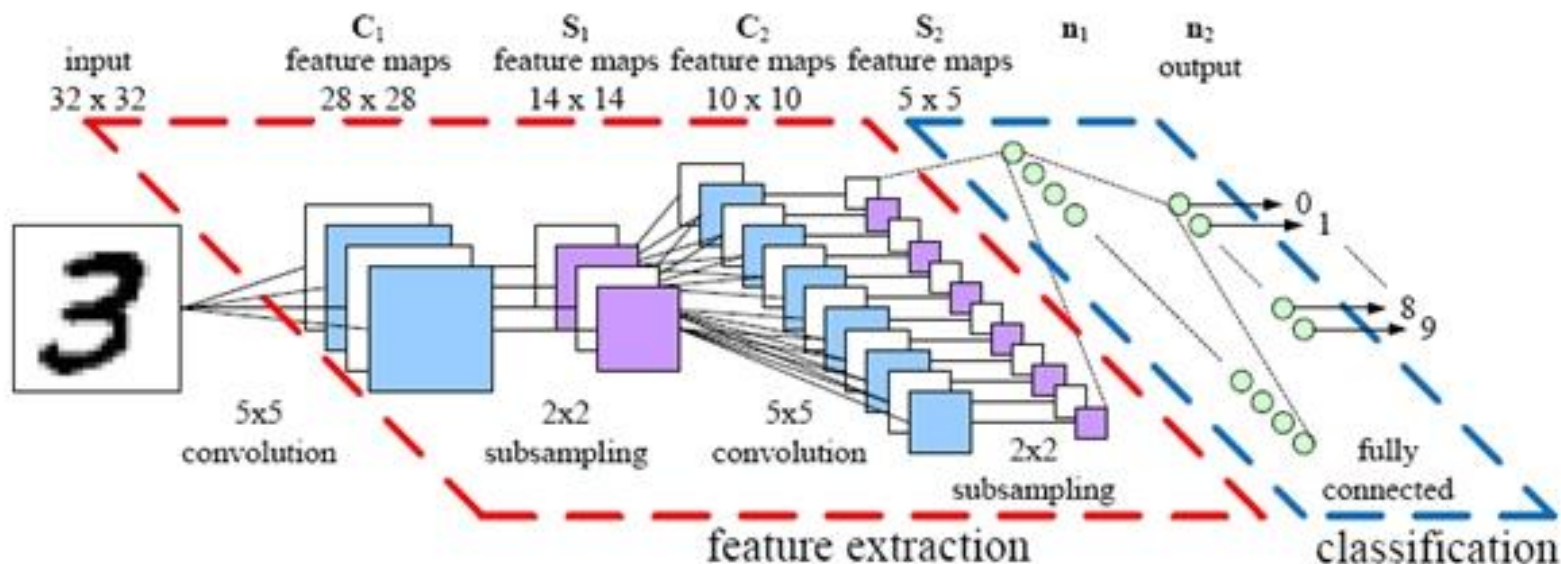


# DEEP LEARNING



# HANDWRITING RECOGNITION

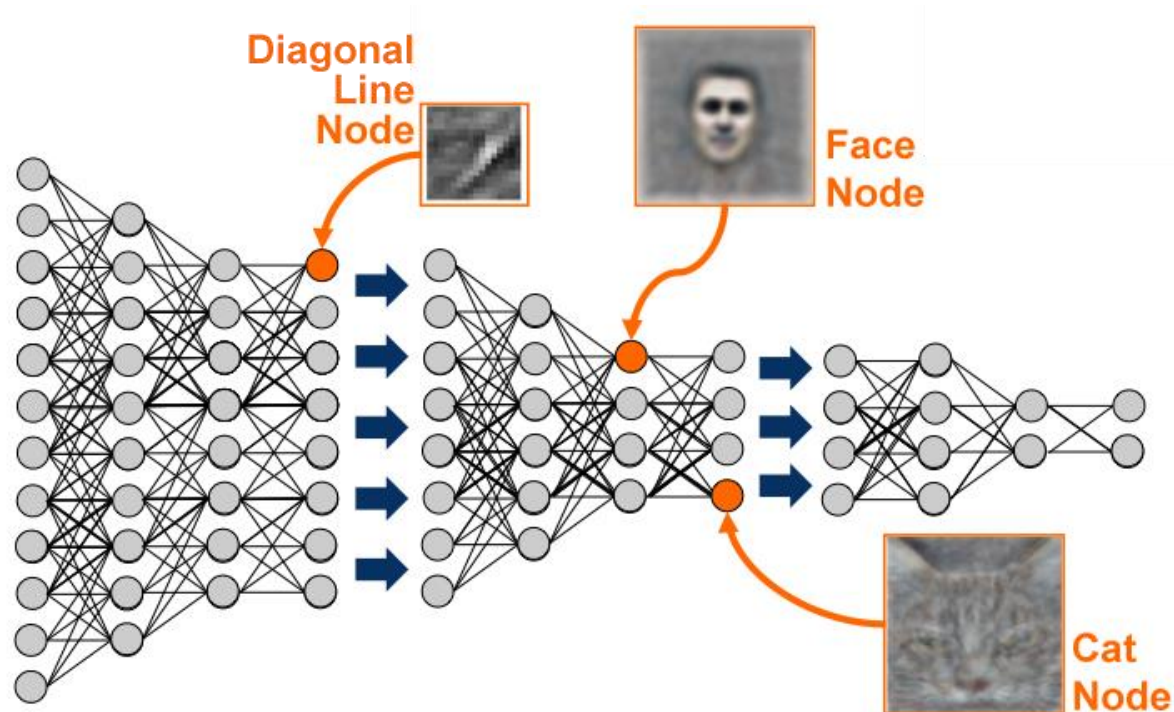


0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9





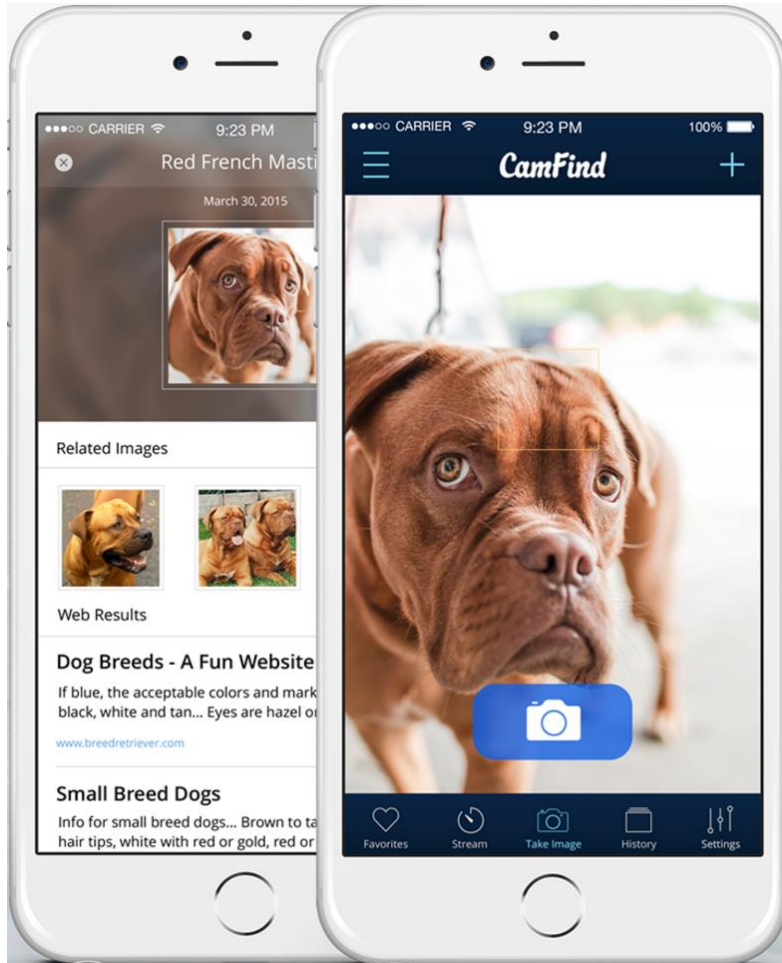
# GOOGLE CAT VIDEOS



Deep Learning  
with GPUs



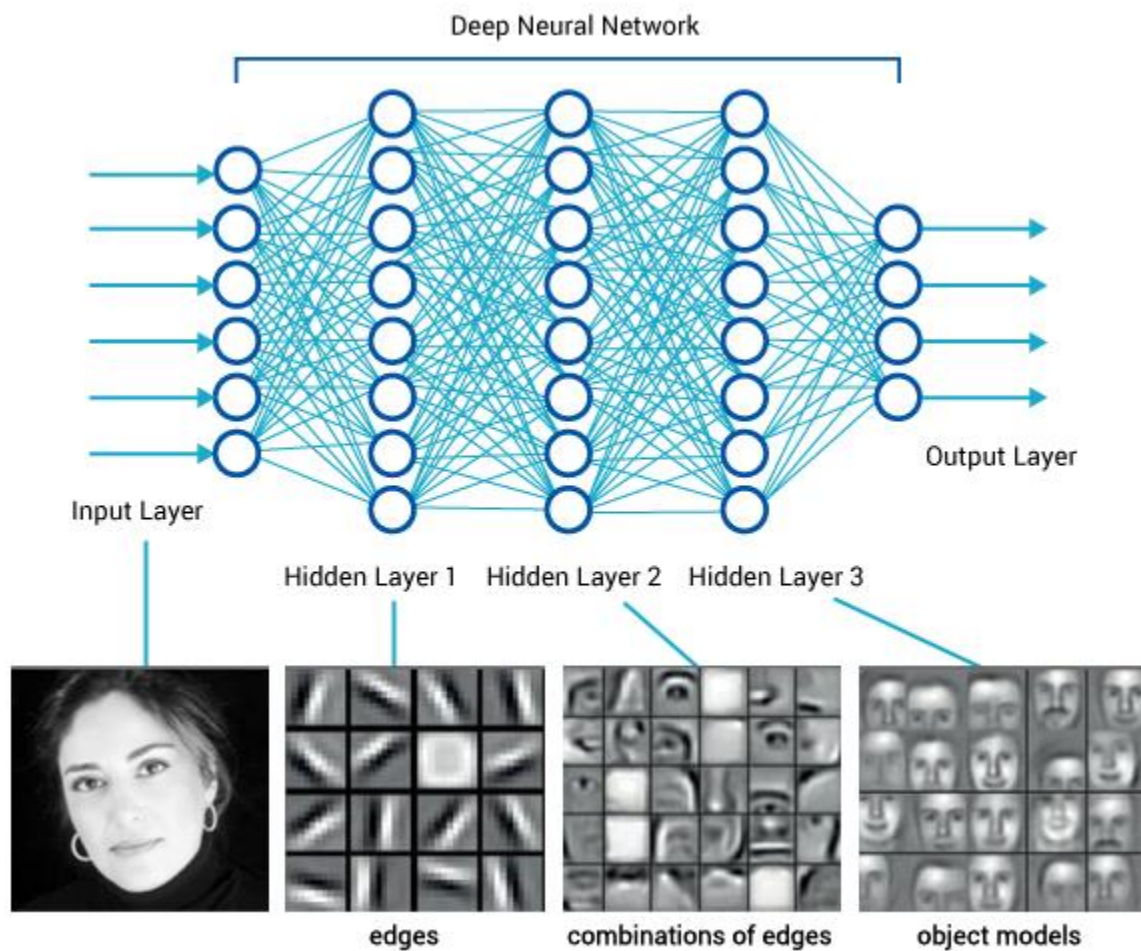
# IMAGE RECOGNITION



**CamFind**  
Visual Search Engine  
(available on iOS, Android)



# FACE RECOGNITION



# SPEECH TRANSLATION



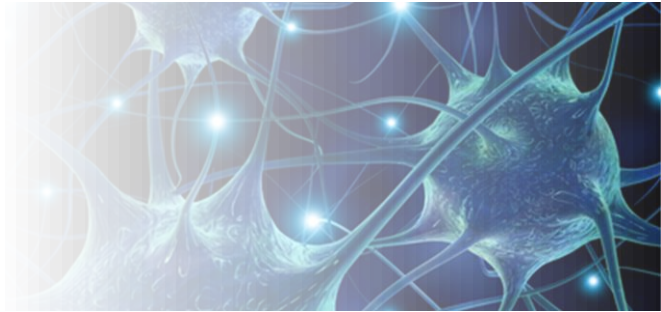
