

Introduction to Deep Learning

- *Natural Language Processing*

*Swayam Mittal - Data Scientist,
Indegene*

)

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.



BECAUSE
YOU'RE A
FAST
WORKER?



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

AGENDA/SCHEDULE

40-50% Theory + 60-50% Hands-On

- 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) Word Embedding
 - i. CBOW, skip-gram, word2vec
- 4) Convolution Neural Networks
 - i. Convolution
 - ii. Sub-sampling, Pooling
 - iii. Dropout
 - iv. Architecture
- 5) Recurrent Neural Network
 - ii. LSTM
- 6) 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 into the side of a mountain, featuring numerous stone terraces, walls, and buildings. Two large, prominent peaks rise behind the city, one of which is partially obscured by thick, billowing clouds. The sky is overcast and gray.

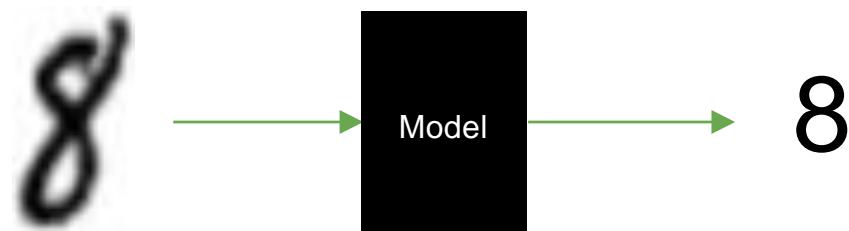
LEARNING

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	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

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

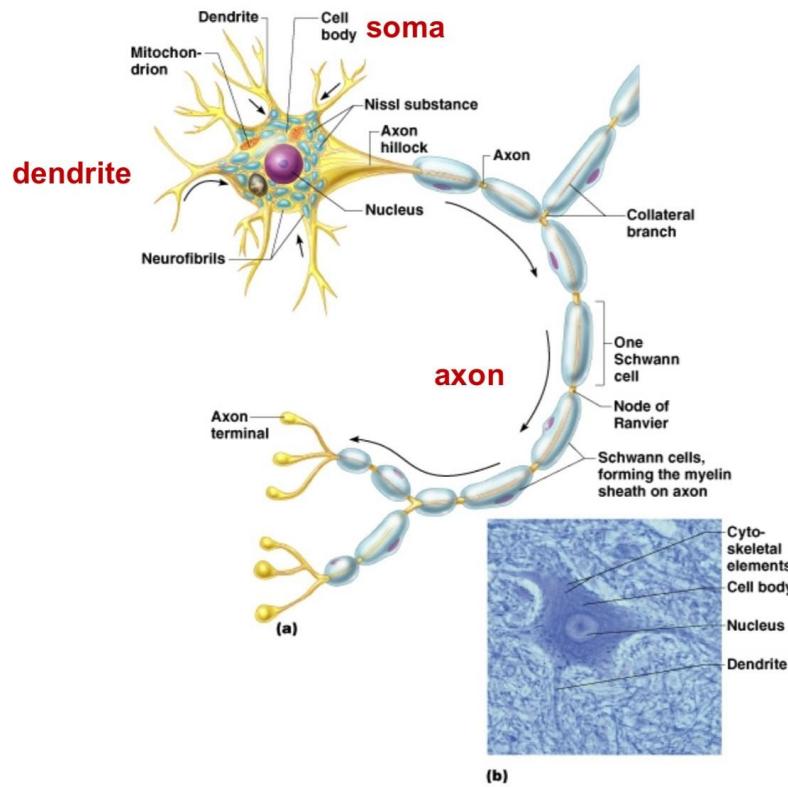
Questions about Learning

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.

Lessons from Biological Learning

Importance of Connectionism

Simple units interacting in a complex network

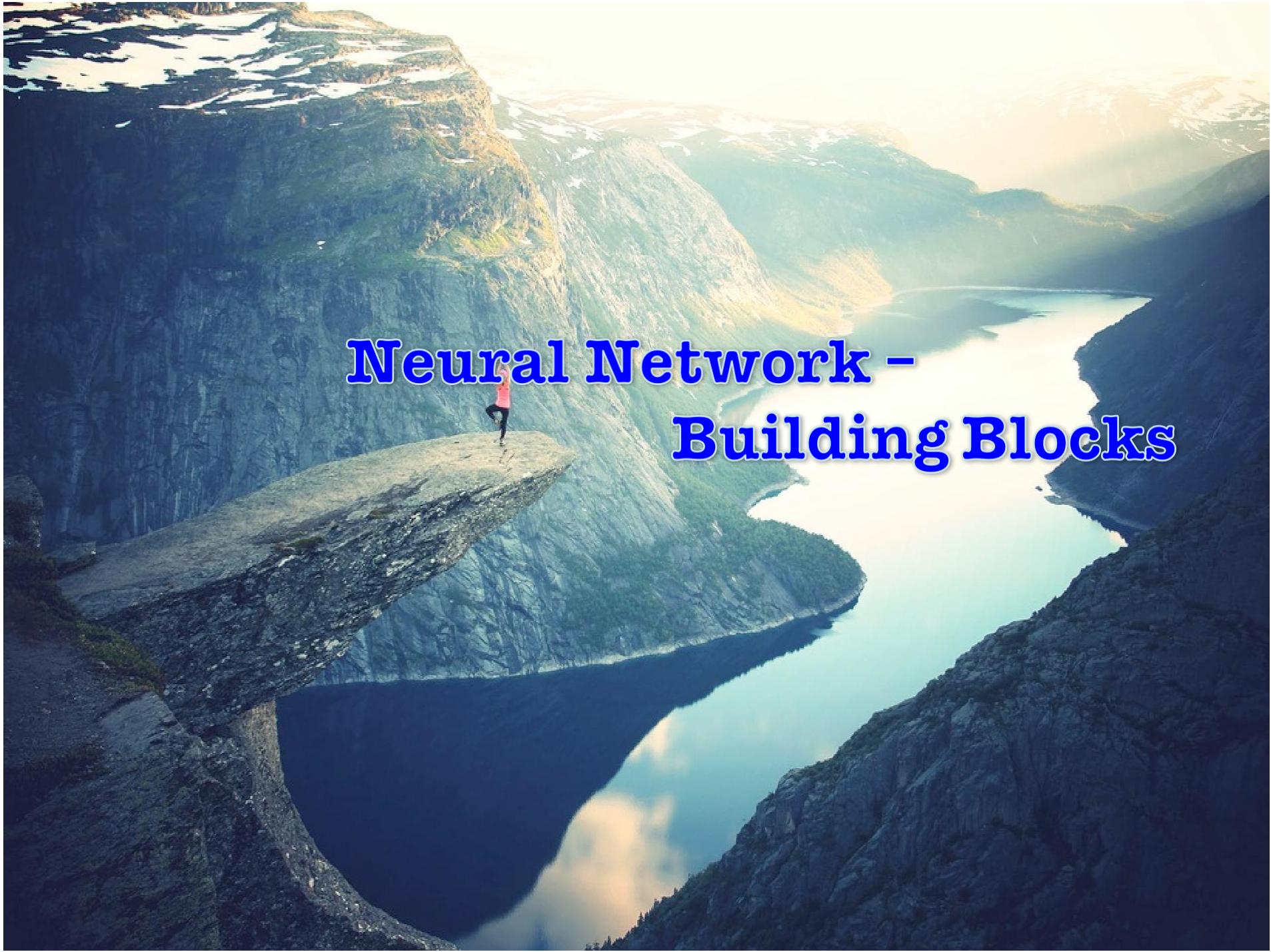
Distributed representation of knowledge

Parallelism

Threshold mechanism for robust classification

Mechanism for learning – adjusting synaptic weights

Comprehend inner structure of observed data

A photograph of a person standing on a small, rocky ledge of a massive cliff face. The cliff is rugged and layered, with patches of snow and green vegetation. Below the cliff, a deep, winding fjord or river cuts through the landscape, its water a vibrant blue-green. The sky is clear and bright, suggesting a sunny day.

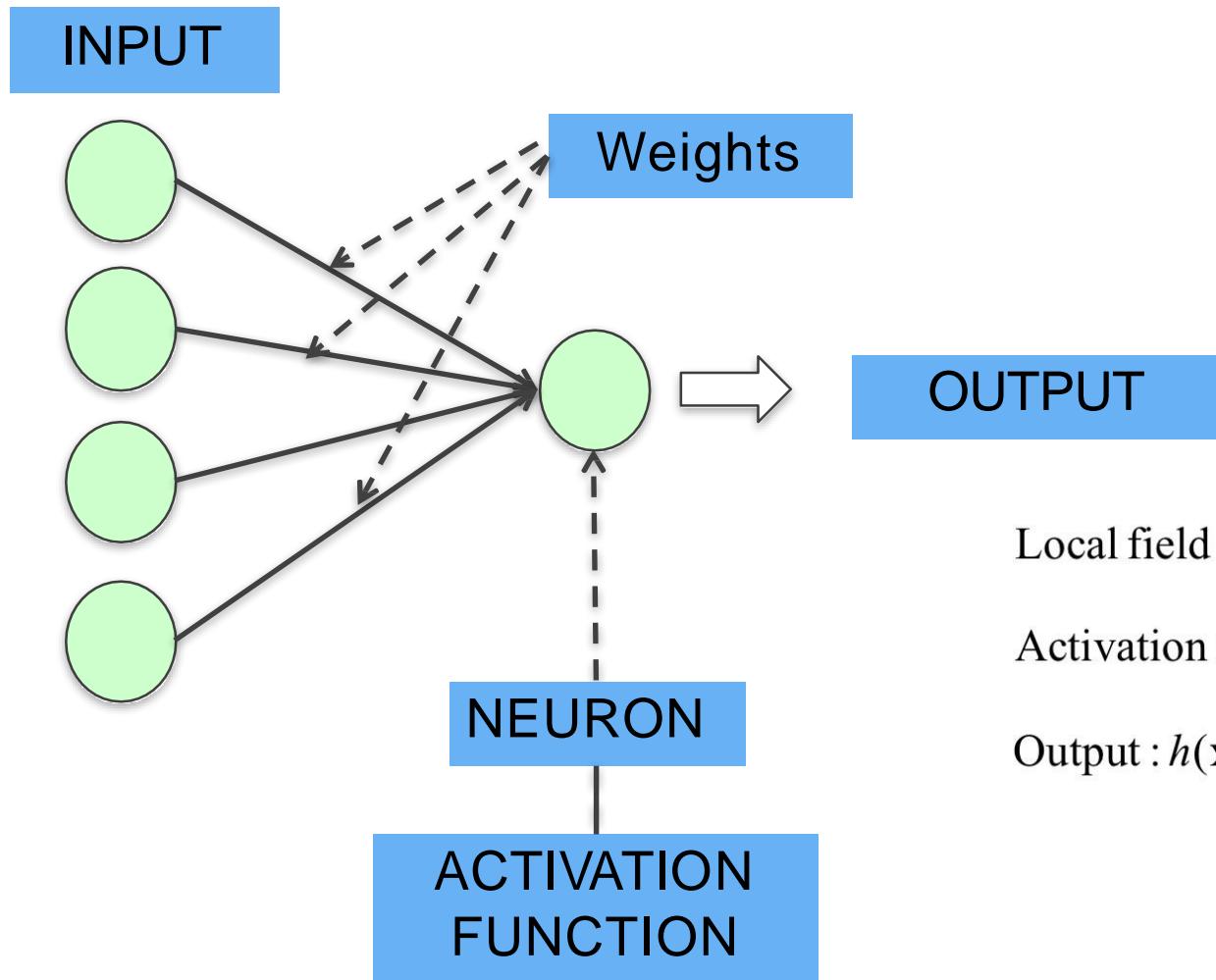
Neural Network – Building Blocks

DID
YOU
KNOW?

