

Introduction to Deep Learning

- Image Recognition

***Bargava Subramanian (@bargava)
Data Scientist, Cisco Systems***

I WANT YOU TO
CREATE ARTIFICIAL
INTELLIGENCE THAT
IS AS SMART AS ME.



Dilbert.com DilbertCartoonist@gmail.com

OKAY. I SHOULD
HAVE THAT BY
LUNCHTIME.



7-4-14 © 2014 Scott Adams, Inc./Dist. by Universal Uclick

BECAUSE
YOU'RE A
FAST
WORKER?



SURE.

AGENDA/SCHEDULE

- 1) Motivation to Machine Learning/Deep Learning
 - i. Biological Motivation, Hierarchical/Representation Learning
- 2) Introduction to Artificial Neural Networks/Deep Learning
 - i. Neuron, Perceptron, Logistic, MLP, Rectified Linear Units
 - ii. Backpropagation Algorithm, Gradient Descent(including SGD), Mini-batch
- 3) Image Recognition: Convolution Neural Networks
 - i. Convolution
 - ii. Sub-sampling, Pooling
 - iii. Dropout
 - iv. Architecture
- 4) Challenges in Deep Learning
 - i. Vanishing Gradients & Local Minima
 - ii. Overfitting

A photograph of the ancient Incan city of Machu Picchu, perched high in the Andes mountains of Peru. The city is built on a series of stone terraces that follow the steep mountain slopes. The architecture is intricate, featuring small stone buildings, walkways, and agricultural terraces. The sky is filled with thick, white clouds and mist, which partially obscure the peaks of the surrounding mountains. The overall atmosphere is mysterious and historical.

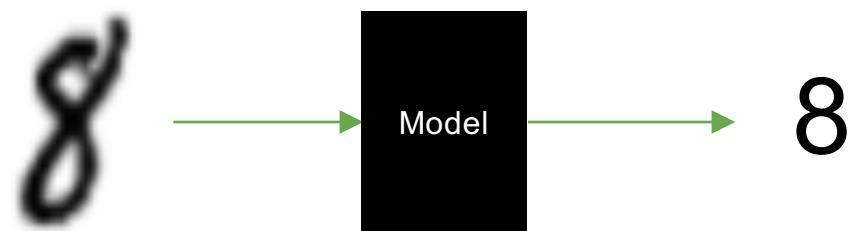
LEARNING

What is Learning?

How do we recognize the digits?

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

Machine Learning Framework



Inputs

Computation

Outputs

Recognizing Digit - An algorithm

How do we recognize the digits?

Use functions that compute relevant information to solve the problem

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

k Nearest-Neighbors

For each image, find “most similar” image. Guess that as the label.

Recognizing Digit – An algorithm

How do we recognize the digits?

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

Difficult to enumerate all possible interactions,
spatial structure, etc. as hand-coded features.

Can we think of another way?

6	7	4	6	8	0	7	8	3	1
---	---	---	---	---	---	---	---	---	---



LEARNING – Biological Inspiration

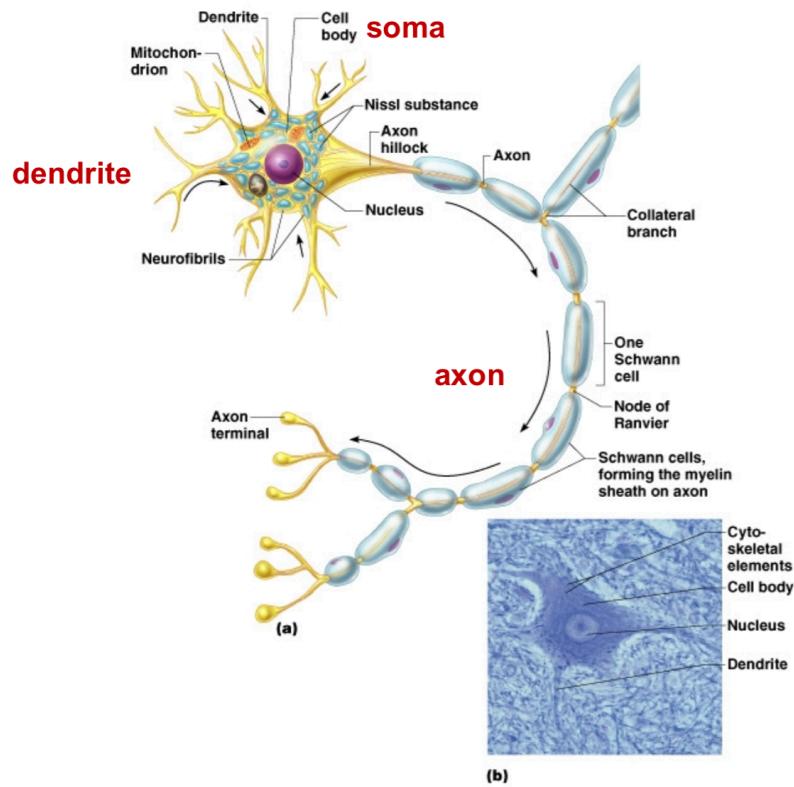
Questions about Learning

How is information detected?

How is it stored?

How does it influence recognition?

Brain



- Connected network of neurons.
 - Communicate by electric and chemical signals
- ~ 10^{11} neurons
~ 1000 synapses per neuron

- Signals come in via dendrites into soma
- Signal goes out via axon to other neurons through synapses

Learnings from Neuro & Cognitive Science



Kids talk grammatically correct sentences even before they are taught formal language.

Kids learn after listening to a lot of sentences

→ Associations and Structural inferences.
Understand context. Eg: Drinking water Vs River Vs Ocean

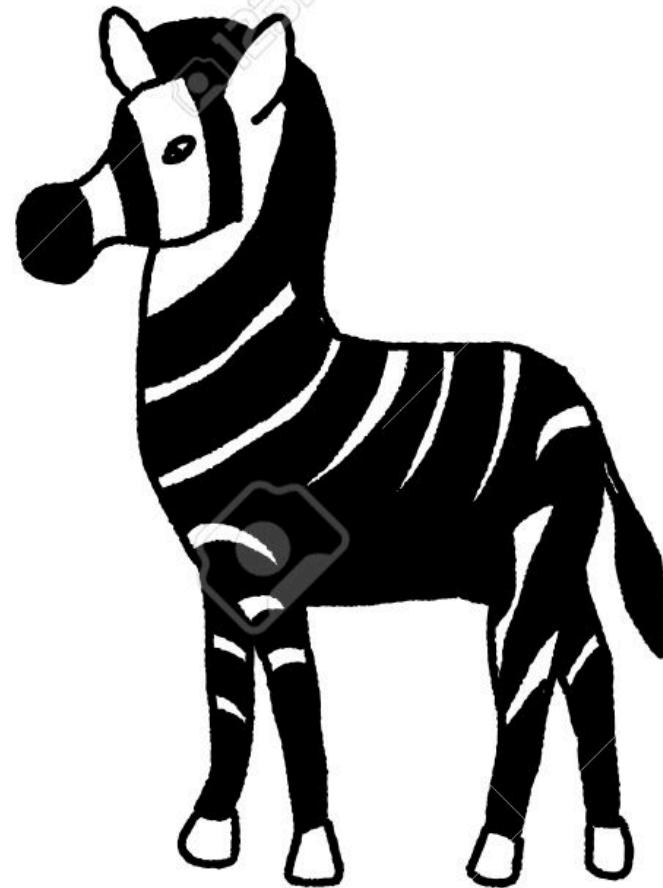
**See/hear/feel first. Assimilate.
Build the context hierarchically.
Recognize. Respond.**

A photograph of a person standing on the edge of a massive, flat-topped rock formation, likely Trolltunga in Norway. The person is wearing a pink top and dark pants, standing in a dynamic pose. Below them is a deep, winding fjord with turquoise water, surrounded by steep, green-covered mountains under a clear sky.

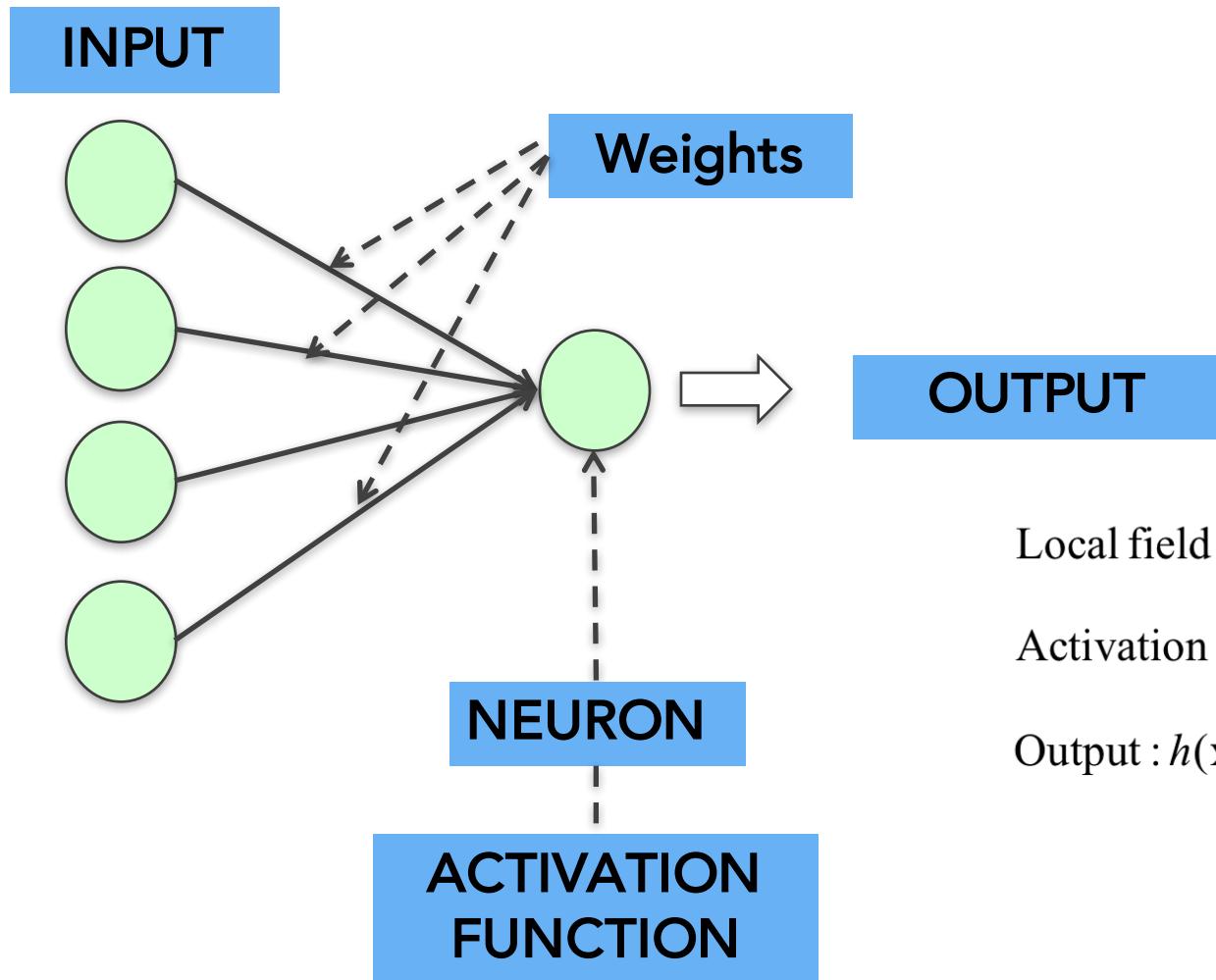
Neural Network – Building Blocks

DID
YOU
KNOW?

ZEBRAS ARE
ACTUALLY BLACK WITH WHITE STRIPES,
NOT WHITE WITH BLACK STRIPES.