From Hidden Markov Models to Recurrent Neural Networks





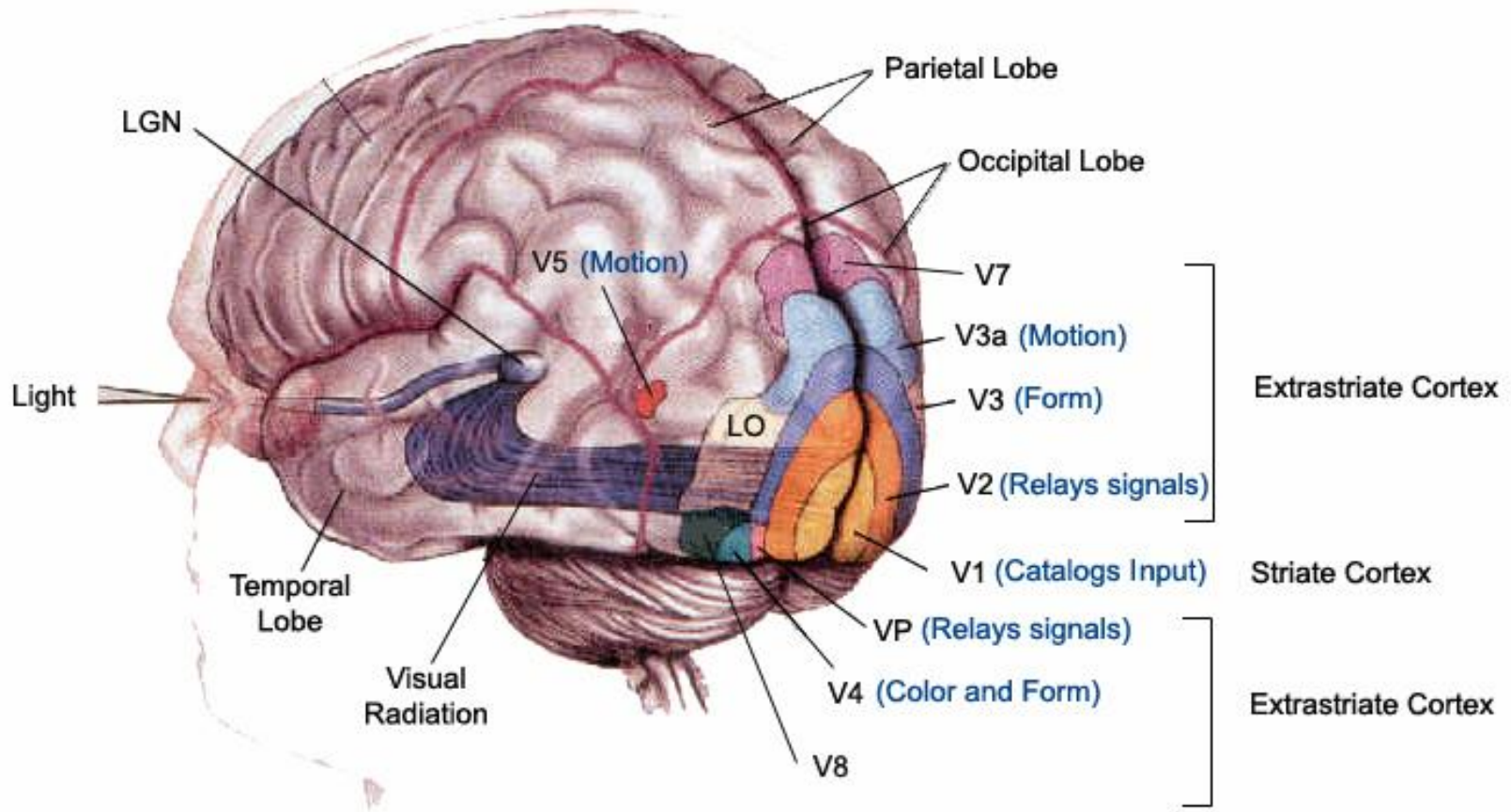


**THEORY**



# WHAT IS DEEP LEARNING?

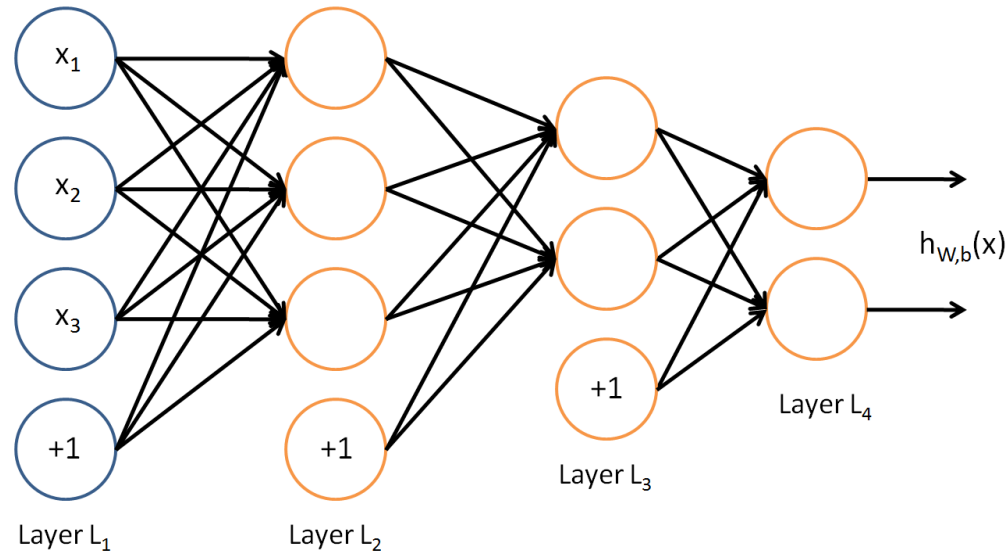
## Visual Cortices





# WHAT IS DEEP LEARNING?

Biologically-inspired multilayer neural networks

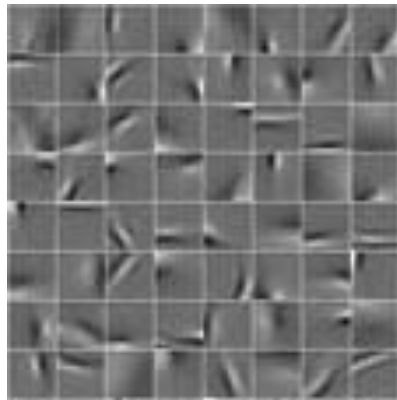


Both supervised and unsupervised

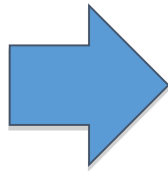


# WHAT IS DEEP LEARNING?

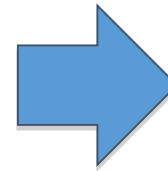
**Example.** Face recognition (Facebook)