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



Neuron, Activation Function

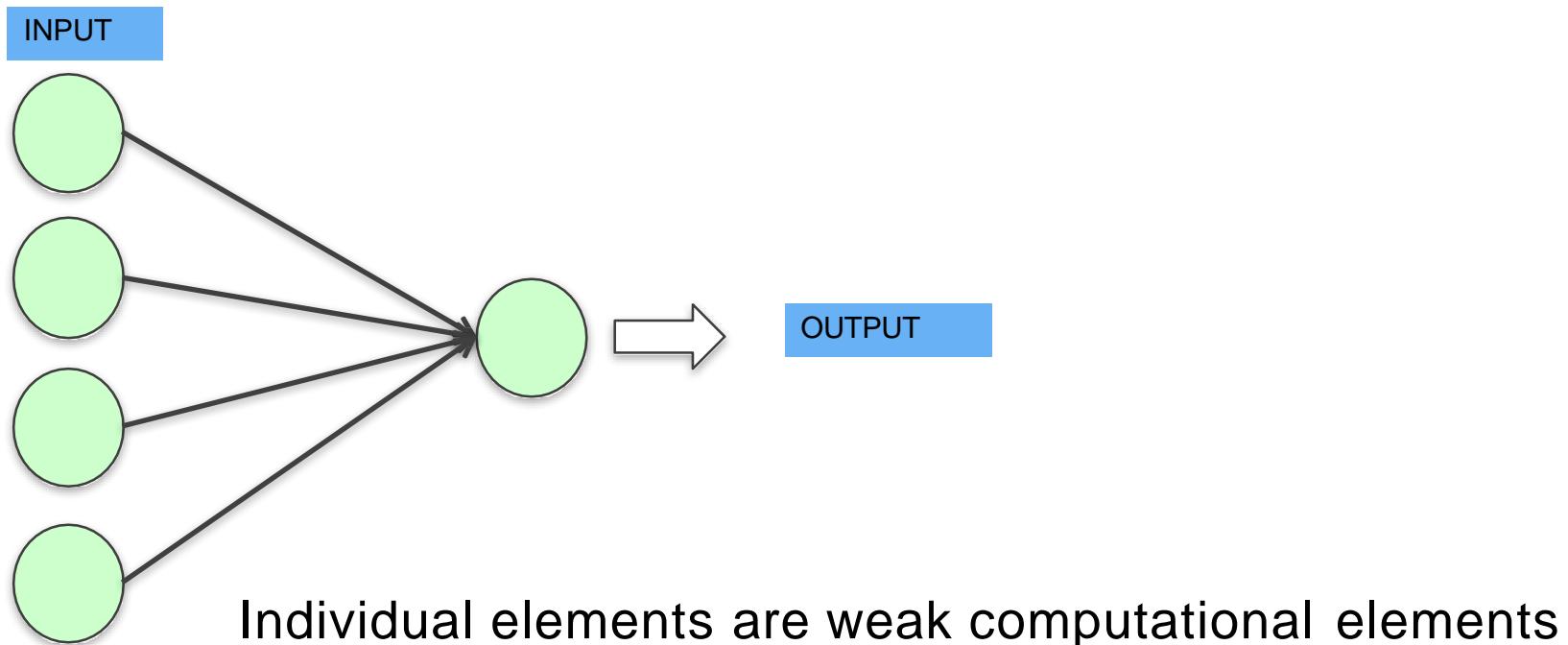


$$\text{Local field} : \sum_{d=0}^D w_d x_d$$

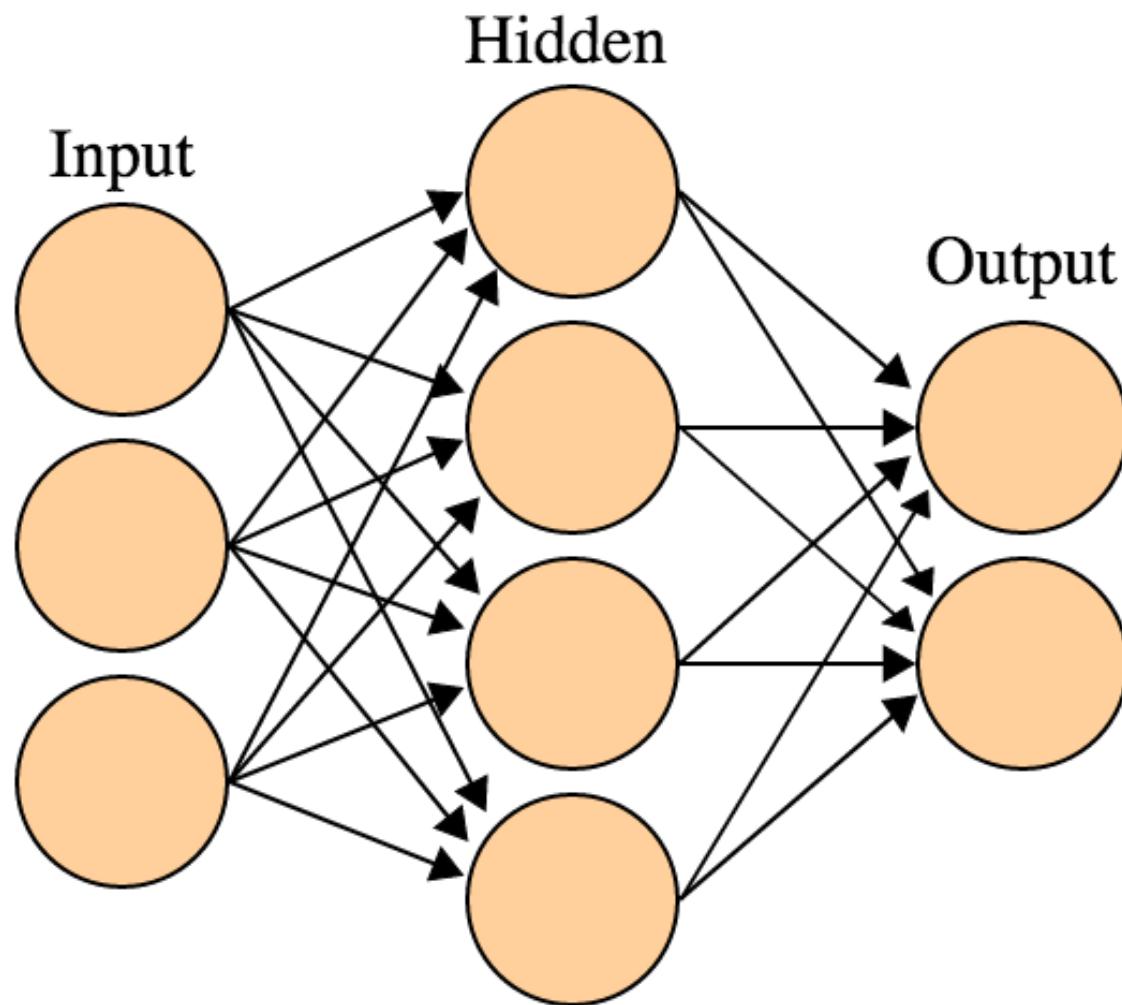
Activation function : $\varphi(\cdot)$

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

Need: Networked Units



A Simple Neural Network

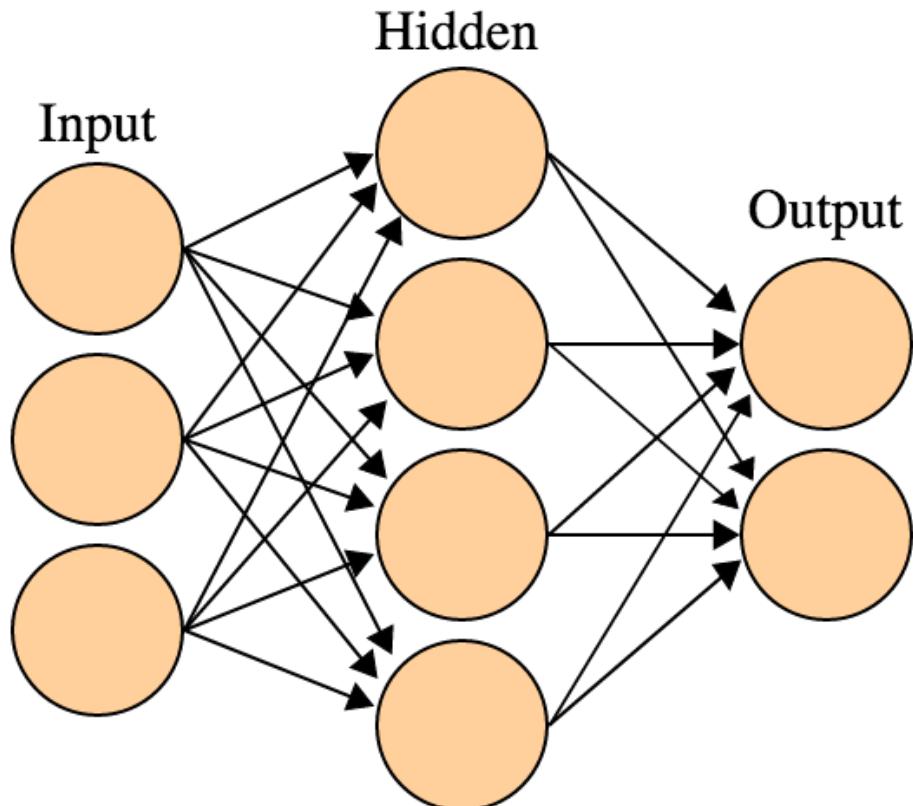


$$\text{Local field} : \sum_{d=0}^D w_d x_d$$

Activation function : $\varphi(\cdot)$

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

Feed-forward Neural Network



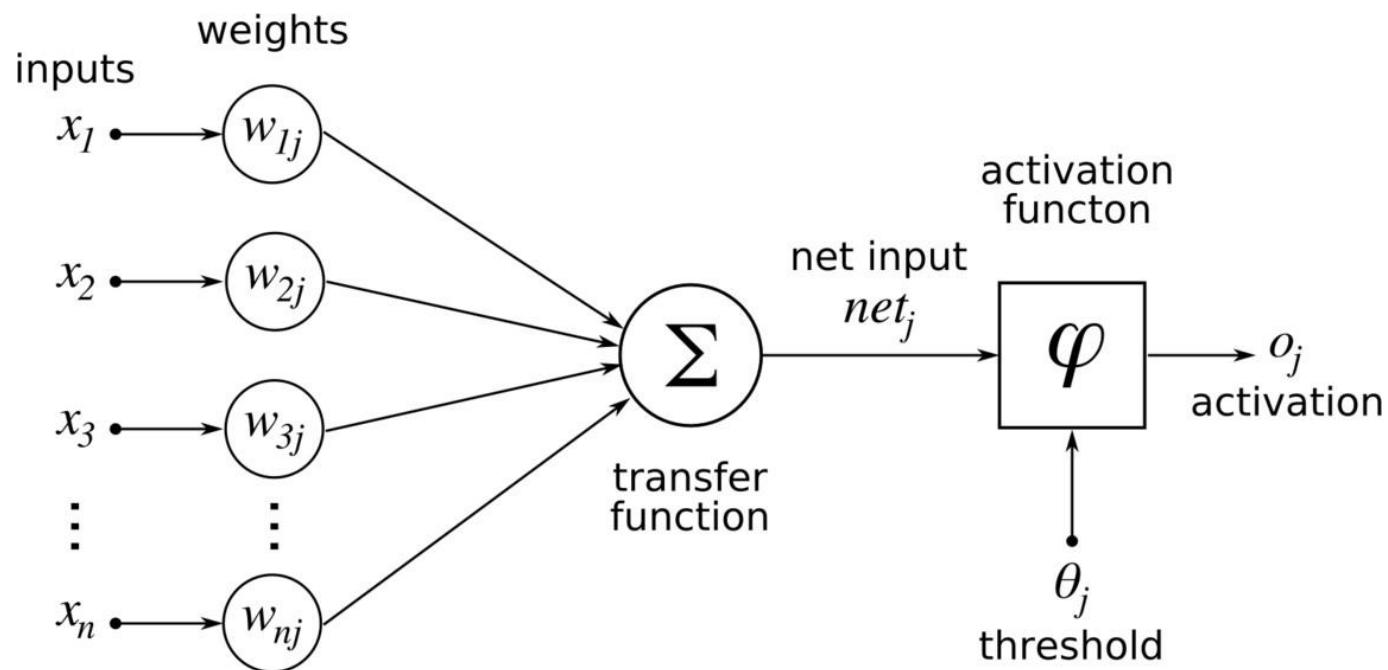
- » The network is formed by
 - Input layer of source nodes
 - One or several hidden layers of processing neurons
 - Output layer of processing neuron(s)
- » Connections only between adjacent layers
- » There are no feedback connections

Activation Functions

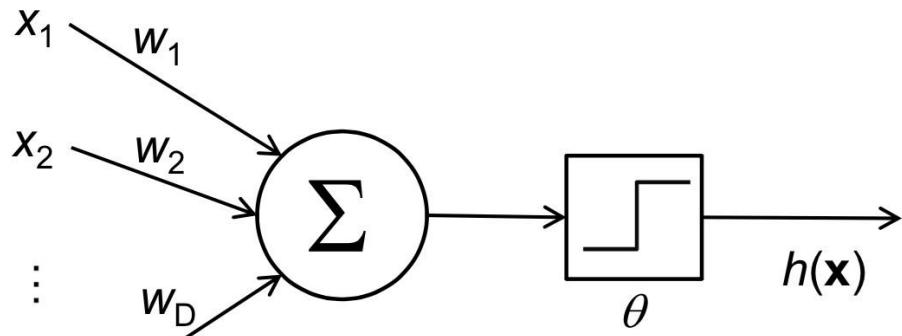


Activation Function

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



Perceptron

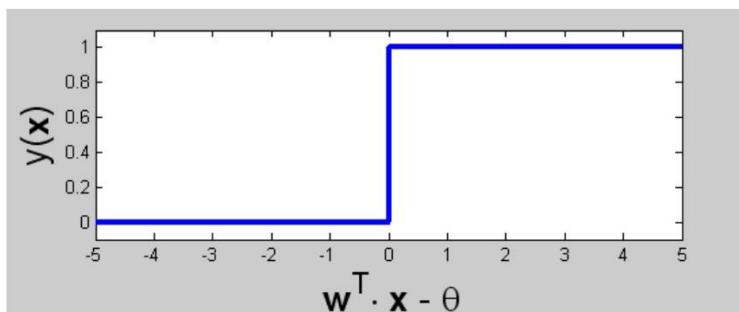


$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } w_1x_1 + w_2x_2 + \dots + w_Dx_D \geq \theta \\ 0 & \text{if } w_1x_1 + w_2x_2 + \dots + w_Dx_D < \theta \end{cases}$$

threshold : θ

In vector notation:

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

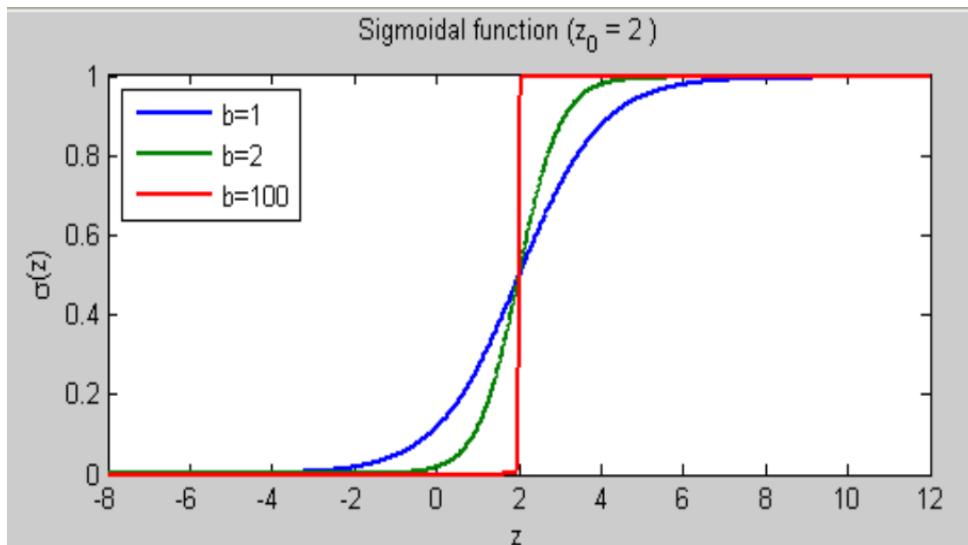


McCulloch, W. S., and Pitts, W. "A Logical Calculus of Ideas Immanent in Nervous Activity." Bulletin of mathematical biophysics, 5, pp. 115-133 (1943). Reprinted in McCulloch, W. S., *Embodiments of mind*. Cambridge, MA: MIT Press.

- 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

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

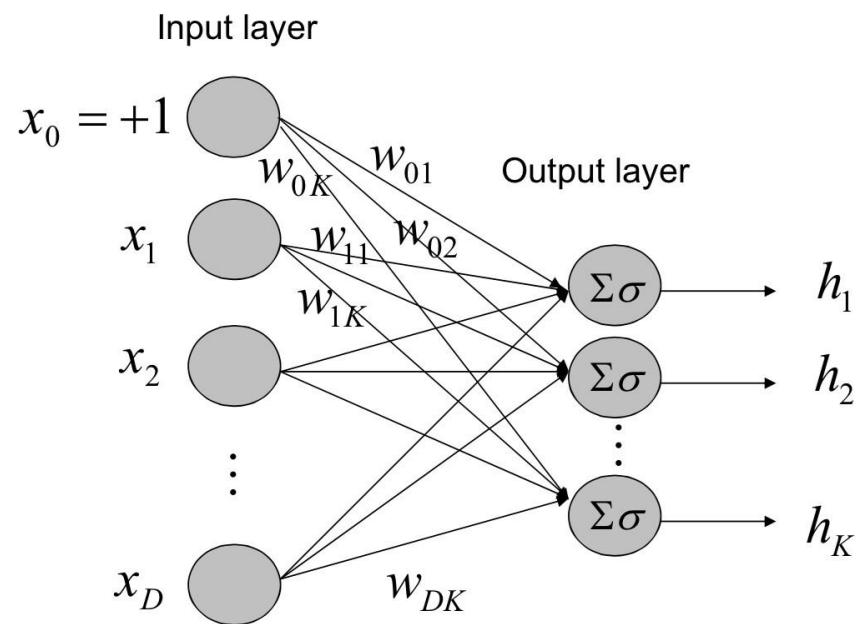
Output : $h(\mathbf{x}) = \sigma(\mathbf{w}^T \cdot \mathbf{x})$

- 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

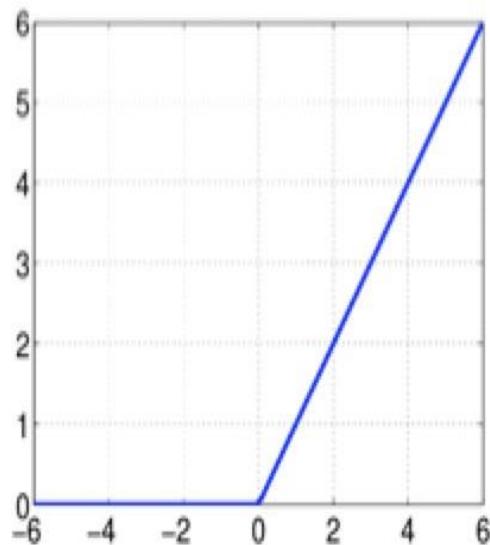
Softmax

Generalization of logistic regression
for Multi-class classification



Rectified Linear Units

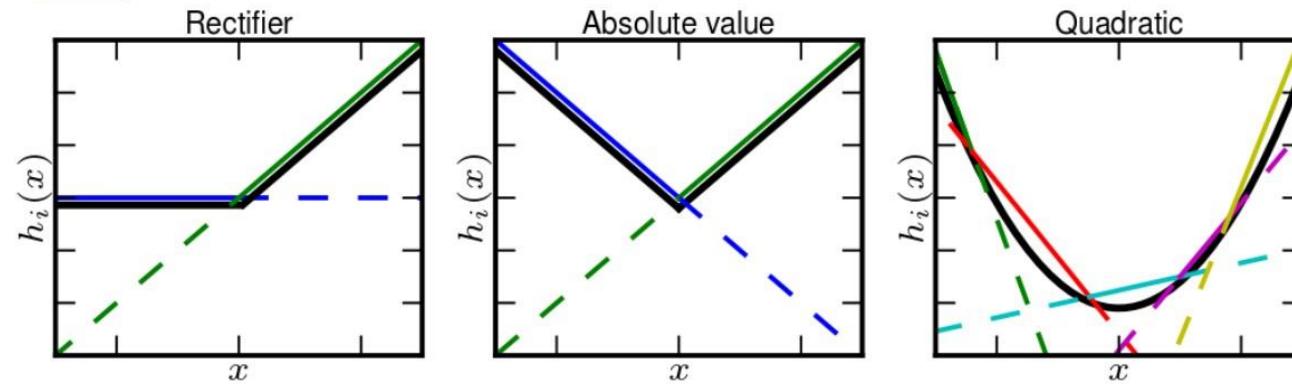
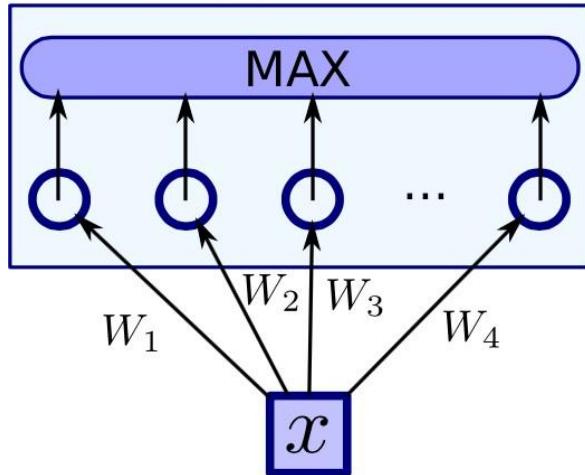
$$f(x) = \max \{0, x\}$$



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

Max-out

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



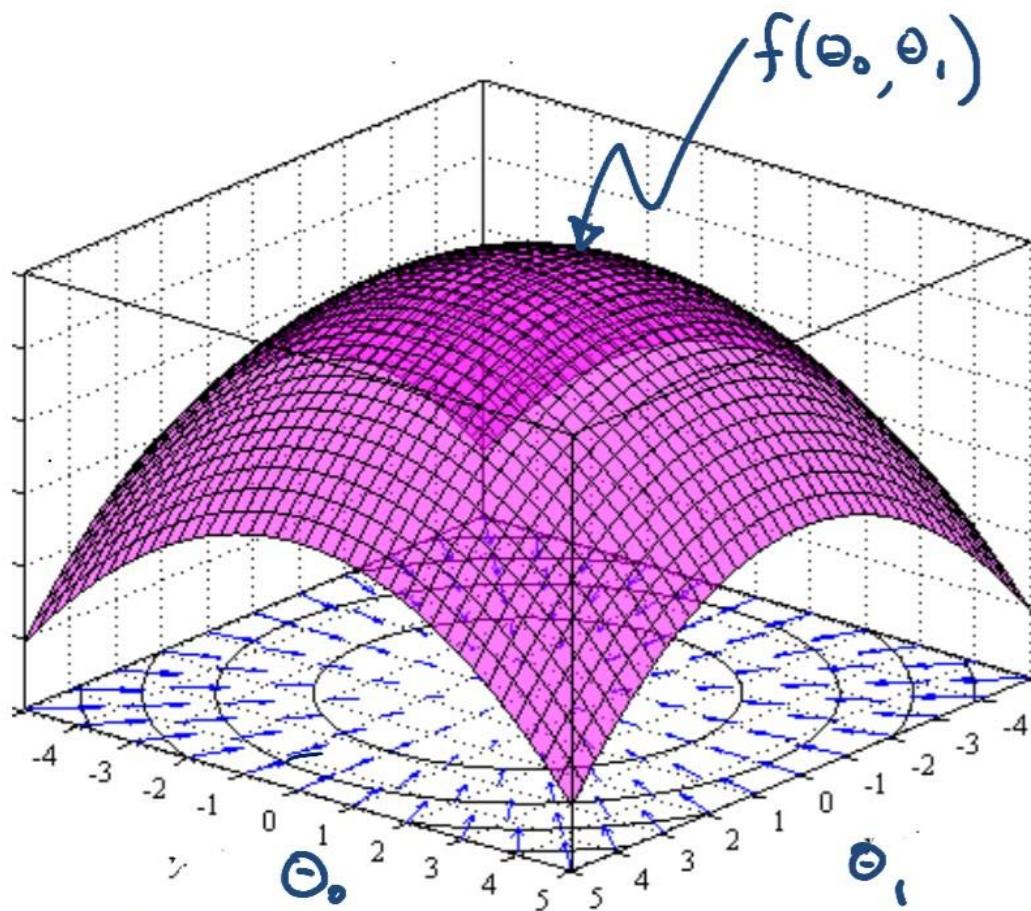


A young woman with dark hair is sitting in a lecture hall, looking directly at the camera. She is wearing a light-colored shirt and is holding a large, greenish-yellow book with both hands, resting it on a white surface. The background consists of rows of blue seats. A black rectangular overlay contains the text "Learning in ANN".