Neuron, Activation Function

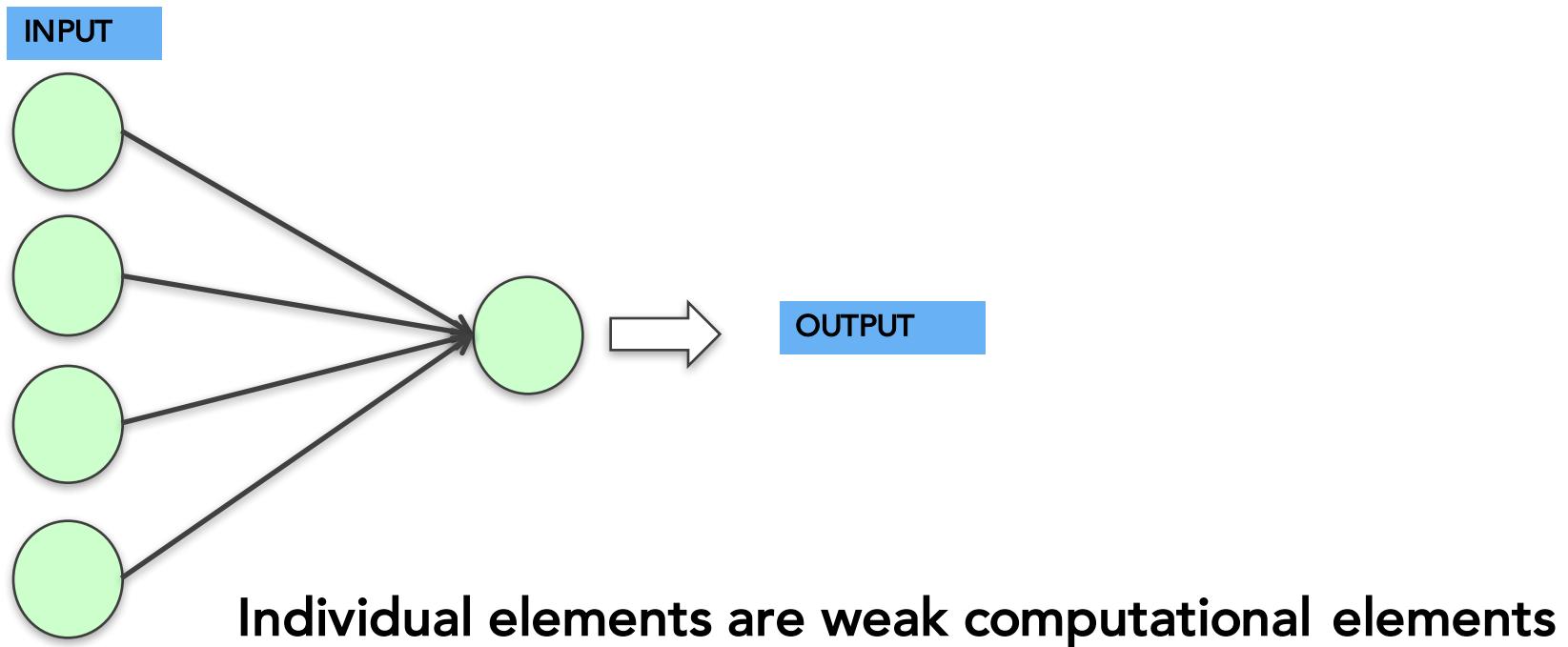


$$\text{Local field} : \sum_{d=0}^D w_d x_d$$

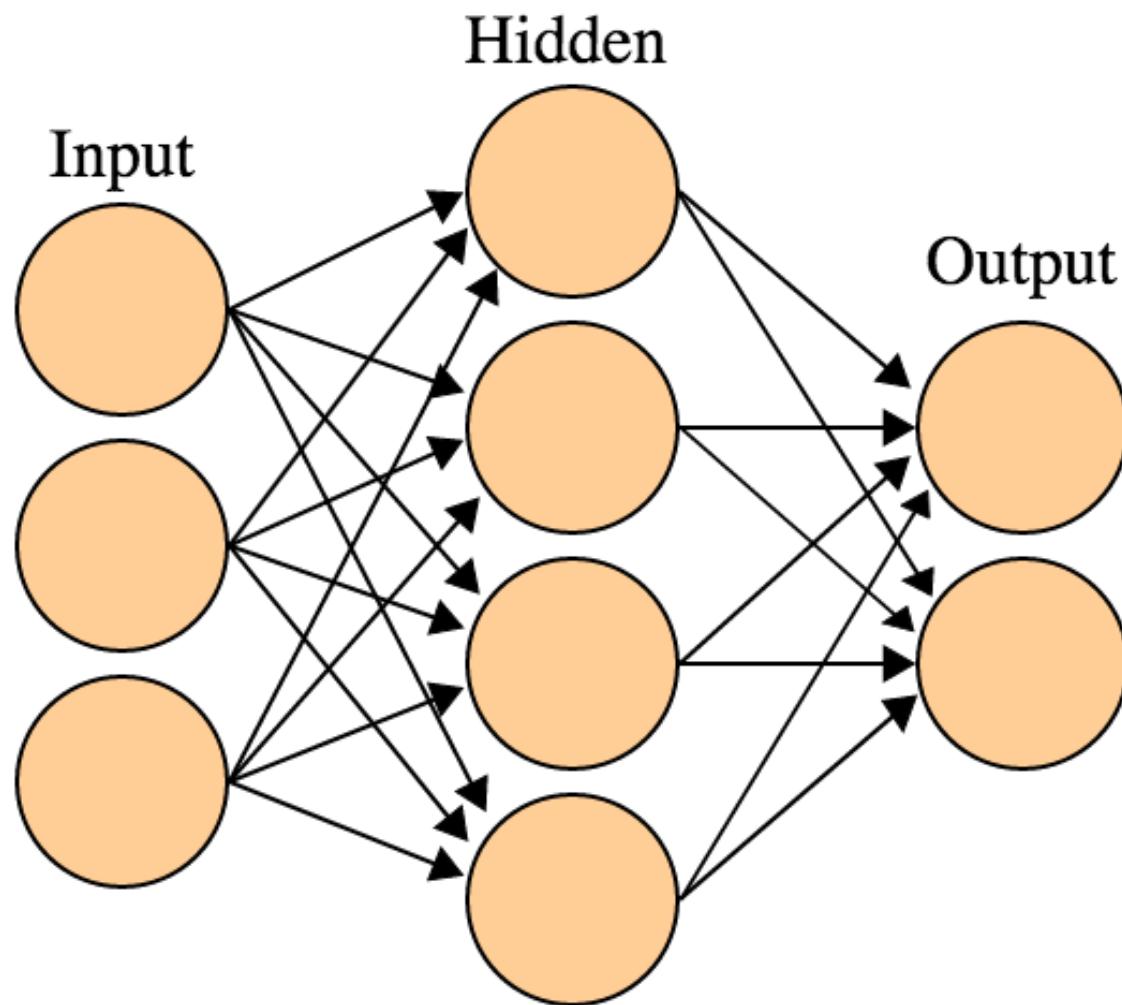
Activation function : $\varphi(\cdot)$

$$\text{Output} : h(\mathbf{x}) = \varphi\left(\sum_{d=0}^D w_d x_d\right)$$

Need: Networked Units



A Simple Neural Network

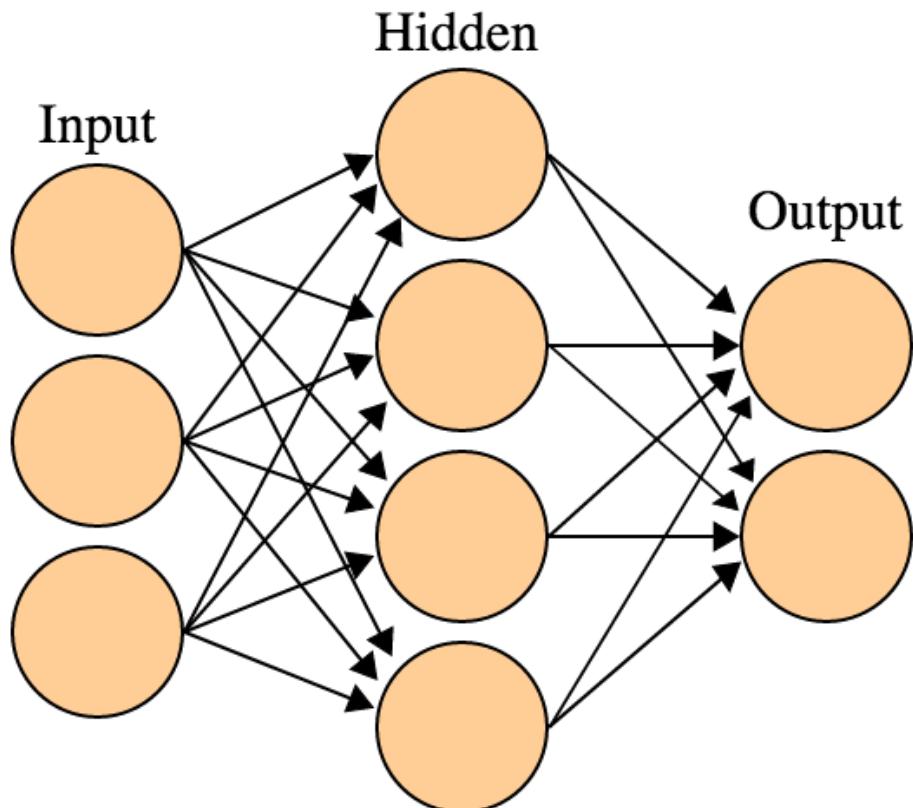


$$\text{Local field} : \sum_{d=0}^D w_d x_d$$

Activation function : $\varphi(\cdot)$

$$\text{Output} : h(\mathbf{x}) = \varphi\left(\sum_{d=0}^D w_d x_d\right)$$

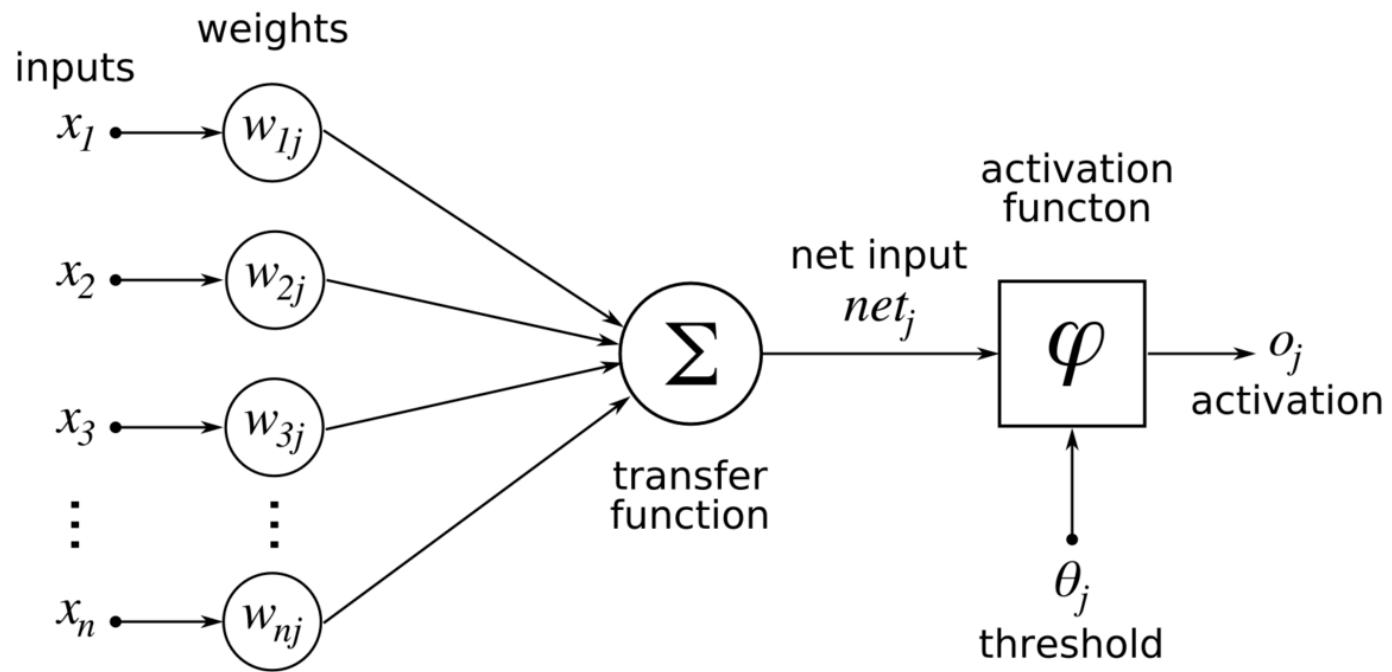
Feed-forward Neural Network



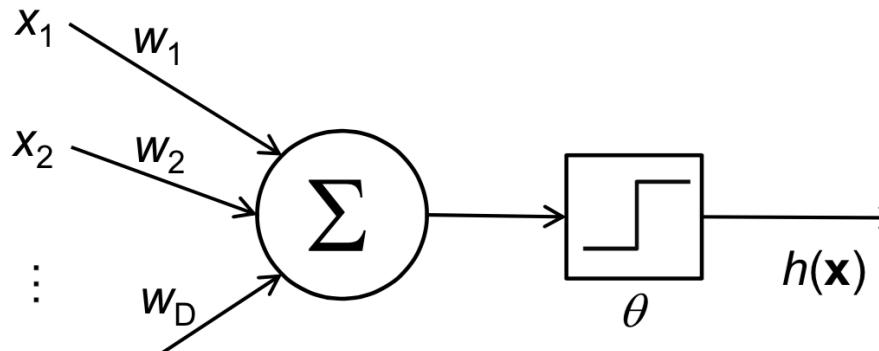
- » The network is formed by
 - **Input layer** of source nodes
 - **One or several hidden layers** of processing neurons
 - **Output layer** of processing neuron(s)
- » Connections only between **adjacent layers**
- » There are **no feedback** connections

Activation Function

activation function of a node defines the output of that node given input(s)



Perceptron

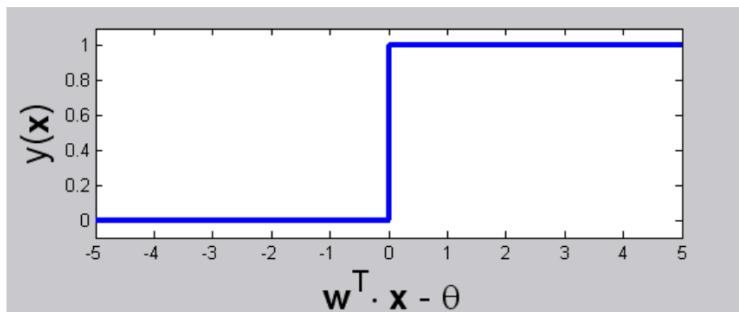


$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } w_1x_1 + w_2x_2 + \dots + w_Dx_D \geq \theta \\ 0 & \text{if } w_1x_1 + w_2x_2 + \dots + w_Dx_D < \theta \end{cases}$$

threshold : θ

In vector notation:

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \cdot \mathbf{x} \geq \theta \\ 0 & \text{if } \mathbf{w}^T \cdot \mathbf{x} < \theta \end{cases}$$

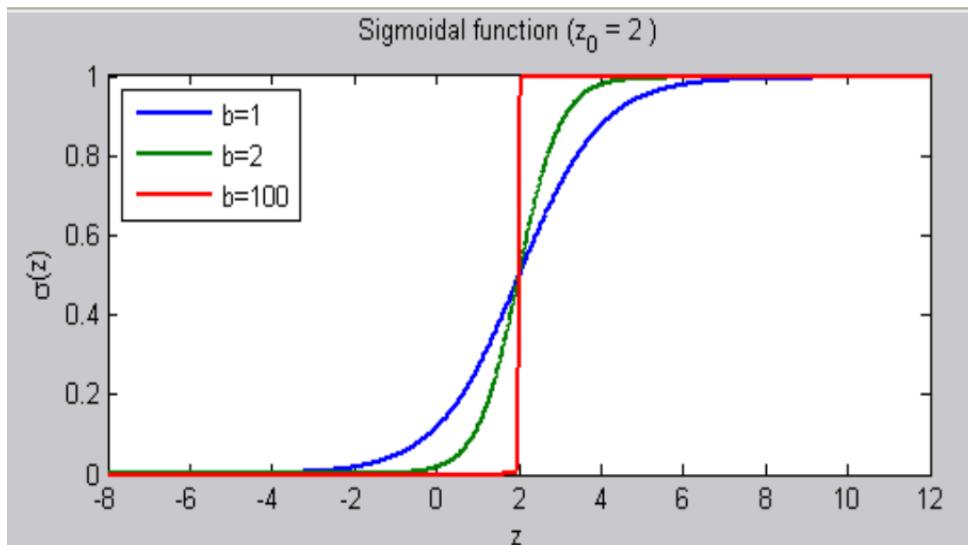


- Invented in 1957.
- Classifies input data into one of the output classes.
- Online learning possible

If the weighted input is more than the threshold, classify as 1.
Else 0

McCulloch, W. S., and Pitts, W. "A Logical Calculus of Ideas Immanent in Nervous Activity." Bulletin of mathematical biophysics, 5, pp. 115-133 (1943). Reprinted in McCulloch, W. S., *Embodiments of mind*. Cambridge, MA: MIT Press.