Edges



Eyes, Noses, Mouths



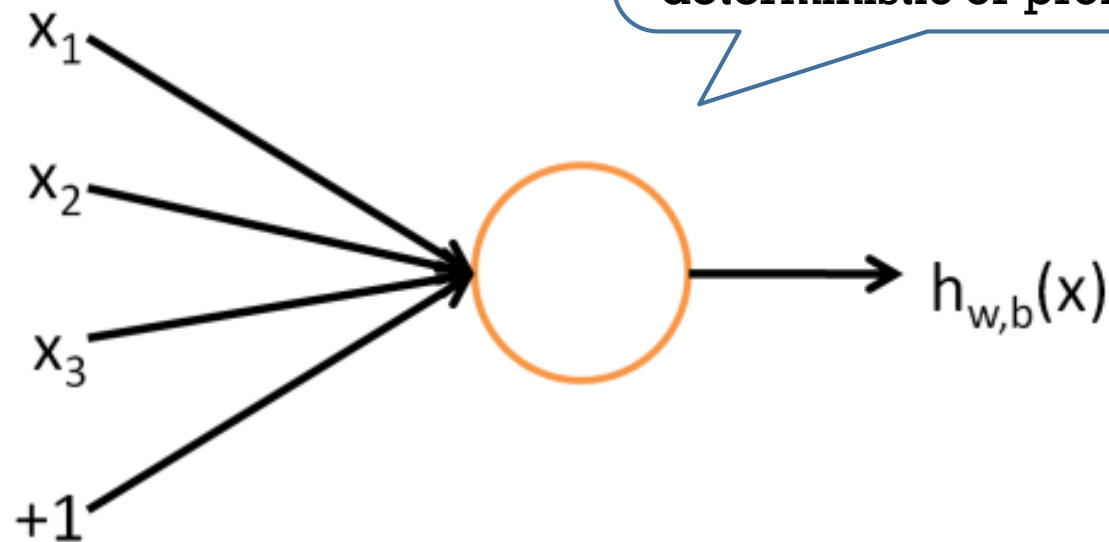
Faces

Deeper layers learn higher-order features



# NEURON

## Perceptron.



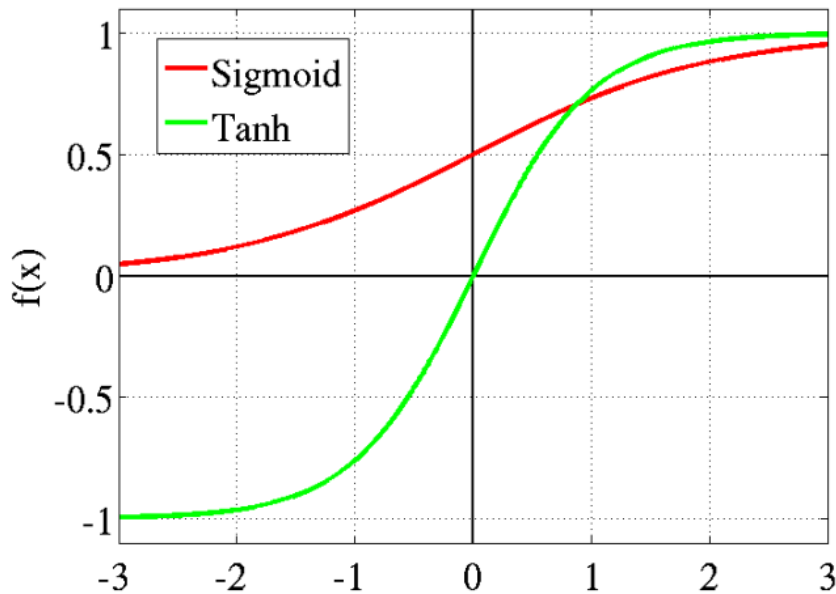
Depending on function  $f$ ,  
neurons can be:  
real-valued or binary-valued;  
deterministic or probabilistic.