Learning in ANN

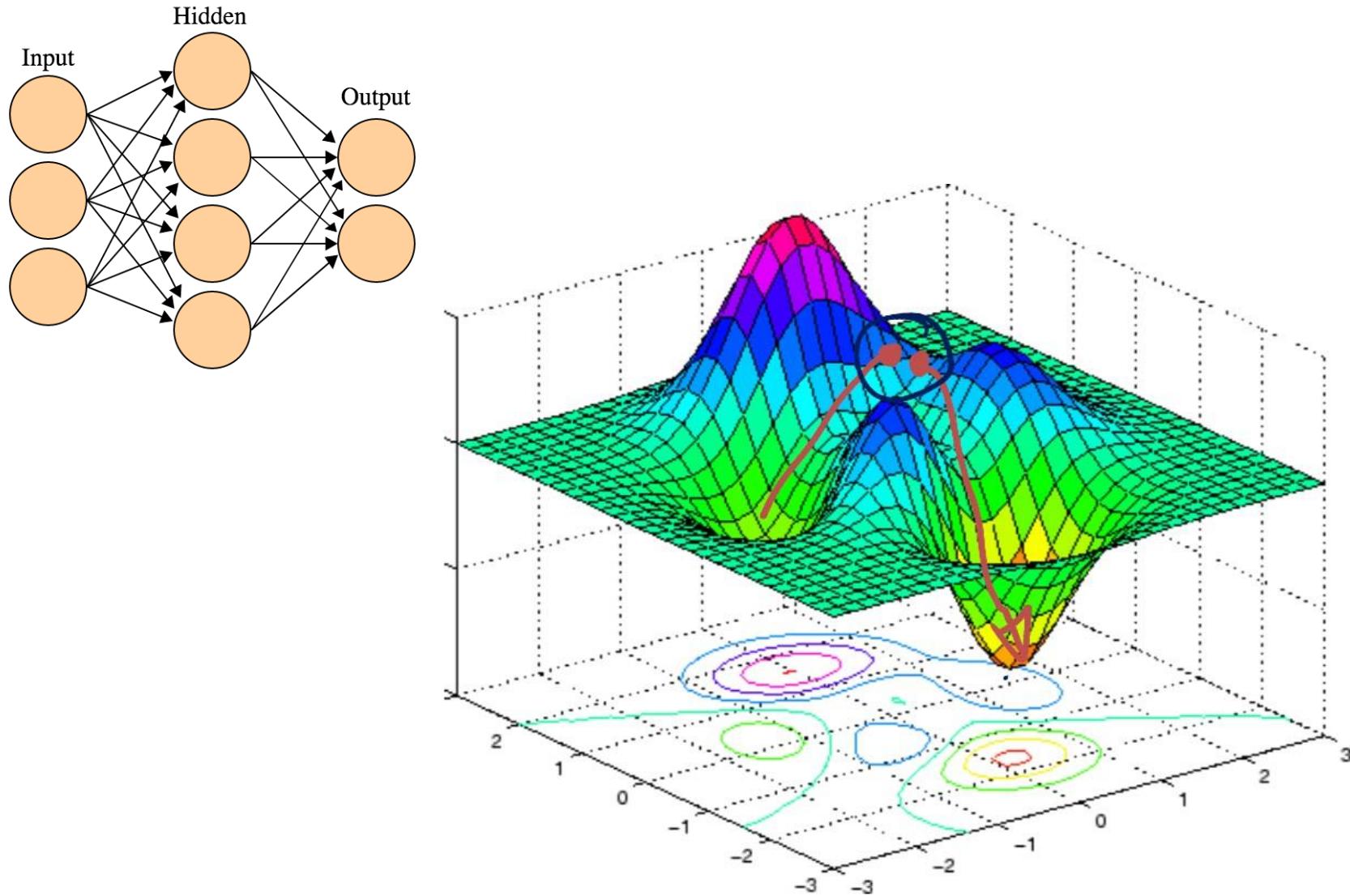
Learning in ANN - Gradient Descent

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 - Pseudo Code

Algorithm 1 Backpropagation

```
for d in data do
```

FORWARDS PASS

Starting from the input layer, forward pass through the network, computing

forward pass through the network, computing

BACKWARDS PASS

Compute the derivatives of the output layer activities

with respect to the output

```
for layer in layers do
```

Compute the derivatives of the upper layer neurons

with respect to the inputs

Compute the derivatives of the weights between the layers

with respect to the weights between the layers

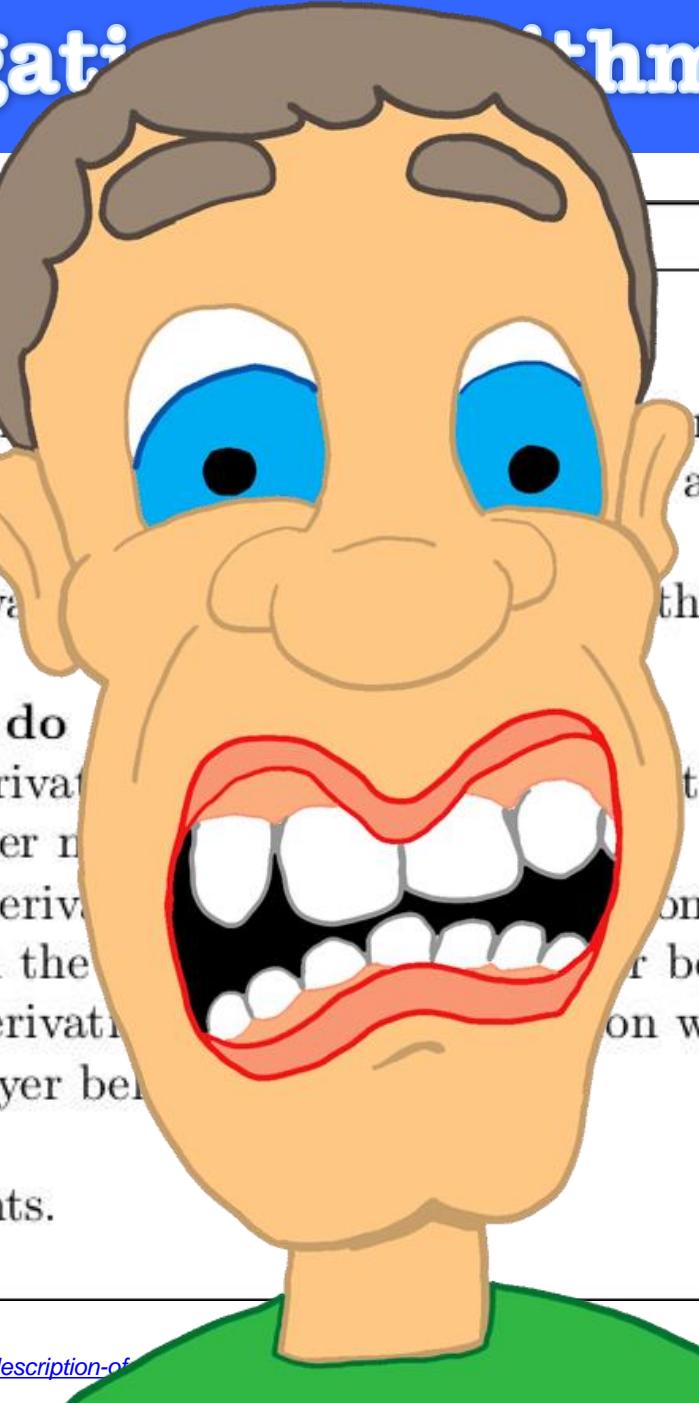
Compute the derivatives of the activities of the layer below

with respect to the activities of the layer below

```
end for
```

Updates the weights.

```
end for
```



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

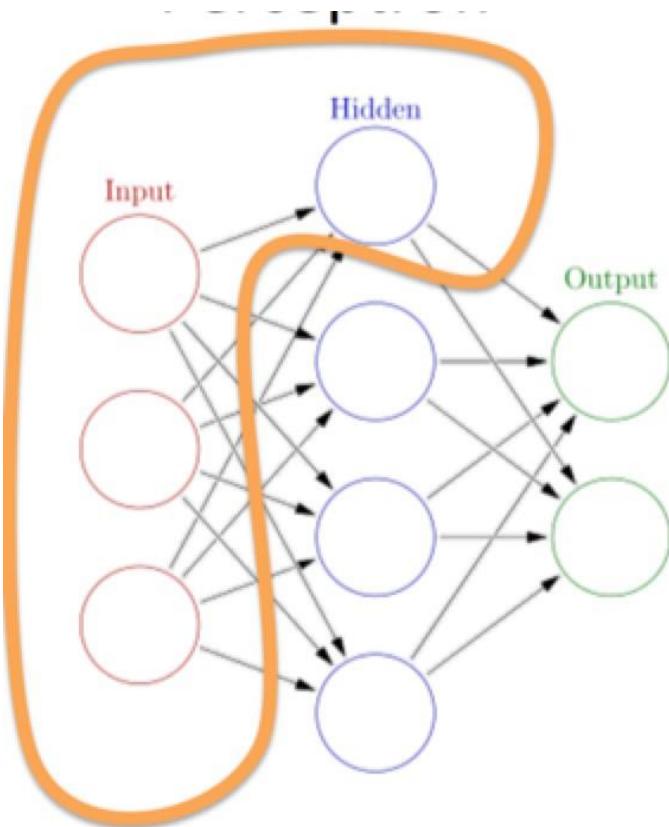
Compute the derivatives of the error function with respect to the activities of the layer below

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

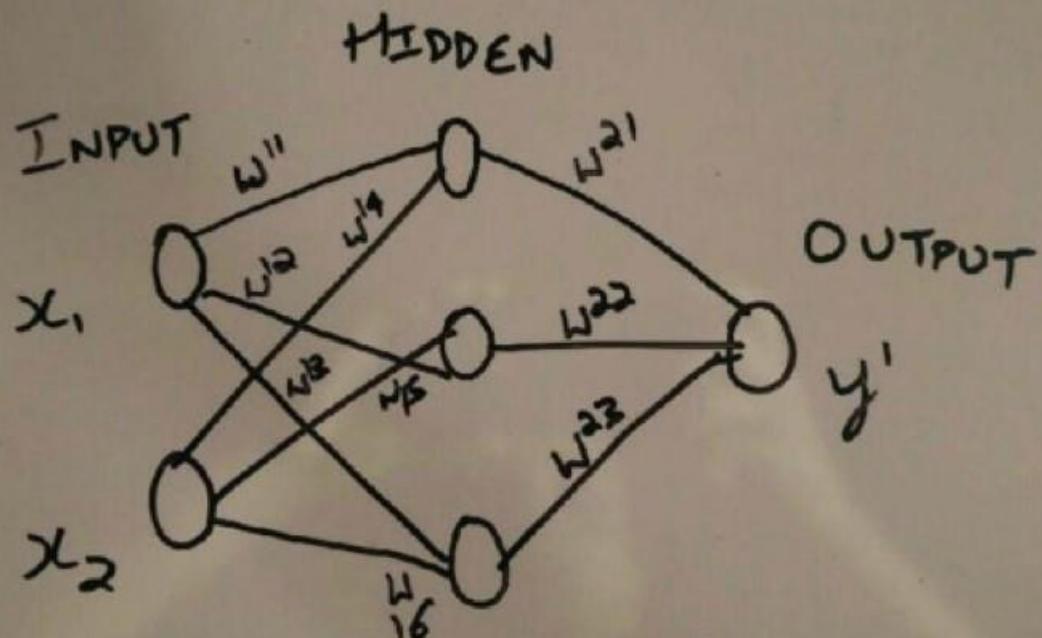
Use gradient descent to update weights

References:

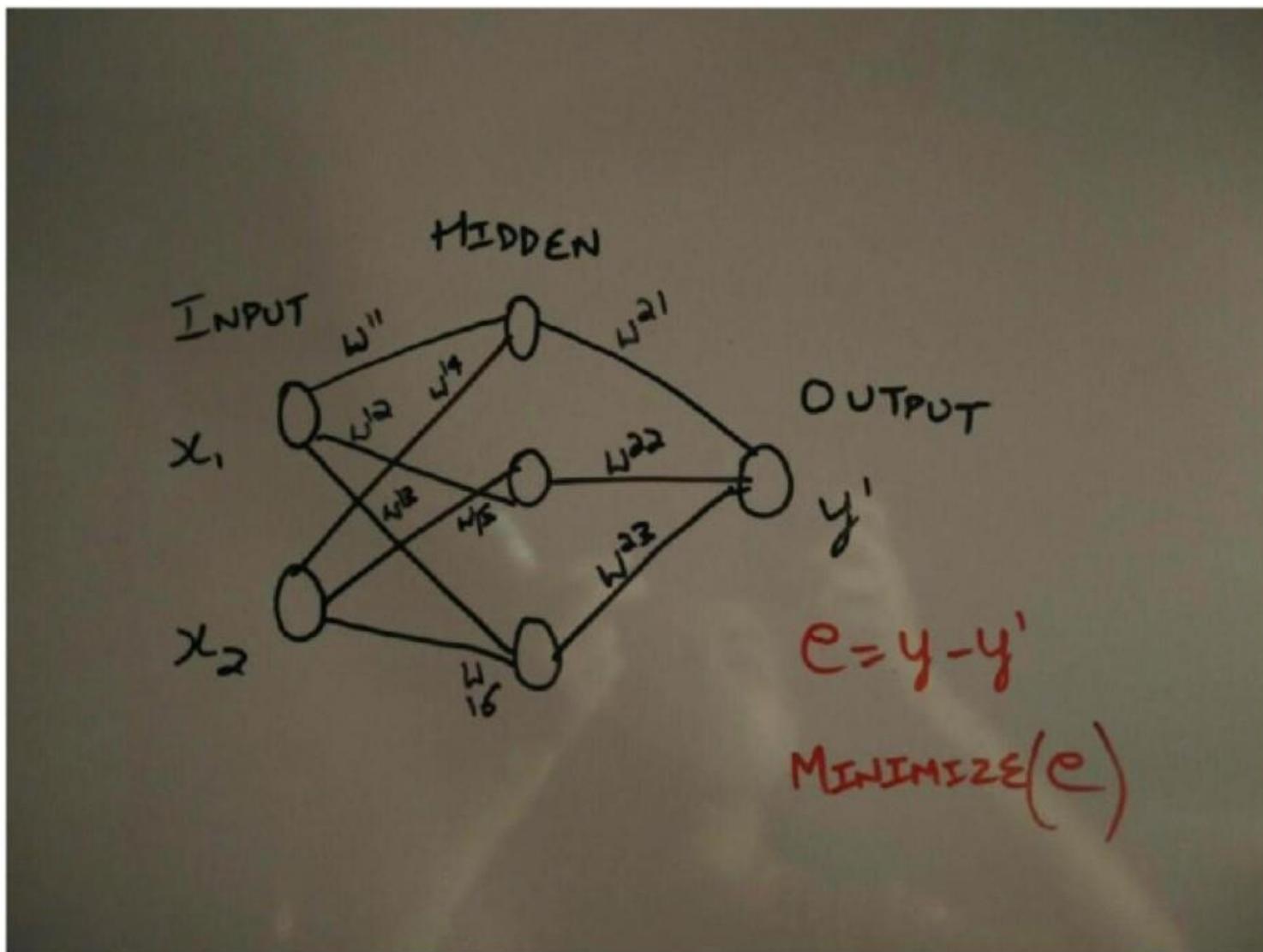
colah.github.io/posts/2015-08-Backprop/

<http://neuralnetworksanddeeplearning.com/chap2.html>

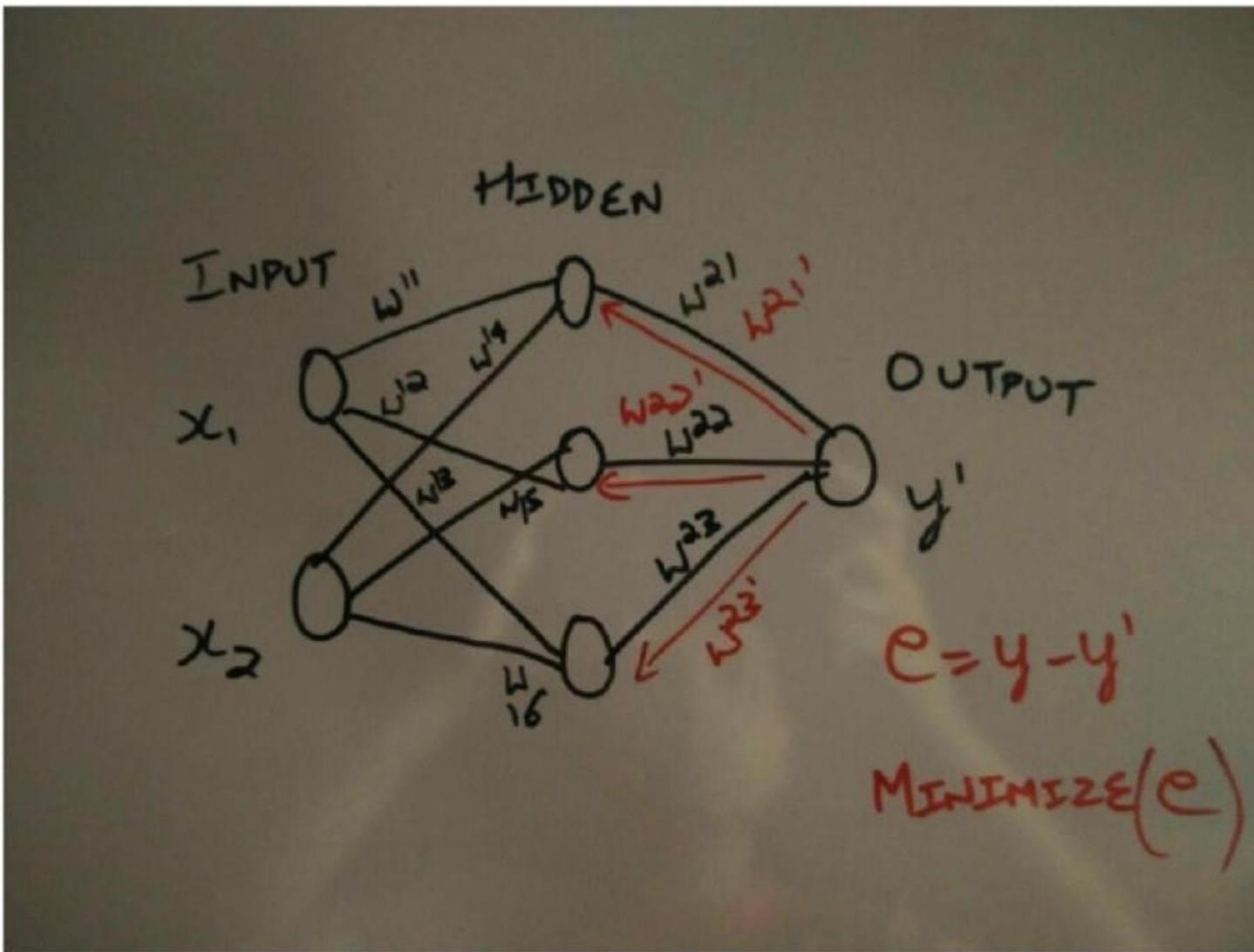
Backpropagation Algorithm



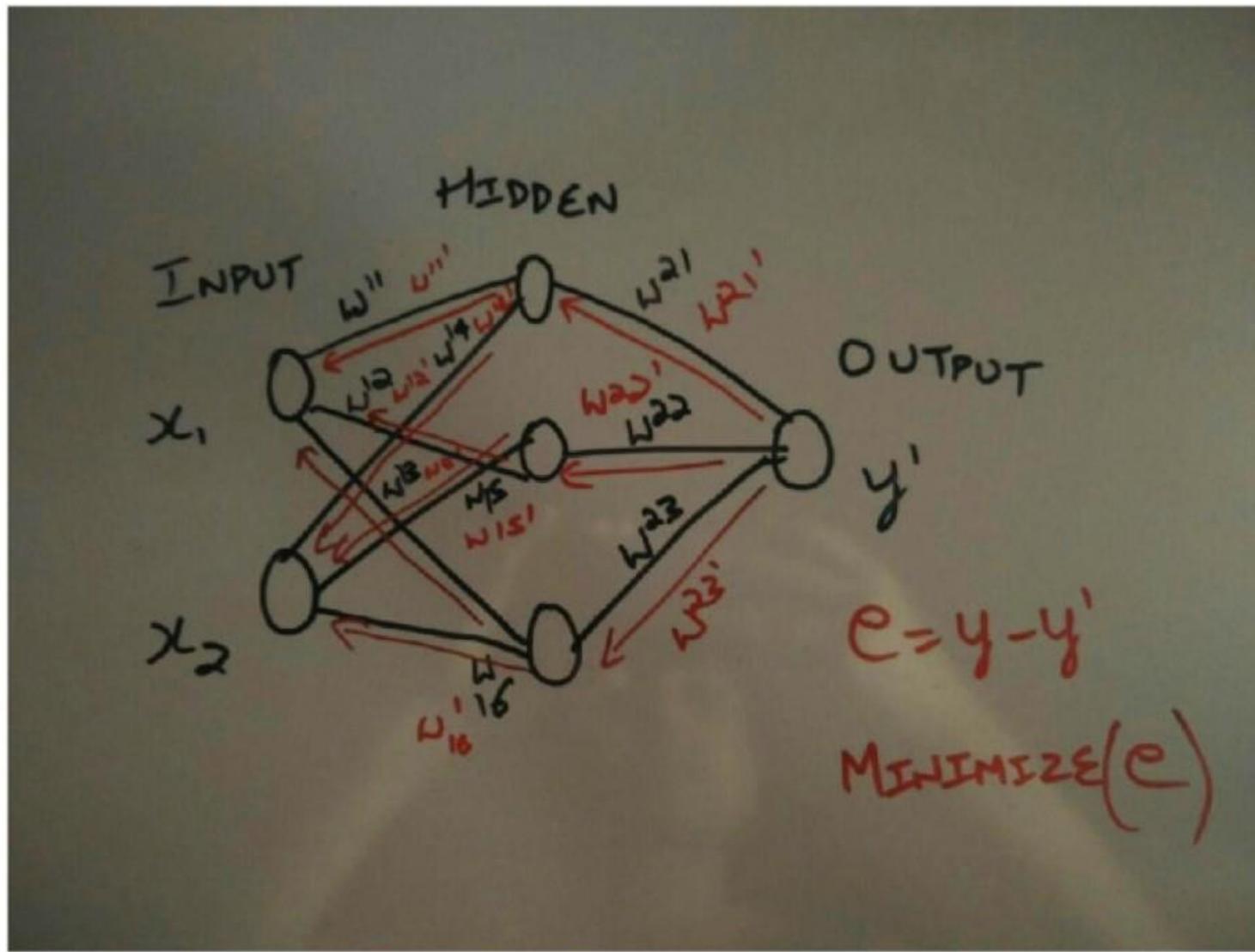
Backpropagation Algorithm



Backpropagation Algorithm



Backpropagation Algorithm



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.
- ❑ There are other variants of the loss function
 - with better empirical results than SGD
(Not covered in the slides)