Sigmoid/Logistic



$$\sigma(z) \equiv \frac{1}{1+e^{-z}}; \quad [\text{logistic sigmoid}]$$

Local field : $\mathbf{w}^T \cdot \mathbf{x} = \sum_{d=0}^D w_d x_d$

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

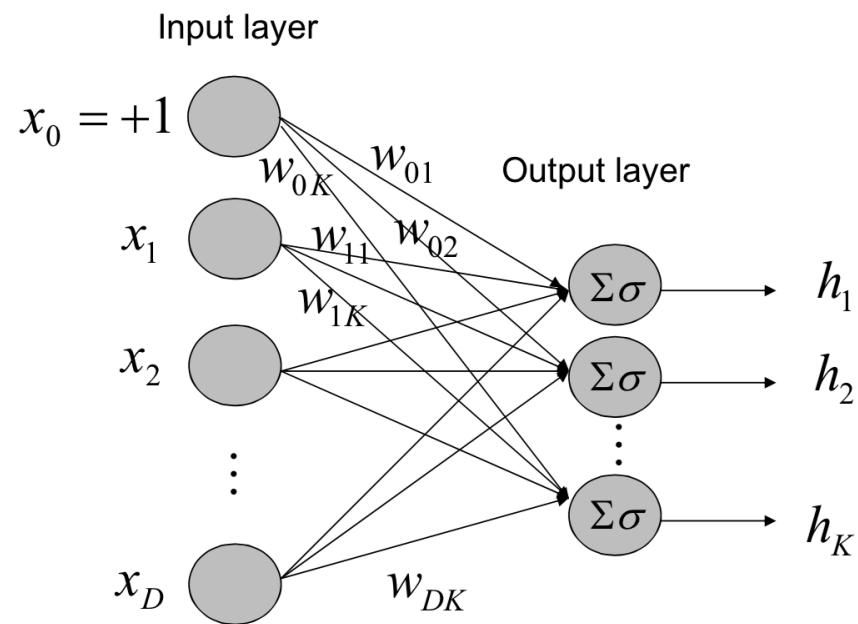
Output : $h(\mathbf{x}) = \sigma(\mathbf{w}^T \cdot \mathbf{x})$

- Output is bounded between 0 & 1
- Domain: Complete set of Real numbers
- Smooth and continuous function.

- Symmetric
- Derivative can be quickly calculated
- Positive, Bounded, strictly positive

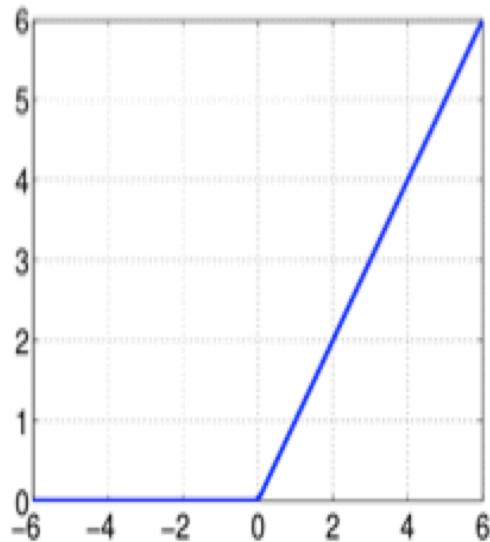
Softmax

Generalization of logistic regression
for Multi-class classification



Rectified Linear Units

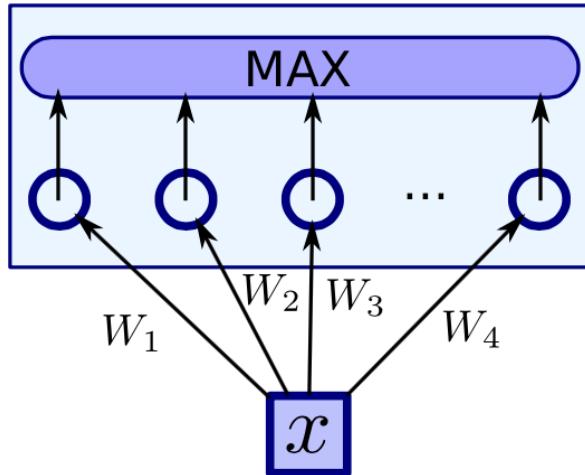
$$f(x) = \max \{0, x\}$$



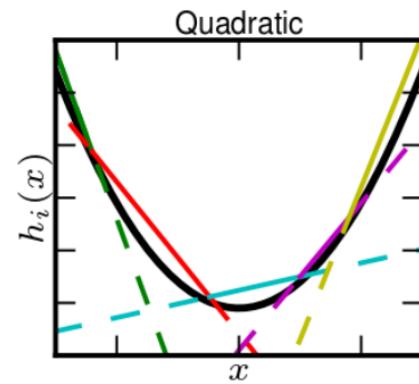
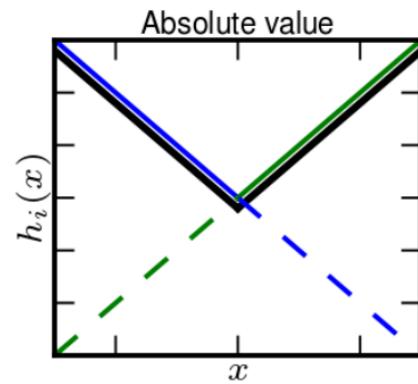
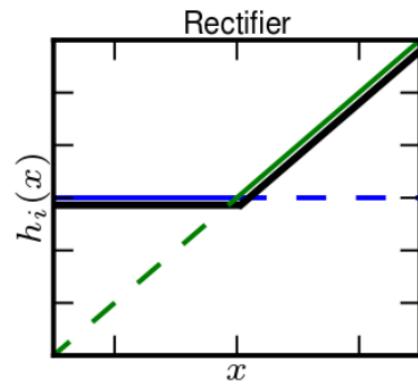
- Cheap to compute (no products/exponentials)
- Faster training
- Sparser networks
- Bounded below 0
- Strictly increasing

Max-out

$$f(x) = \max_{i=1}^k (x^T W_i + b_i)$$

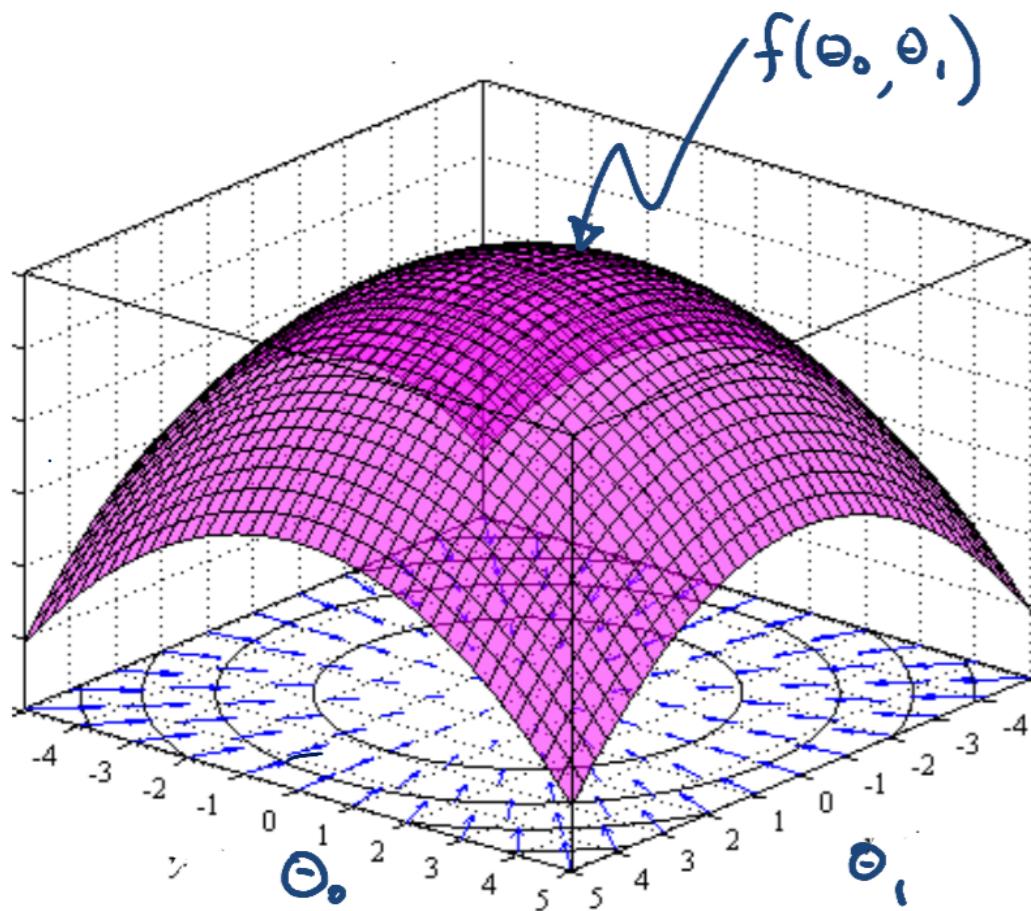


- Generalization of Rectified Linear
- Max of k linear functions -> piece-wise linear
- At large k, can approximate a non-linear function



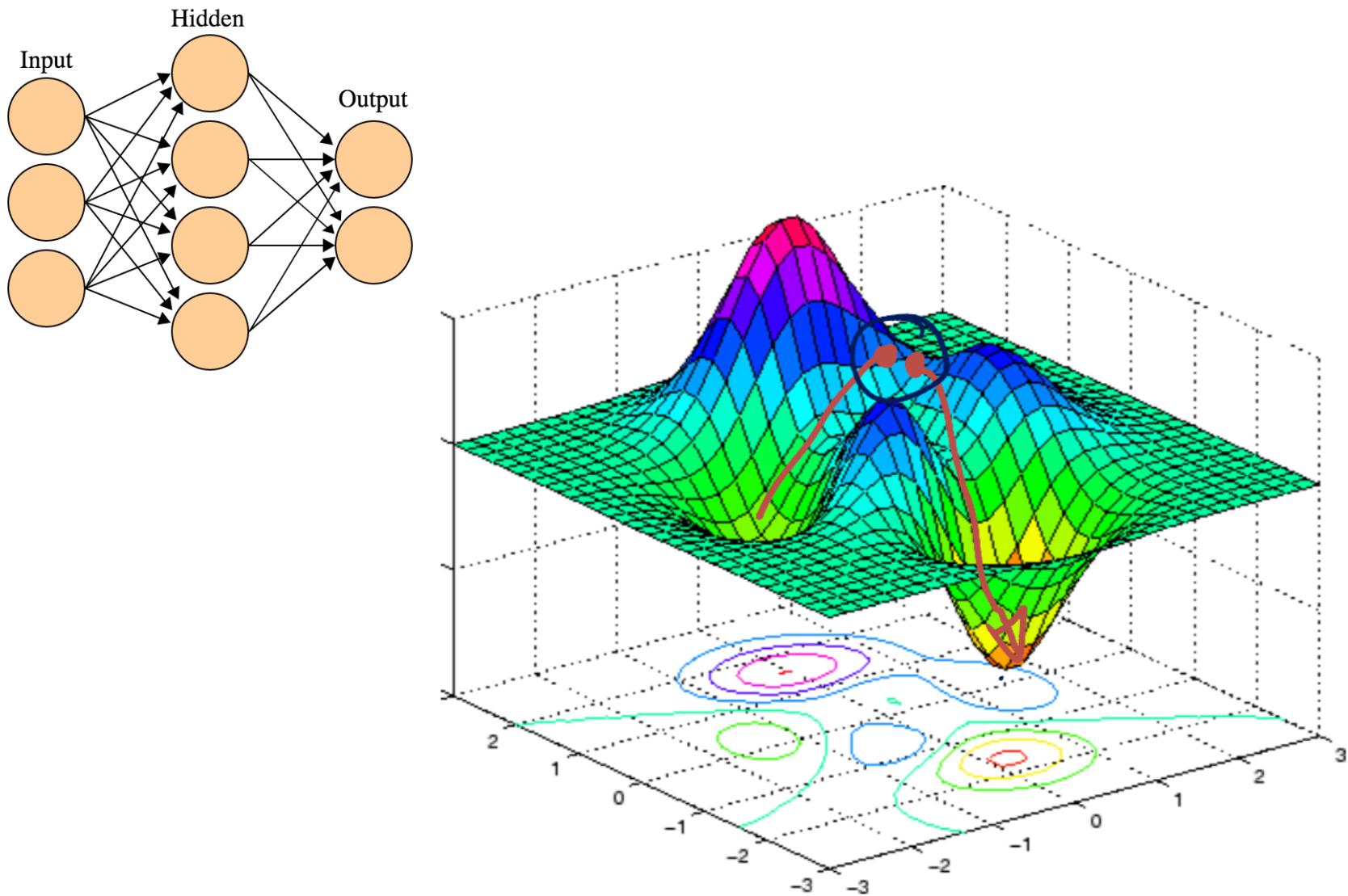
Learning in ANN - Gradient Descent

Goal: To find minimum of the loss function (minimize error of the model)



Repeat until convergence {
 $\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$
}

Er.... How to compute gradient?



Backpropagation Algorithm – Pseduo Code

Algorithm 1 Backpropagation learning algorithm

for d in data **do**

FORWARDS PASS

Starting from the input layer, use eq. 1 to do a forward pass trough the network, computing the activities of the neurons at each layer.

BACKWARDS PASS

Compute the derivatives of the error function with respect to the output layer activities

for layer in layers **do**

Compute the derivatives of the error function with respect to the inputs of the upper layer neurons

Compute the derivatives of the error function with respect to the weights between the outer layer and the layer below

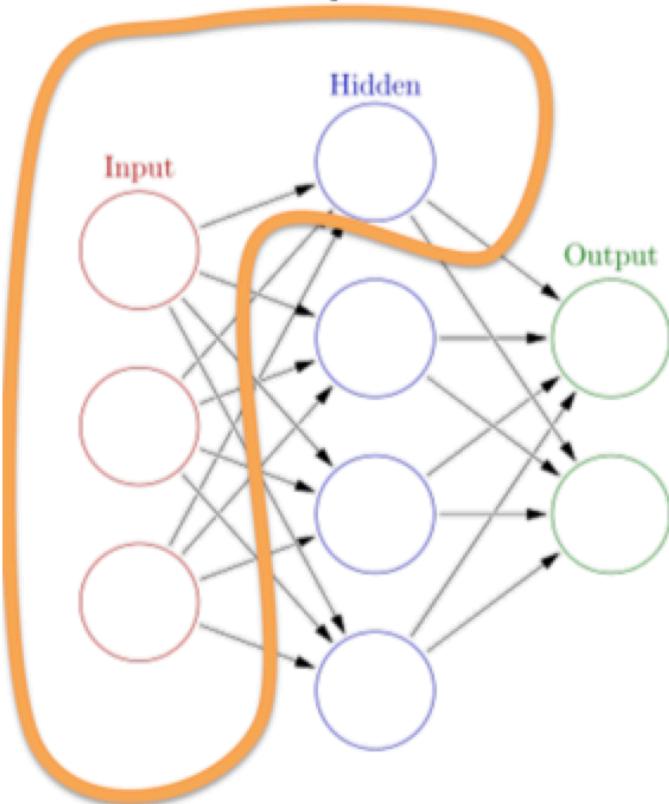
Compute the derivatives of the error function with respect to the activities of the layer below

end for

Updates the weights.

end for

Backpropagation Algorithm



Computes gradient of the loss function w.r. to the weights.

