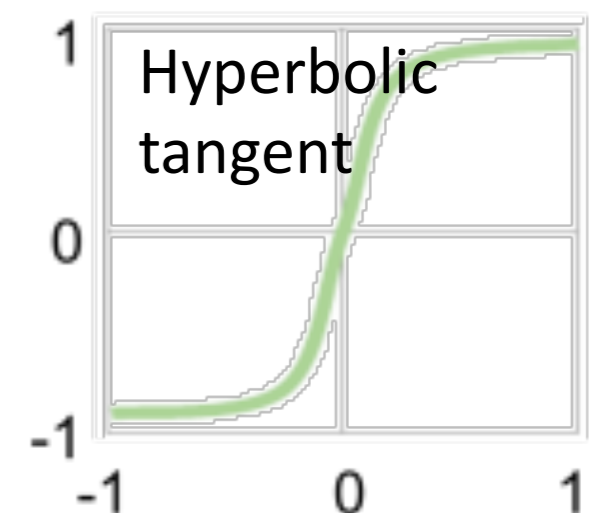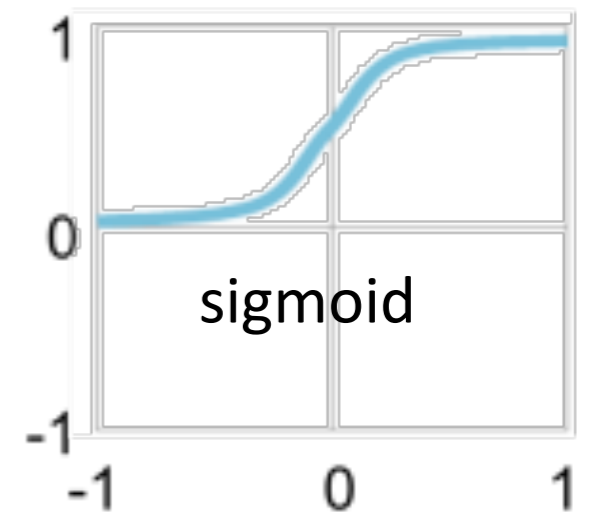-Most popular choice these days
-Fragile during training and neurons can "die off"…
be careful about learning rates
-"Noisy" or "leaky" variants

•Softplus $g(x) = \log(1+\exp(x))$
-Smooth approximation to rectifier activation



sigmoid



Hyperbolic tangent



ReLU

# General neural networks

- Layers and layers and layers of linear models and non-linear transformations

- Around for about 50 years
  - Fell in "disfavor" in 90s

- In last few years, big resurgence
  - Impressive accuracy on several benchmark problems
  - Powered by huge datasets, GPUs, & modeling/learning algorithm improvements

# Overfitting

Are NNs likely to overfit?
- *Yes*, they can represent arbitrary functions!!!

Avoiding overfitting?
- More training data
- Fewer hidden nodes / better topology
  - **Rule of thumb:** 3-layer NNs outperform 2-layer NNs, but going deeper rarely helps (different story for convolutional networks!)
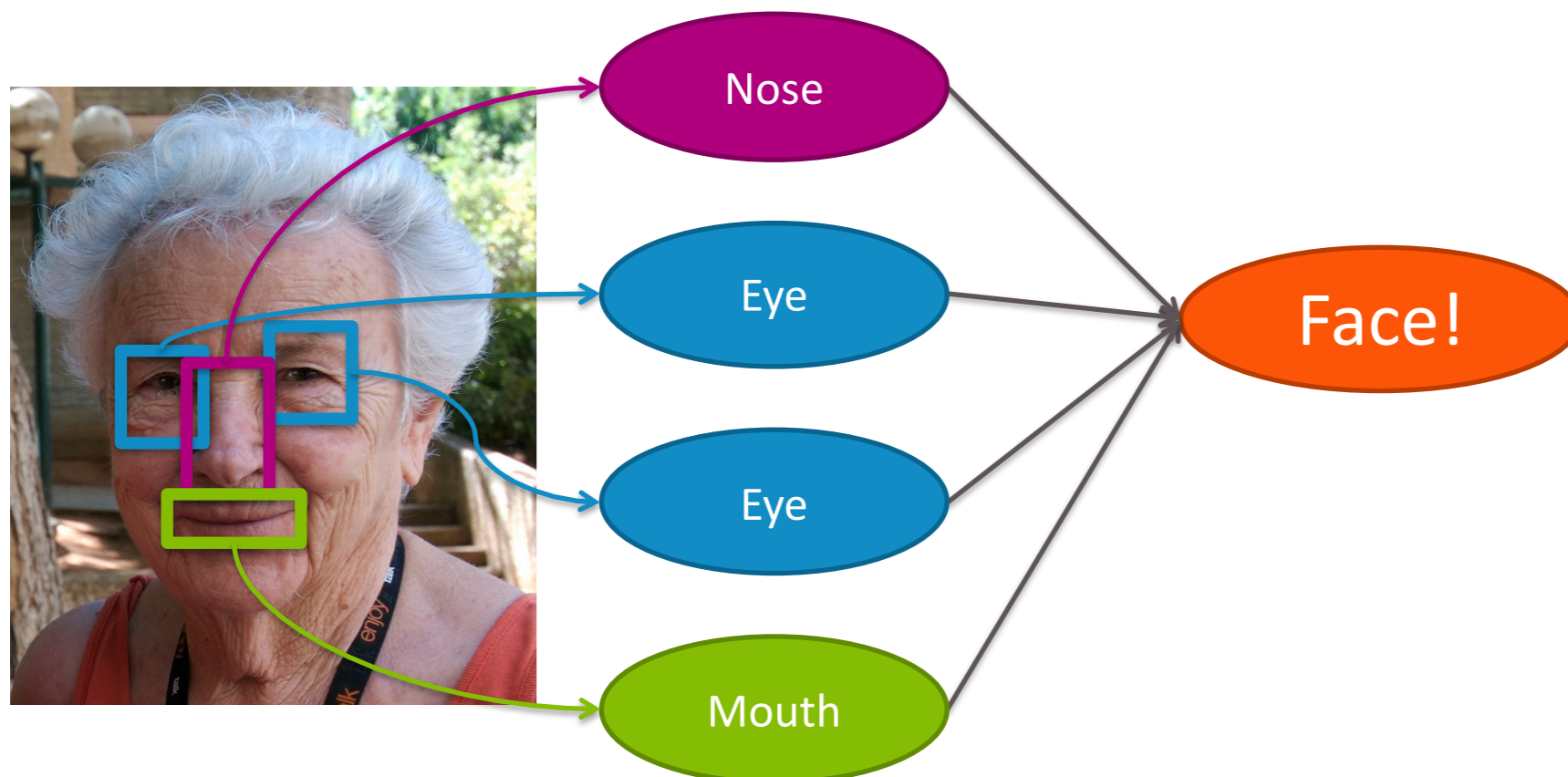- Regularization
- Early stopping

# Applications to vision problems

- Classical image processing manually extracts features

Features = local detectors
- Combined to make prediction
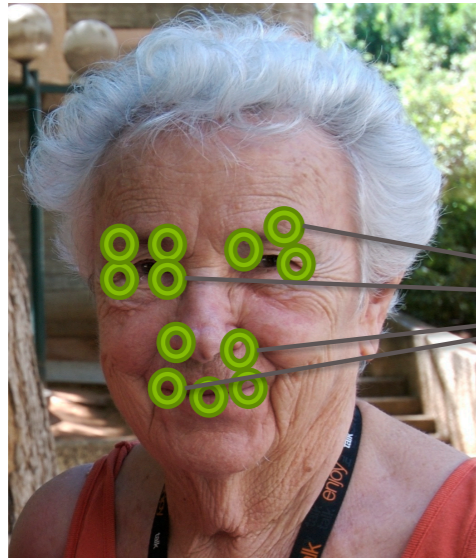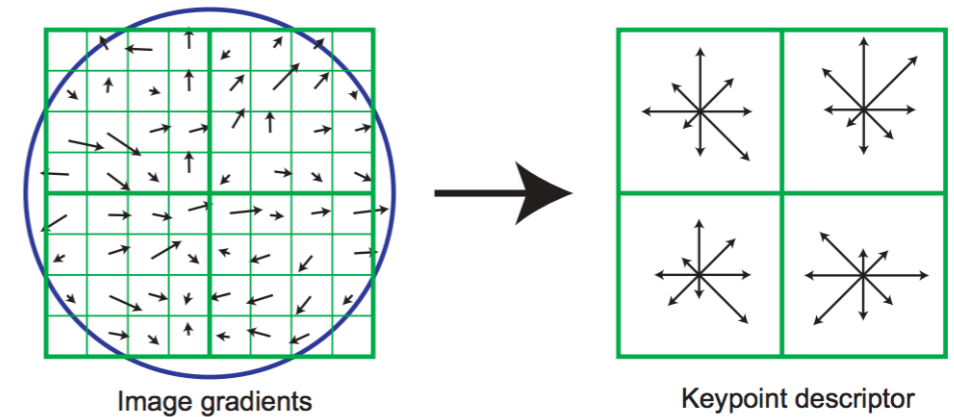- (in reality, features are more low-level)

Typical local detectors look for locally "interesting points" in image

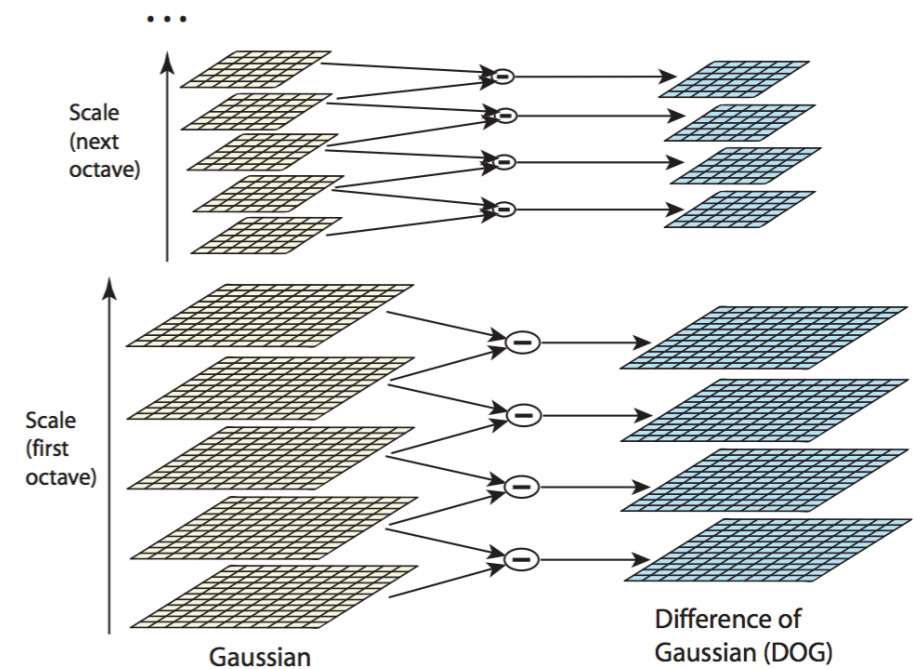*Image features:* collections of locally interesting points
−Combined to build classifiers



Face!