SGD/Mini-Batch/Online

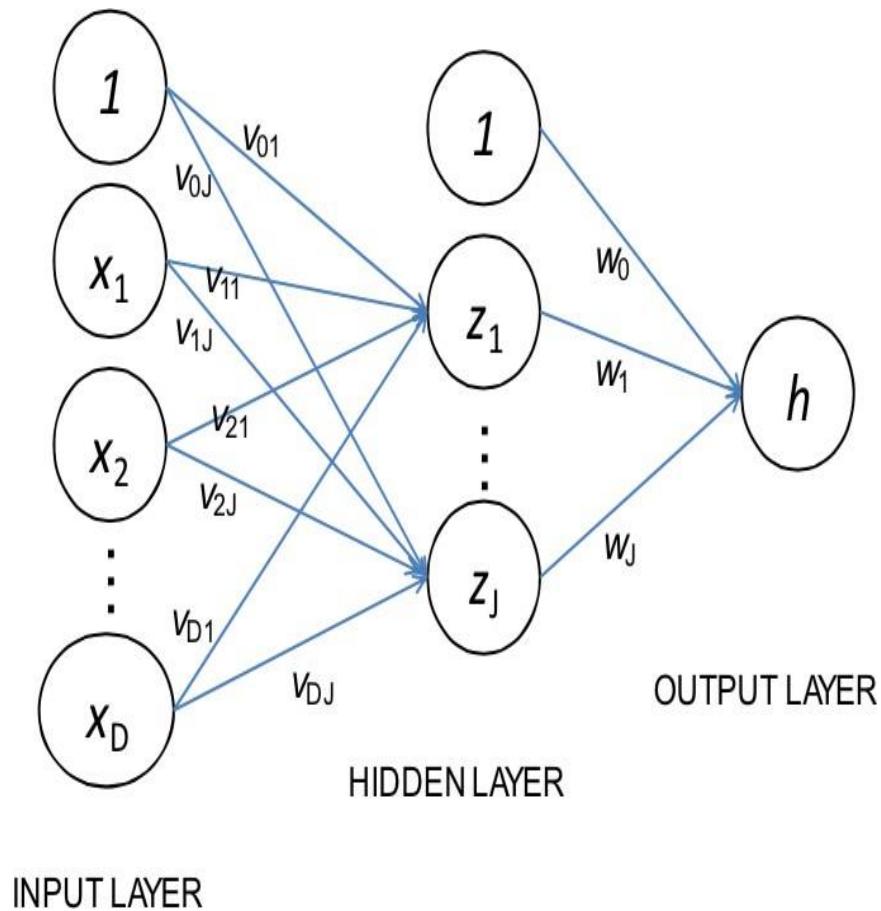
Algorithm 2 Minibatch Stochastic Gradient Descent Training

```
1: Input: Function  $f(\mathbf{x}; \theta)$  parameterized with parameters  $\theta$ .  
2: Input: Training set of inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$  and outputs  $\mathbf{y}_1, \dots, \mathbf{y}_n$ .  
3: Input: Loss function  $L$ .  
4: while stopping criteria not met do  
5:   Sample a minibatch of  $m$  examples  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_m, \mathbf{y}_m)\}$   
6:    $\hat{\mathbf{g}} \leftarrow 0$   
7:   for  $i = 1$  to  $m$  do  
8:     Compute the loss  $L(f(\mathbf{x}_i; \theta), \mathbf{y}_i)$   
9:      $\hat{\mathbf{g}} \leftarrow \hat{\mathbf{g}} + \text{gradients of } \frac{1}{m}L(f(\mathbf{x}_i; \theta), \mathbf{y}_i) \text{ w.r.t } \theta$   
10:     $\theta \leftarrow \theta + \eta_k \hat{\mathbf{g}}$   
11: return  $\theta$ 
```



Our First Architecture

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



Machine Translation



source:<http://blog.webcertain.com/machine-translation-technology-the-search-engine-takeover/18/02/2015/>

Video Classification

Google's neural net learns just by watching youtube videos

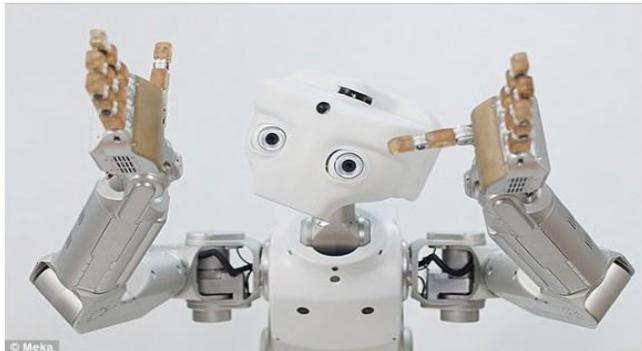
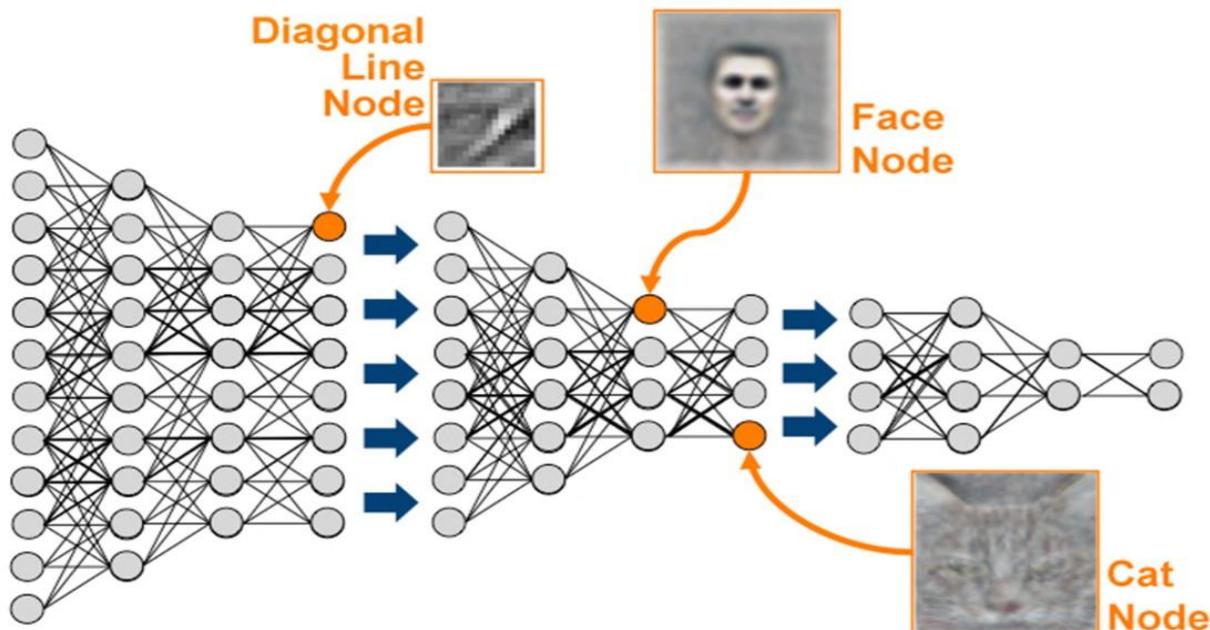


Image Recognition

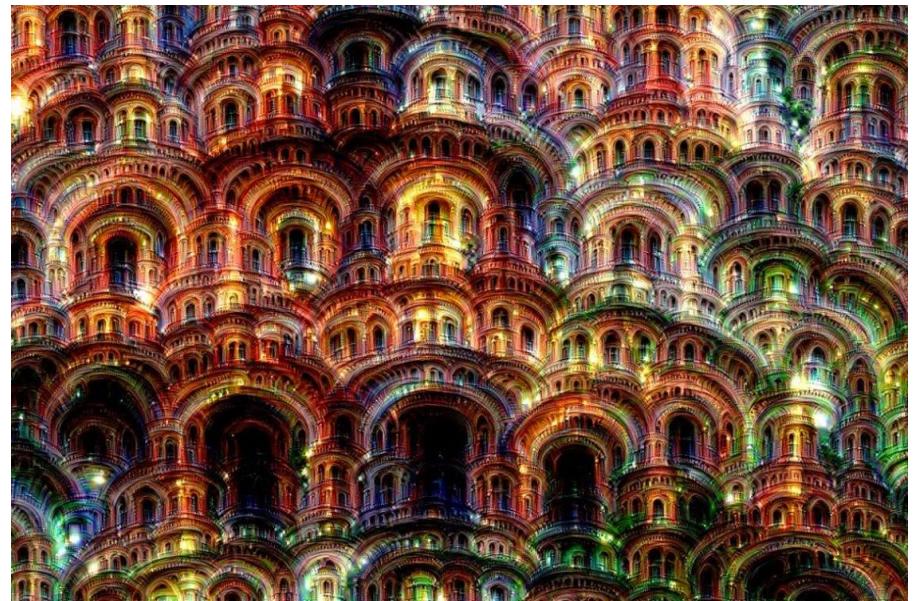


Automatic Generation of Cooking Recipe

```
1  MMMMM
2
3  MMMMM----- Recipe via Meal-Master (tm) v8.05
4
5      Title: BARBECUE RIBS
6  Categories: Chinese, Appetizers
7  Yield: 4 Servings
8
9      1 pk Seasoned rice
10     1   Beer -- cut into
11         -cubes
12     1 ts Sugar
13     3/4 c Water
14         Chopped finels,
15         -up to 4 tb1sp of chopped
16     2 pk Yeast Bread/over
17
18  MMMMM-----FILLING-----
19     2 c Pineapple, chopped
20     1/3 c Milk
21     1/2 c Pecans
22         Cream of each
23     2 tb Balsamic cocoa
24     2 tb Flour
25     2 ts Lemon juice
26         Granulated sugar
27     2 tb Orange juice
28     1 c Sherry wheated curdup
29     1 Onion; sliced
30     1 ts Salt
31     2 c Sugar
32     1/4 ts Salt
33     1/2 ts White pepper, freshly ground
34         Sesame seeds
35     1 c Sugar
36     1/4 c Shredded coconut
37     1/4 ts Cumin seeds
38
39     Preheat oven to 350. In a medium bowl, combine milk, the sugar, vanilla and seasoned
40     flour and water and then cornstarch. add tomatoes, oregano, and
41     nutmeg; serve.
```



Image Generation – Google Inceptionism



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

Natural Language Processing

As ducemos lo esto: digamos el mejor.
Lo q caye en precio: del seo confessor.
Hizo se aduzir: este aigo lajada.
A la casa del monge: de suso ementado.
Ca creye bien afirmes: estua syntado.
Sarie desta coyta: por ell, terminado.
Hizo fue ala pueria: de sanz fabustan.
Non quis el mesqno: pedir uno ni p'm.
As dice ay padres: por senor san...
Te prend
Pare allá
E sei tu
Enno: yo non podria: partarme des' legar.
Tu me m'indes: ser etor...
Padre delos fajados: dena ne usitar.
Y en sobre mi tu mano: signa me del polgar.
Solo q yo pridiesse la tu mano: besar.
De aya esta coyta: cuydaria sacar
Al padre beneyto: bien entre do estudi.
A los apellidos: q este aigo dava.
Caro: e preguntó le: q'l cosa demandava.
Iro el a humne: ca el non cobriana.
Venme s'a comisa: un enemigo.

Echol con el ylopo: del agua salada.
O onsigno li los ojos con la cruz consagrad.
La dolor: la coyta fue luego amansada.
La humore q pdiera: fue toda recobrada.
Entenderio pridieses: amigos i señores.
Quie muchos males: de diuersos colores.
V nos de ceguedad: al de gues dolores.
Mas de todo bién sano: rendicados lodos.

Cuanta q no pagas: no tagal tolha.
Ca sera po tu tido: faces recadia.
Muchos son los mirados: q des' pad' sabemos.
El p'mel us' uino o ov'mos: los otros q leemos.
Endubda: ei portavas: en q'l empazaremos.
No q has aq'l parte: q sei: adechar aruremos.
Saco Desta saxon los est: qero los fer esquios.
en su D esir uno i mebie vos: mette fuerdes uno.
mida. Como gano la grá: q saca los catinos.
Por ond de huengas tierras: le enbia bodigos.
Ena en essi tiempo: los moros muy uegros.

Some NLP Tasks

- Sentiment Analysis
- Machine Translation
- Document Classification
- Language Modeling
- Event Detection
- Question-Answering
- Chat Bots (Response Generation)

PLEASE

RELAX ~~STEAL~~ DANCE TOUCH

LIRT ~~SME~~ FEEL

Input Features

USE ~~EAT~~ SING LISTEN TALK

OUCH ~~NEON~~ LOOK COMMUNIC

EACH OTHER CAMERA ~~FLA~~

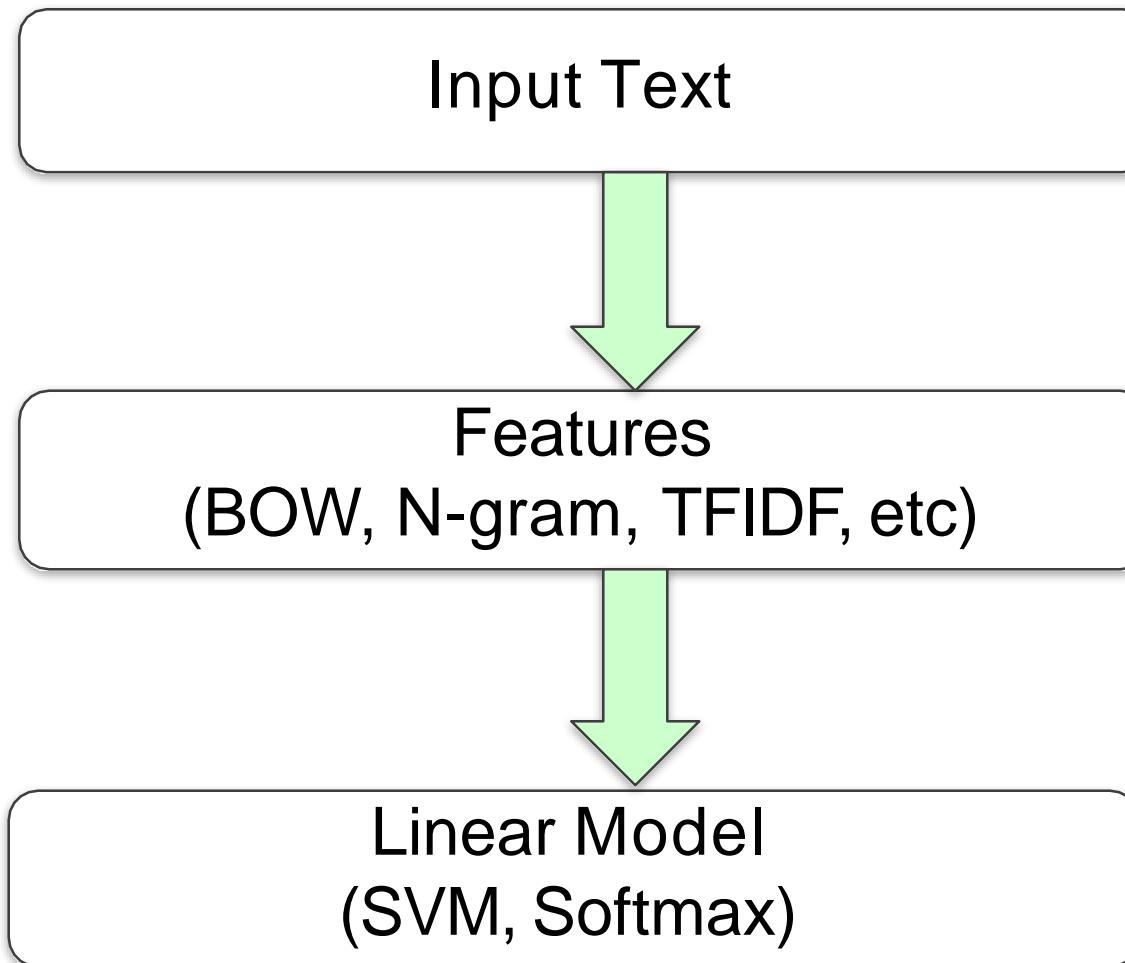
How to create input features?

- ❑ Bag of Words

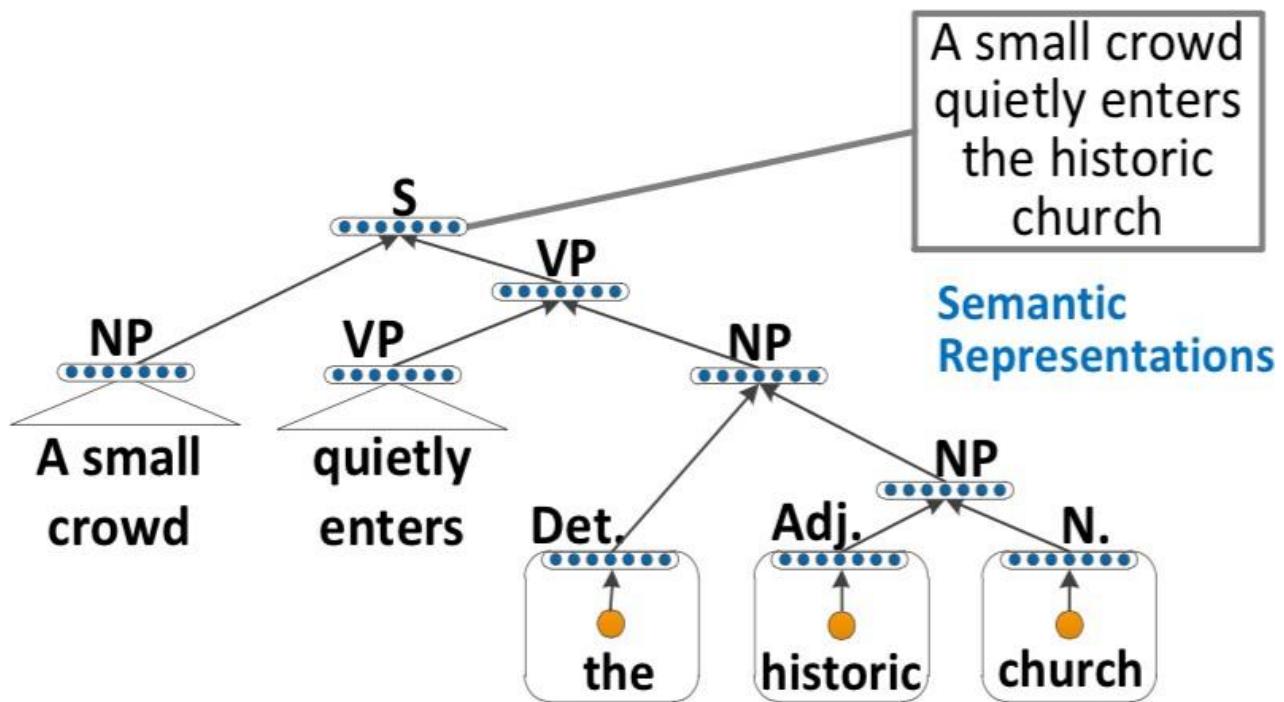
- ❑ N-gram

- ❑ TF-IDF

Existing ML framework for NLP



A Sentence



Structure is very important

A small crowd quietly
enters the historic church



A Historic crowd enters
the small quietly church

Structure important for:
Humor
Sarcasm

Representing structure is hard

n-grams
– quickly explodes

Limitation of the architectures so far

Fixed Size Input
(Eg: Image)

Fixed computational steps
(Eg: Number of layers)

Fixed Size Ouput
(Eg: probabilities of different classes)

Limitation of the architectures so far

Hierarchy captured. But
context, structure ?
– mostly NOT !

NLP -Deep Learning Classification System

The general structure for an NLP classification system based on a feed-forward neural network is thus:

1. Extract a set of core linguistic features f_1, \dots, f_k that are relevant for predicting the output class.
2. For each feature f_i of interest, retrieve the corresponding vector $v(f_i)$.
3. Combine the vectors (either by concatenation, summation or a combination of both) into an input vector \mathbf{x} .
4. Feed \mathbf{x} into a non-linear classifier (feed-forward neural network).