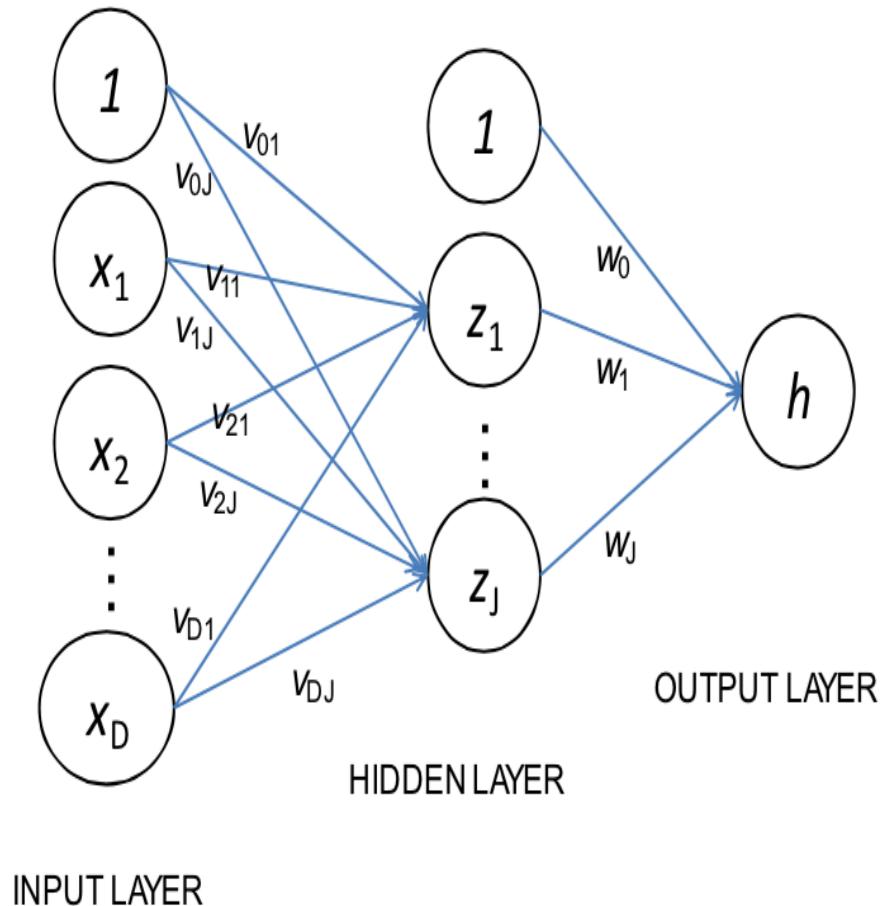
Backpropagate training error to generate deltas of all the neurons from hidden layers to output layer

Use gradient descent to update weights

SGD/Mini-Batch/Online

- ❑ Stochastic Gradient Descent: Instead of using all of the training data, train iteratively on “mini-batches”
- ❑ Online Learning: Mini-batch size is 1. Weights are adjusted for every single data point.

Multi-Layer Perceptron



- Feedforward ANN
- Activation function:
Mostly sigmoid
- Improvement over
basic perceptron: Can
classify data that aren't
linearly separable

SOME DEEP LEARNING EXAMPLES



Video Classification

Google's neural net learns just by watching youtube videos

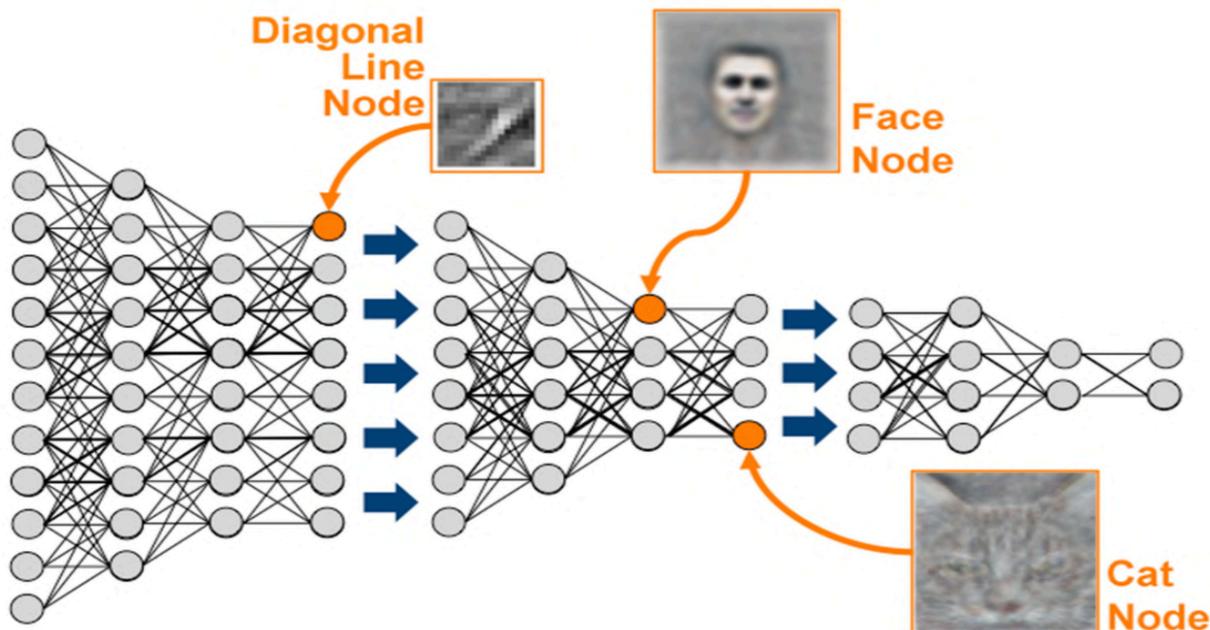


Image Recognition



Image Generation – Google Inceptionism



source: <http://googleresearch.blogspot.in/2015/06/inceptionism-going-deeper-into-neural.html>



HANDS-ON

CONVOLUTIONAL NEURAL NETWORKS

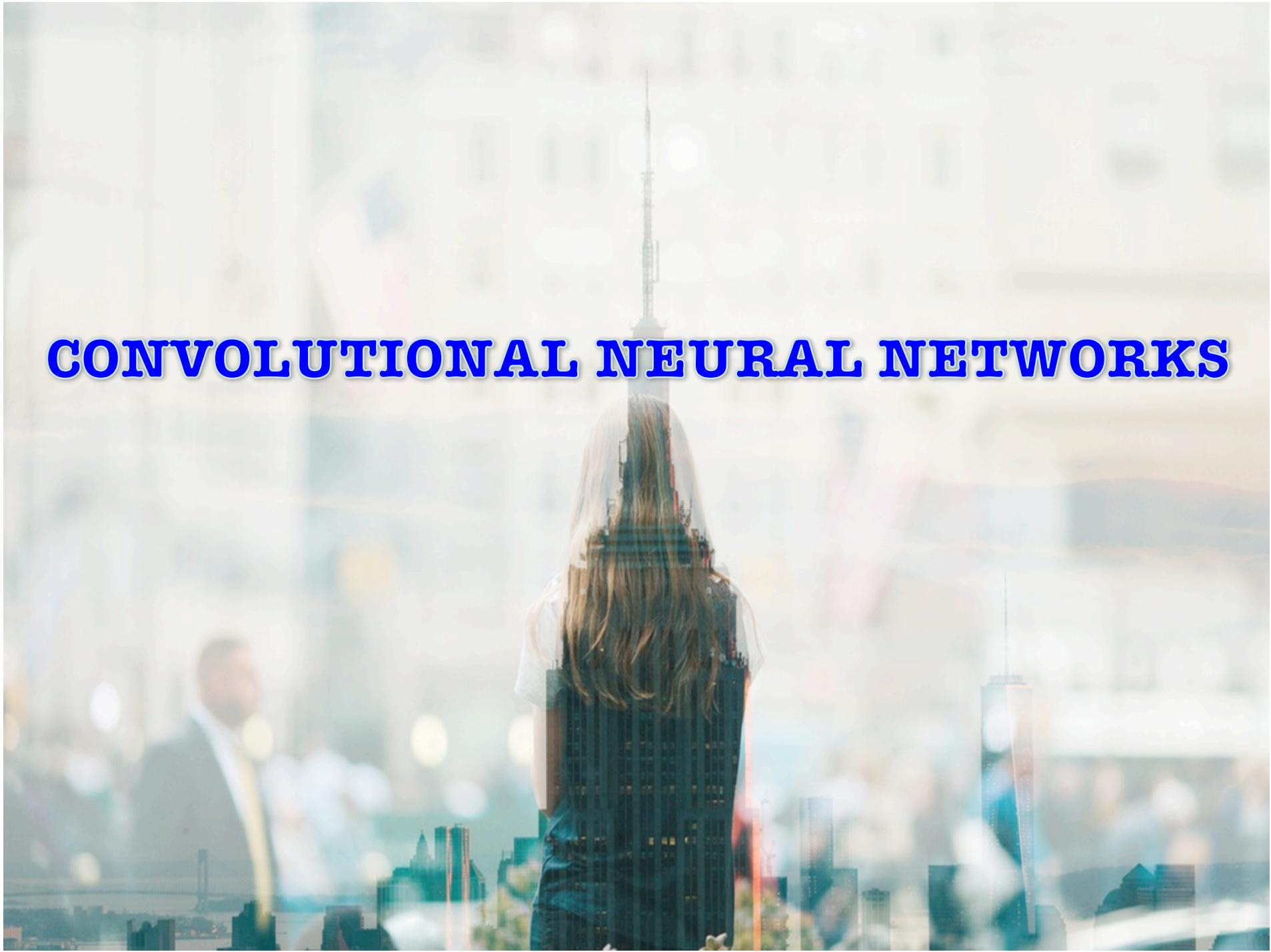
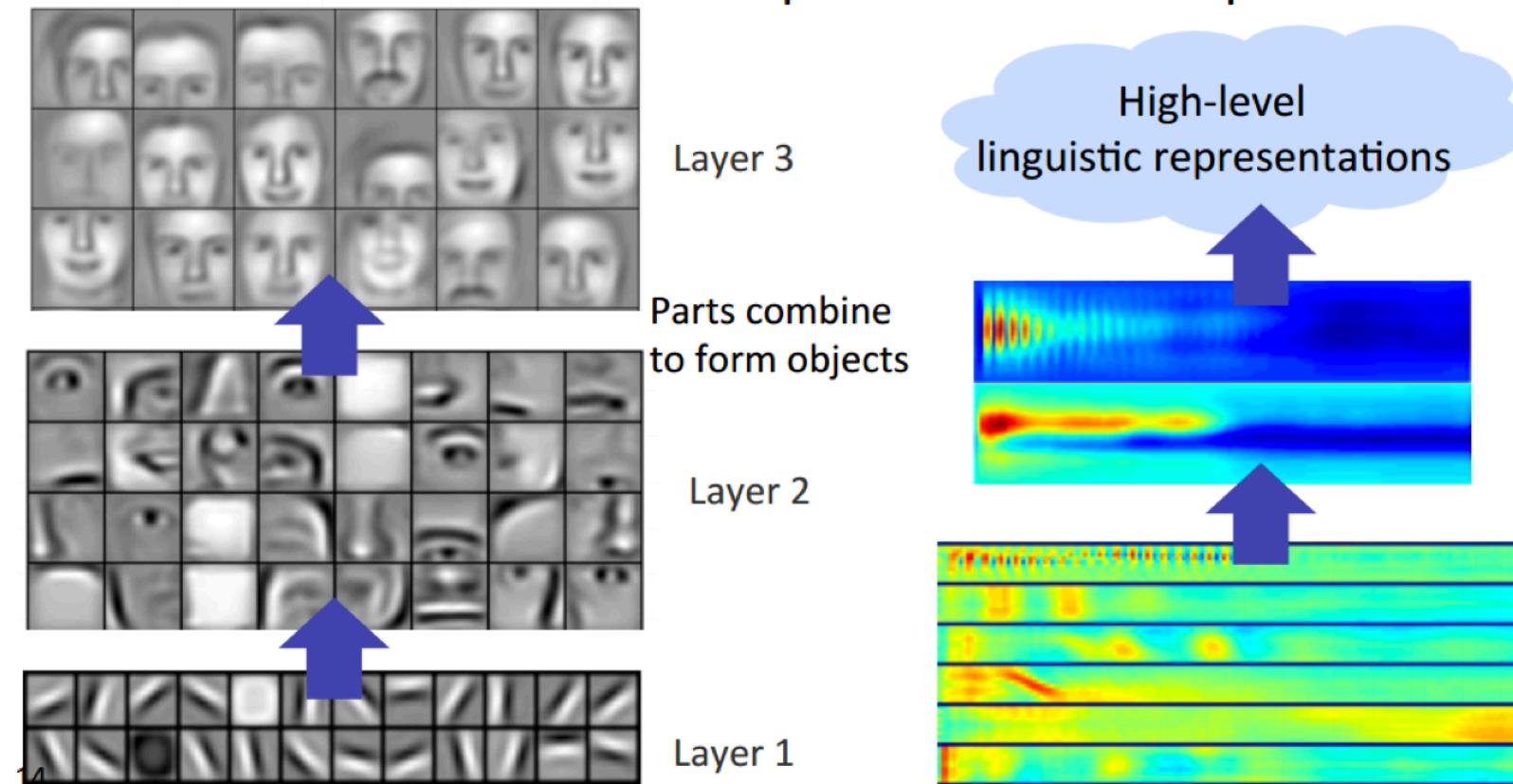


Image Recognition/Classification

Successive model layers learn deeper intermediate representations



Prior: underlying factors & concepts compactly expressed w/ multiple levels of abstraction

Convolution

Image Processing Technique to change intensities of a pixel to reflect the intensities of the surrounding pixels.
Eg: image effects like blur, sharpen, and edge detection