Image gradients

Keypoint descriptor

...

Many hand created features exist
 for finding interest points…

Scale
(next
octave)

Scale
(first
octave)
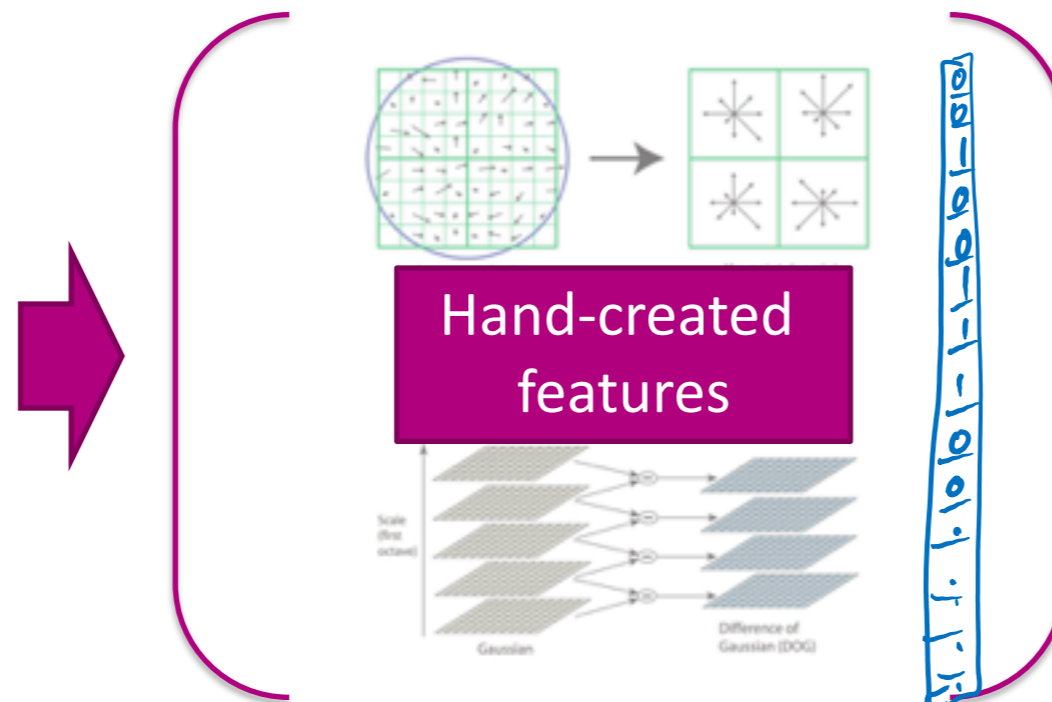
Gaussian

Difference of
Gaussian (DOG)

*SIFT* [Lowe '99]
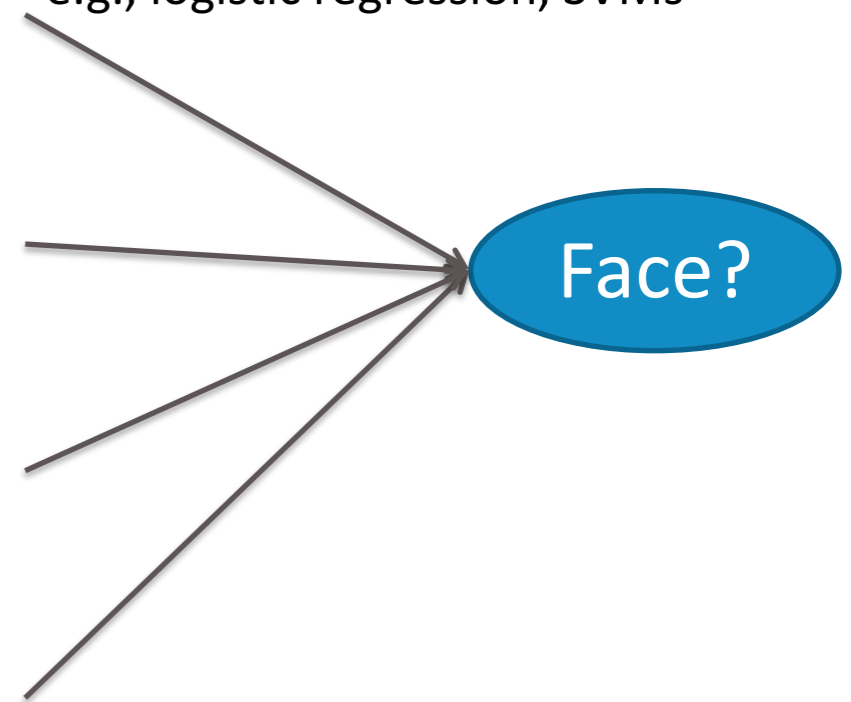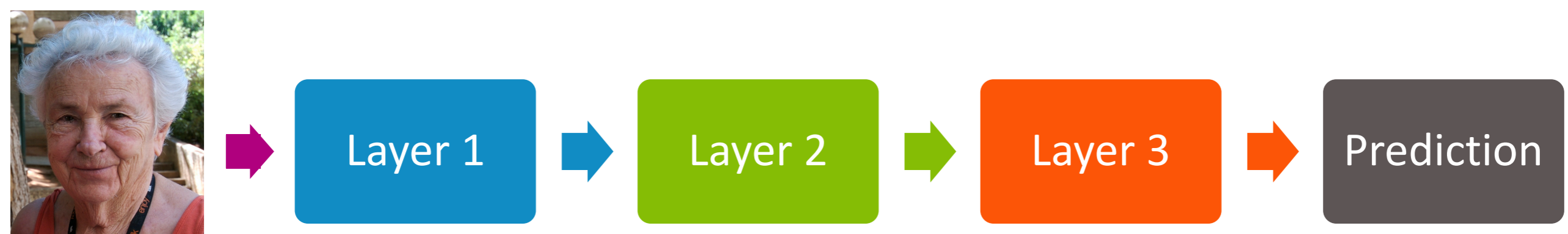
# Classical image classification

Input        Extract features        Use simple classifier
e.g., logistic regression, SVMs



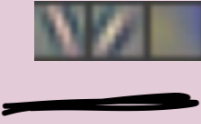Hand-created features

Face?

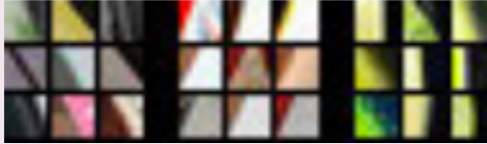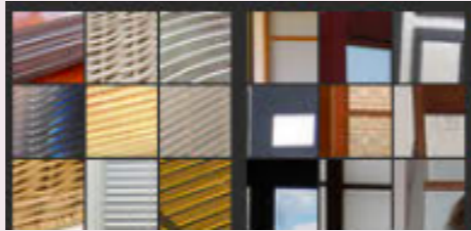- Critically relies on having good features manually chosen

# Instead, neural network (implicitly) discovers those features from data



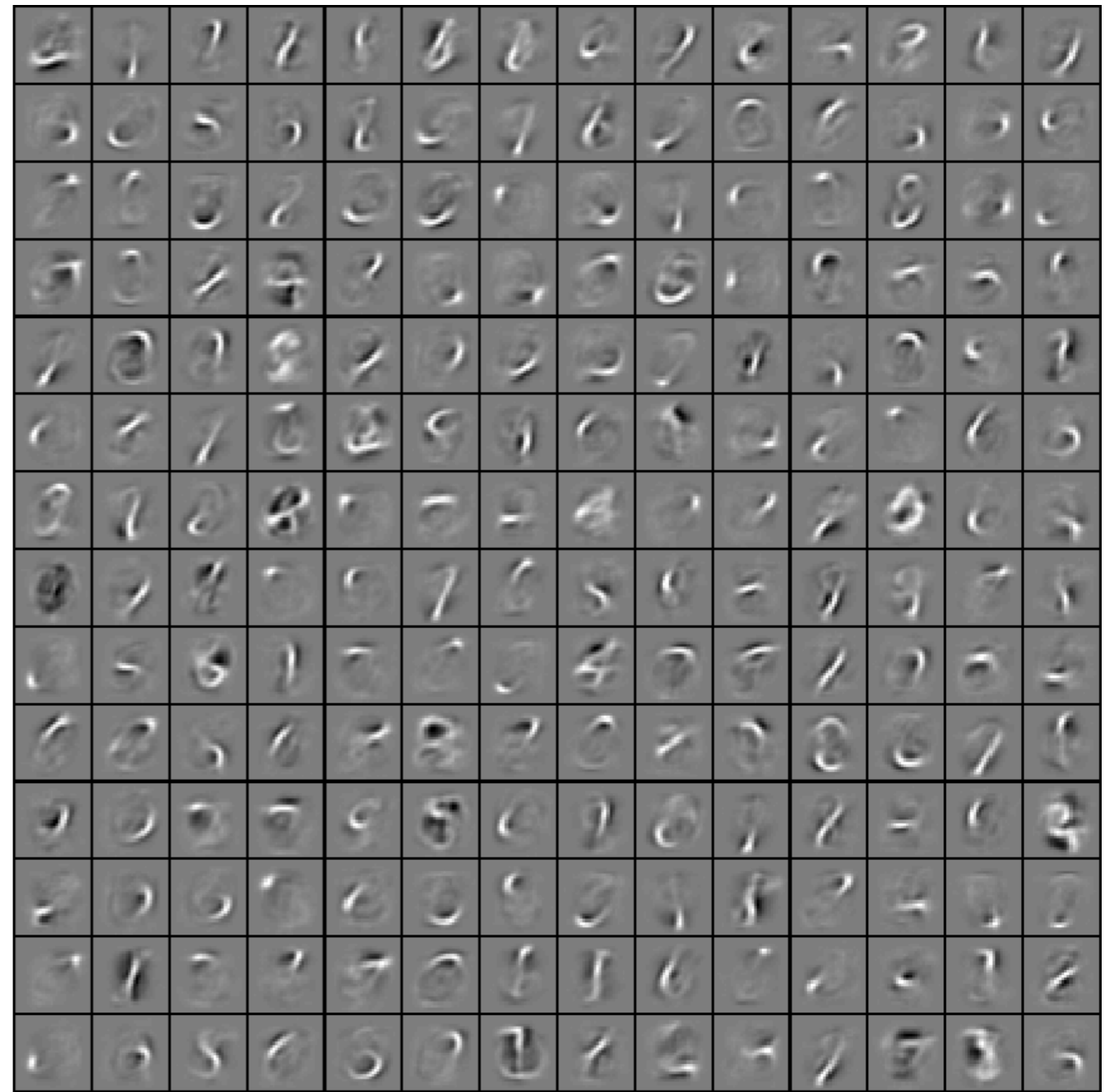| | Layer 1 | Layer 2 | Layer 3 | Prediction |
|---|---|---|---|---|
| Example detectors learned | | | | |
| Example interest points detected | | | | |