Word Embeddings

- ❑ Vectorization
- ❑ Dense-representations on fixed-size vectors
(low dimensional space)

Word Embeddings

Goal

- “learning high-quality word vectors from huge data sets with billions of words, and with millions of words in the vocabulary”
- Resulting word representations
 - Similar words tend to be close to each other
 - Words can have multiple degrees of similarity

Word Embeddings

“You shall know a word
by the company it keeps”

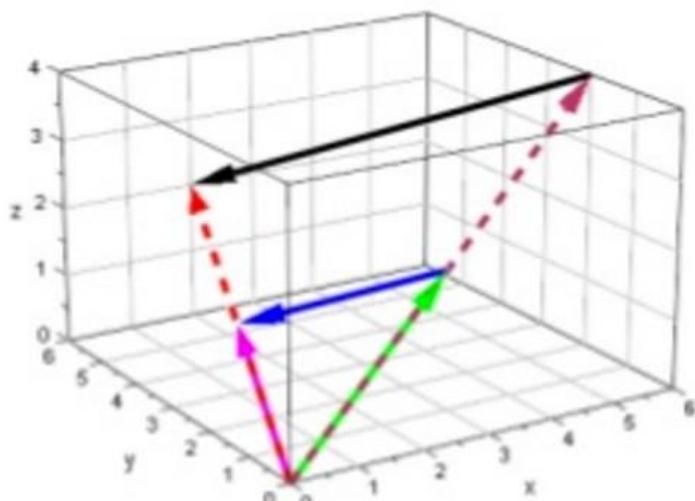
One of the most successful ideas of modern statistical NLP

government debt problems turning into banking crises as has happened in
saying that Europe needs unified banking regulation to replace the hodgepodge

☛ These words will represent *banking* ☛

Quoted after Socher

Word Embeddings



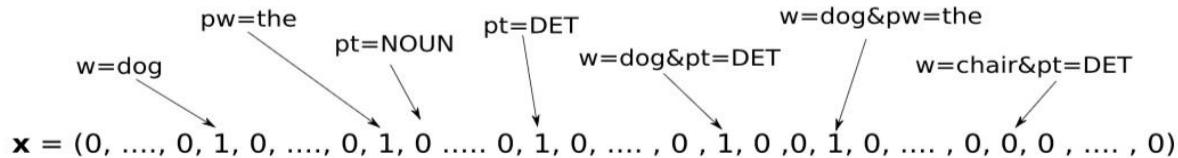
$$\text{linguistics} = \begin{bmatrix} 0.286 \\ 0.792 \\ -0.177 \\ -0.107 \\ 0.109 \\ -0.542 \\ 0.349 \\ 0.271 \end{bmatrix}$$

Quoted after Socher

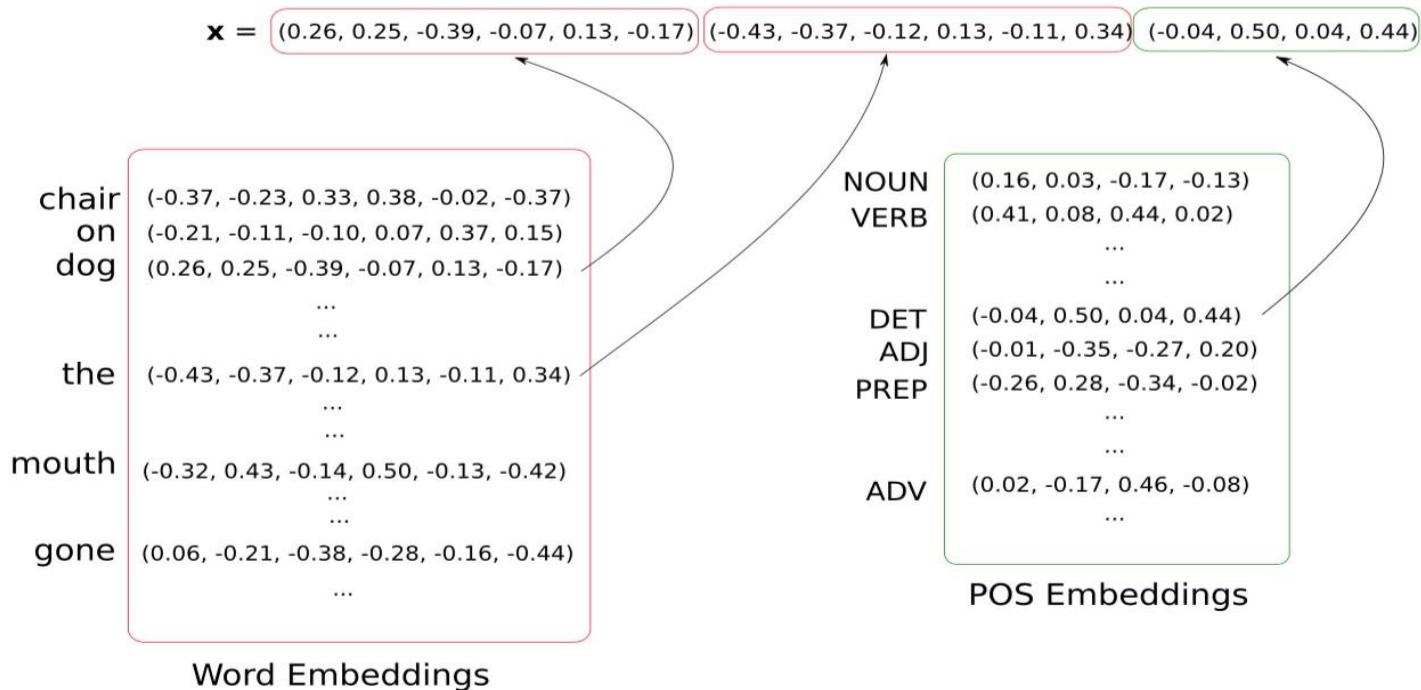
Sparse Vs Dense Feature Representation:

Example: “the dog”

(a)



(b)

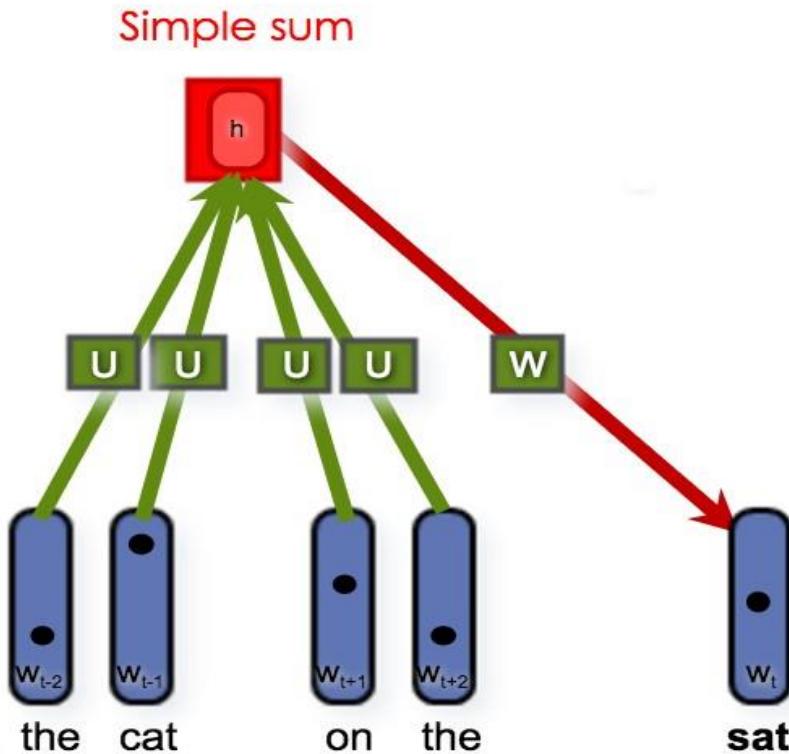


CBOW

word embedding space \mathbb{R}^D
in dimension $D=100$ to 300

Word embedding matrices

discrete word space $\{1, \dots, V\}$
 $V > 100k$ words



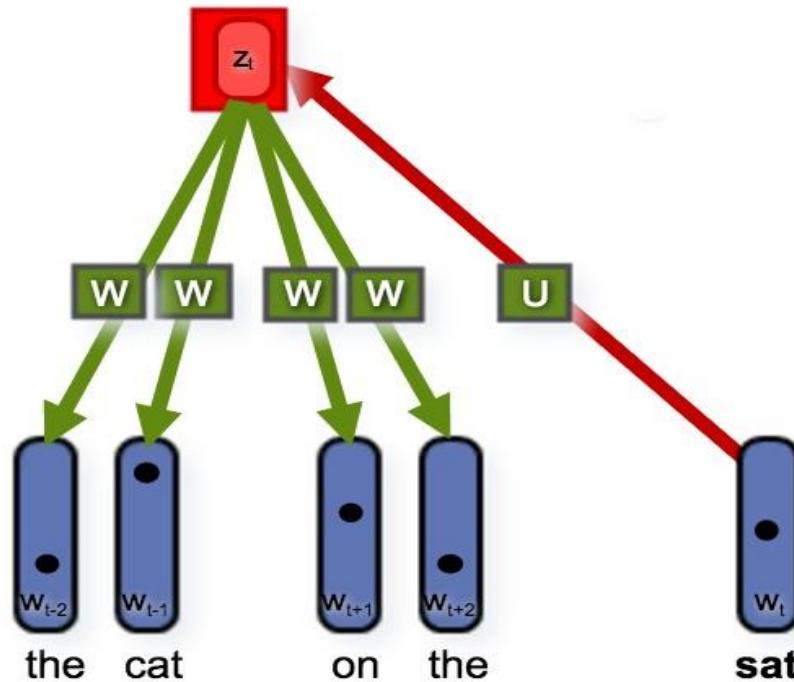
[Mikolov et al, 2013a; Mnih & Kavukcuoglu, 2013;
<http://code.google.com/p/word2vec>]

Skip Grams

word embedding space \mathbb{R}^D
in dimension
 $D=100$ to 1000

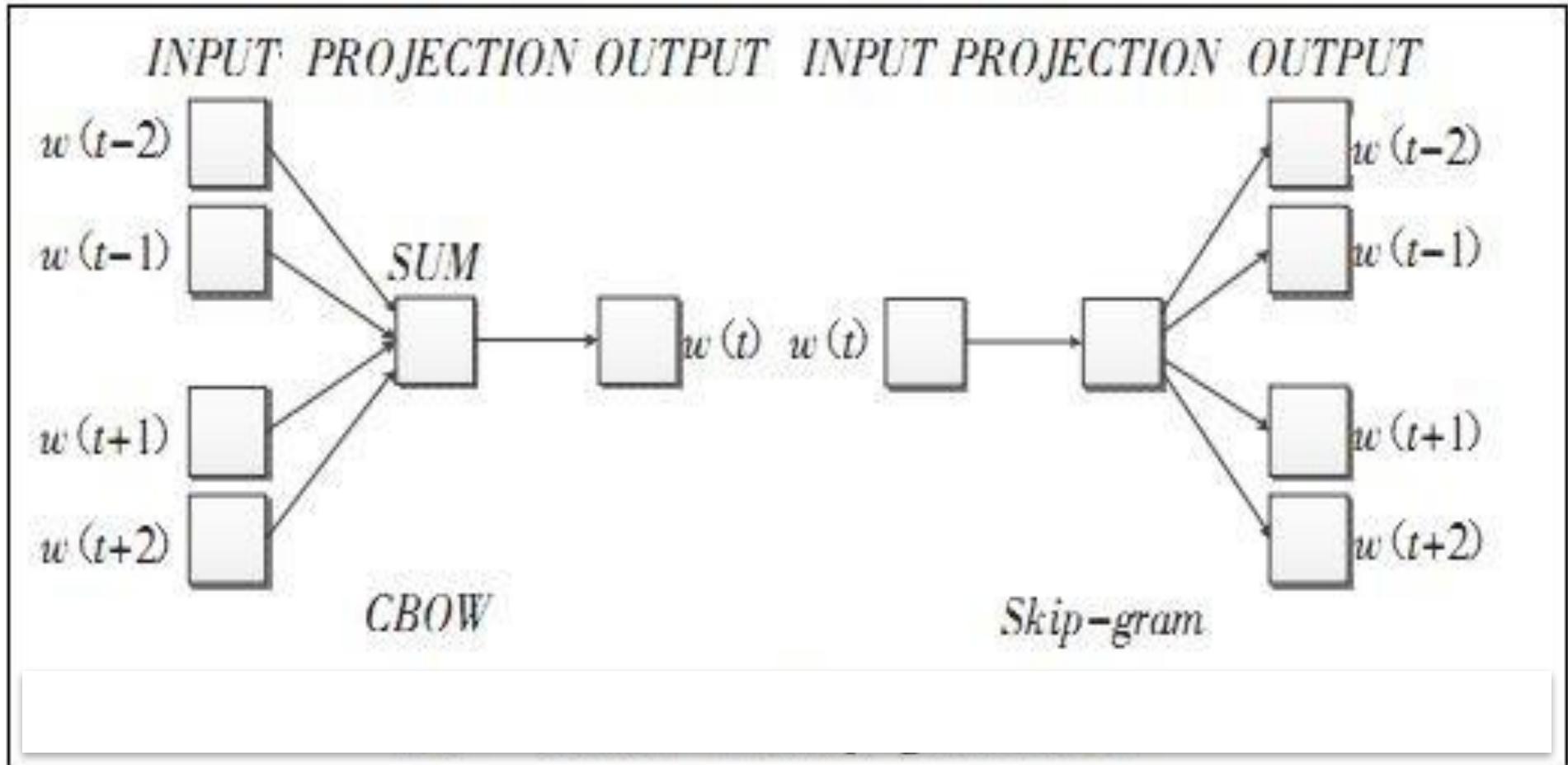
Word embedding matrices

discrete word space $\{1, \dots, V\}$
 $V > 100k$ words



[Mikolov et al, 2013a, 2013b; Mnih & Kavukcuoglu, 2013;
<http://code.google.com/p/word2vec>]

CBOW Vs Skip-gram



source: <http://blog.csdn.net/eastmount/article/details/50637476>

word2vec (Google)

Use documents to train a neural network model maximizing the conditional probability of context given the word

Apply the trained model to each word to get its corresponding vector

Calculate the vector of sentences by averaging the vector of their words

Construct the similarity matrix between sentences

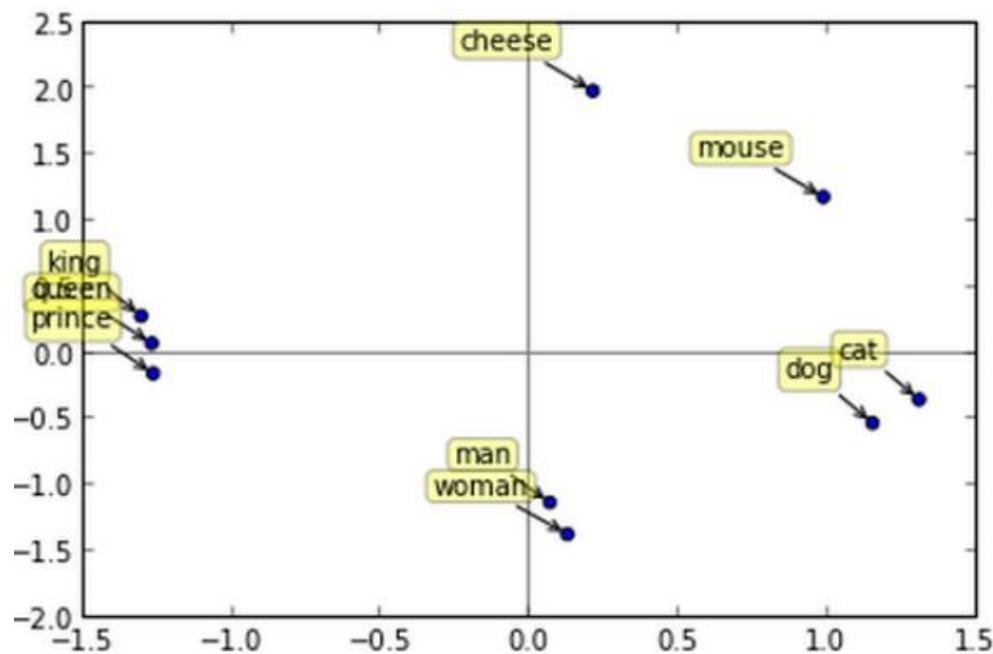
Use pagerank to score the sentences in graph

word2vec

word2vec associates words to points in space

Similar words are closer together

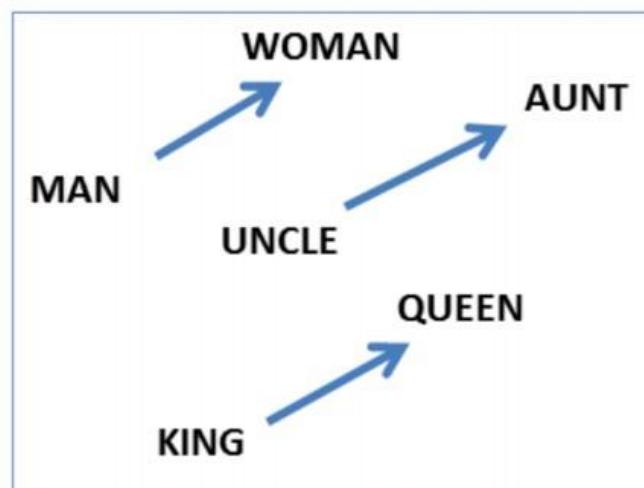
Developed by Mikolov, Sutskever, Chen, Corrado and Dean in 2013 at Google Research



source: http://files.meetup.com/12426342/5_An_overview_of_word2vec.pdf

word2vec

Word relationships are displacements



Source: *Linguistic Regularities in Continuous Space Word Representations*, Mikolov et al, 2013

$$\text{vec[queen]} - \text{vec[king]} = \text{vec[woman]} - \text{vec[man]}$$

source: http://files.meetup.com/12426342/5_An_overview_of_word2vec.pdf

A silhouette of a person in mid-air, performing a dynamic pose, set against a backdrop of a mountain range at sunset. The sky is filled with dramatic, colorful clouds transitioning from blue to orange and yellow. The mountains are dark silhouettes against the bright horizon.