Original

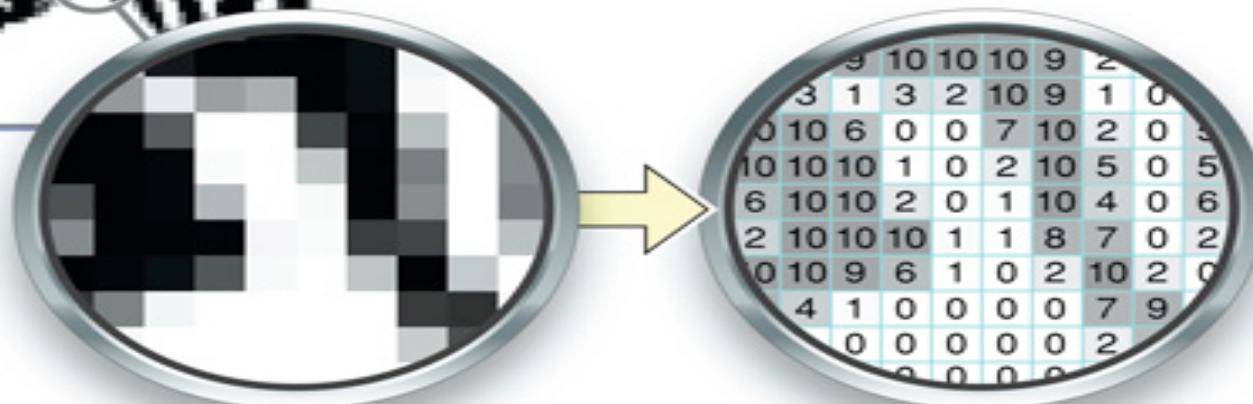


Emboss

Convolution

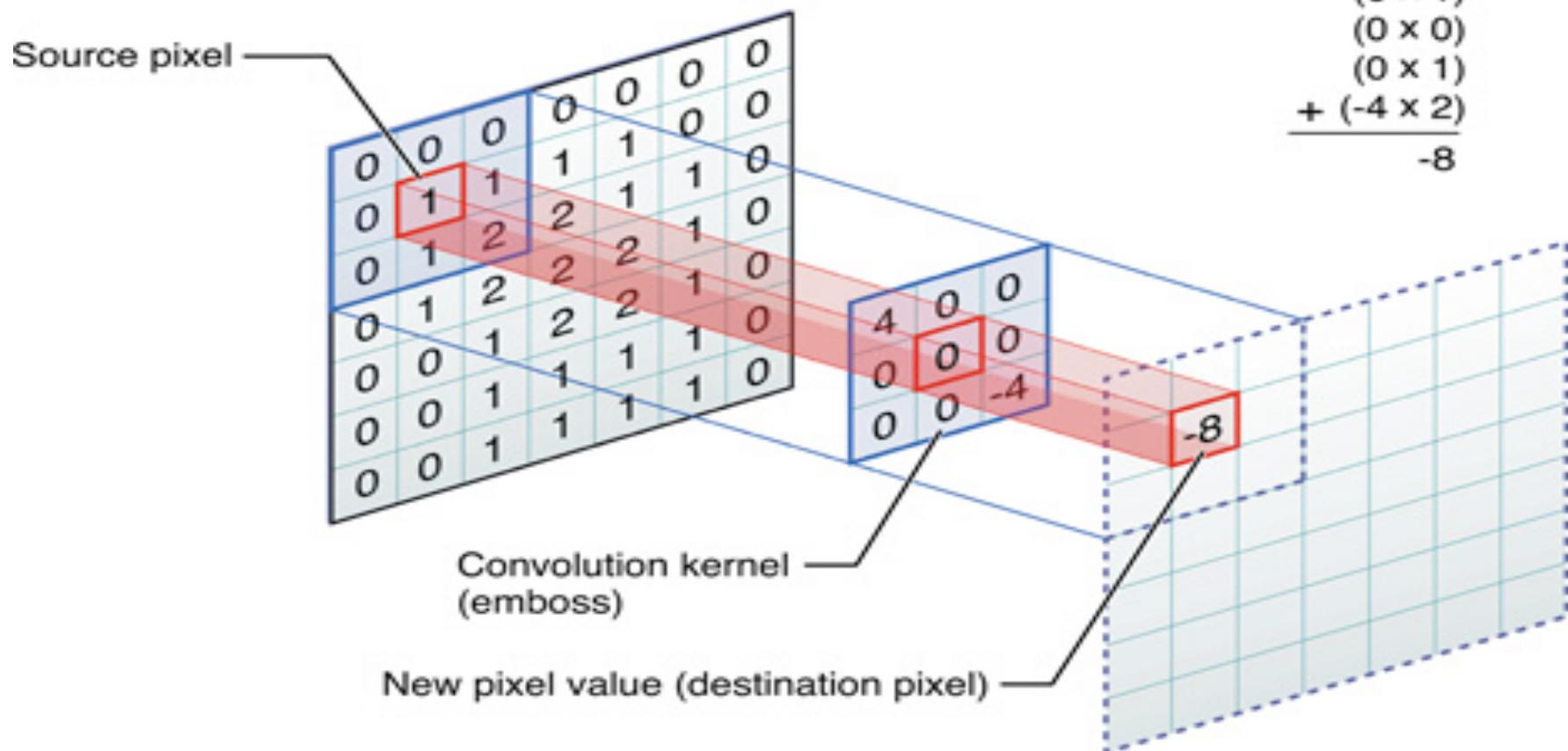


Pixels depicted by a grid of numbers representing intensity



Convolution

Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.

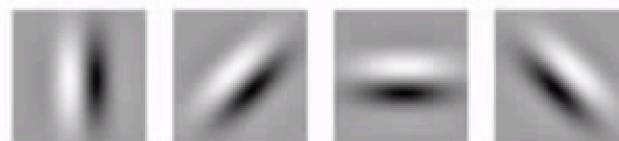


Convolution – One more example

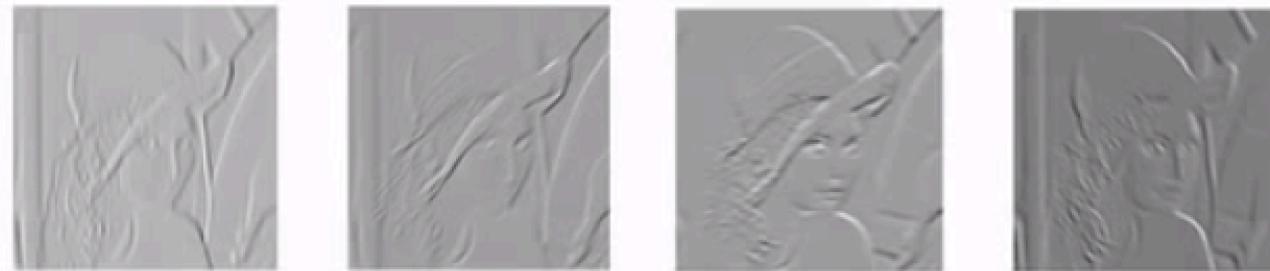
Lena



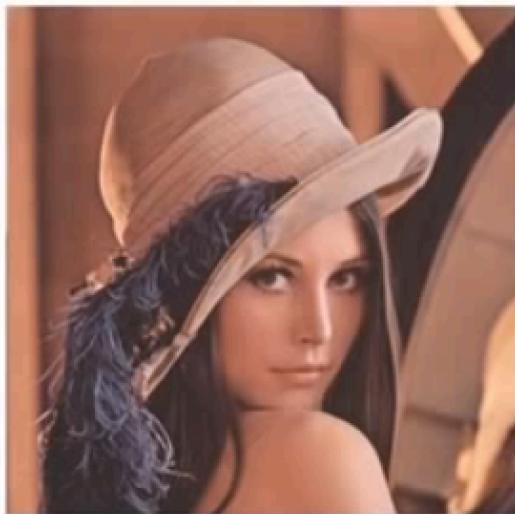
Gabor filters



Convolved by Gabor filters lena



Sub-sampling



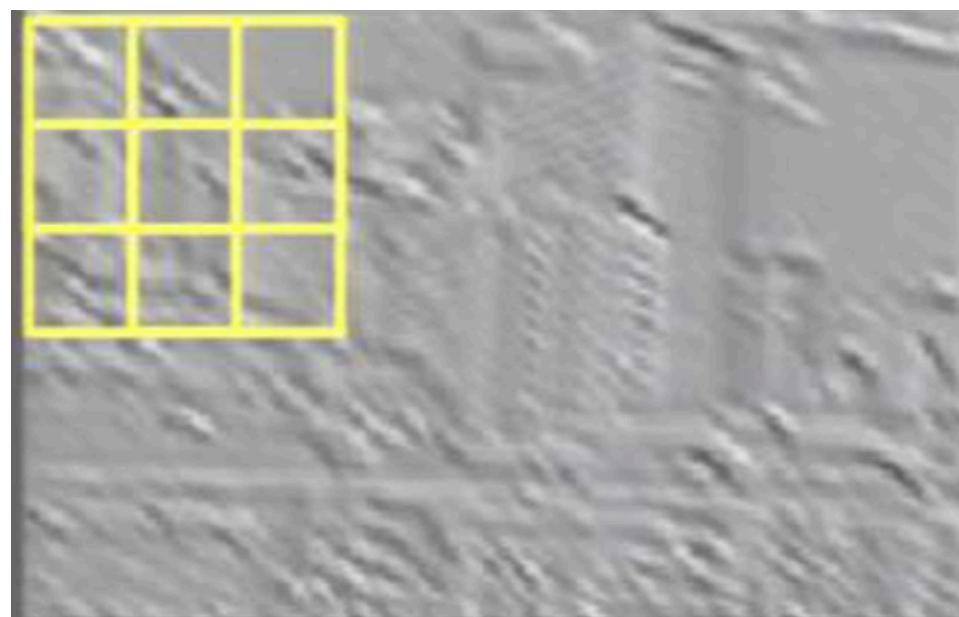
Input 508x508 pixels Lena

down-sampled by 2x2
=====>



Input 254x254 pixels Lena

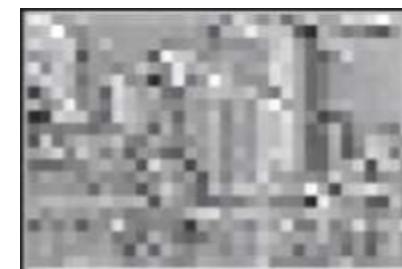
Pooling



Max

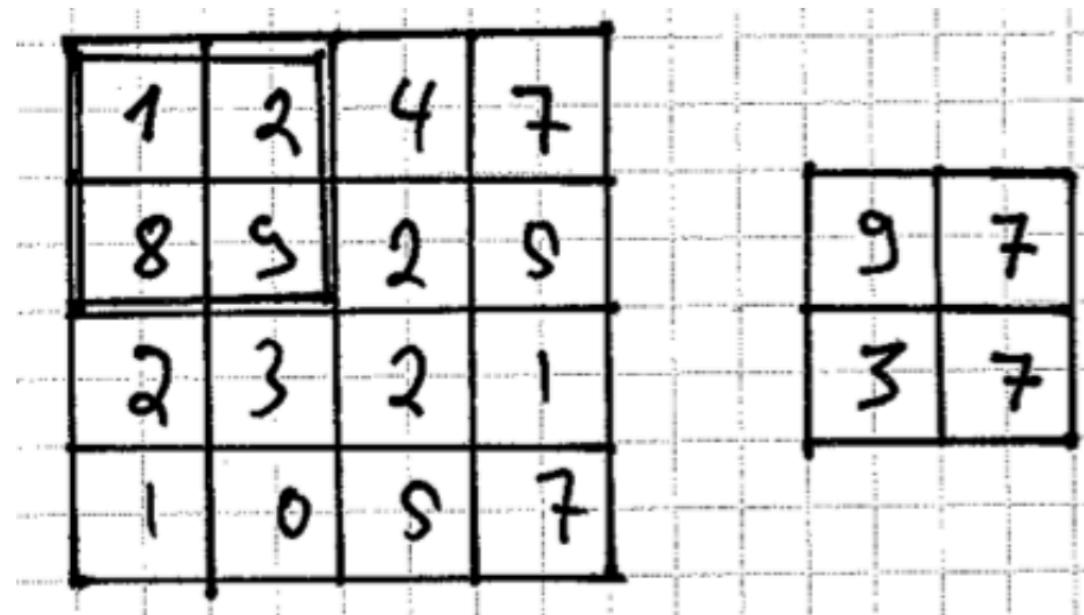


Sum



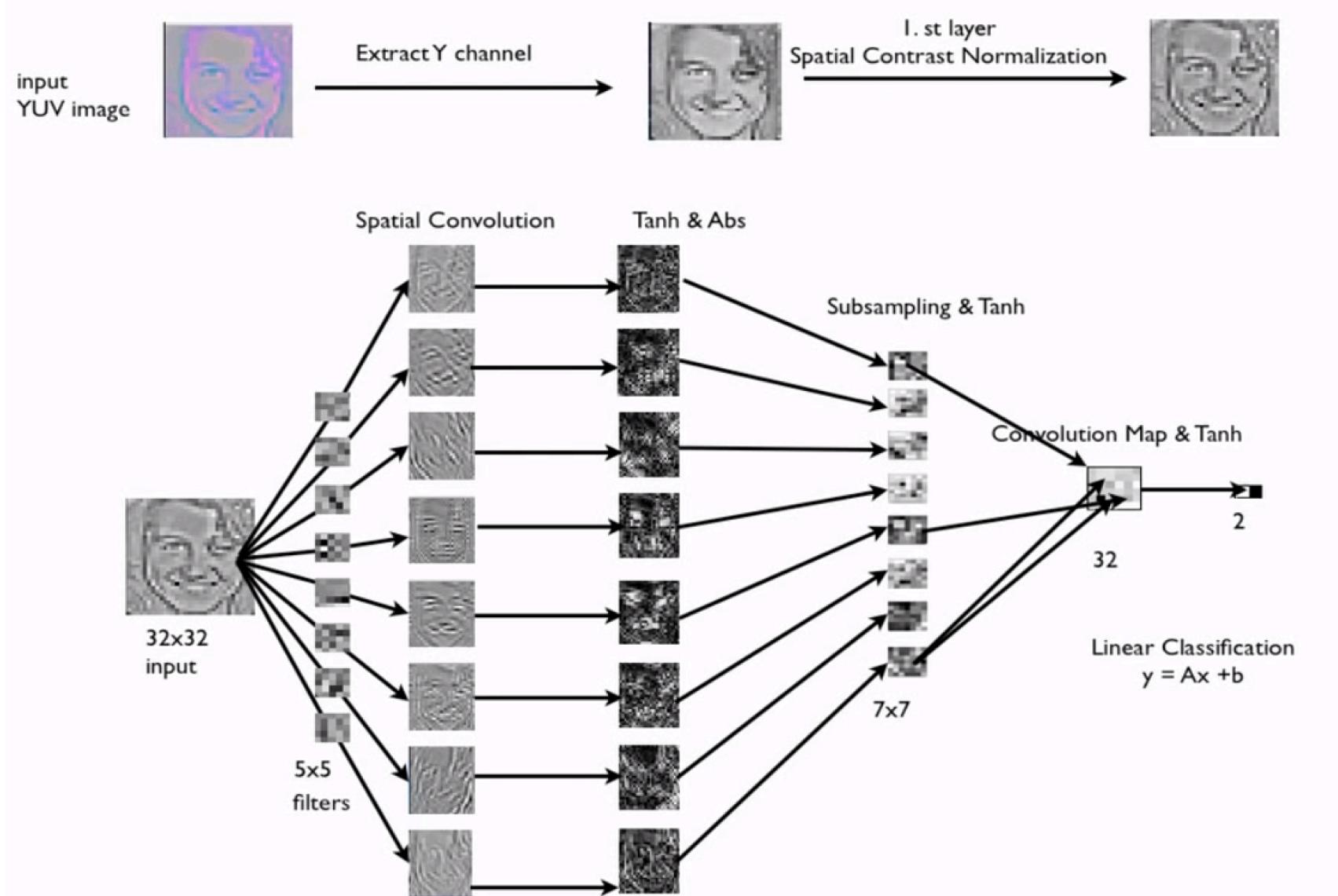
Max Pooling

Hinton: The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster"