$$h_{w,b}(x) = f(w^T x) = f\left(\sum_{i=1}^d w_i x_i + b\right)$$



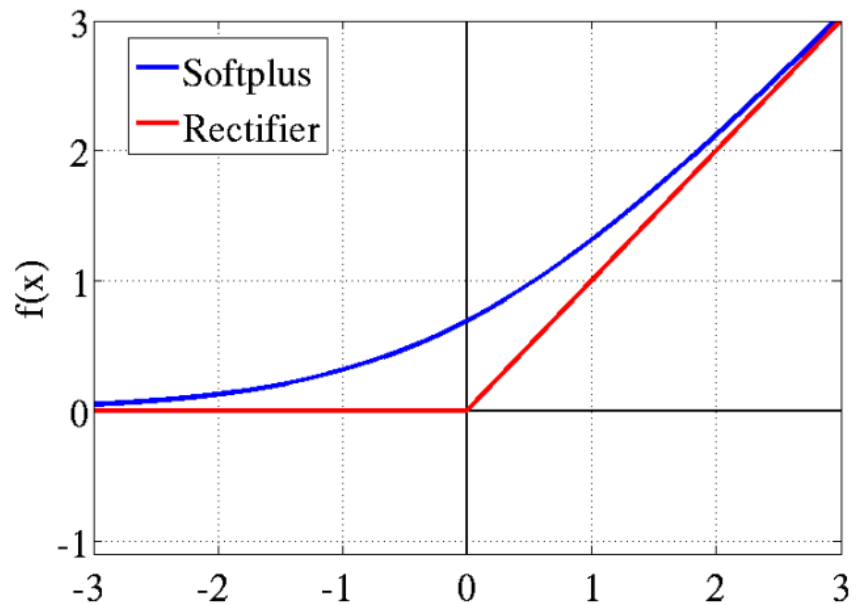


# ACTIVATION FUNCTIONS



$$\text{sigmoid } f(z) = \frac{1}{1+e^{-z}}$$

$$\text{tanh } f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

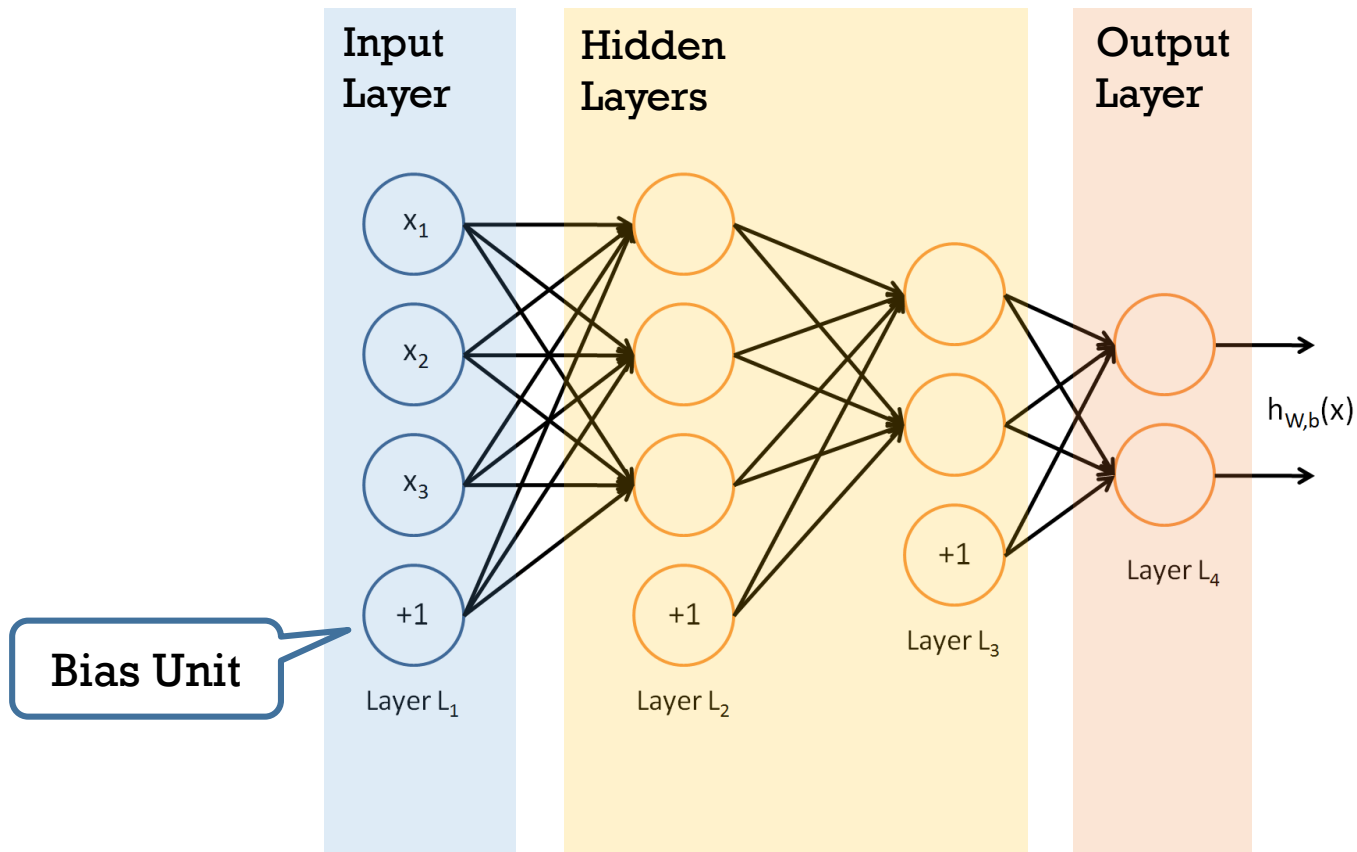


$$\text{softplus } f(z) = \ln(1 + e^{-z})$$

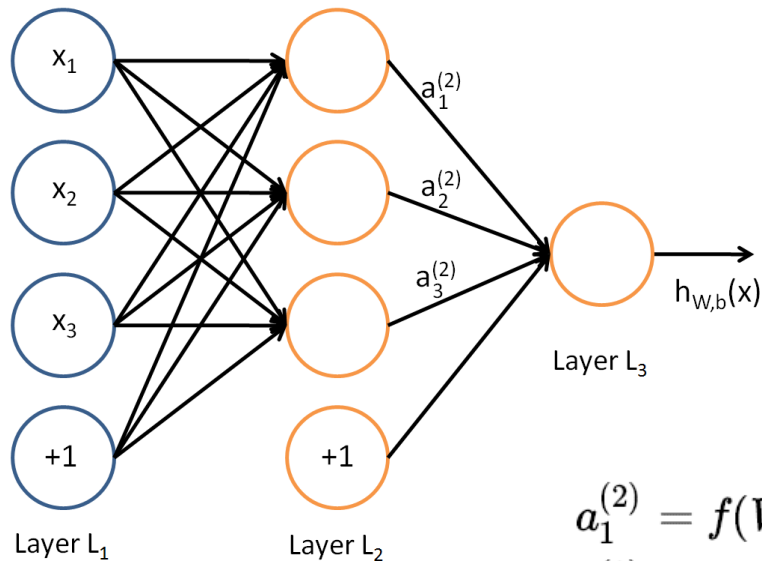
$$\text{rectified linear unit (ReLU)} \quad f(z) = \max(0, z)$$



# MULTI-LAYER NEURAL NETWORK



# MULTI-LAYER NEURAL NETWORK

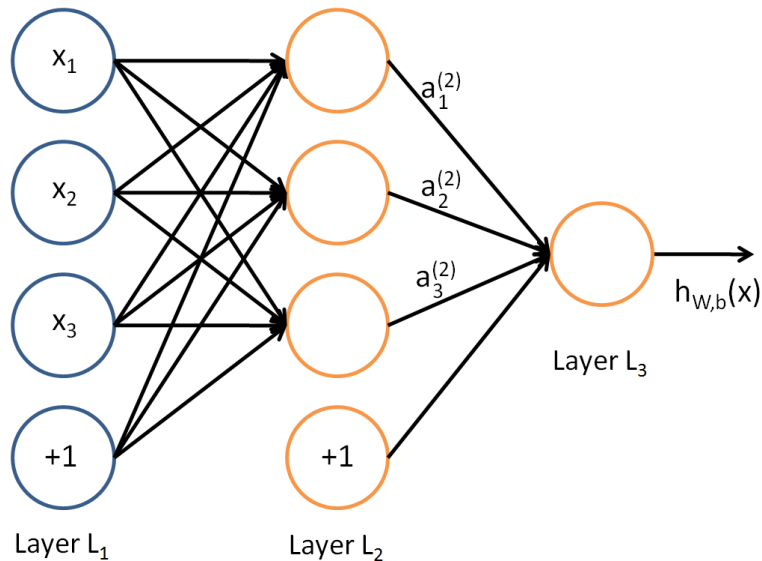


$$\begin{aligned}a_1^{(2)} &= f(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + W_{13}^{(1)} x_3 + b_1^{(1)}) \\a_2^{(2)} &= f(W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{23}^{(1)} x_3 + b_2^{(1)}) \\a_3^{(2)} &= f(W_{31}^{(1)} x_1 + W_{32}^{(1)} x_2 + W_{33}^{(1)} x_3 + b_3^{(1)}) \\h_{W,b}(x) &= a_1^{(3)} = f(W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + W_{13}^{(2)} a_3^{(2)} + b_1^{(2)})\end{aligned}$$





# MULTI-LAYER NEURAL NETWORK



## Forward Propagation.

$$z^{(2)} = W^{(1)} x + b^{(1)}$$

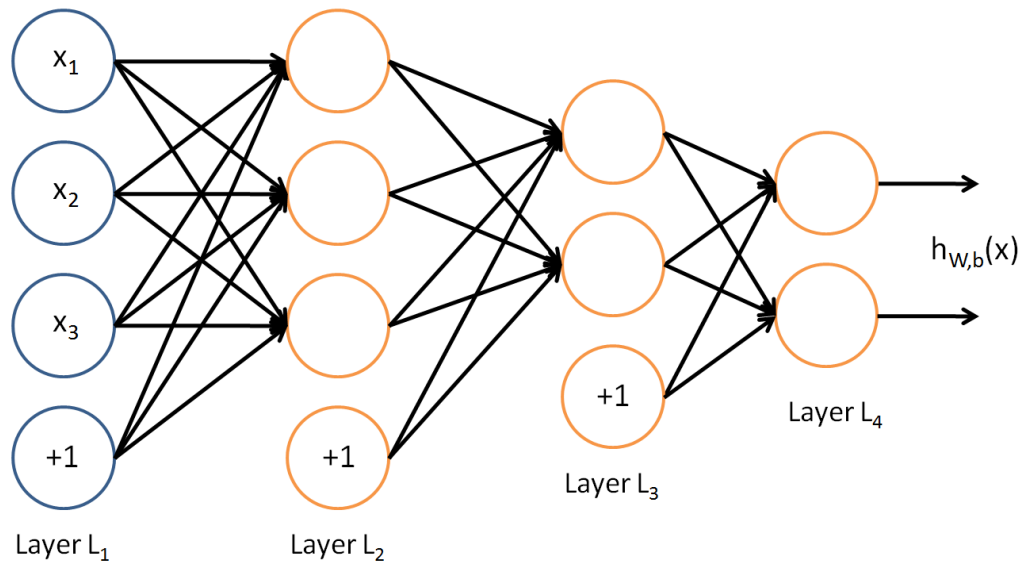
$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)} a^{(2)} + b^{(2)}$$

$$h_{W,b}(x) = a^{(3)} = f(z^{(3)})$$



# MULTI-LAYER NEURAL NETWORK



**Feed-Forward Neural Network.**

**Neural Network Architecture.**

Arrangement of neurons and their connections.

$$z^{(l+1)} = W^{(l)} a^{(l)} + b^{(l)}$$
$$a^{(l+1)} = f(z^{(l+1)})$$

Activation



# COST FUNCTION

Error for one data point

For binary neurons, other loss functions are used.

$$J(W, b; x, y) = \frac{1}{2} \|h_{W,b}(x) - y\|^2$$

Error for training set + **Weight decay** regularization

$$\begin{aligned} J(W, b) &= \left[ \frac{1}{m} \sum_{i=1}^m J(W, b; x^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left( W_{ji}^{(l)} \right)^2 \\ &= \left[ \frac{1}{m} \sum_{i=1}^m \left( \frac{1}{2} \|h_{W,b}(x^{(i)}) - y^{(i)}\|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^{n_l-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} \left( W_{ji}^{(l)} \right)^2 \end{aligned}$$

$\lambda$  weight decay parameter





# BACK PROPAGATION

## Chain Rule for Neural Networks.

$$\begin{aligned}\frac{\partial}{\partial W_{ij}^{(l)}} \left( \frac{1}{2} \|a^{(n_l)} - y\|^2 \right) &= (a^{(n_l)} - y) \frac{\partial}{\partial W_{ij}^{(l)}} f(z^{(n_l)}) \\&= (a^{(n_l)} - y) f'(z^{(n_l)}) \frac{\partial}{\partial W_{ij}^{(l)}} (W^{(n_l-1)} a^{(n_l-1)} + b^{(n_l-1)}) \\&= (a^{(n_l)} - y) f'(z^{(n_l)}) W^{(n_l-1)} \frac{\partial}{\partial W_{ij}^{(l)}} a^{(n_l-1)} \\&= \underbrace{(a^{(n_l)} - y) f'(z^{(n_l)}) W^{(n_l-1)} f'(z^{(n_l-1)})}_{\delta^{(n_l)}} \frac{\partial}{\partial W_{ij}^{(l)}} z^{(n_l-1)} \\&\quad \underbrace{\hspace{10em}}_{\delta^{(n_l-1)}}\end{aligned}$$