HANDS-ON

CONVOLUTIONAL NEURAL NETWORKS

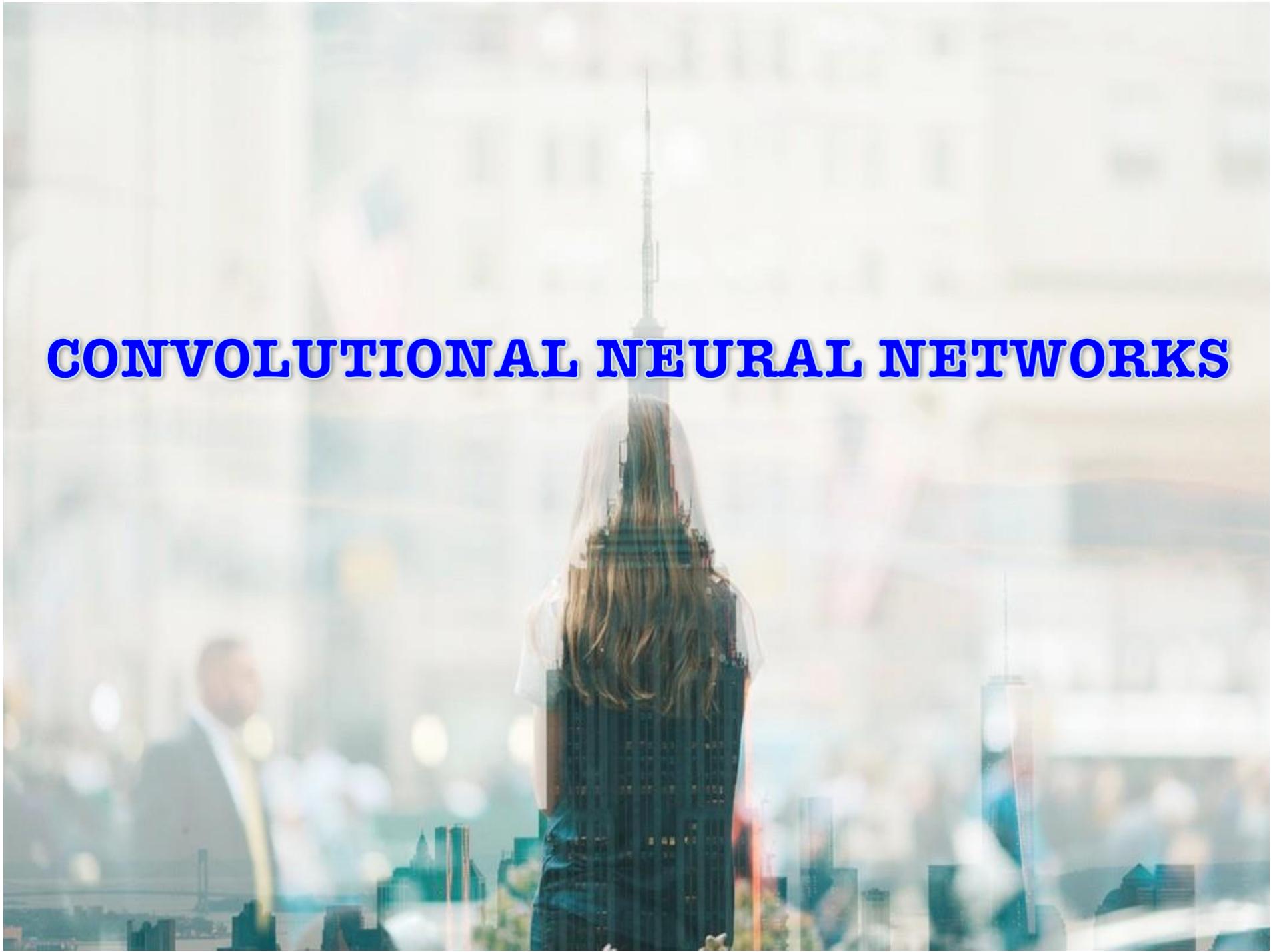
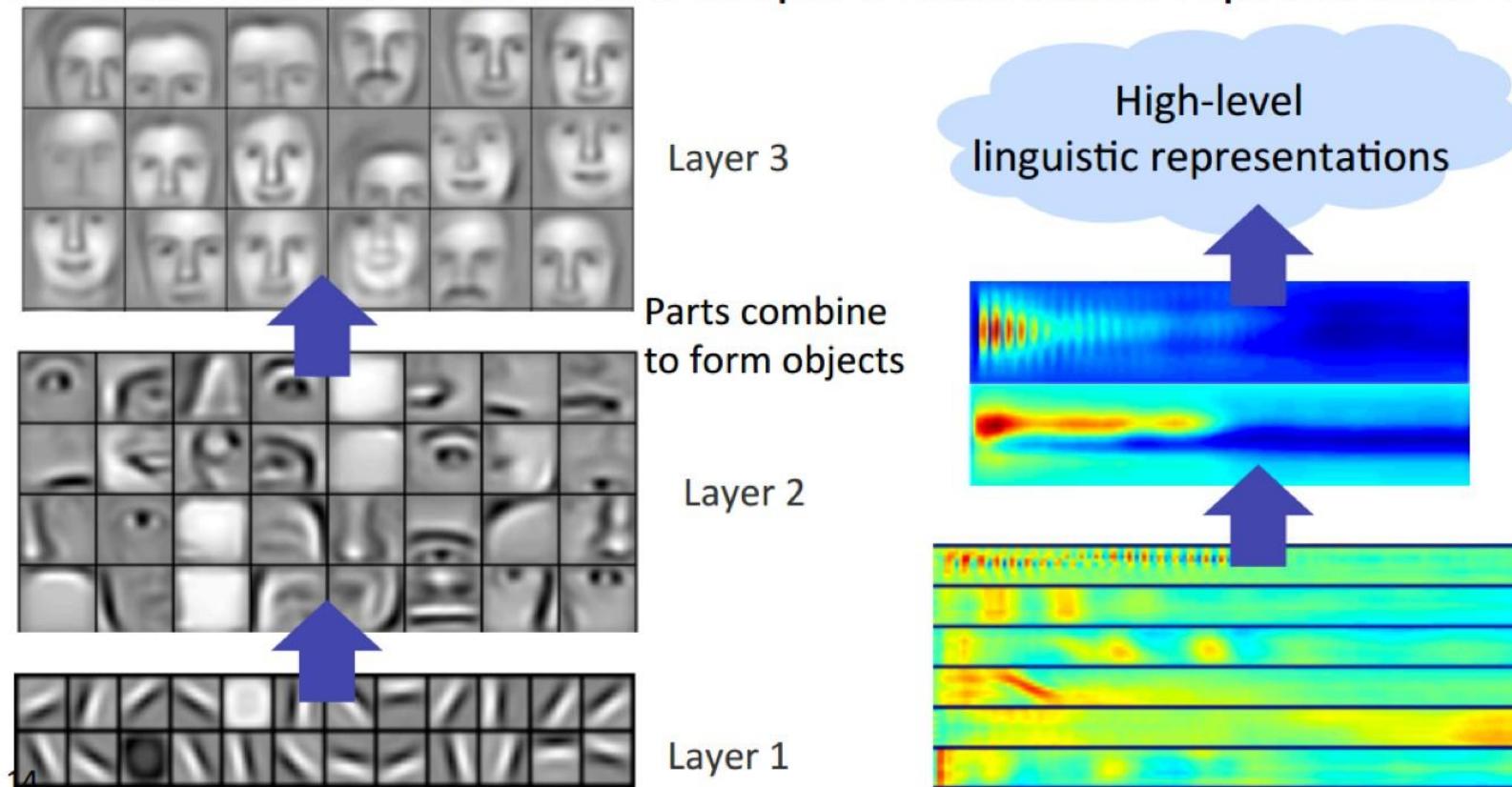


Image Recognition/Classification

Successive model layers learn deeper intermediate representations



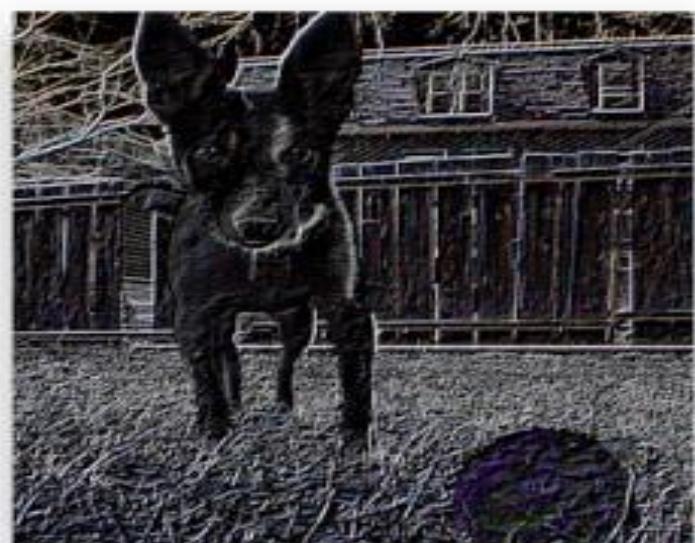
Prior: underlying factors & concepts compactly expressed w/ multiple levels of abstraction

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

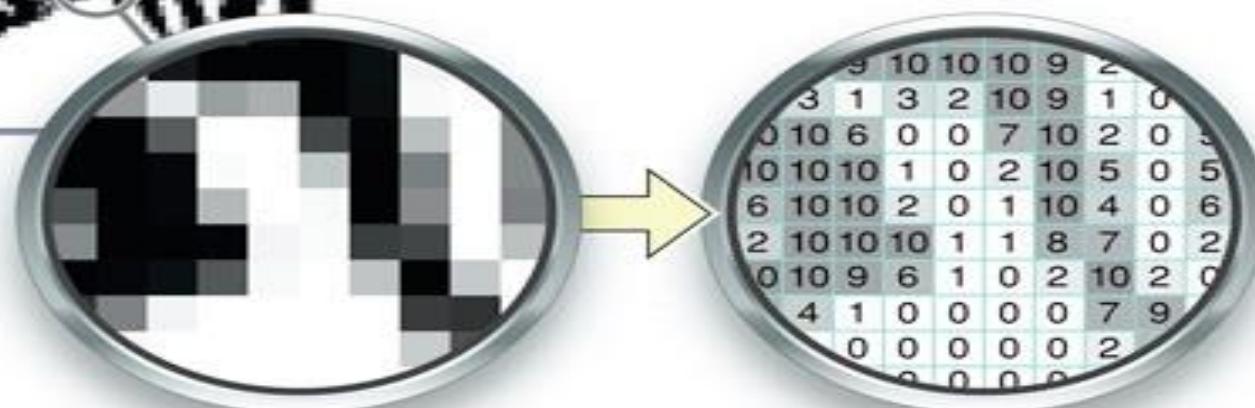


Emboss

Convolution

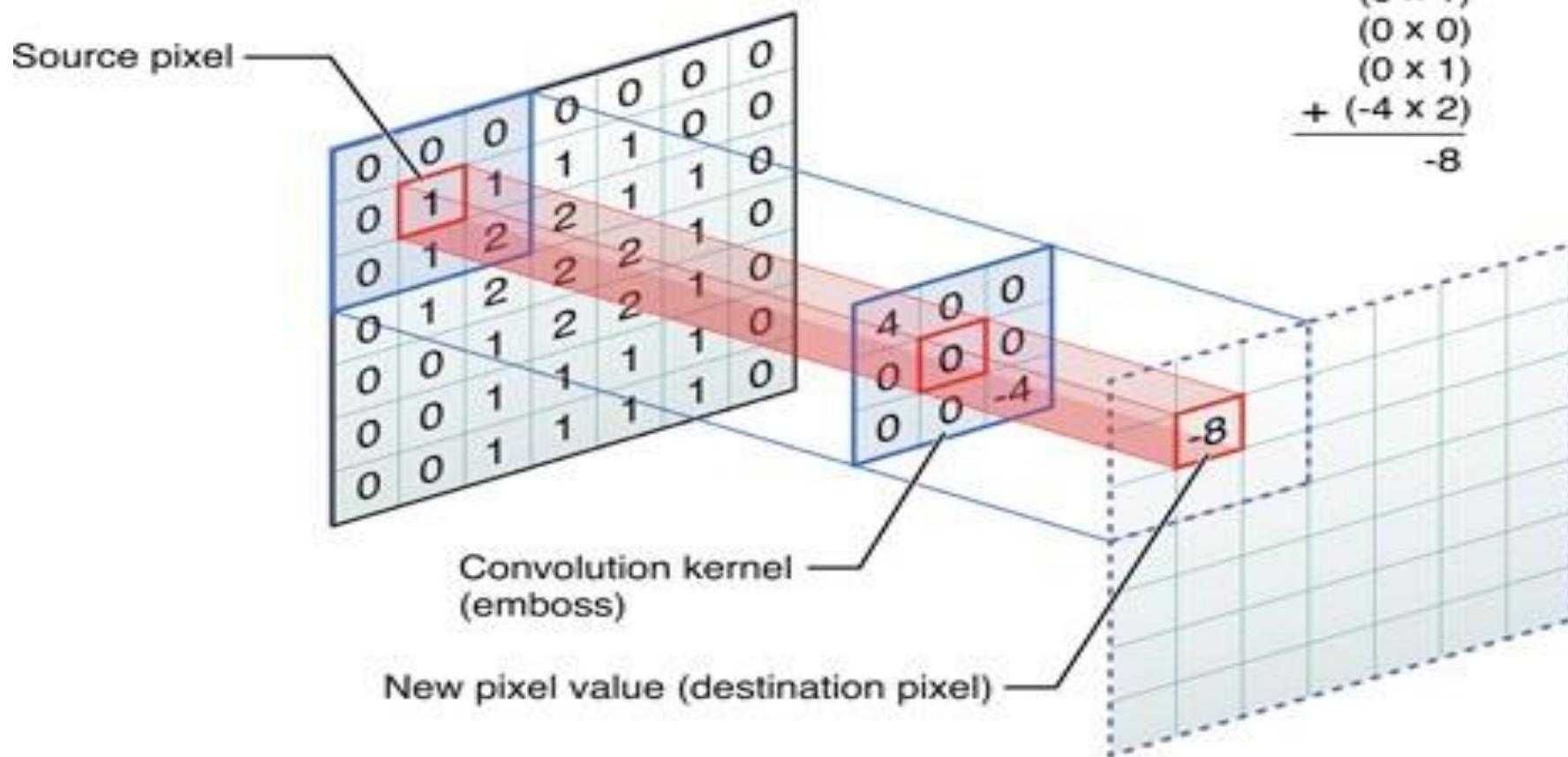


Pixels depicted by a grid of numbers representing intensity



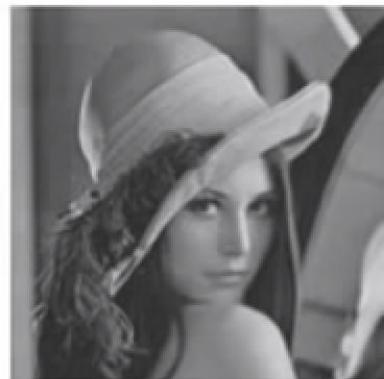
Convolution

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.