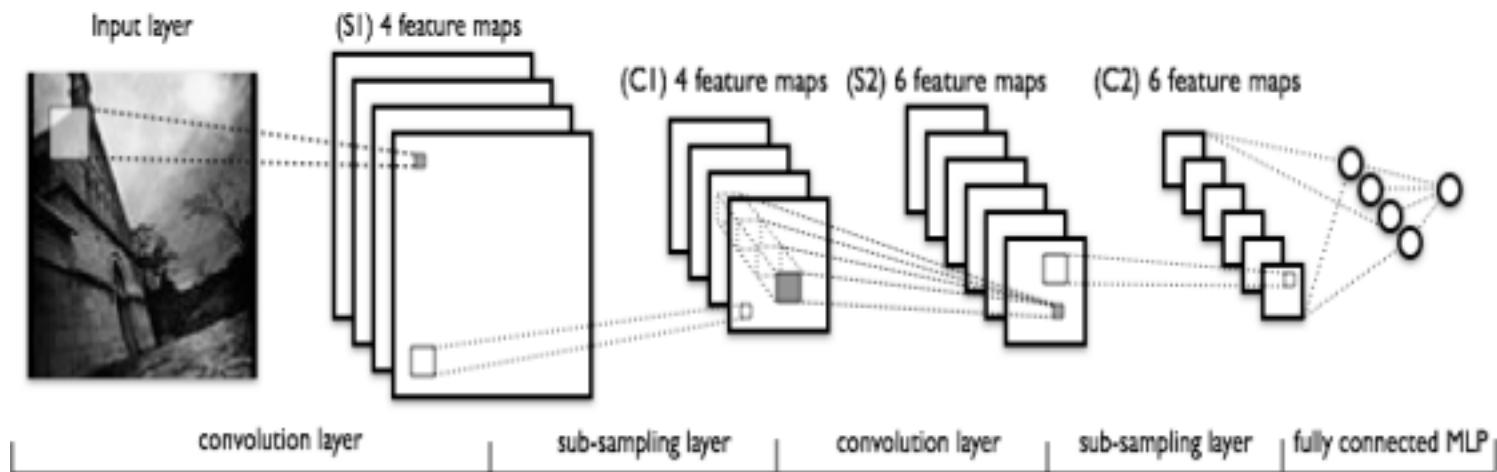


Simply join e.g. 2x2 adjacent pixels in one.

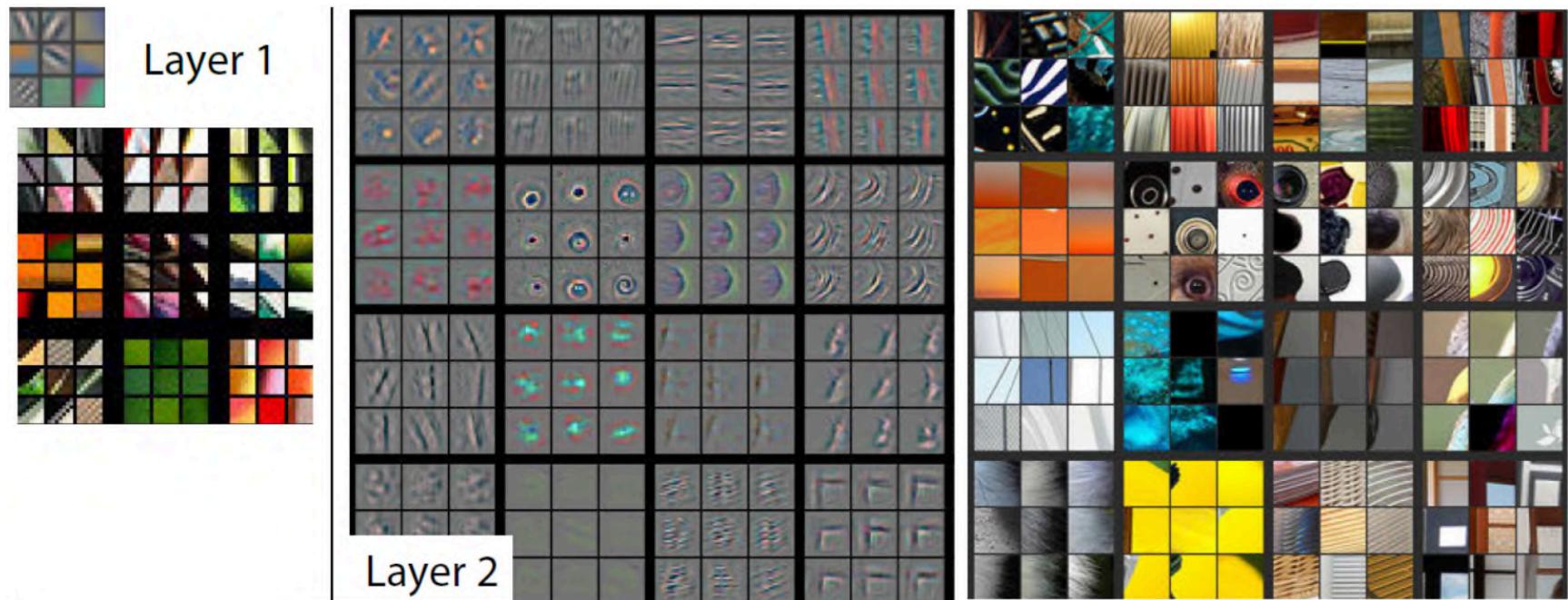
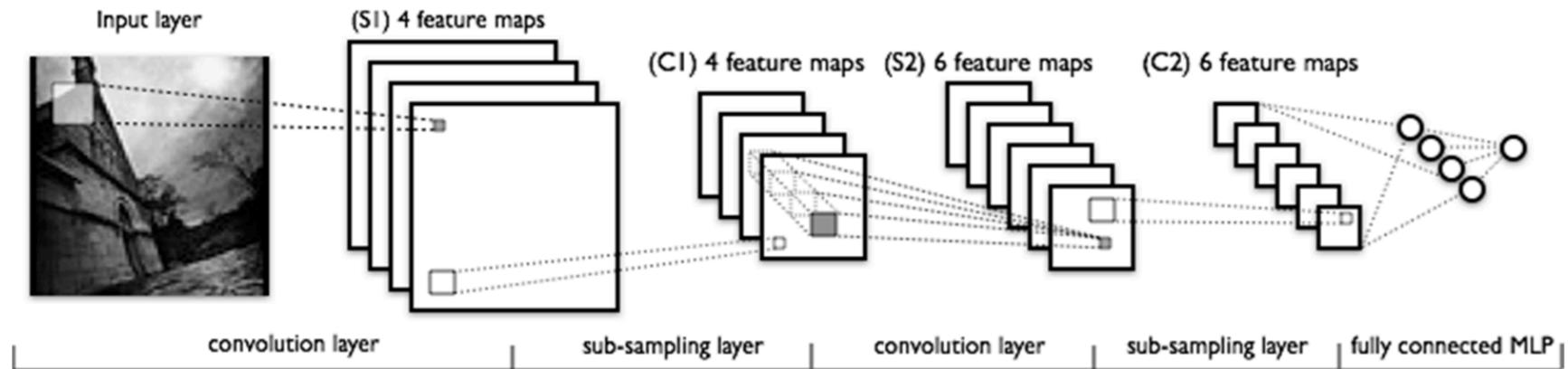
Basic Architecture



CNN Architecture

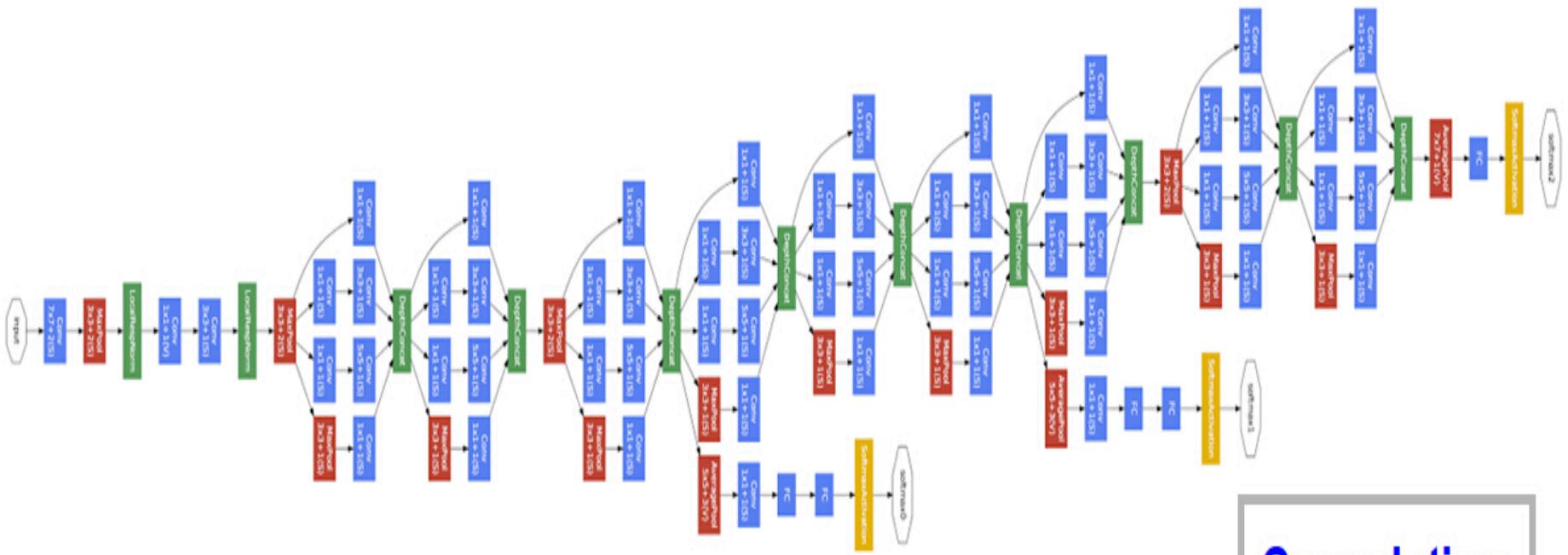


CNN Architecture



[Matthew Zeiler & Rob Fergus]

Example of a CNN model



Convolution
Pooling
Softmax
Other

A ... model? - GoogLeNet

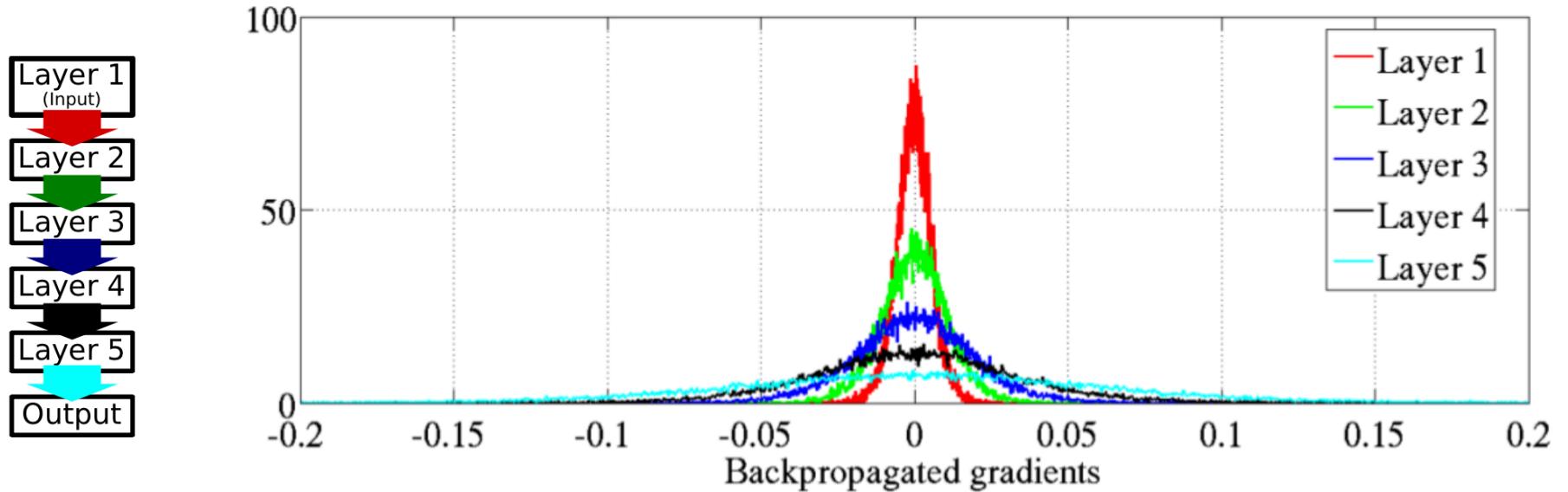
CHALLENGES IN DEEP LEARNING



People who are more likeable and attractive are generally perceived as more intelligent. This is called the "Halo Effect."