# BACK PROPAGATION

1. Perform a feed-forward pass, computing the activations layer by layer.

2. For the output layer (layer  $n_l$ ), set

$$\delta^{(n_l)} = -(y - a^{(n_l)}) \bullet f'(z^{(n_l)})$$

3. For  $l = n_l - 1, n_l - 2, n_l - 3, \dots, 2$ , set

$$\delta^{(l)} = \left( (W^{(l)})^T \delta^{(l+1)} \right) \bullet f'(z^{(l)})$$

4. Compute the desired partial derivatives:

$$\nabla_{W^{(l)}} J(W, b; x, y) = \delta^{(l+1)} (a^{(l)})^T + \lambda W^{(l)}$$

$$\nabla_{b^{(l)}} J(W, b; x, y) = \delta^{(l+1)}.$$

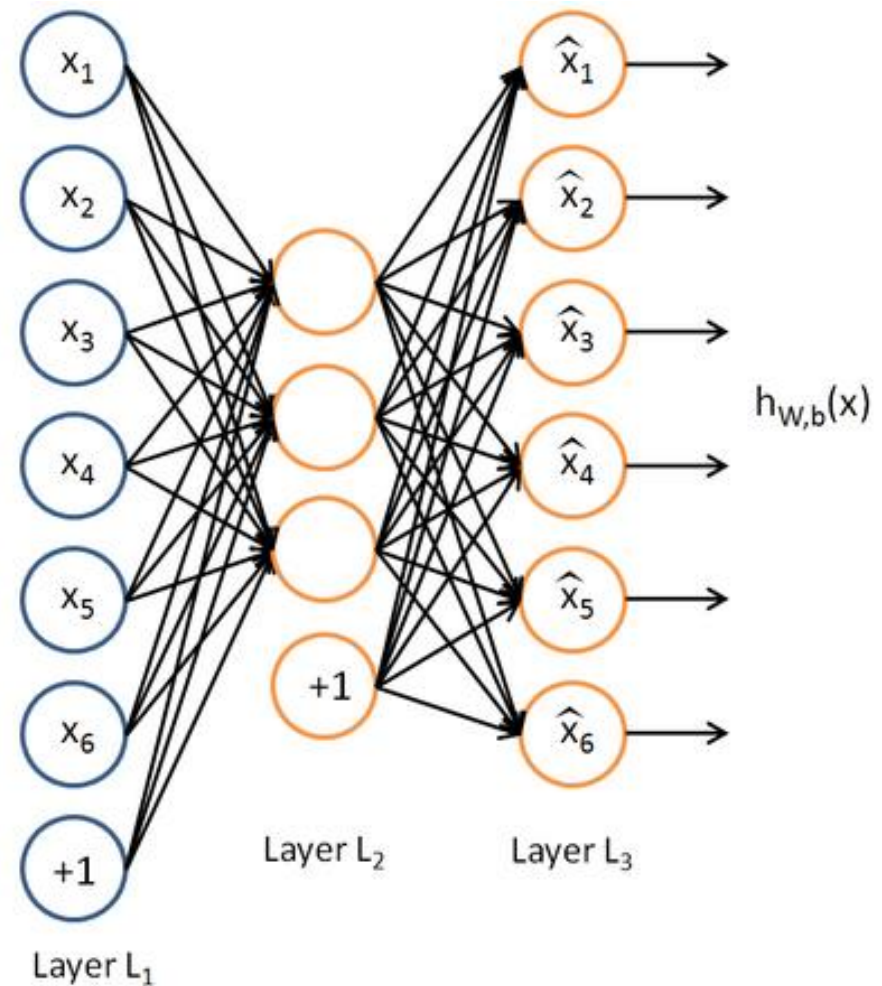


# AUTOENCODERS

Training a multilayer neural network to reconstruct the input from a **reduced representation**.

## Strategies for Dimensionality Reduction

- Few hidden neurons
- Sparse activations





# Sparse Autoencoder

## Sparsity Penalty.

Average activation  $\hat{\rho}_j = \frac{1}{m} \sum_{i=1}^m \left[ a_j^{(2)}(x^{(i)}) \right]$

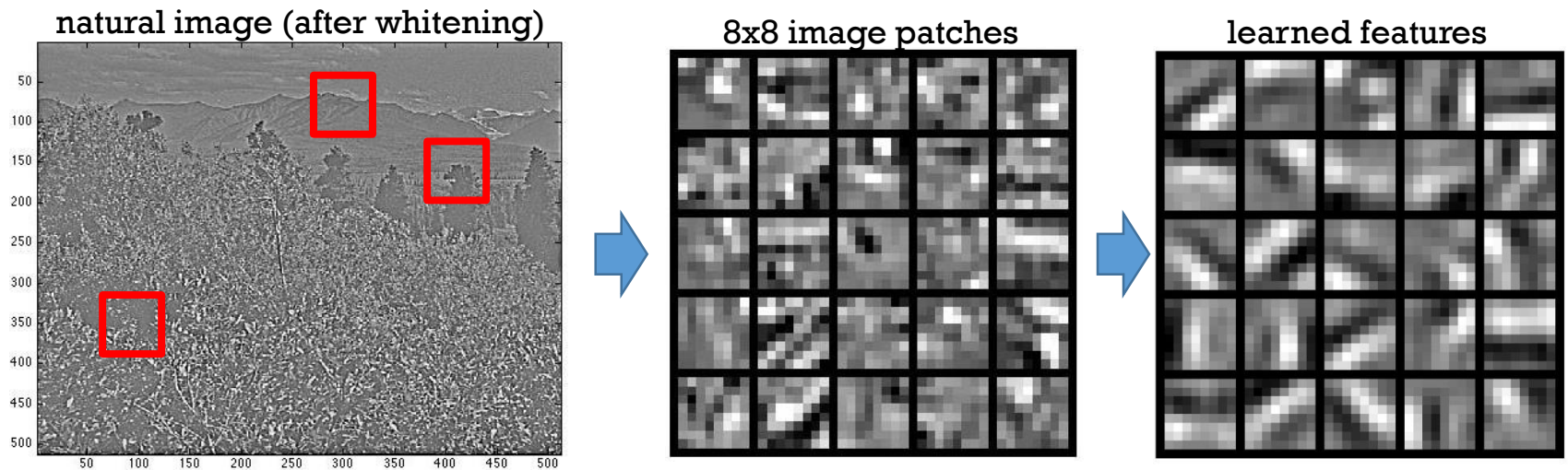
$$\text{KL}(\rho || \hat{\rho}_j) = \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j}$$

$$J_{\text{sparse}}(W, b) = J(W, b) + \beta \sum_{j=1}^{s_2} \text{KL}(\rho || \hat{\rho}_j),$$

$\beta$  sparsity parameter



# NATURAL IMAGES



Edge features similar to those from neuroscience experiments (see Hubel & Wiesel Cat Experiment, 1959).

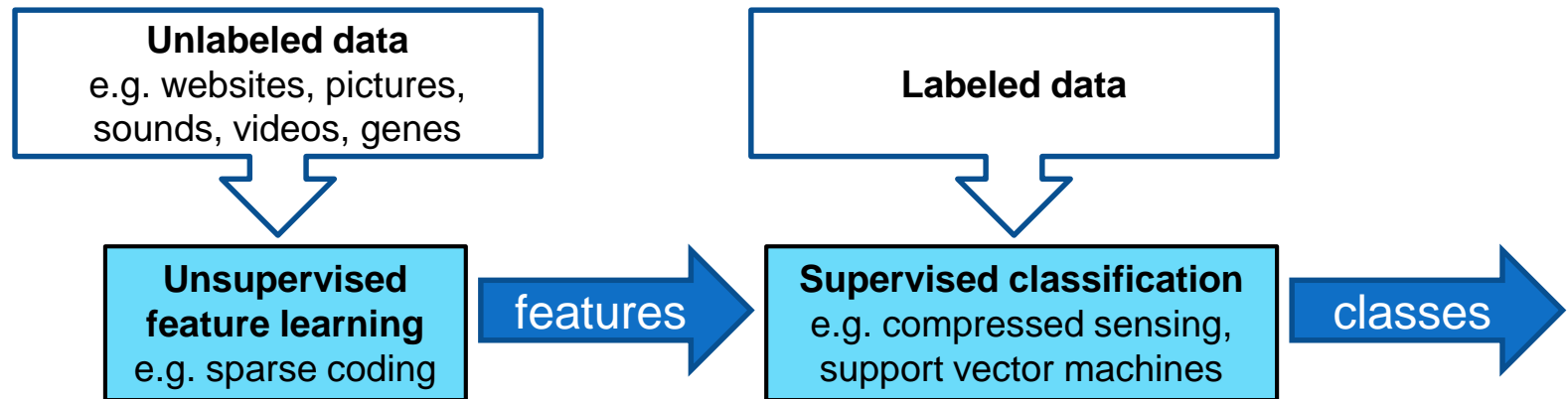


# GREEDY INITIALIZATION OF LAYERS

- Initialization of neural weights
  1. Train an autoencoder on input data.
  2. Feed-forward data to hidden layer.
  3. Train a second autoencoder on hidden activations.
  4. Feed-forward data to second hidden layer.
  5. ... and so on.
- Fine tuning of neural weights
  - Backpropagation