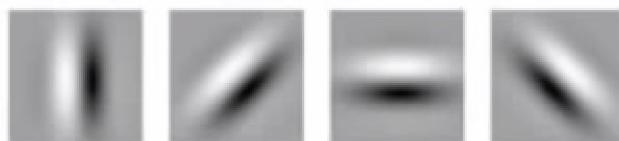


Convolution – One more example

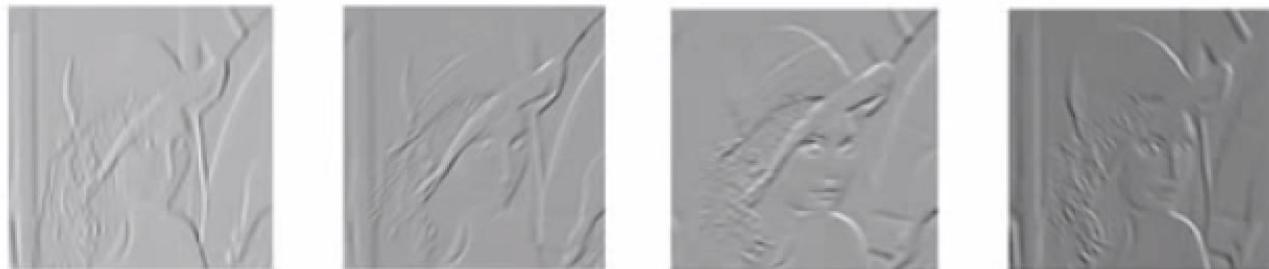
Lena



Gabor filters



Convolved by Gabor filters lena



Sub-sampling



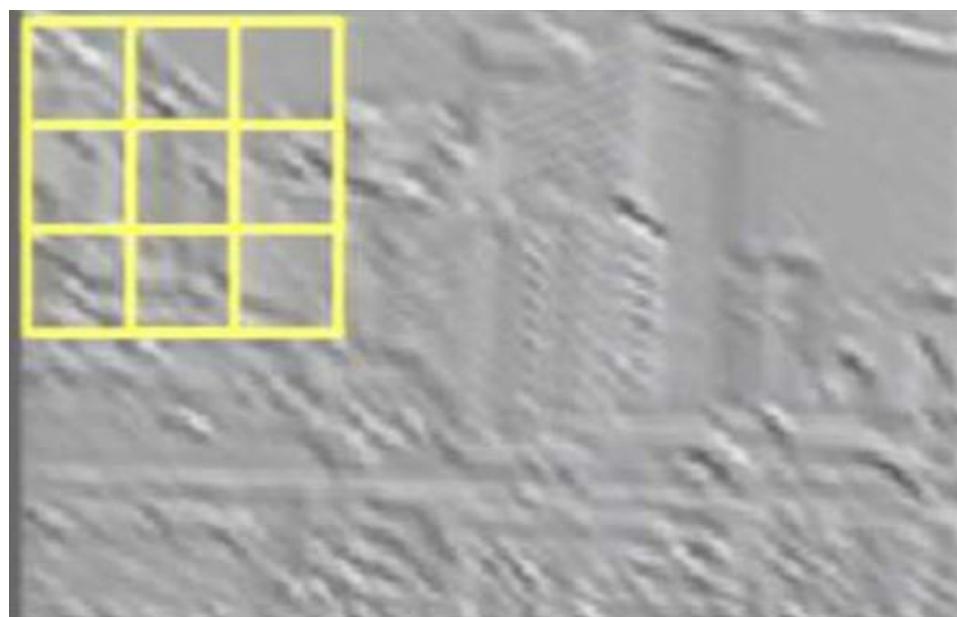
Input 508x508 pixels Lena

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



Input 254x254 pixels Lena

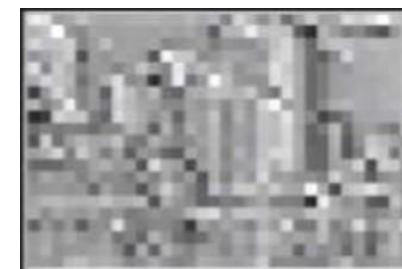
Pooling



Max

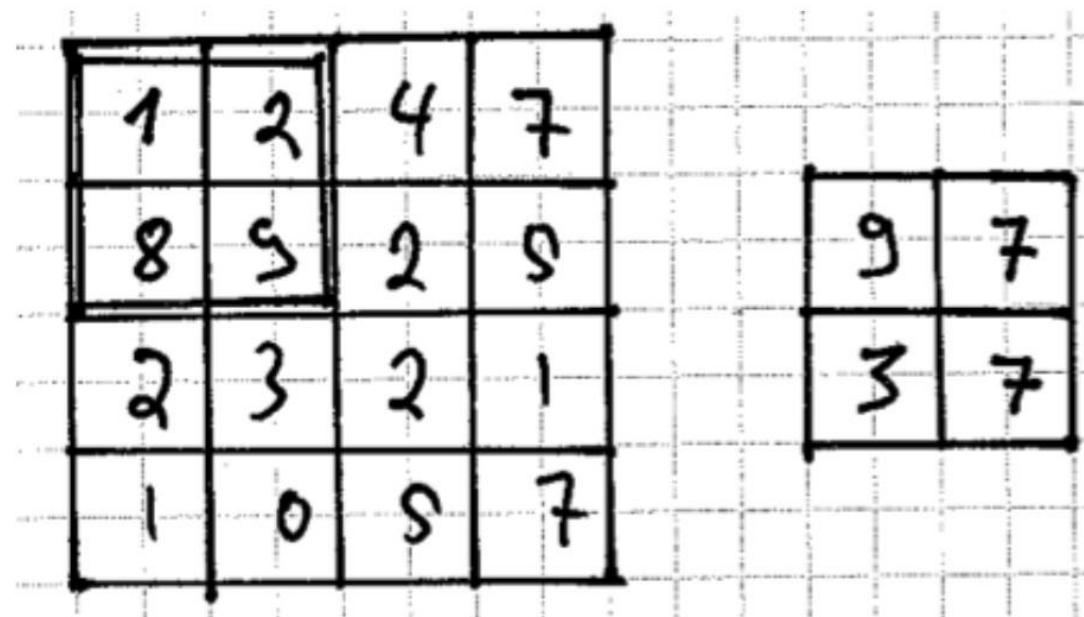


Sum



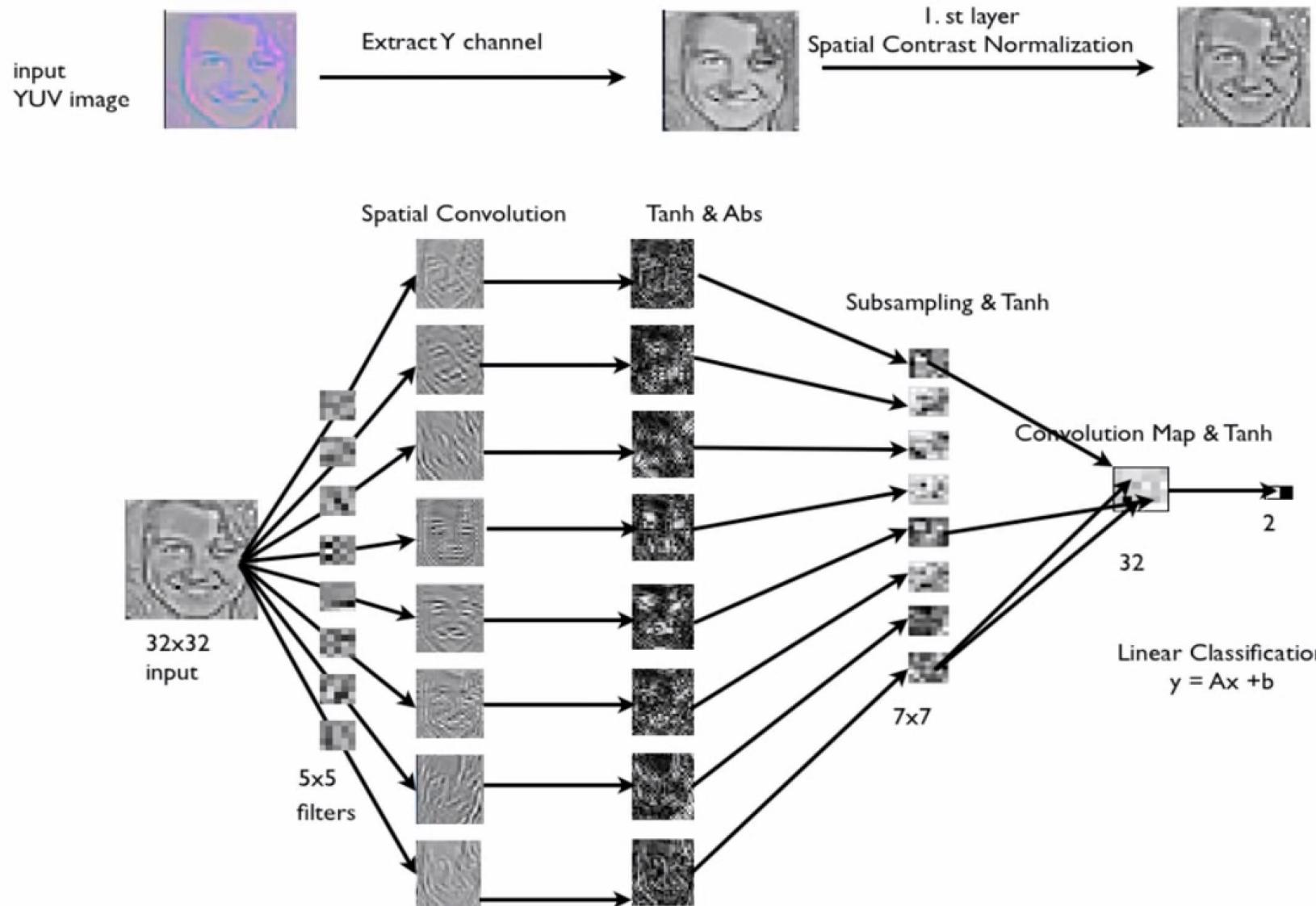
Max Pooling

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

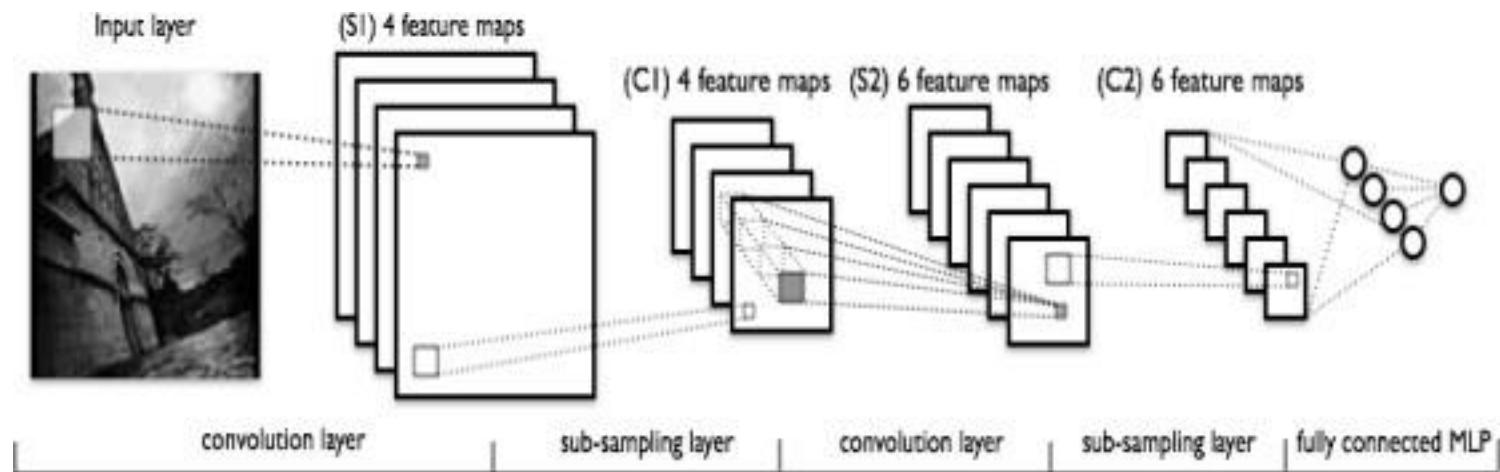


Simply join e.g. 2x2 adjacent pixels in one.

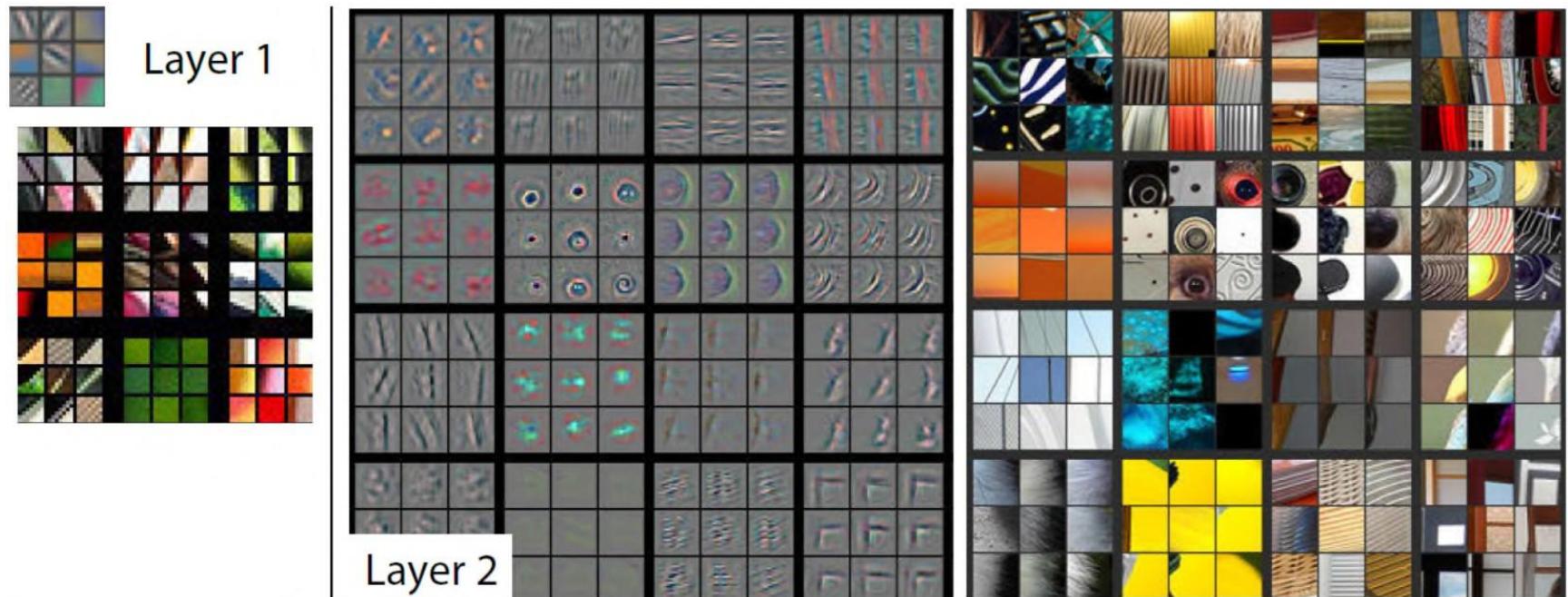
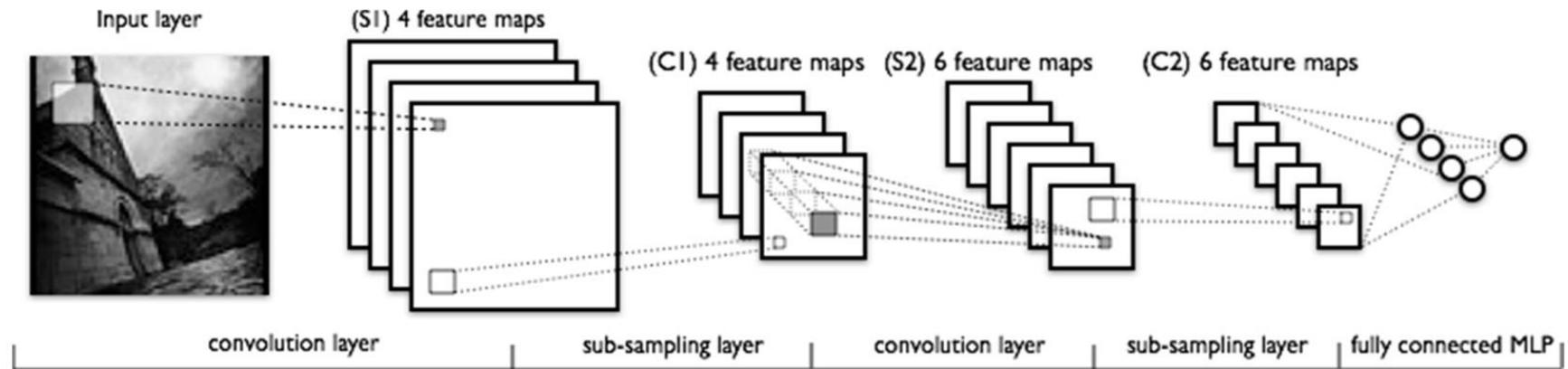
Basic Architecture



CNN Architecture

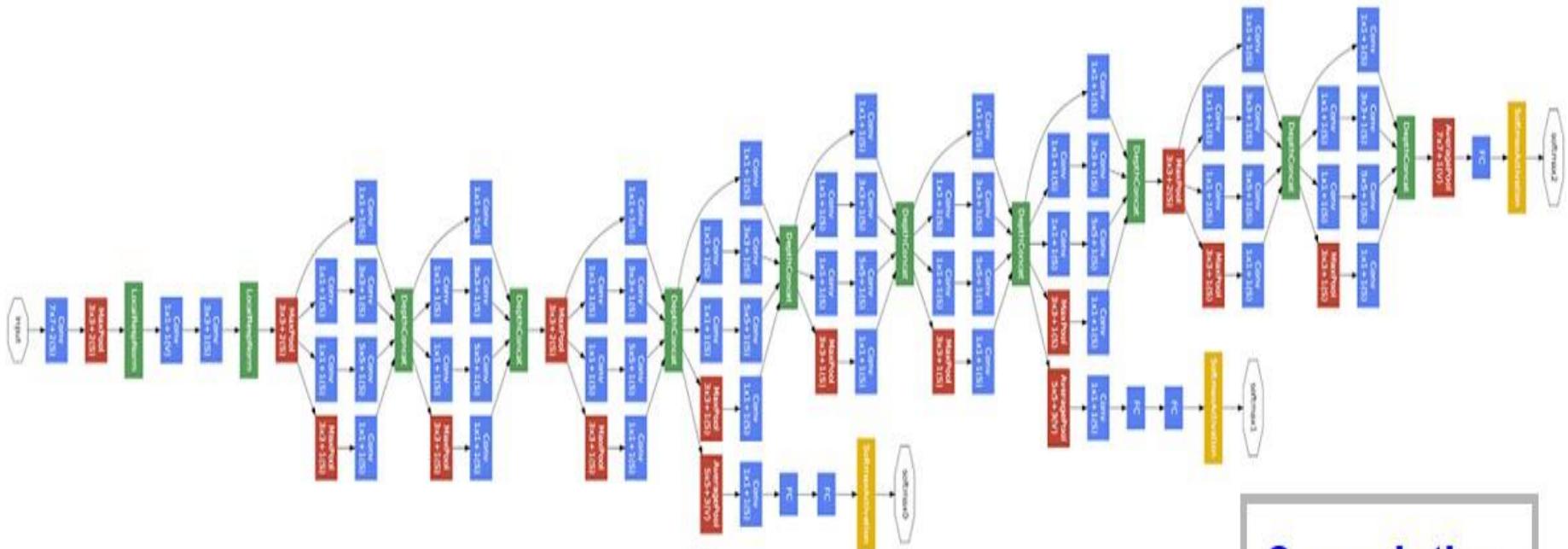


CNN Architecture



[Matthew Zeiler & Rob Fergus]