Vanishing Gradient

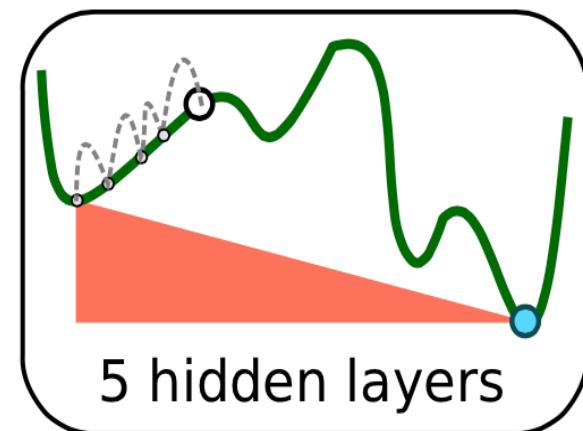
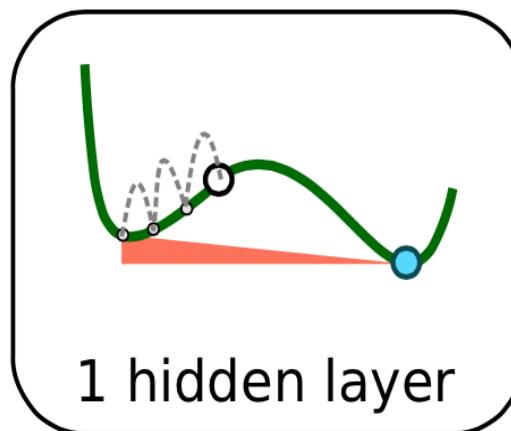
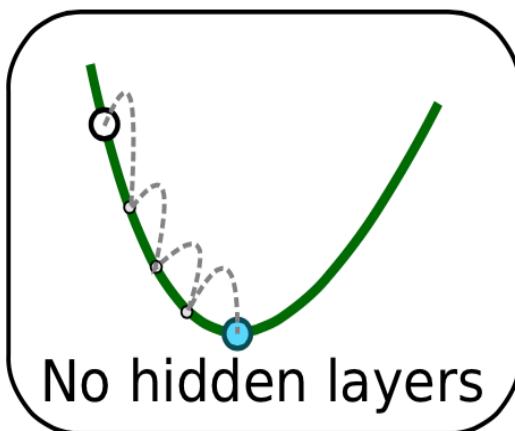


- In backpropagation, the gradient becomes **zero** in the layers nearest to the input.
- **Intuition** : layers near to the output have larger influence in the error rate.
- Training deep layers successfully is terribly slow.
- Until deep layers learn some sensible weights, subsequent layers are trying to learn mostly from noise.

(*) Glorot & Bengio. Understanding the difficulty in training deep feedforward neural networks.

High Non-Linearity – Local Minima

- Backpropagation is a gradient descent method.
- Only guarantees finding a local minimum of the error function, which can be worse than the global minimum .
- When introducing more non-linearities in the network (more layers), local minima multiply .



- It is hard to find good weight configurations for a deep network.

Alleviating vanishing gradients & local minima

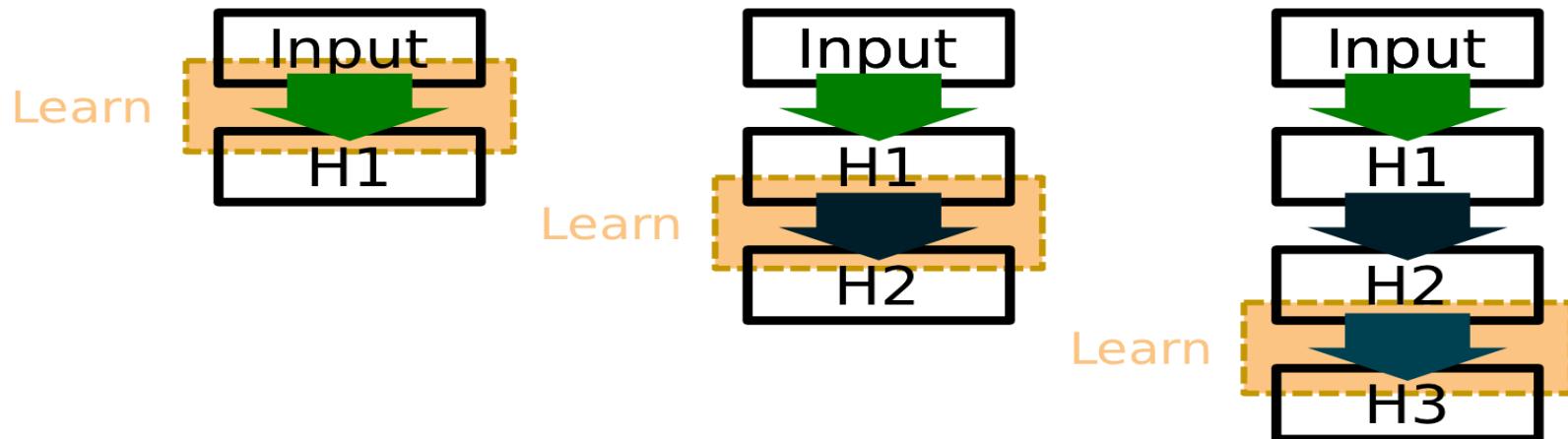
- Smart weight initialization
- Use RL and Maxout units
- Train for a longer time .
- Use restarts .

Incremental Pre-training

Instead of learning all the network weights jointly, build each layer iteratively .

Incremental Pre-training

- Start learning the layer nearest to inputs.
- proceed adding layers until last hidden layer
- When adding a new layer, weights from previous layers are kept fixed .
- One layer is learned at a time, vanishing gradients are avoided .



Incremental Pre-training

Two options for output

- Keep pre-trained weights fixed .
- Deep network ≈ feature generator : more useful features are obtained from the inputs.
- The output layer is trained as a perceptron , inputs are the features obtained.
- Discriminative fine-tuning . Optimize jointly over all network weights (e.g. using BP).
- Pre-trained weights seem to be a better choice than random initialization .
- Slower Fine tuning

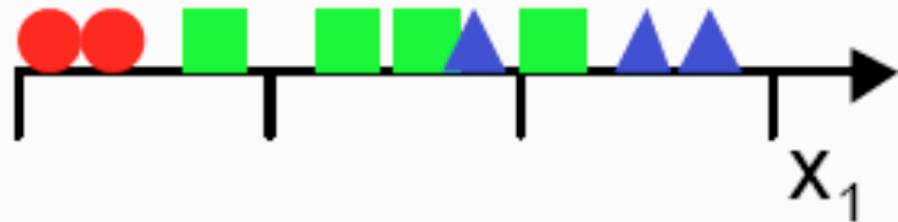


UNSUPERVISED LEARNING

Curse of Dimensionality

Single feature

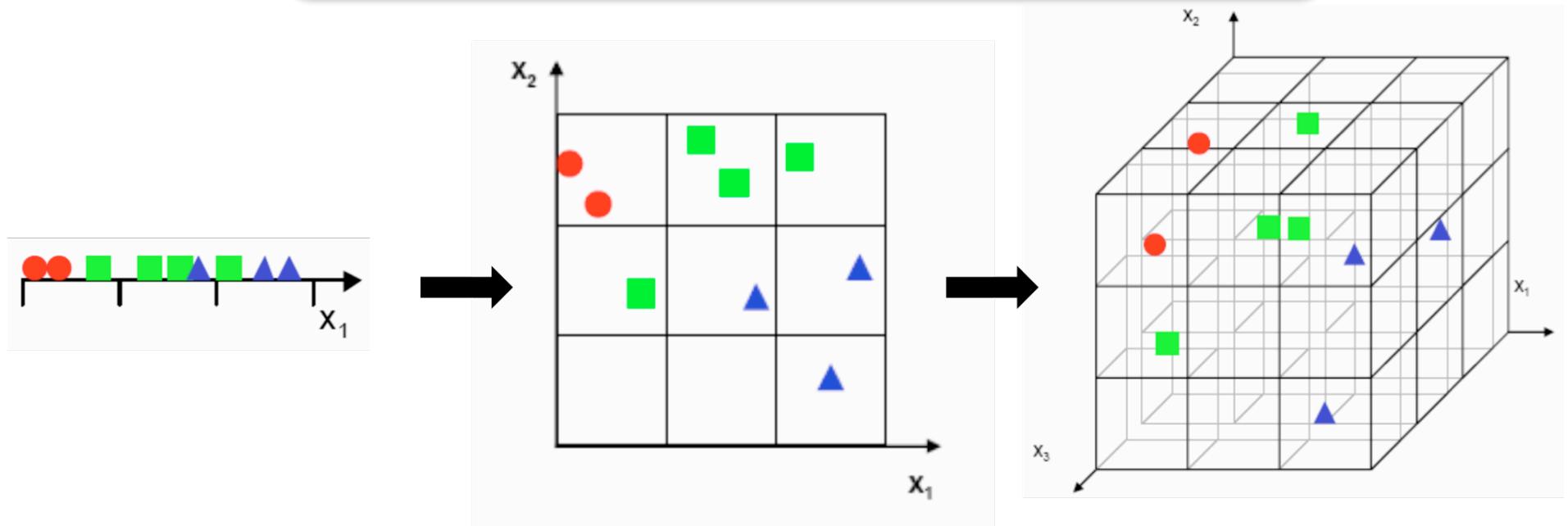
Divide feature space into 3 simple bins



Too simplistic –
Too much overlap between classes.

Curse of Dimensionality

But by adding features to improve separability, the feature space explodes.



Feature Selection

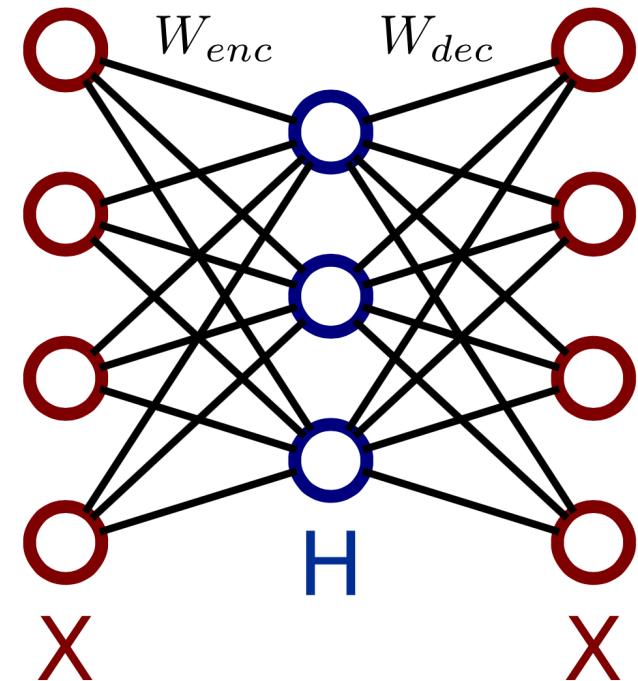
Need: Automated feature selection

Autoencoders

GOAL: Efficient Reconstruction of Input Data

Use input as input and output of the network and train (Don't use the target)

Learn the weights using a standard neural network method (e.g. backpropagation).

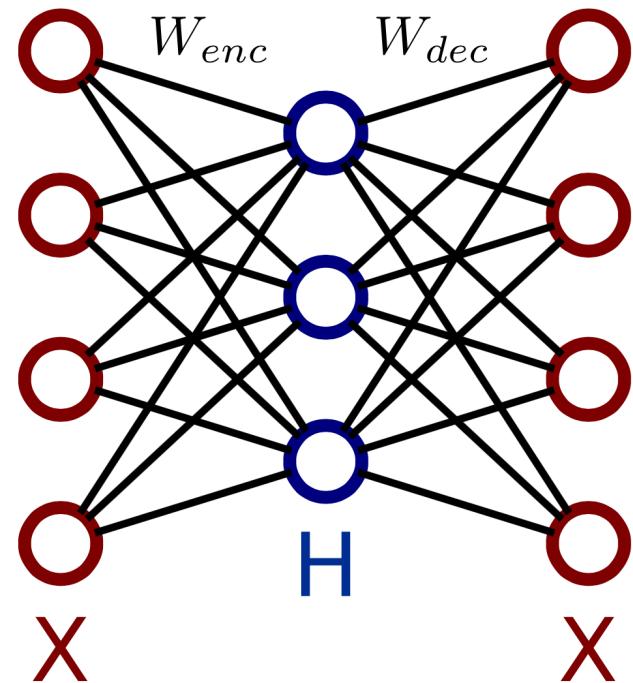


Autoencoders - Output

After learning,

W_{enc} : compresses the information in the data into the hidden units.

W_{dec} decompresses the info in the hidden units back to its original form (with some loss).

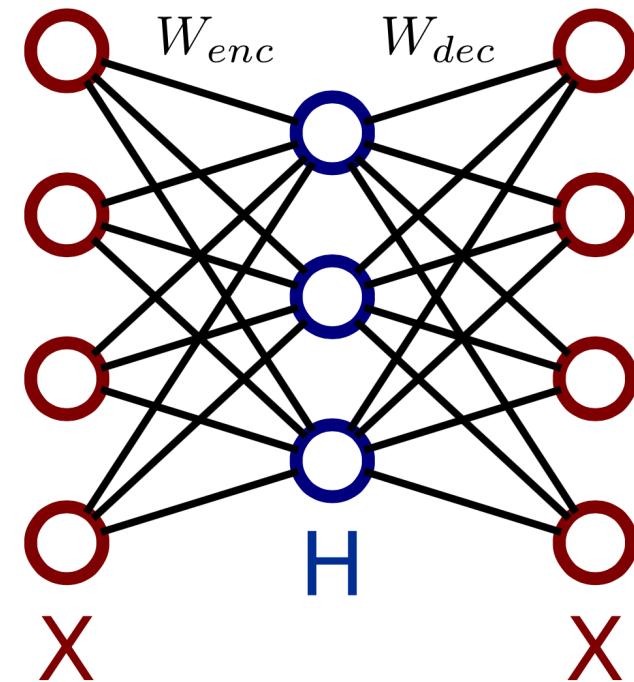


Autoencoders – To make them effective

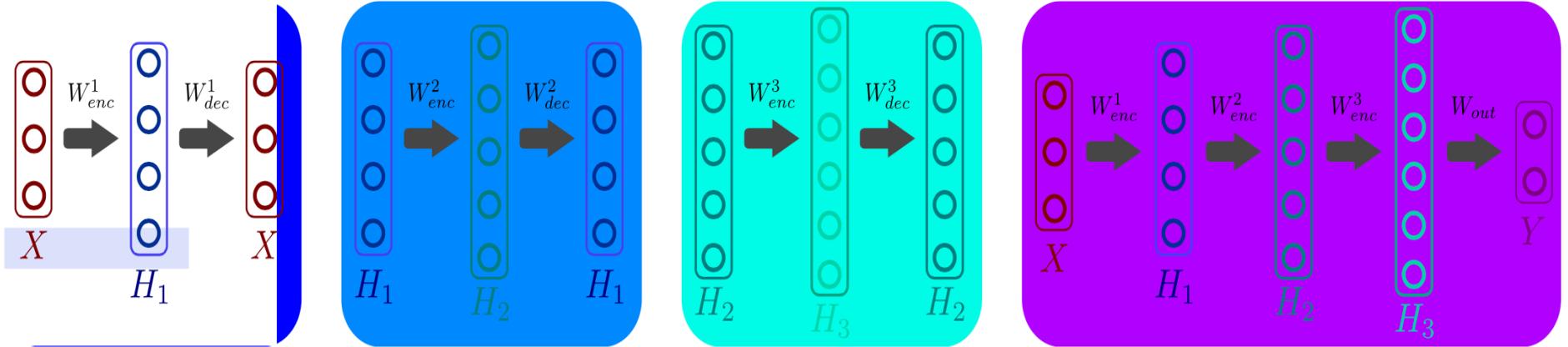
Enforce sparsity .