[Zeiler & Fergus '13]

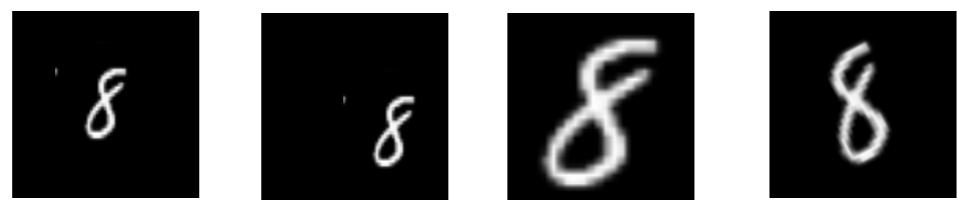- Each layer learns increasingly complex features, as we go higher in the layers
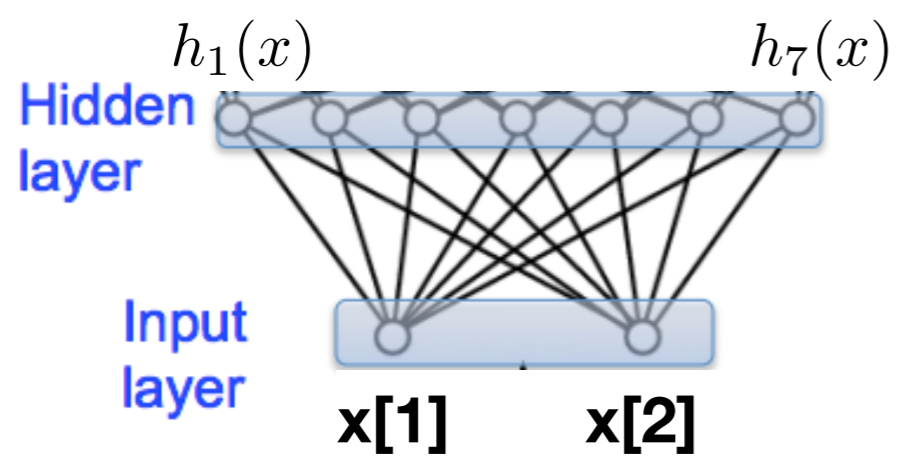
# Convolutional neural networks

# The challenge of applying regular neural networks (multilayer perceptrons) to images

- Interesting images are very high-dimensional
- And images have particular structures
  - Invariance to shift, scale, rotation



## Convolutional Neural Networks (CNN)

**Main building block of NN: fully connected layer**

**Main building block of CNN: convolutional layer**



$h_1(x)$     $h_7(x)$

Hidden layer

Input layer

x[1]    x[2]

$$h_i(x) = g(\, w_{i0} + w_{i1}x[1] + w_{i2}x[2]\,)$$

**Activation function of choice**

$h_{1,1,1}(x)$

x[1,1,1] →

Depth = # filters

x[32,32,3]    output volume

$h_{32,32,24}(x)$

**Both input and output are 3-dimensional arrays called tensors**

# Convolution

- Consider a one-dimensional signal (a sequence or a vector)
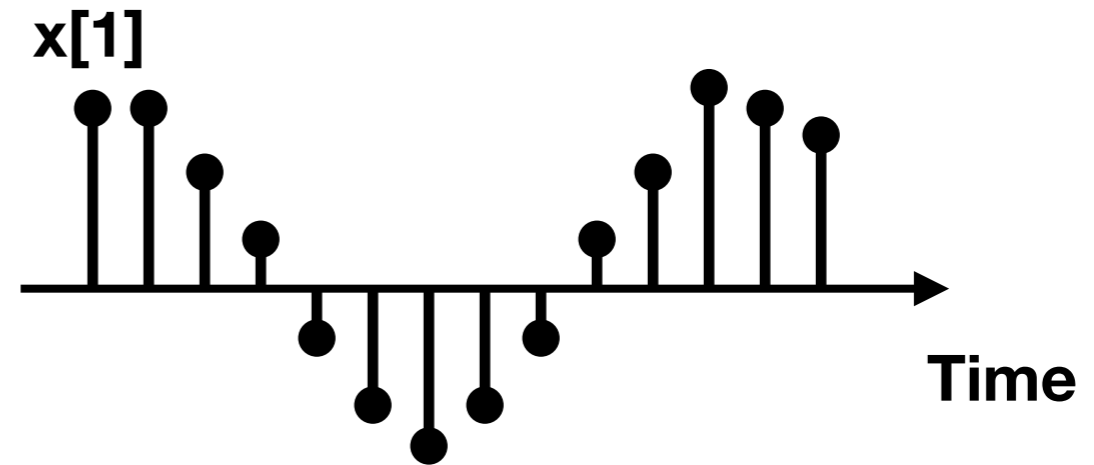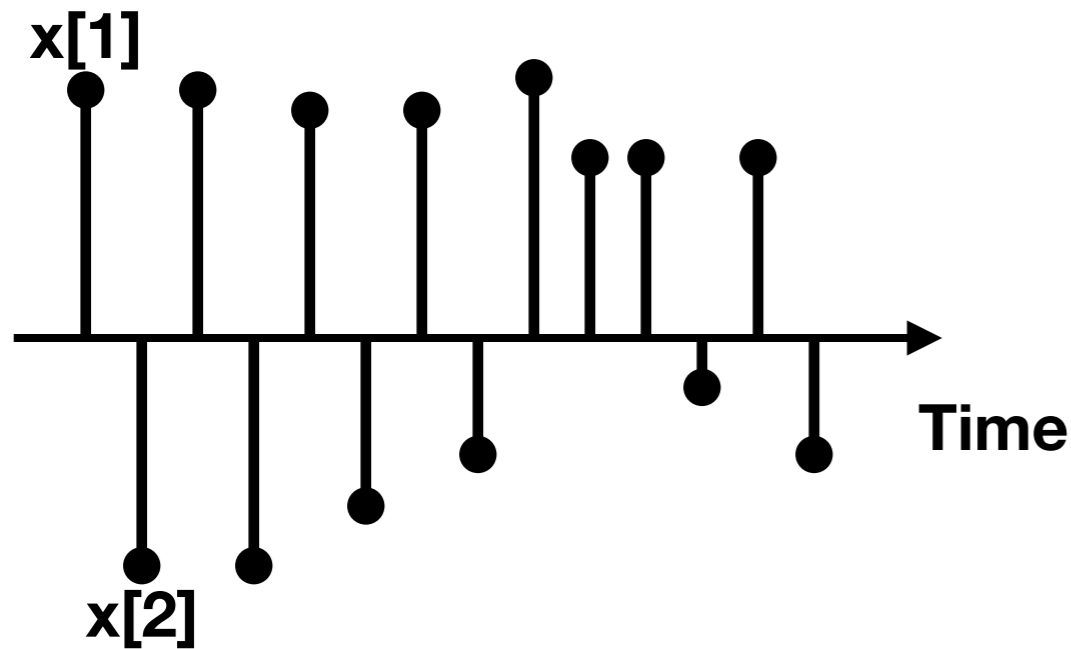  - For example, speech recognition



- A popular method for extracting features from a sequential data is convolution

# Convolution

- Consider the task of classifying whether the signal **x** is high pitch (high frequency) or low pitch (low frequency)



- We use a filter **w** and convolve **x** and **w** to get **x** • **w**

Example of length 2 filter

w=(w[1],w[2])

- Convolving high pass filter with a high pitch signal
- Slide the filter from left to right and compute the inner-product (entrywise product and sum)

high pitch signal $x$

high pass filter $w$

$x \bullet w$

high pitch signal $x$

high pass filter $w$

$x \bullet w$

- Convolution

high pitch signal $x$
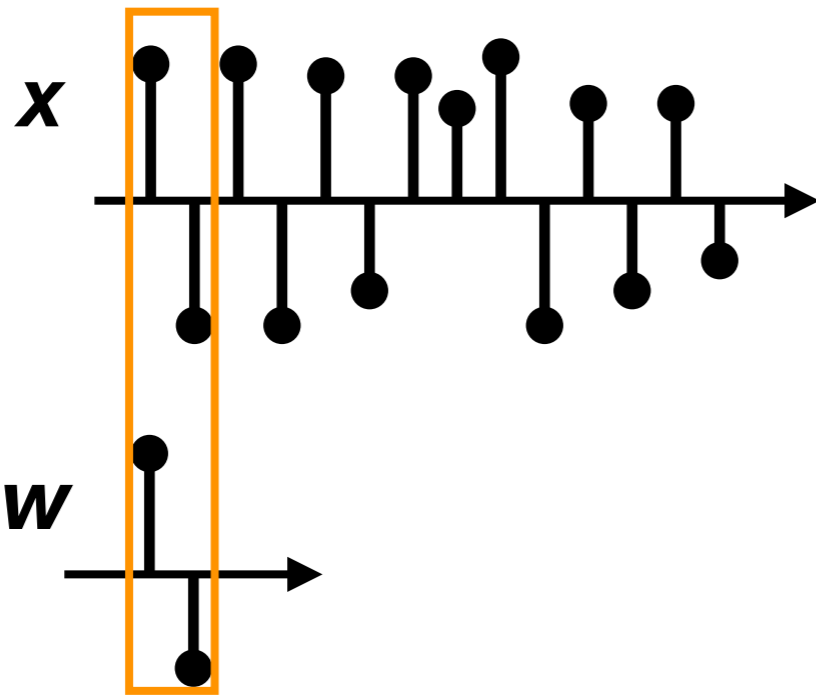
high pass filter $w$