# SEMI-SUPERVISED LEARNING



# WHY DOES DEEP LEARNING WORK?

*It's not really about the autoencoders...*

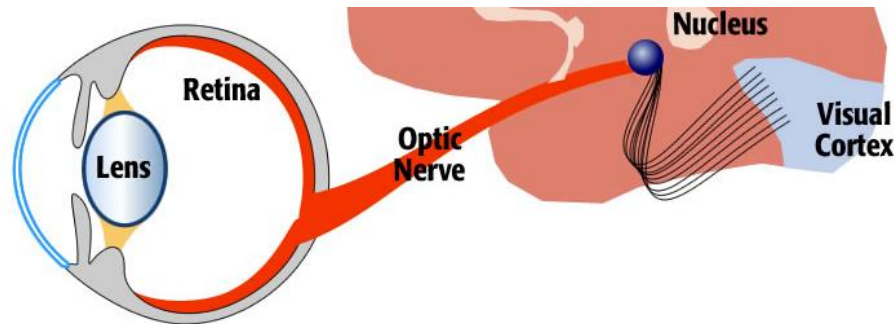
- Greedy layer-wise initialization can be any of the following:
  - Contrastive divergence
  - Sparse autoencoder
  - Sparse coding
  - K-means clustering
  - Random data vectors





# WHY DOES DEEP LEARNING WORK?

*... but it has something to do with sparsity.*



## Compressed Sensing

- If input is sparse in some known basis, then it can be reconstructed from almost any random projections.

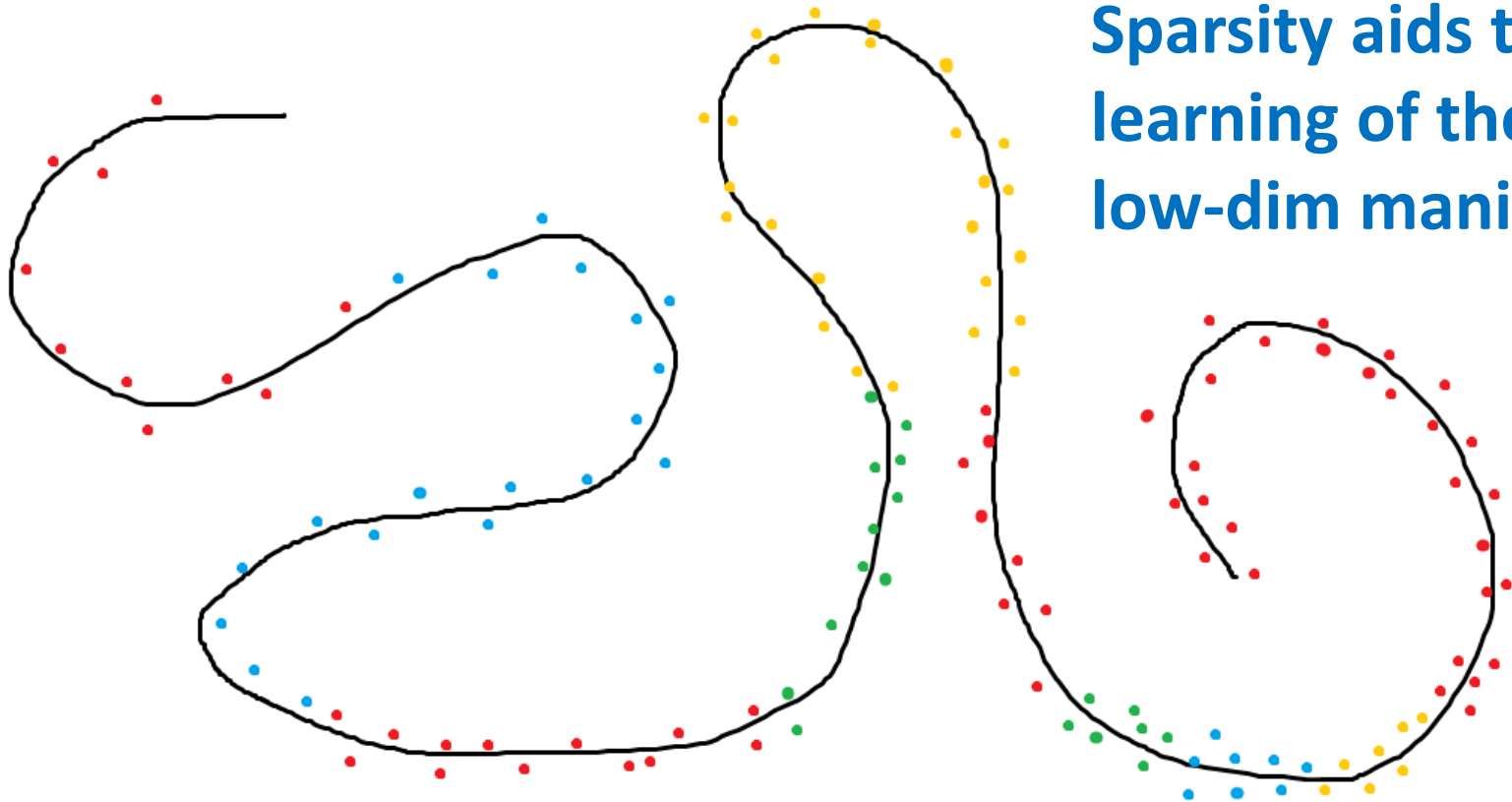


# WHY DOES DEEP LEARNING WORK?

Data is often near  
low-dim manifold  
in high-dim space



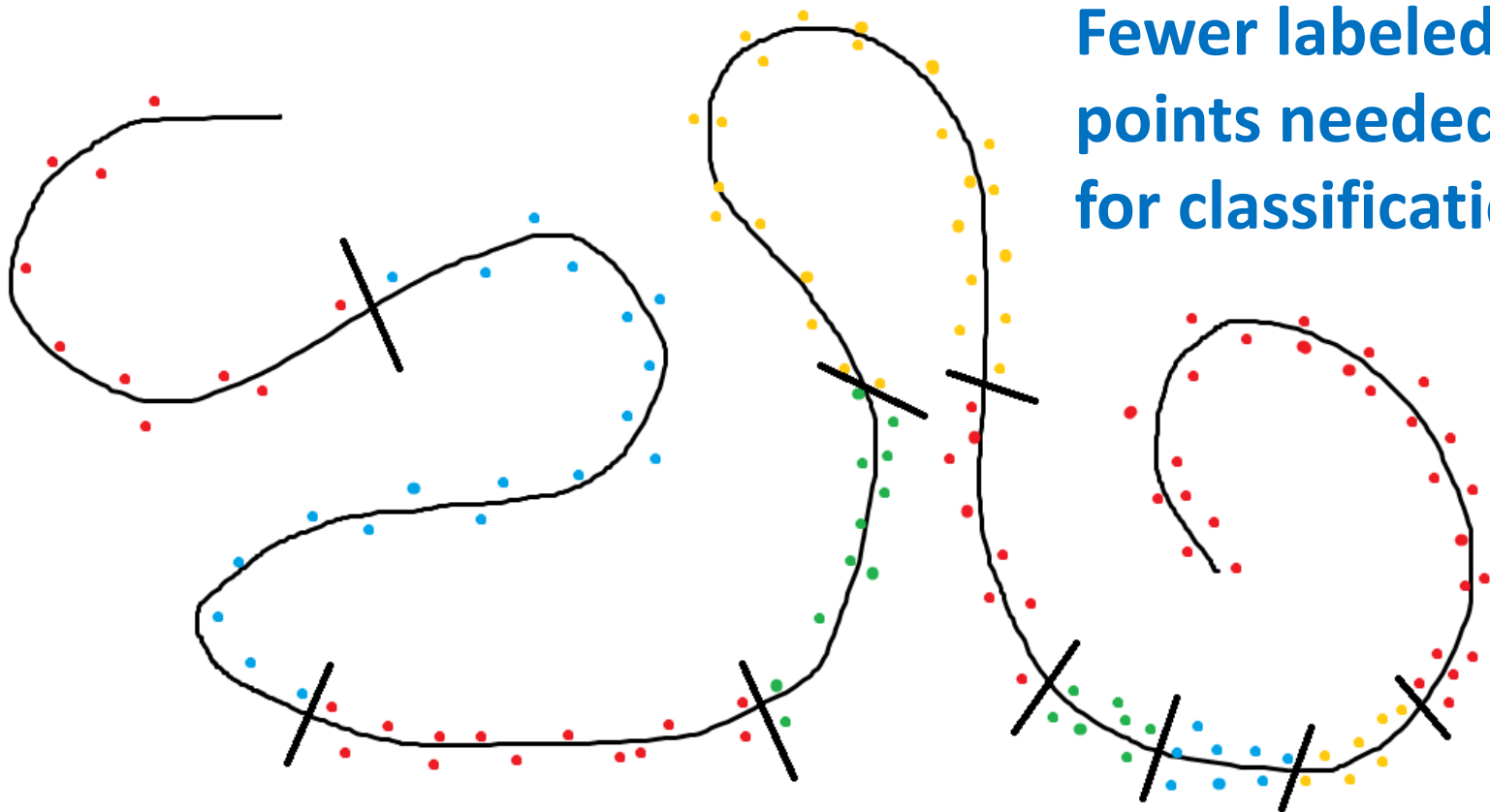
# WHY DOES DEEP LEARNING WORK?