(A network is sparse if only a few
of its units >0 simultaneously)

Introduce random noise in the
input patterns when training the
network. (Only in the inputs!)



Deep Autoencoders



Train autoencoder to reconstruct the training data .

Proceed iteratively, building new autoencoders reconstructing the value of the hidden units of the previous stage.

Create a feedforward network using the encoding weights W_{enc} from all the trained autoencoders. As output layer, use the training labels .

Fine-tuning : train the last layer or the full network using standard backpropagation.

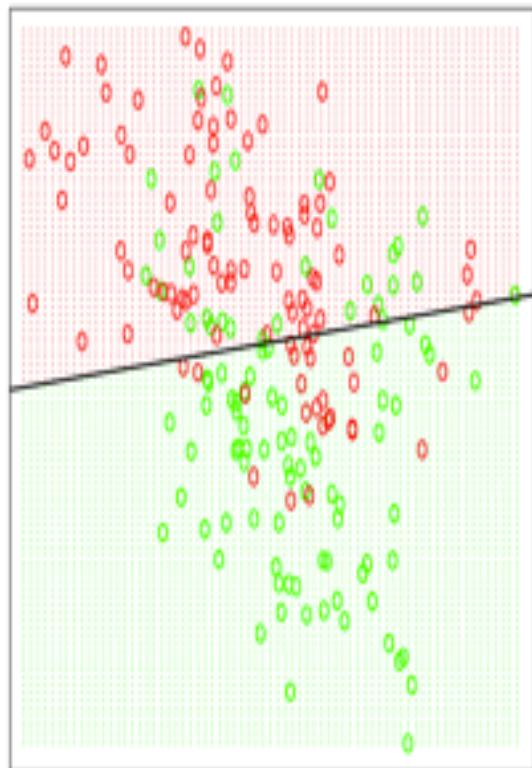
The background of the slide is a photograph of a large, tranquil lake or sea. The water is a deep blue, with slight ripples across its surface. In the distance, a dark, horizontal line of trees or a distant shoreline is visible against a bright, clear sky. The overall atmosphere is peaceful and expansive.

CHALLENGES IN DEEP LEARNING

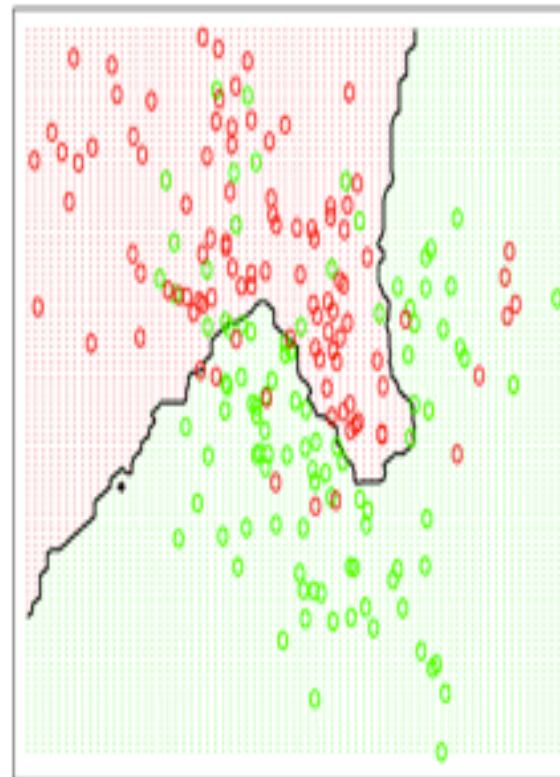
- Continued

Overfitting

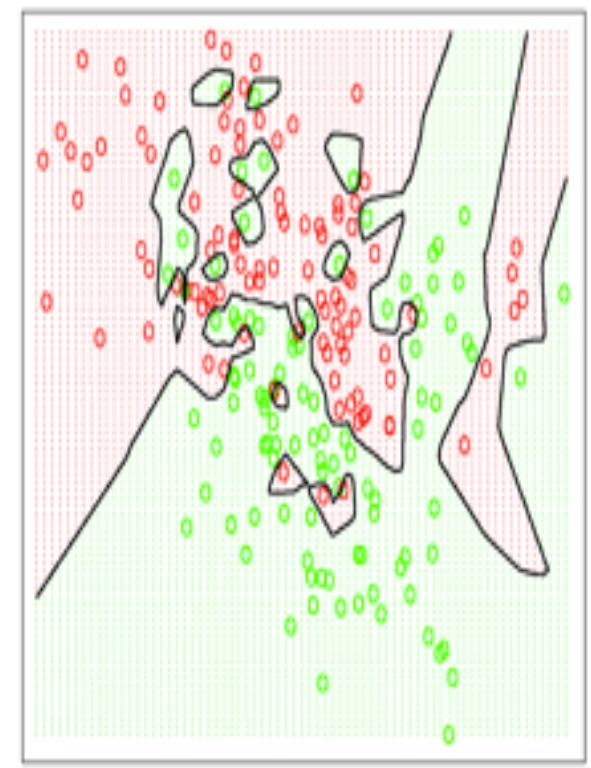
**Underfitted
Model**



**Good
Model**



**Overfitted
Model**

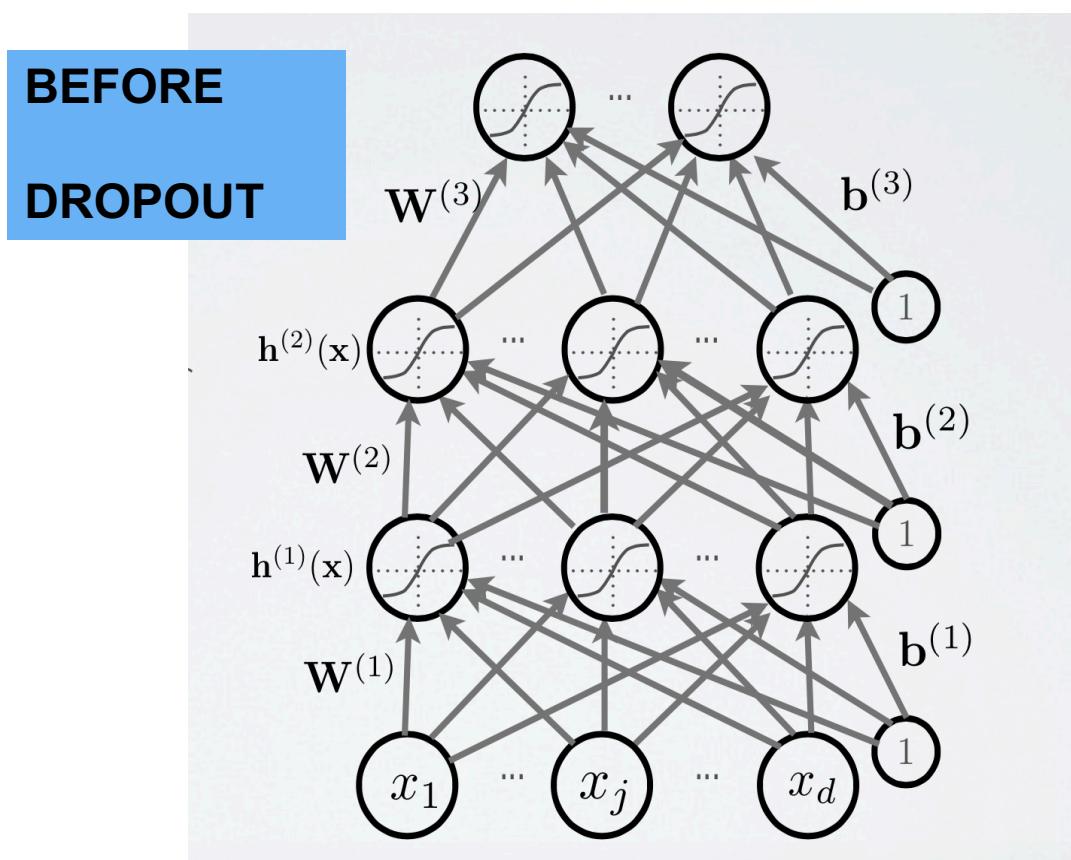


Some ways to address Overfitting

- Weight Decay
- L1/L2 Regularization
- Suitable Model Architectures (depth and width of the layers)
- Unsupervised Pre-training
- Dropout
- Data Augmentation

Dropout

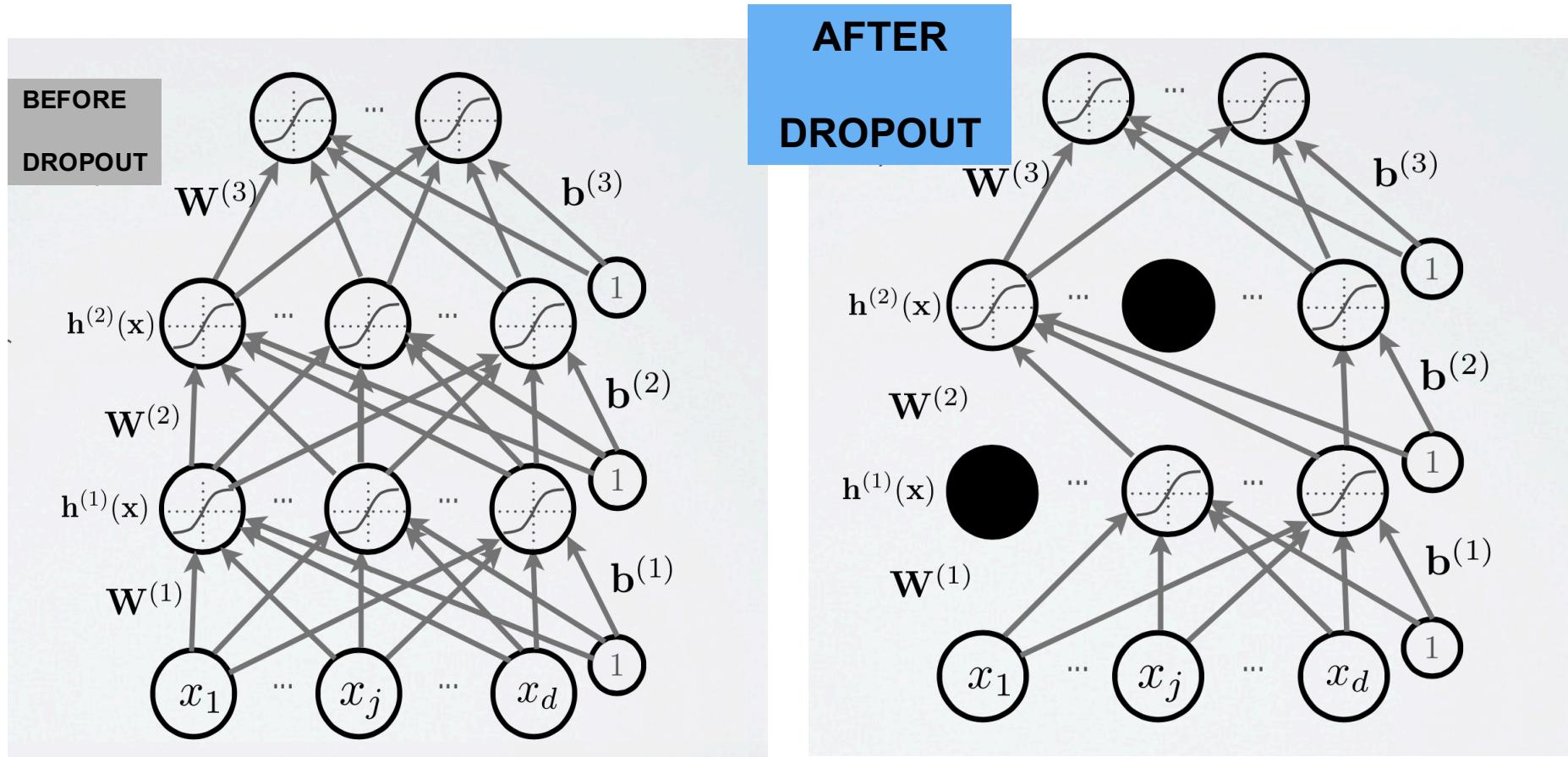
- Cripple the network by removing hidden units stochastically
- In practice, probability of 0.5 works well.



source: http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html

Dropout

- Cripple the network by removing hidden units stochastically
- In practice, probability of 0.5 works well.



source:http://info.usherbrooke.ca/hlarochelle/neural_networks/content.html

Data Augmentation

Some ways to augment data:

- **Rotation:** random angle between 0° and 360°
- **Translation:** random shift between -10 and 10 pixels
- **Rescaling:** random scale factor between 1/1.6 and 1.6
- **Flipping:** yes or no
- **Shearing:** random angle between -20° and 20°
- **Stretching:** random with stretch factor between 1/1.3 and 1.3



HANDS-ON



Graphical Processing Unit

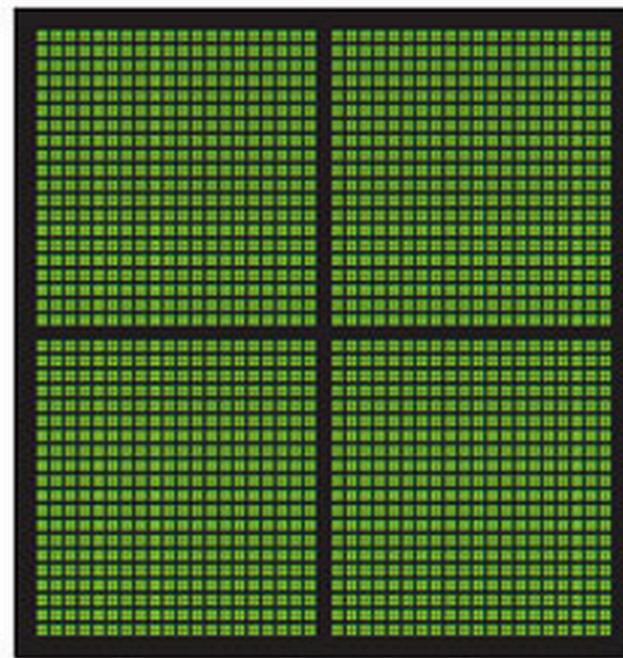
Impact of GPU

Compared to CPUs 20x speedups are typical

GPUs have thousands of cores to process parallel workloads efficiently



CPU
MULTIPLE CORES



GPU
THOUSANDS OF CORES

Impact of GPU

- Accelerated computations on float32 data
- Matrix multiplication, convolution, and large element-wise operations can be accelerated a lot (5-50x)
- Difficult to parallelize dense neural networks on multiple GPU efficiently (Active area of research)
- Convolutional neural networks – unlike dense neural networks – can be run very efficiently on multiple GPUs. Their use of weight sharing makes data parallelism very efficient
- Copying of large quantities of data to and from a device is relatively slow.
- CUDA has released cuDNN.

Summary