$x \bullet w$

- How high frequency is $x$?
- **Pooling** operation aggregates the data
  - max-pooling: $max(x \bullet w)$
  - Average pooling: $(1/N)( |(x \bullet w)[1]| + \ldots + |(x \bullet w)[N]| )$

- Convolved and Pooled value will be large for high-frequency data

# high-freq. signal vs. low-freq. signal

# Two-dimensional convolution

- Consider a task of classifying 0's and 1's



- One manual way is to use some 2-d filters



**W1=**

| -1 | -1 | -1 | -1 | -1 |
|----|----|----|----|----|
| 1  | 1  | 1  | 1  | 1  |

**Convolve** → **Pooling**

**W2=**

| -1 | 1 |
|----|---|
| -1 | 1 |
| -1 | 1 |
| -1 | 1 |
| -1 | 1 |

**Convolve** → **Pooling**

**Declare "Zero" if small or "One" if large**

# Example of convolutional layer with 3x3 filter (9 parameters)

**To understand the convolution of 3-dimensional arrays (tensors), let's consider the convolution of 2-dimensional arrays (matrices)**

**Input image**

**Output image**



output volume

Depth = # filters

**In this example, we consider a 3x3 convolution**

**Sub-region in input image**

**Filter represented by a matrix $W$**
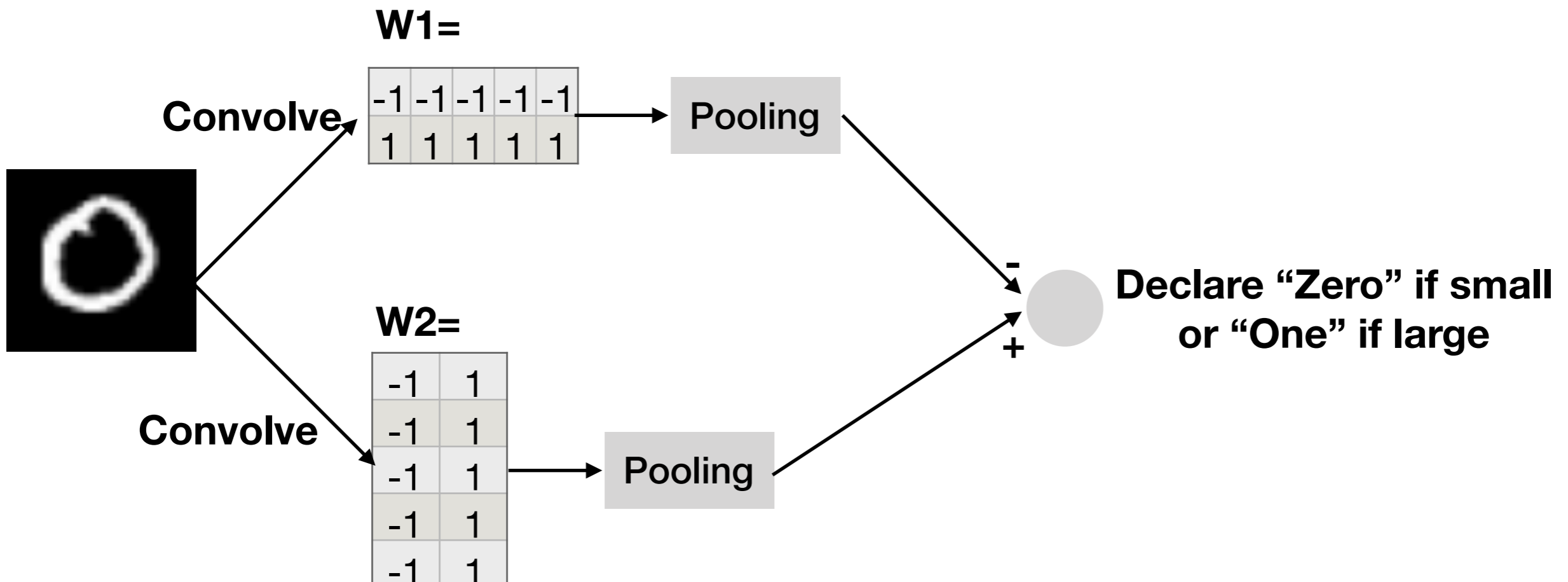
**Output Image (pixel)**

| 3 | 3 | 2 |
|---|---|---|
| 0 | 0 | 1 |
| 3 | 1 | 2 |

●

$W =$

| 0 | 1 | 2 |
|---|---|---|
| 2 | 2 | 0 |
| 0 | 1 | 2 |

➡

| 12 |
|----|

= inner product of two matrices of the same size

| 0 | 0 | 1 |
|---|---|---|
| 3 | 1 | 2 |
| 2 | 0 | 0 |

●

| 0 | 1 | 2 |
|---|---|---|
| 2 | 2 | 0 |
| 0 | 1 | 2 |

➡

**Key aspect of convolutional layer is that it applies the same filter (with the same weights) to all sub-regions [also known as weight sharing]**

➡ **This gives efficiency + shift invariance**

# Convolution

- What is the output image?

W = 
| 0 | 1 |
|---|---|
| 2 | 0 |

Filter (2x2)

X = 
| 1 | 1 | 0 |
|---|---|---|
| 0 | 1 | 1 |
| 1 | 1 | 1 |

Input Image (3x3)

Option A

| 0 | 2 | 2 | 0 |
|---|---|---|---|
| 1 | 1 | 2 | 2 |
| 0 | 3 | 3 | 2 |
| 1 | 1 | 1 | 0 |

Option B

| 3 | 3 |
|---|---|
| 3 | 4 |

Option C

| 1 | 2 |
|---|---|
| 3 | 3 |

Option D

| 0 | 7 |
|---|---|
| 14 | 0 |

# In practice,

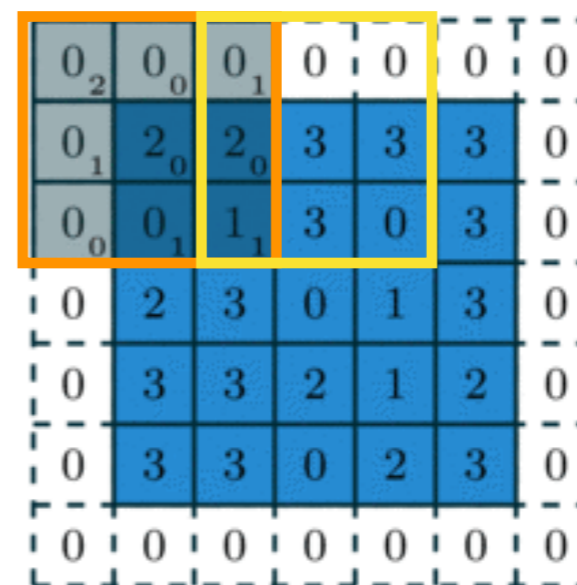- We commonly use convolution with **zero-padding** and **stride**

- **Zero-padding:**

- Pad zeros around the boundary to preserve information and avoid boundary effect



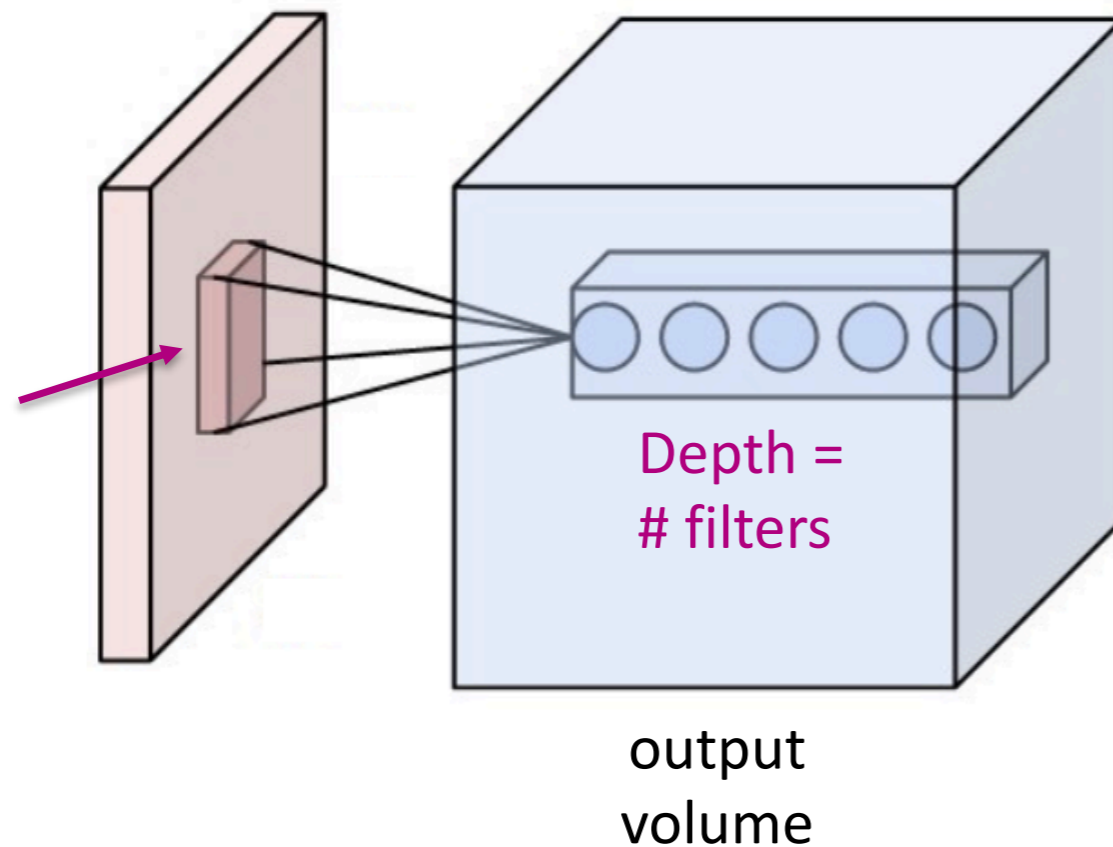If we have 7x7 filter, how many zeros do we need to pad on one row on one side?