Sparsity aids the learning of the low-dim manifold

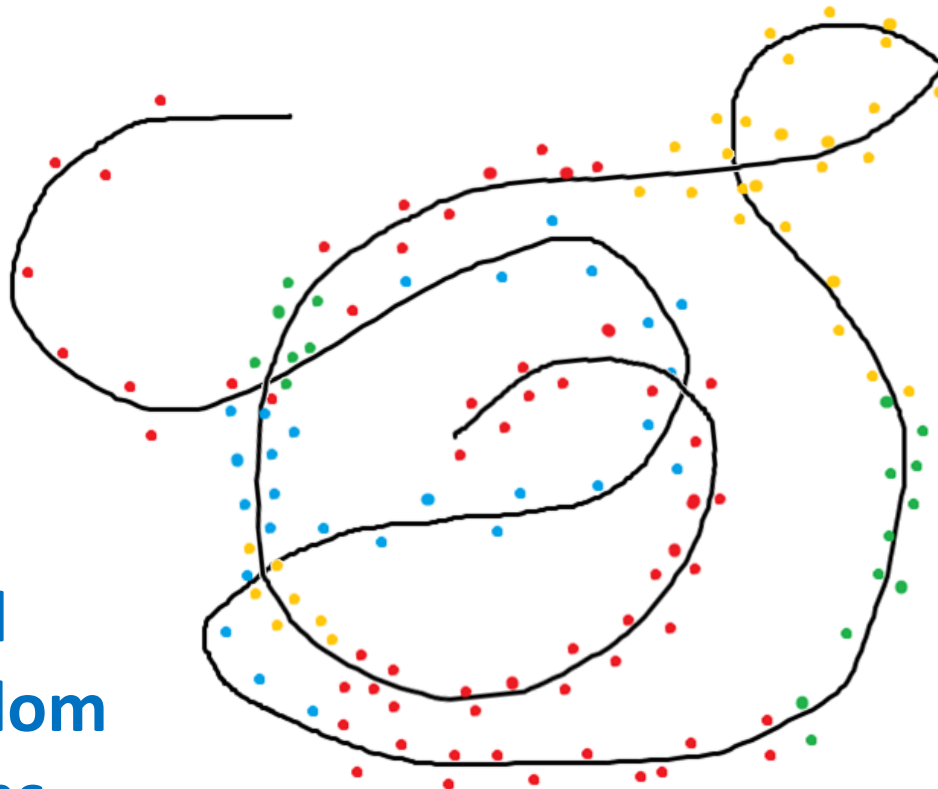


# WHY DOES DEEP LEARNING WORK?



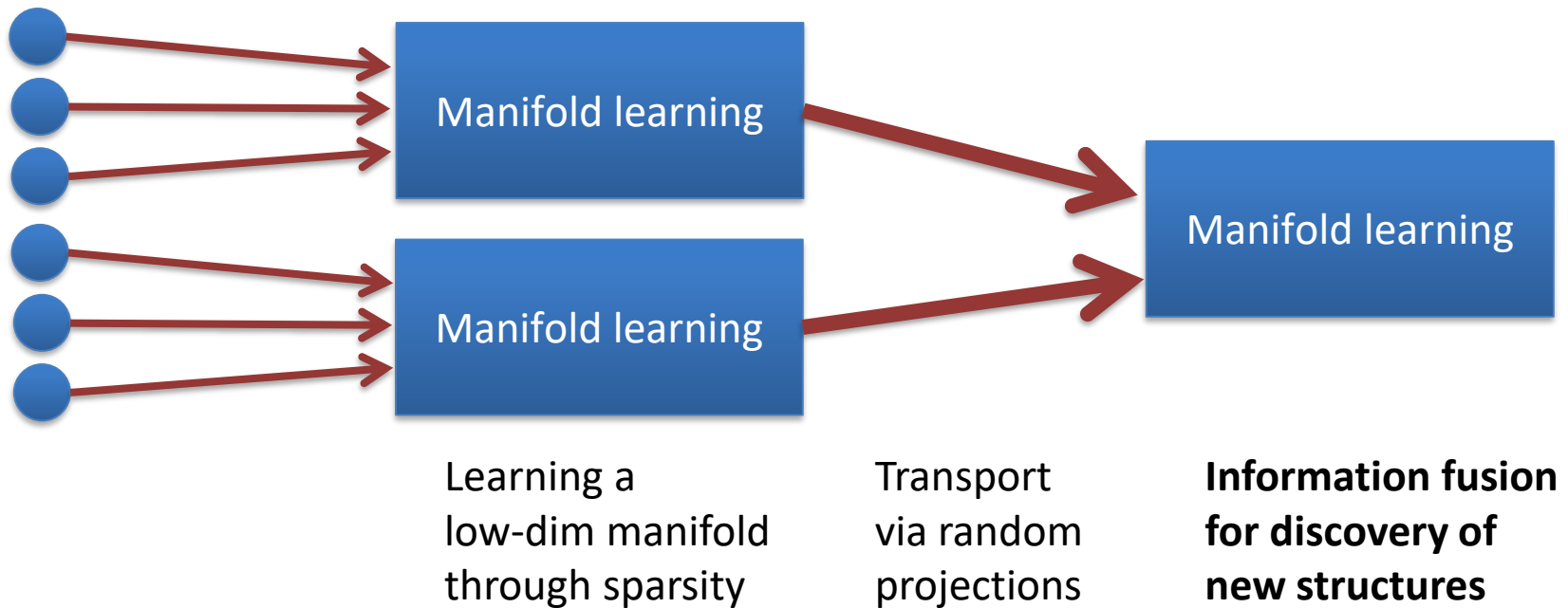
# WHY DOES DEEP LEARNING WORK?

Structure  
preserved  
after random  
projections





# HOW DOES DEEP LEARNING WORK?



**Example.** Speech recognition – from hidden markov models (one layer) to deep learning (multiple layers)

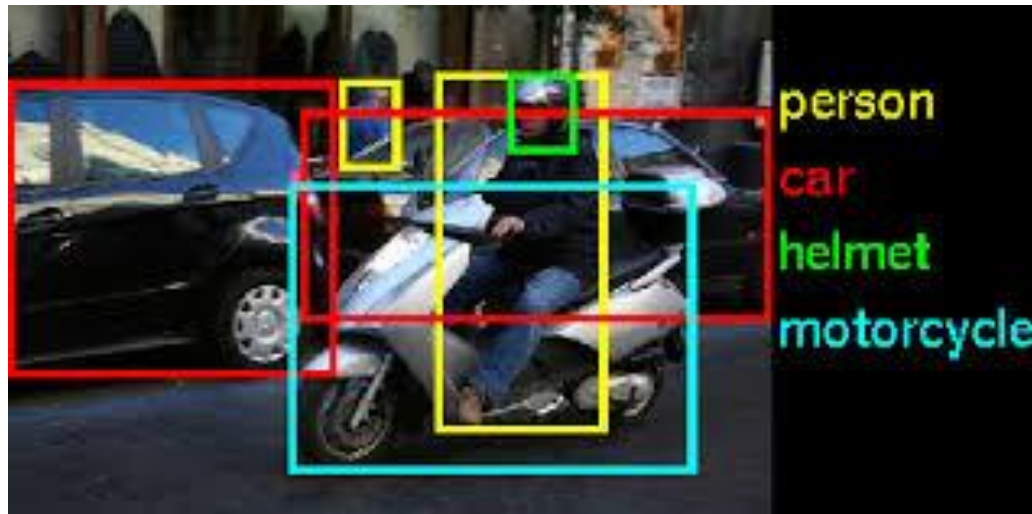




# CONVOLUTIONAL NEURAL NETWORKS



# IMAGE RECOGNITION



**Too many parameters.**

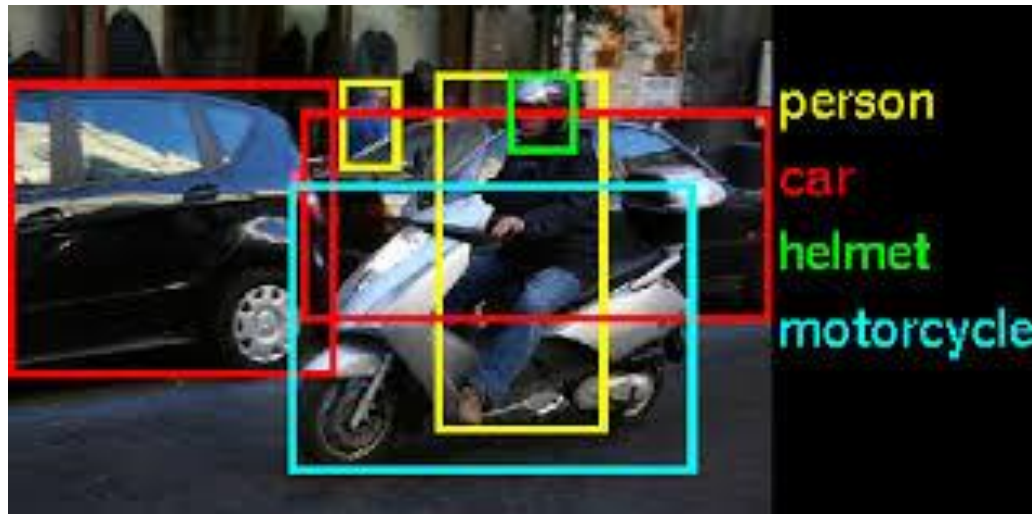
- $(\text{MNIST } 28 \times 28 \text{ pixels}) * (100 \text{ features}) = 78,400 \text{ params}$

**Translation invariance.**

- Statistics of a patch (usually) does not depend on its location.