- **Stride:**

- skip patches periodically to reduce redundancy and increase efficiency, and capture different resolutions
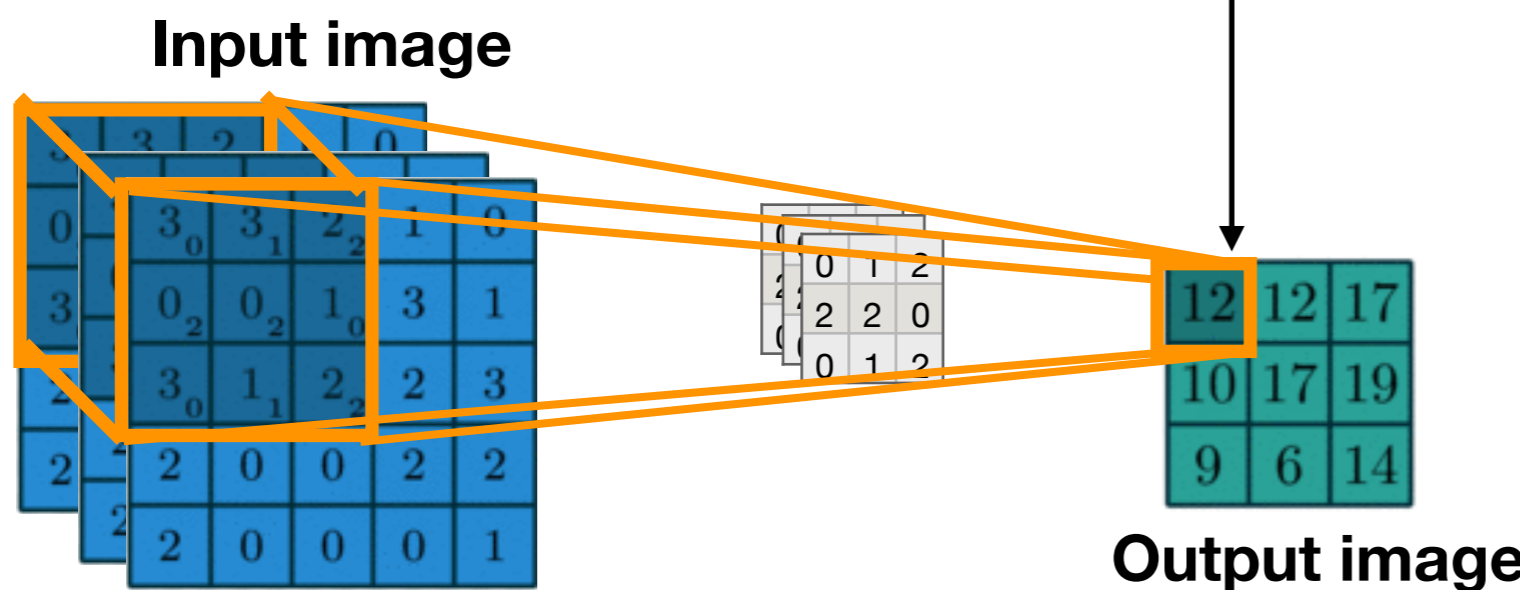


This is an example with 3x3 filter and **stride 2**

# Component in CNN: Convolutional layer
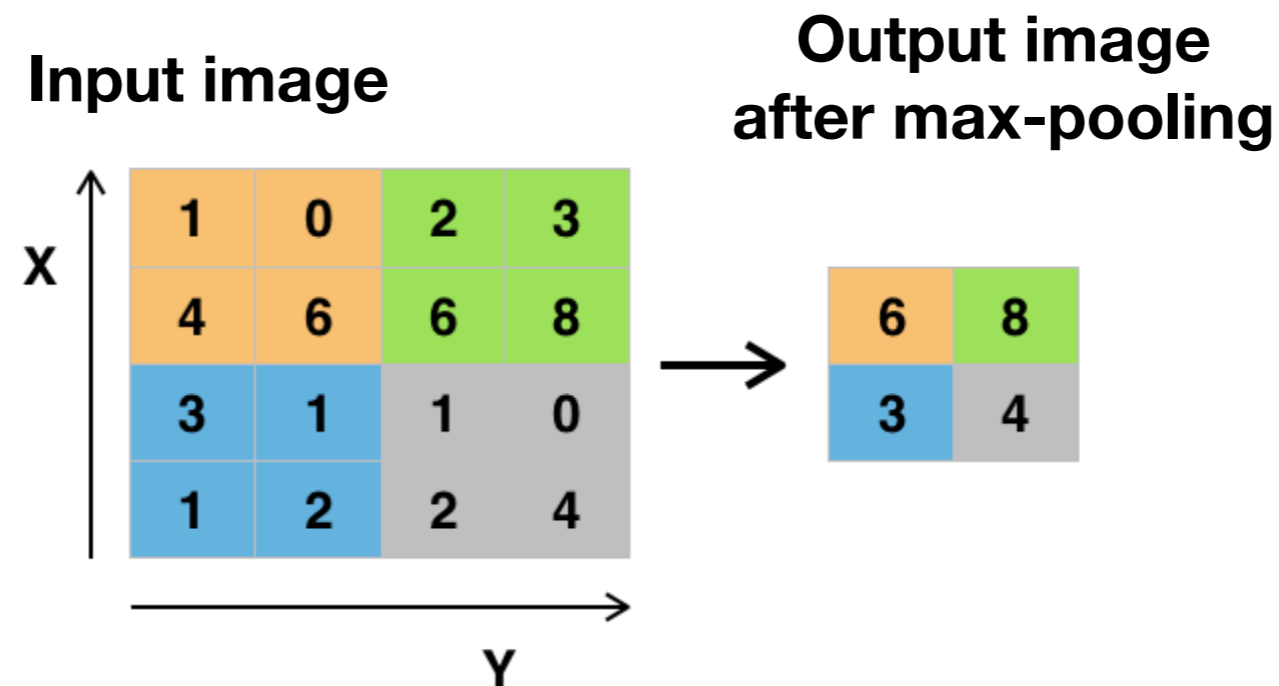


output
volume

In image processing, convolution is typically an operation over three-dimensional arrays

each scalar output = inner product of two
tensors of the same size $(3\times3\times3)$

**Input image**

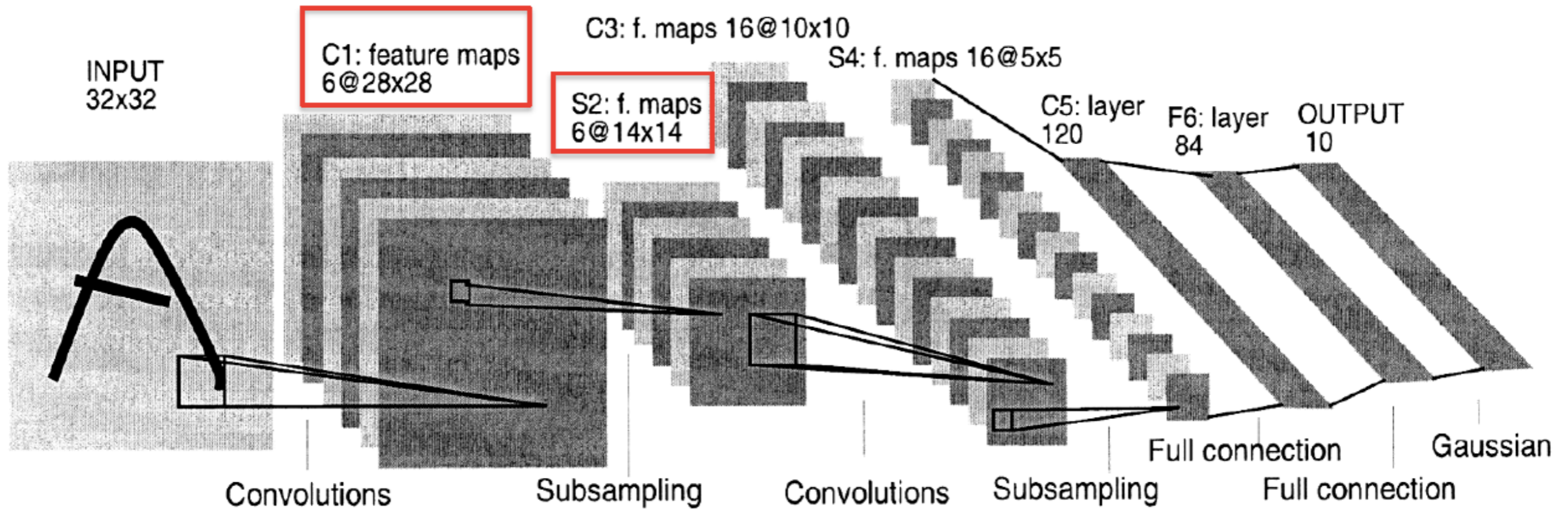**Output image**

# Component in CNN: Pooling layer

- Downsampling the spatial dimensions

- Common to insert between successive **conv** layers

- Typically, **max pooling of size 2x2 with stride 2**
  - Applied separately to each depth slice
  - Tends to work better than average pooling

**Input image**

**Output image after max-pooling**

# Performance of deep learning

- LeNet, 1990's

- 82 error made by LeNet on MNIST

35 error made by Ciresan et al.

further, most of the time the true answer is in the top-2 prediction

idea: train with transformed samples [data augmentation]

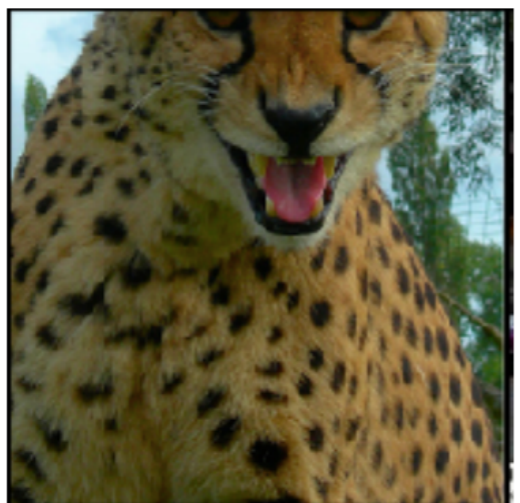# ImageNet 2012 competition: 1.2M training images

Challenging dataset:

High-dimensional data from previous 28 x28 grey-scale to now 256x256 color
10 classes to 1,000 classes
multiple objects
natural 3-d scene

# ImageNet 2012 competition:
# 1.2M training images, 1000 categories

Winning entry:  SuperVision
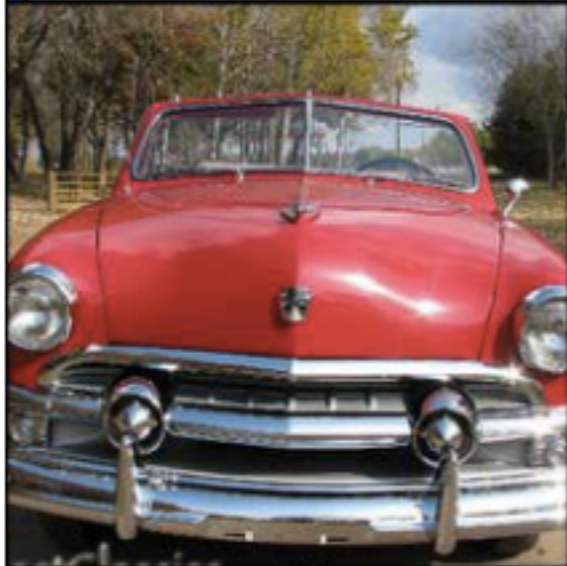8 layers,  60M parameters [Krizhevsky et al. '12]

**mite**

| | |
|---|---|
| ▬ | mite |
| ▮ | black widow |
| ▮ | cockroach |
| ▮ | tick |
| ▮ | starfish |

**container ship**

| | |
|---|---|
| ▬ | container ship |
| ▮ | lifeboat |
| ▮ | amphibian |
| ▮ | fireboat |
| ▮ | drilling platform |

**motor scooter**

| | |
|---|---|
| ▬ | motor scooter |
| ▮ | go-kart |
| ▮ | moped |
| ▮ | bumper car |
| ▮ | golfcart |

**leopard**

| | |
|---|---|
| ▬ | leopard |
| ▮ | jaguar |
| ▮ | cheetah |
| ▮ | snow leopard |
| ▮ | Egyptian cat |

**grille**

| | |
|---|---|
| ▮ | convertible |
| ▬ | grille |
| ▮ | pickup |
| ▮ | beach wagon |
| ▮ | fire engine |

**mushroom**

| | |
|---|---|
| ▮ | agaric |
| ▬ | mushroom |
| ▮ | jelly fungus |
| ▮ | gill fungus |
| ▮ | dead-man's-fingers |

**cherry**

| | |
|---|---|
| ▮ | dalmatian |
| ▮ | grape |
| ▮ | elderberry |
| ▮ | ffordshire bullterrier |
| ▮ | currant |

**Madagascar cat**

| | |
|---|---|
| ▮ | squirrel monkey |
| ▮ | spider monkey |
| ▮ | titi |
| ▮ | indri |
| ▮ | howler monkey |

47

# Going even deeper…



Won 2014 ImageNet challenge with 6.66% top-5 error rate

GoogLeNet, 2014

Huge CNN depth has proven helpful in recognition systems… Maybe because images contain hierarchical structure (faces contain eyes contain edges, etc.)
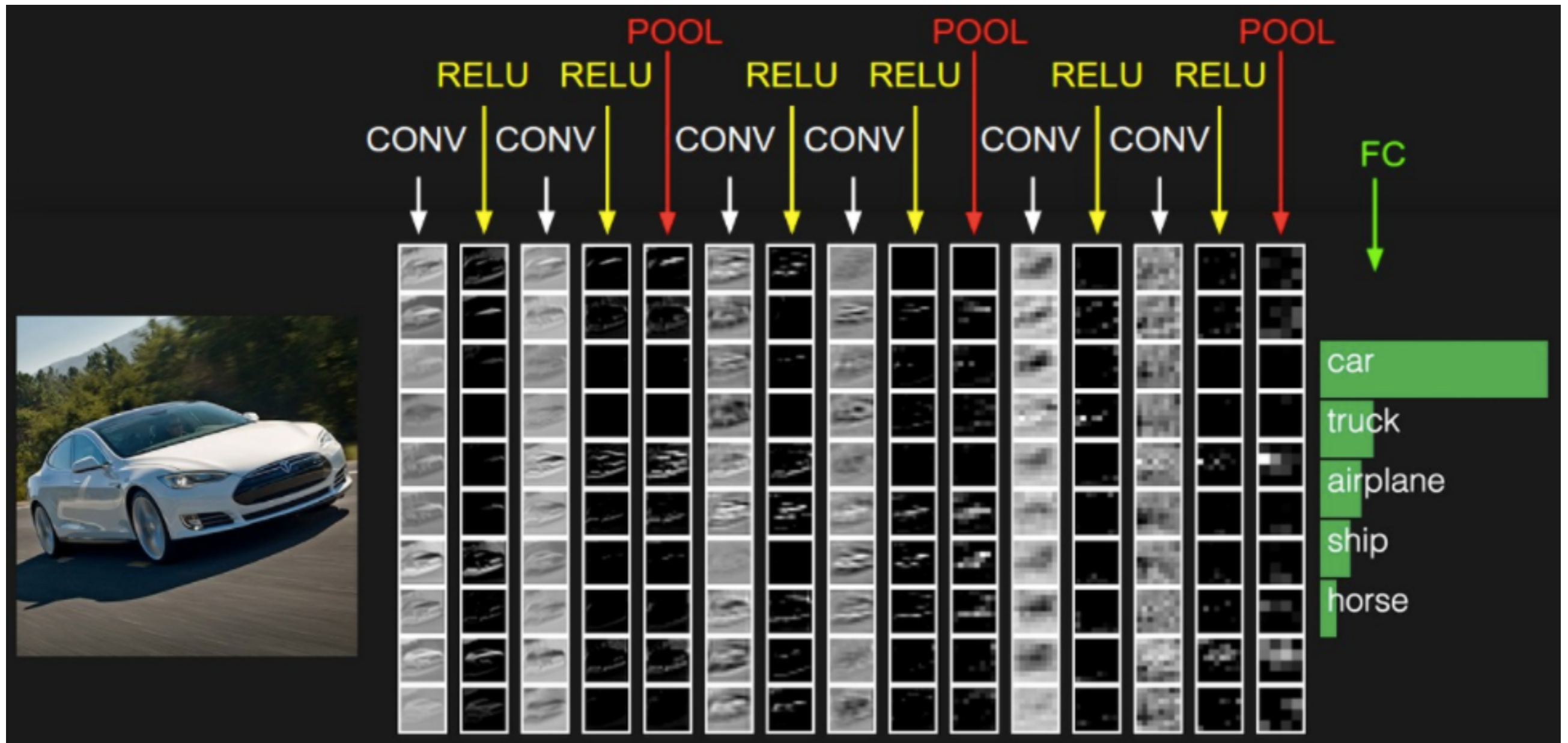
48

# What happens when we train convolutional neural networks?

- First convolutional layer trained on natural images looks like the following



- Simple geometric patterns are "detected" or "matched" in the first layer

POOL    POOL    POOL

RELU  RELU    RELU  RELU    RELU  RELU

CONV  CONV    CONV  CONV    CONV  CONV    FC

car

truck

airplane

ship

horse

50

# Returning to our example…
# "Detectors" are the learned filters



| | Layer 1 | Layer 2 | Layer 3 |
|---|---|---|---|
| Example detectors learned | | | |
| Example interest points detected | | | |

[Zeiler & Fergus '13]

**Filter *W* at layer 1**          **Feature map at layer 3 shown on image space**

# Performance of deep learning



German traffic sign
recognition benchmark

   – 99.5% accuracy (IDSIA team)