# IMAGE RECOGNITION



## **Solution.**

- Learn good features from small image patches.
- Use features to compute better representation of image.
- Down-sample image to reduce dimensionality.
- Repeat.



# CONVOLUTION

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved  
Feature





# FEATURE MAP



Image



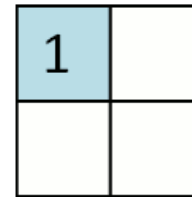
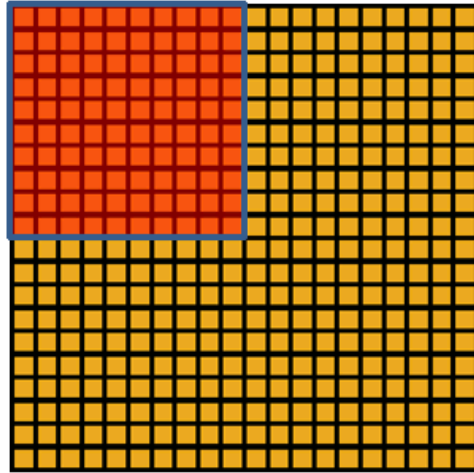
Feature Map

-1	-1	-1
2	2	2
-1	-1	-1

Filter



# POOLING



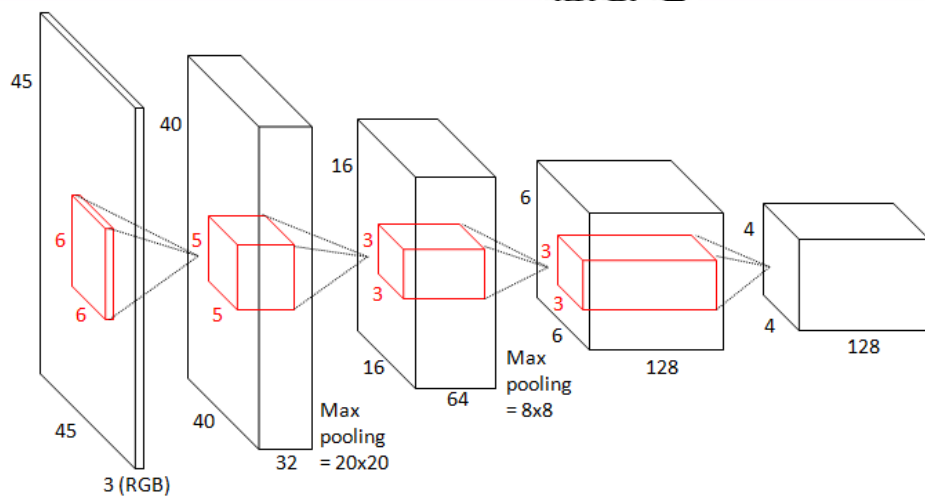
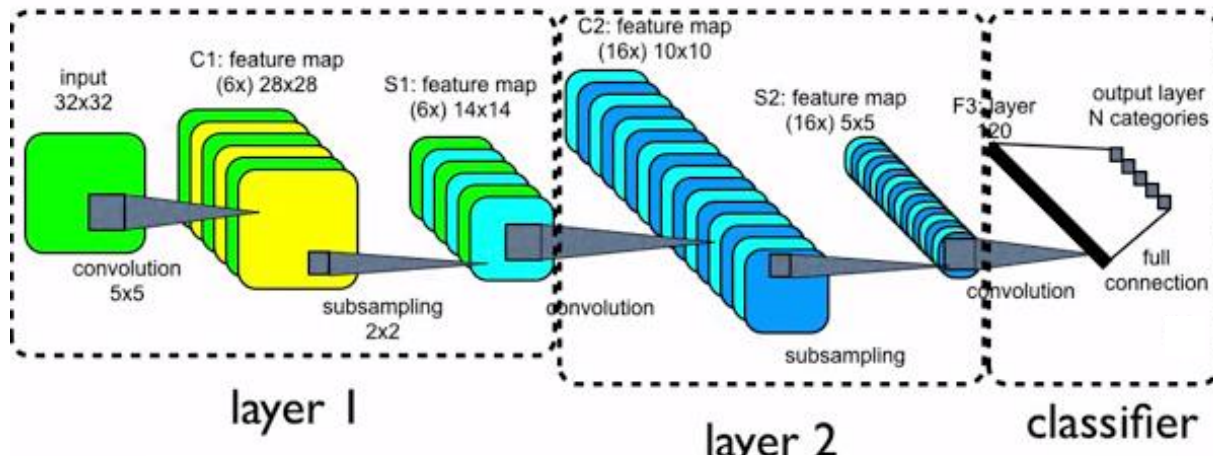
Convolved  
feature

Pooled  
feature

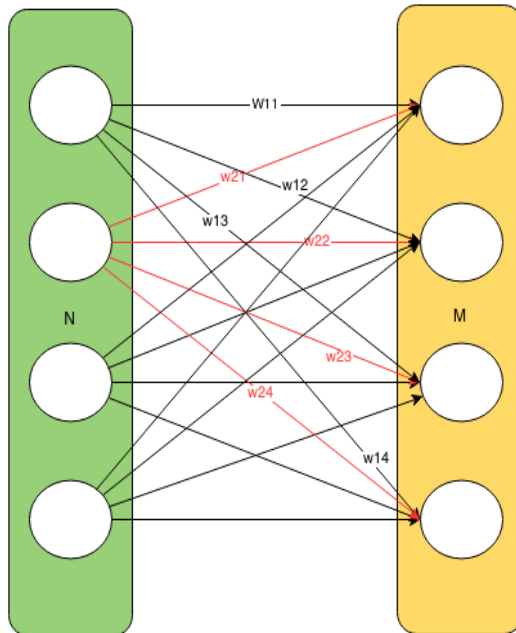
**Examples.** Mean Pooling, Max Pooling.



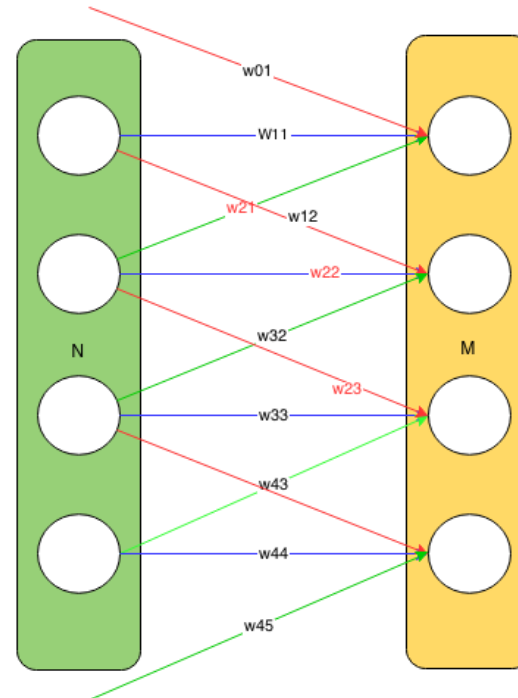
# CONVOLUTIONAL PYRAMID



# SHARED WEIGHTS



**Fully-Connected  
(Dense)**



**Locally-Connected  
with Shared Weights**



# DISCUSSION

## **Question 1.**

How is backpropagation and fine-tuning performed?

## **Question 2.**

How to introduce rotation, skew, scale invariance?







# ONLINE TUTORIAL

## Basic Usage.

- [https://www.tensorflow.org/versions/r0.11/get\\_started/basic\\_usage.html](https://www.tensorflow.org/versions/r0.11/get_started/basic_usage.html)

## Simple MNIST Example.

- <https://www.tensorflow.org/versions/r0.11/tutorials/mnist/beginners/index.html#mnist-for-ml-beginners>

## Convolutional Neural Networks.

- <https://www.tensorflow.org/versions/r0.11/tutorials/mnist/pros/index.html#deep-mnist-for-experts>



# SUMMARY

- Multilayer Neural Networks
  - Neuron
  - Activation Function
  - Forward Propagation
- Convolutional Neural Networks
  - Convolution
  - Pooling
  - Shared Weights
- Learning Algorithm
  - Cost Function
  - Backpropagation
  - Sparse Autoencoder
  - Greedy Initialization