House number recognition

   – 97.8% accuracy per character
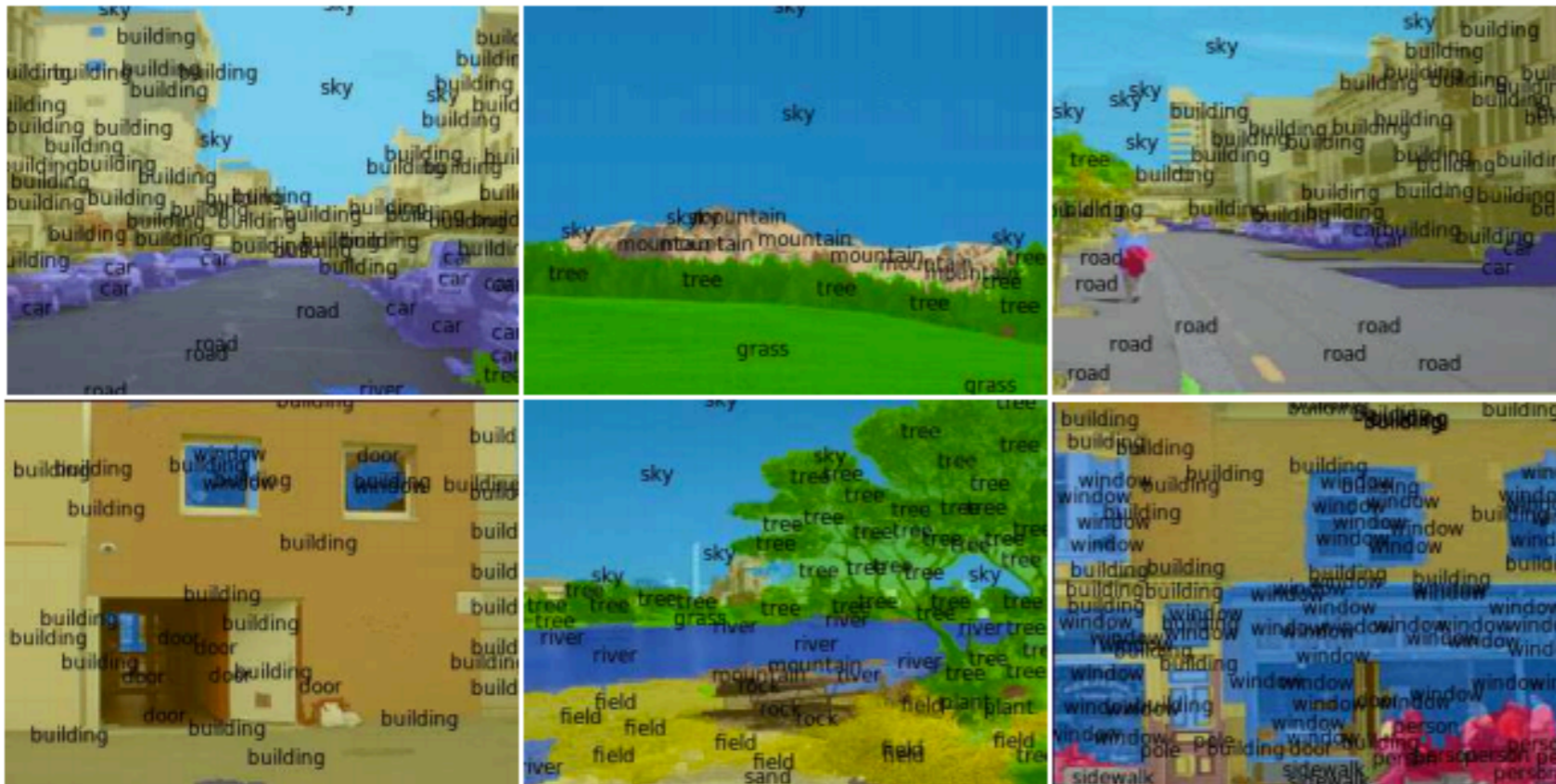     [Goodfellow et al. '13]
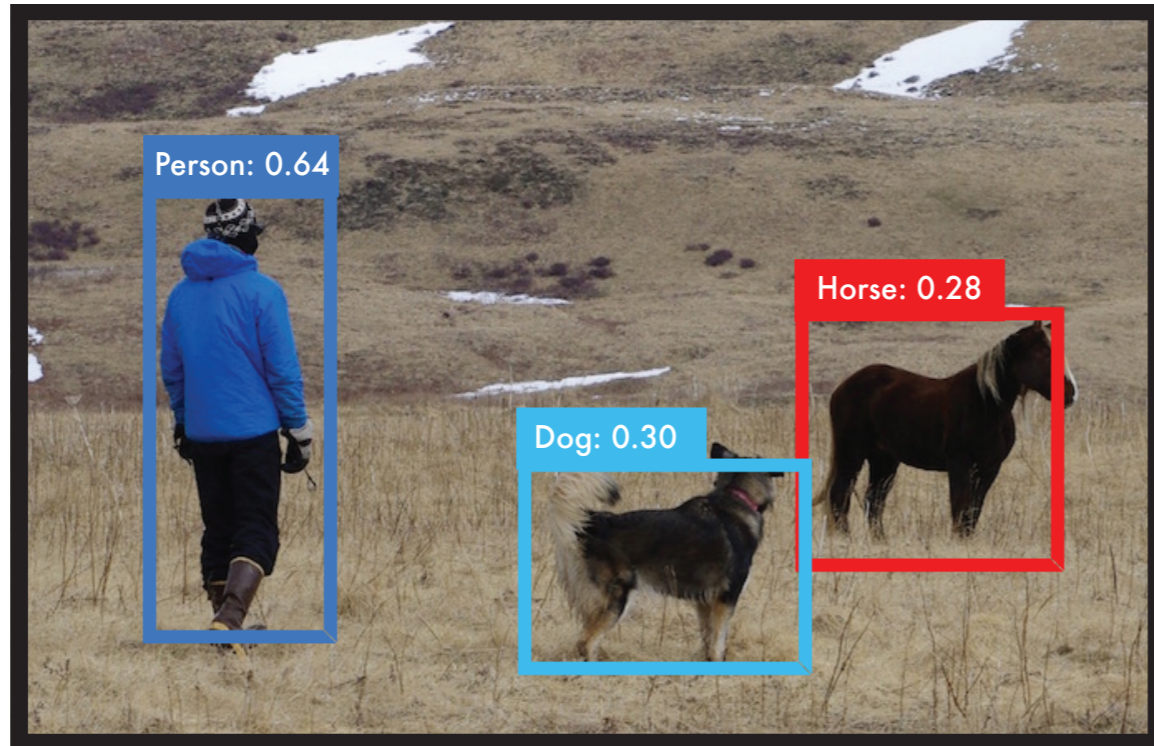
- Image classification



Input: **x**
**Image pixels**

Output: y
Predicted object

- Scene parsing

[Farabet et al. '13]

- Object detection



Redmon et al. 2015
http://pjreddie.com/yolo/

- Retrieving similar objects
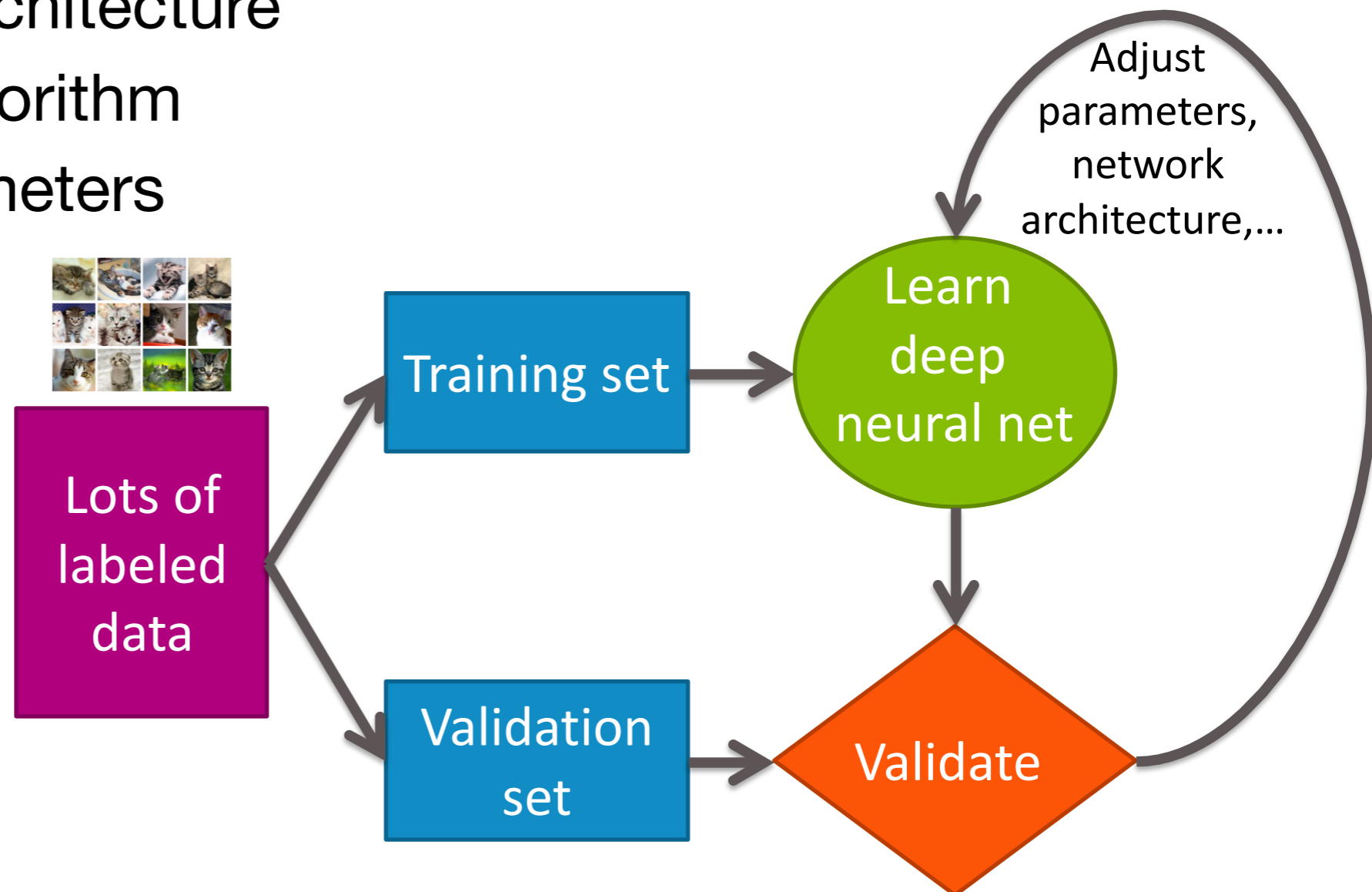
Input Image     Nearest neighbors

# Deep Learning practice

- Pros
  - Instead of manually engineering features, enable automated learning of features
  - Impressive performance gains in practice
    - Image processing
    - Natural language processing
    - Speech recognition
  - Making huge impacts in many applications in many fields
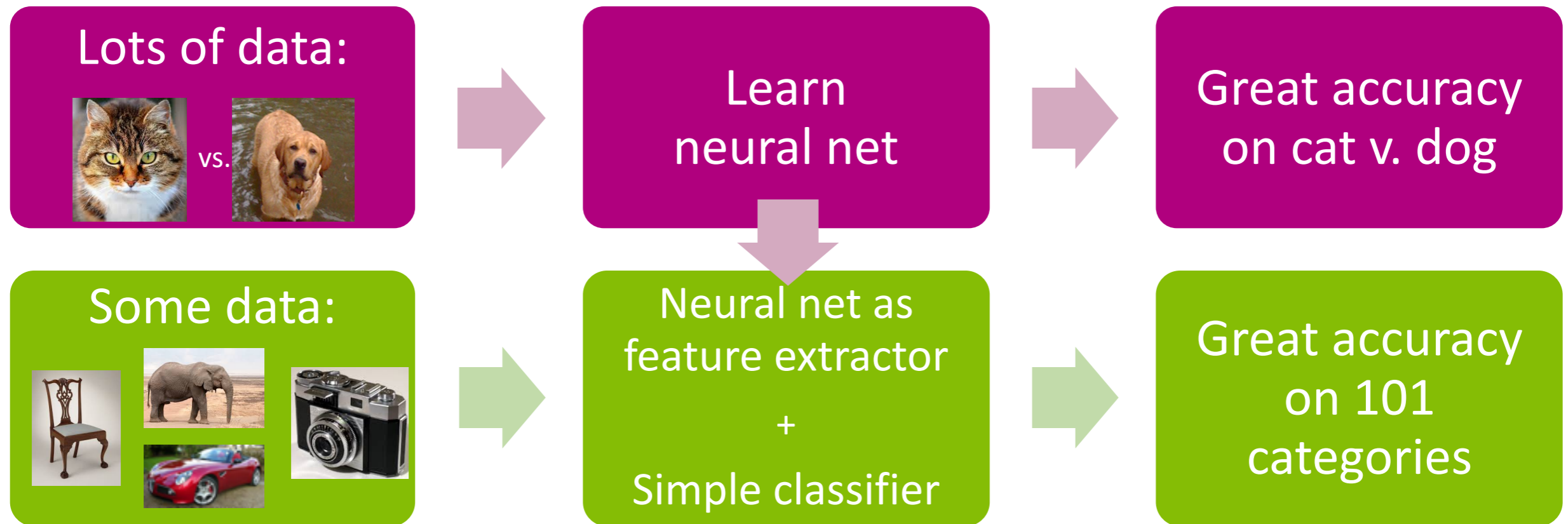
# Deep Learning practice

- Cons
  - Requires a lot of data
  - Computationally really expensive
  - Hard to tune hyper-parameters
    - Choice of architecture
    - Learning algorithm
    - Hyper-parameters

# Transfer Learning

# Transfer Learning

- Transfer Learning
  - Use data from one task to help learn on another task
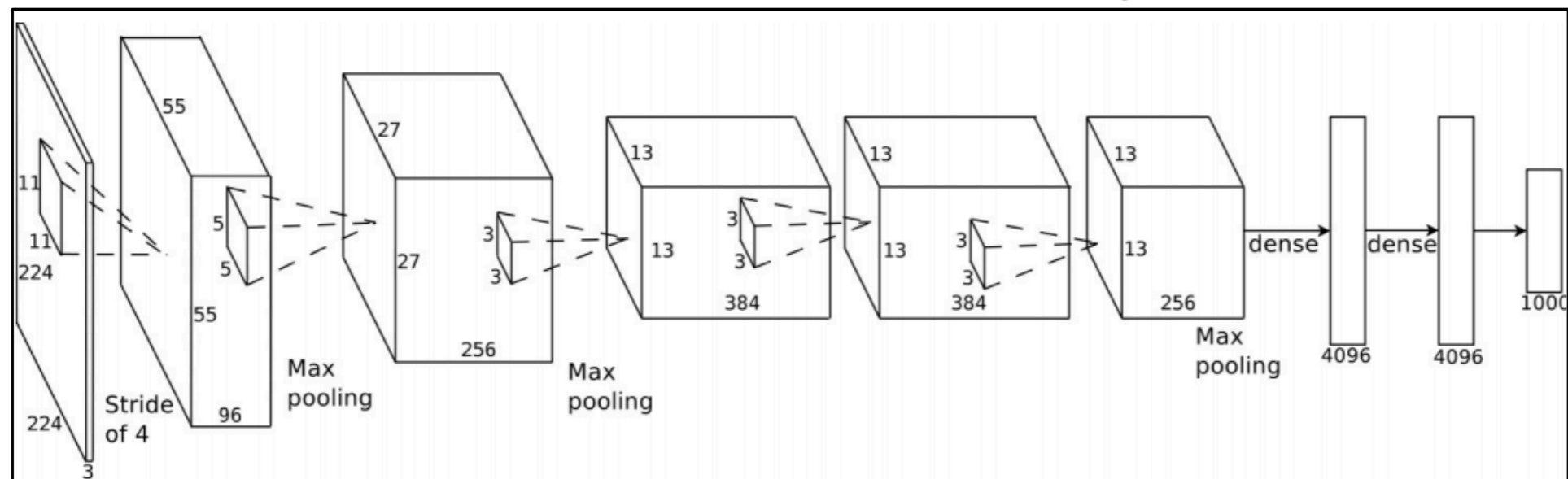  - Old idea, explored for deep learning by Donahue et al. '14 & others
  -

# What is learned in a neural networks

- Initial layers are not too sensitive/specific to the task at training

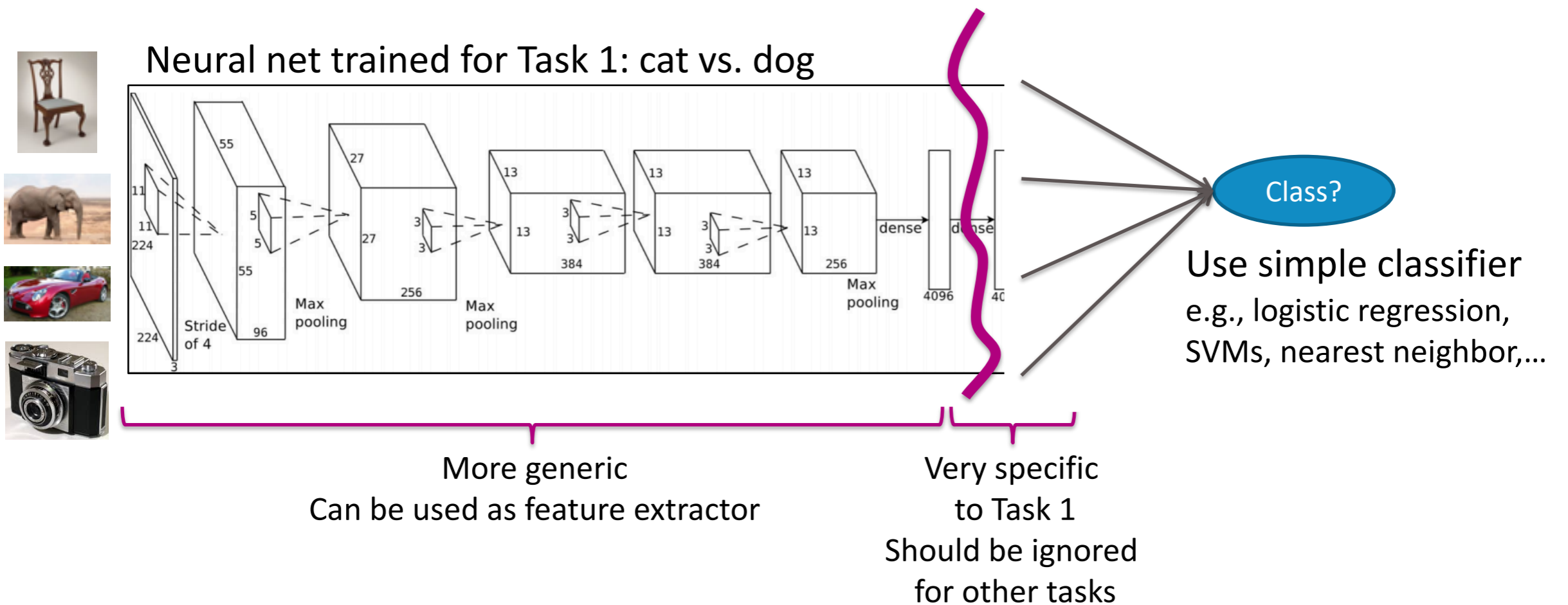Neural net trained for Task 1: cat vs. dog



vs.

More generic
Can be used as feature extractor

Very specific
to Task 1
Should be ignored
for other tasks

# Transfer learning

- For the second task of predicting 101 categories, (re)-train only the last layer of the neural network
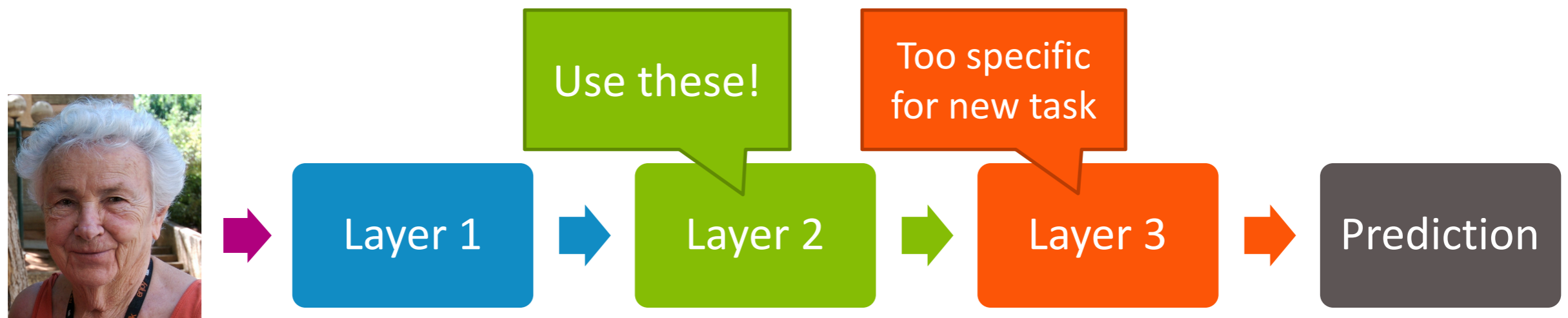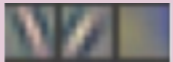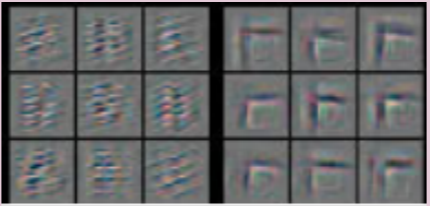
Neural net trained for Task 1: cat vs. dog



Class?

Use simple classifier
e.g., logistic regression, SVMs, nearest neighbor,…

More generic
Can be used as feature extractor

Very specific
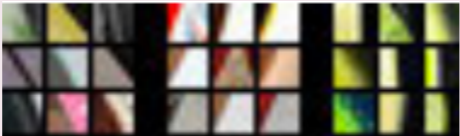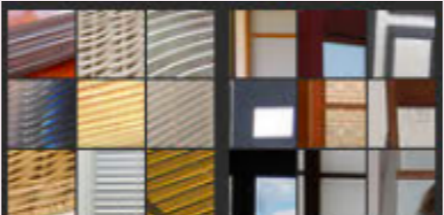to Task 1
Should be ignored
for other tasks

Keep weights fixed!

Re-train

# Transfer learning

- Need to be careful about where you cut, as latter layers may be too task specific



[Zeiler & Fergus '13]