Example of a CNN model

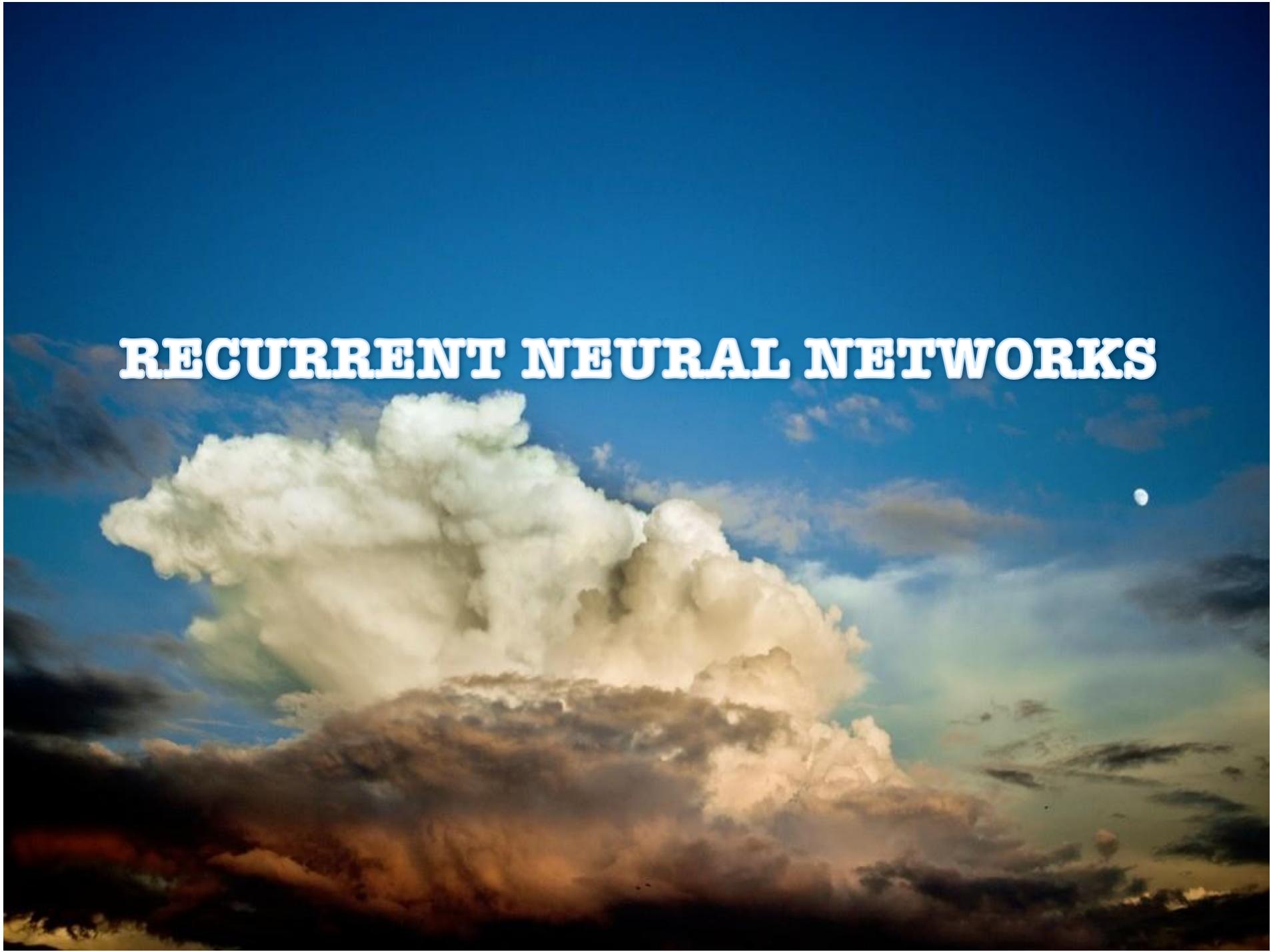


Convolution
Pooling
Softmax
Other

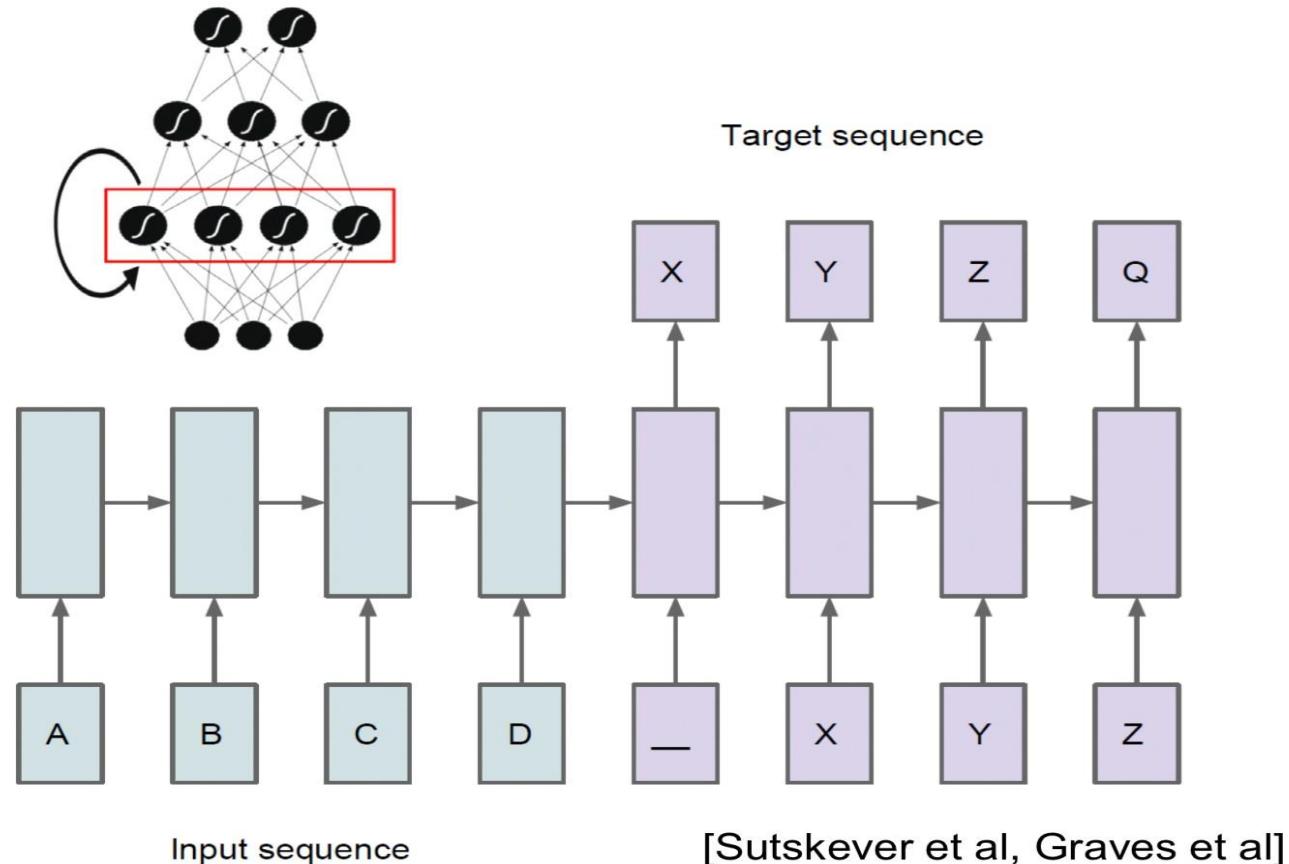
A ... model? - GoogLeNet

source: arXiv:1409.4842v1 [cs.CV] 17 Sep 2014

RECURRENT NEURAL NETWORKS

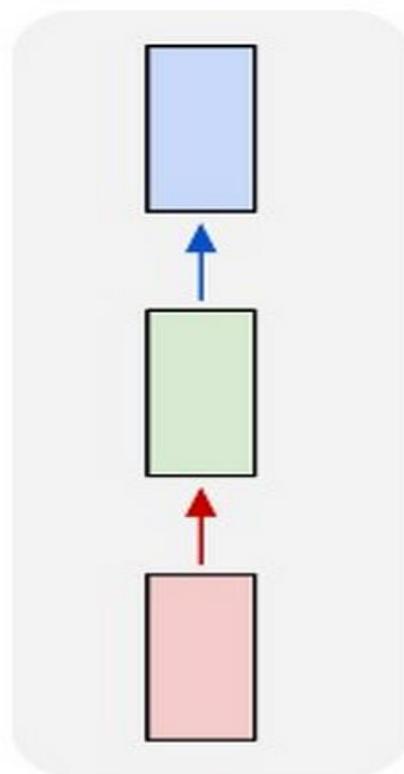


Recurrent Neural Network



One-to-One Sequence

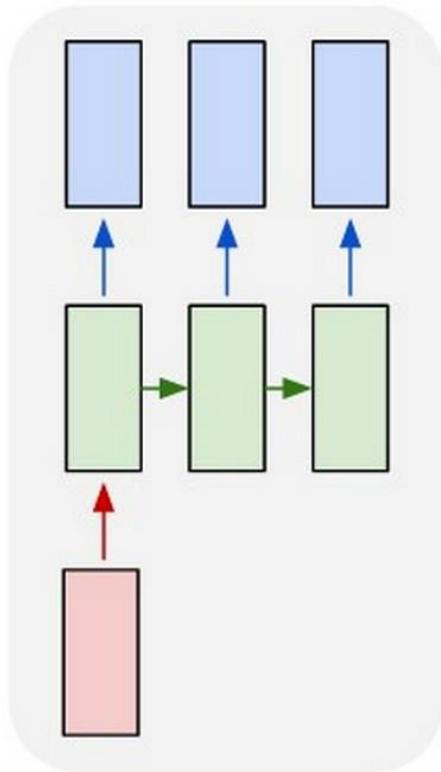
one to one



Vanilla-model. No RNN

One-to-Many Sequence

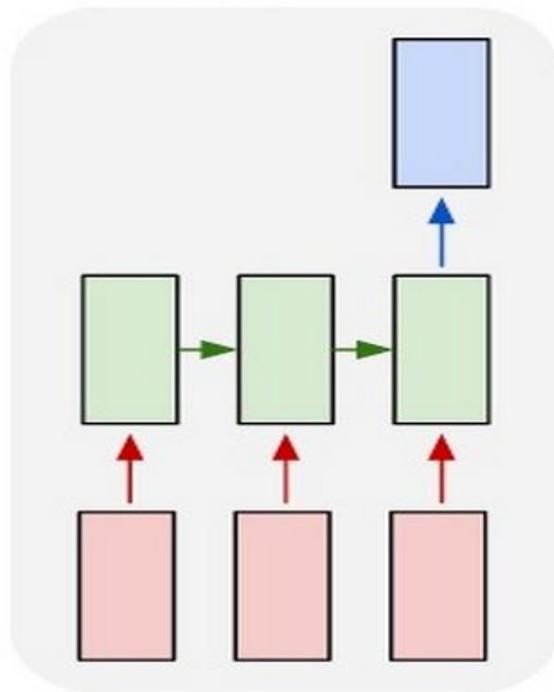
one to many



Sequence output
(Eg:Recognize image
and explain it in words)

Many-to-One Sequence

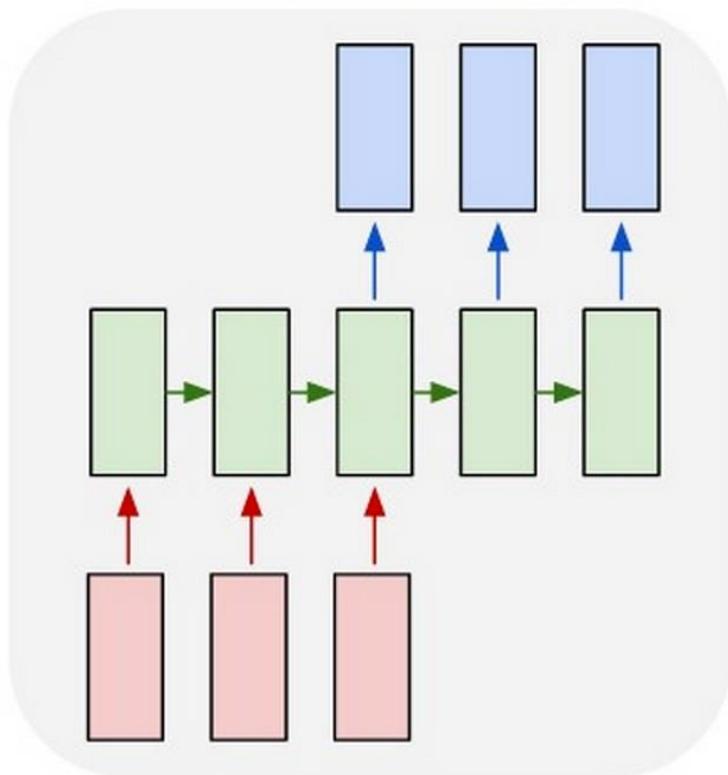
many to one



Sequence input
(Eg:Sentiment analysis
of a sentence)

Many-to-Many Sequence

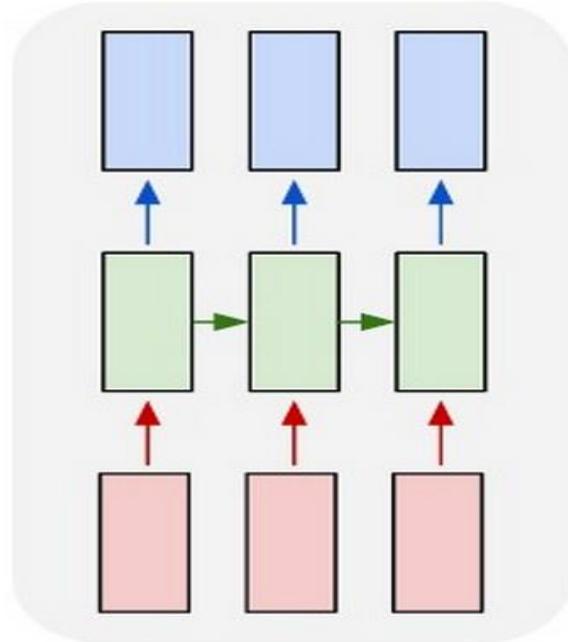
many to many



Sequence input
(Eg:Machine Translation
– English to Spanish)

Many-to-Many Sequence

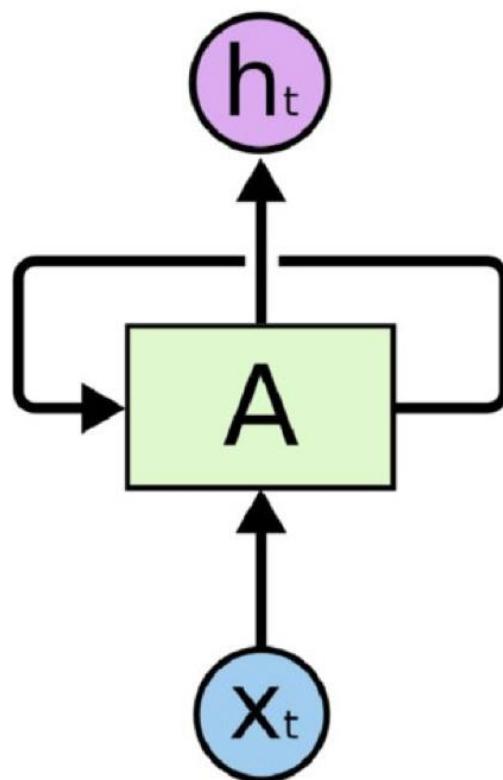
many to many



Synced Sequence
input & Output
(Eg:Video classification
of each frame)

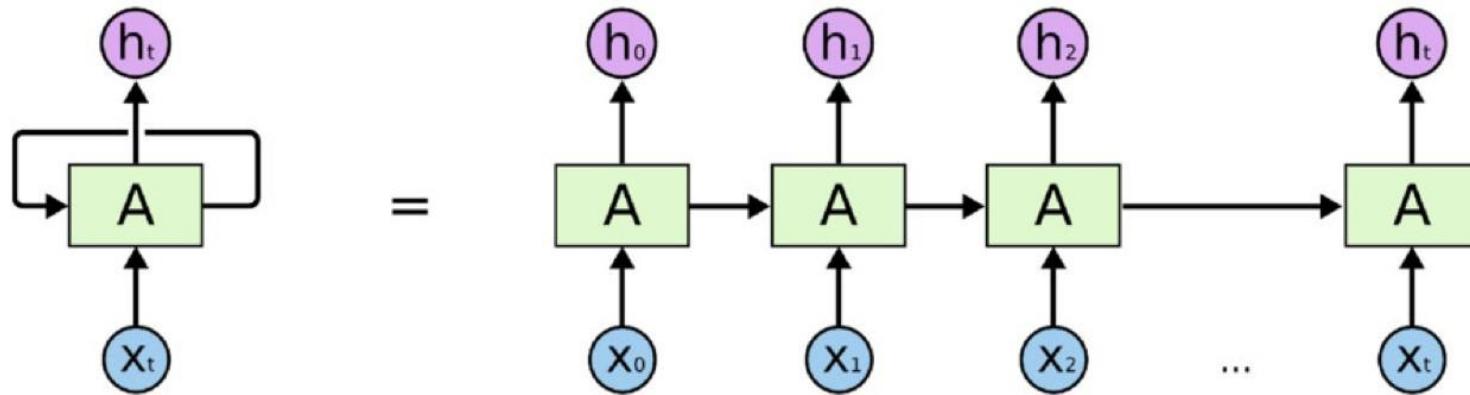
Many-to-Many Sequence

RNN



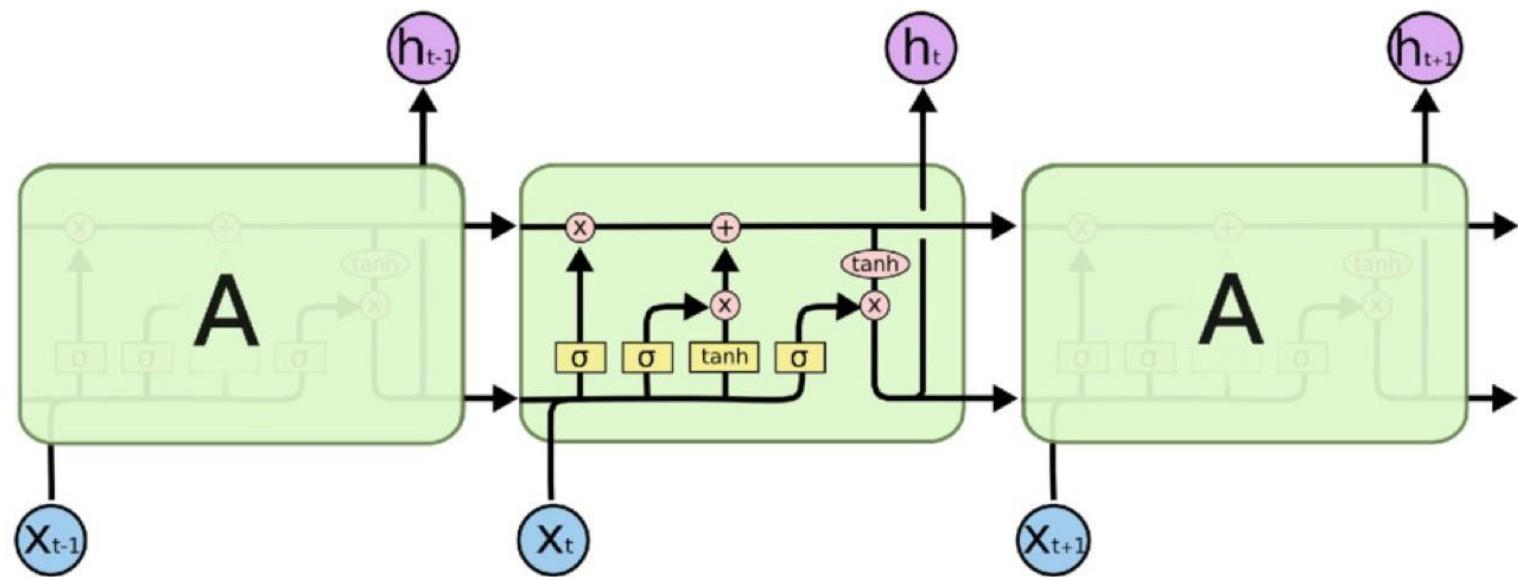
source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Many-to-Many Sequence



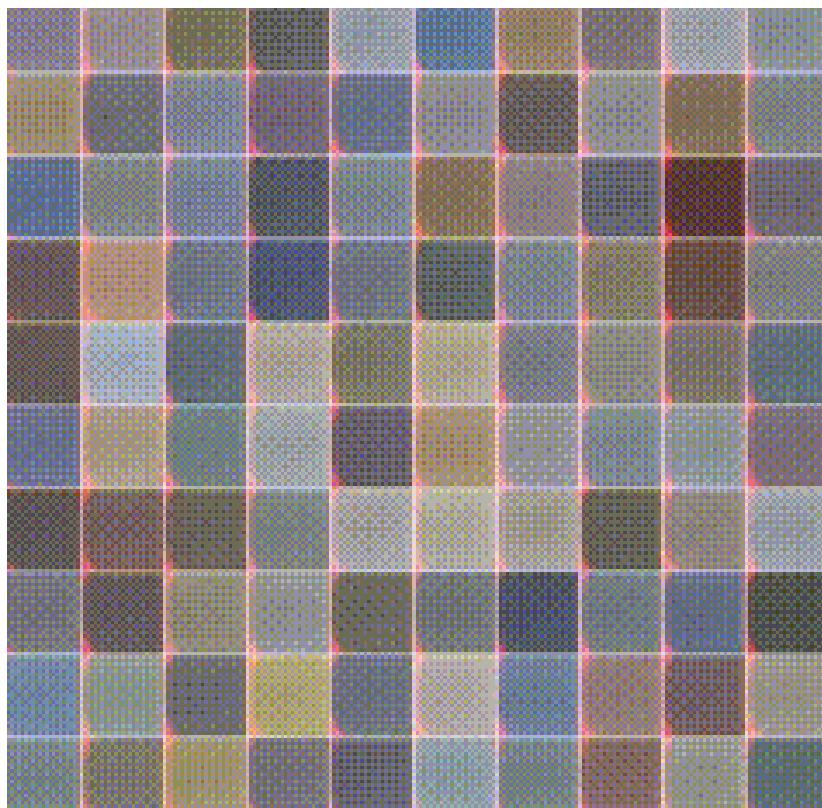
source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Long Short-Term Memory



source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

RNN - Generating Images



Source: <http://arxiv.org/abs/1502.04623c>

Effect of Dataset Size

- RNNs have poor generalization properties on small datasets.
 - 1K labeled examples 25-50% worse than linear model...
- RNNs have better generalization properties on large datasets.
 - 1M labeled examples 0-30% better than linear model.
- Crossovers between 10K and 1M examples
 - Depends on dataset.

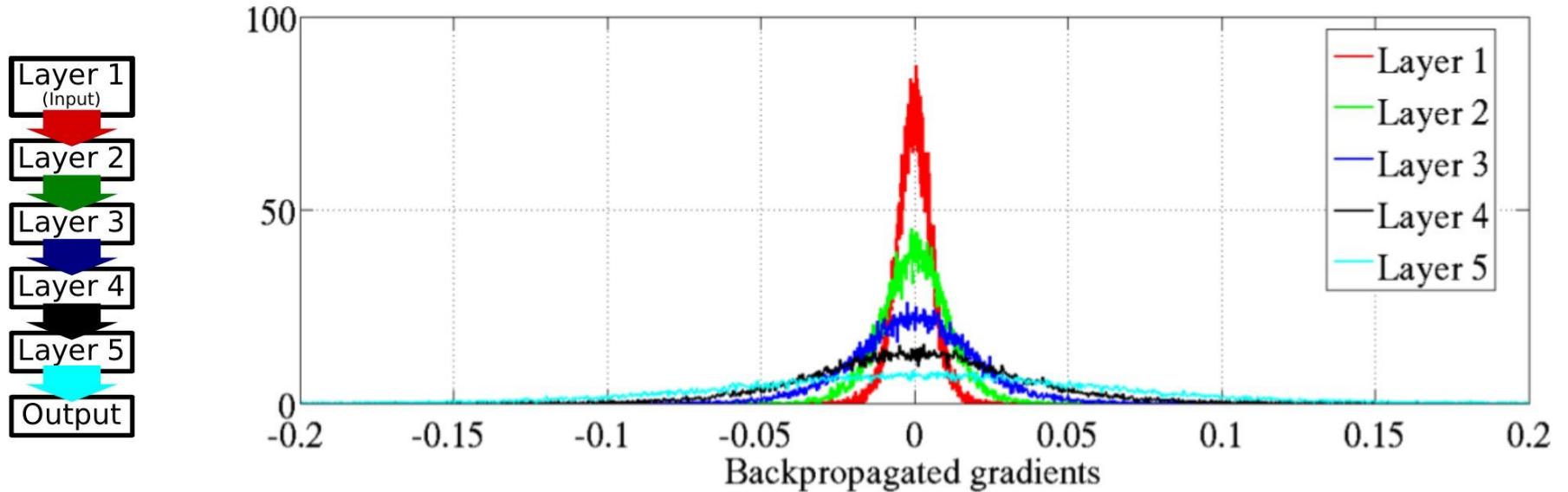
CHALLENGES IN DEEP LEARNING



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



Vanishing Gradient

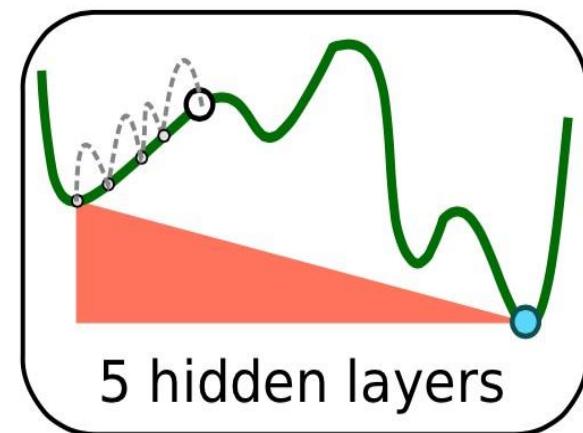
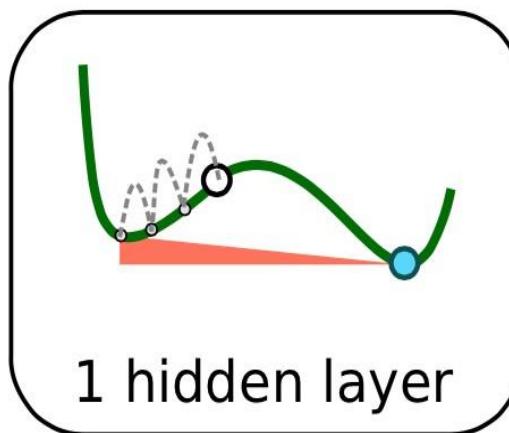
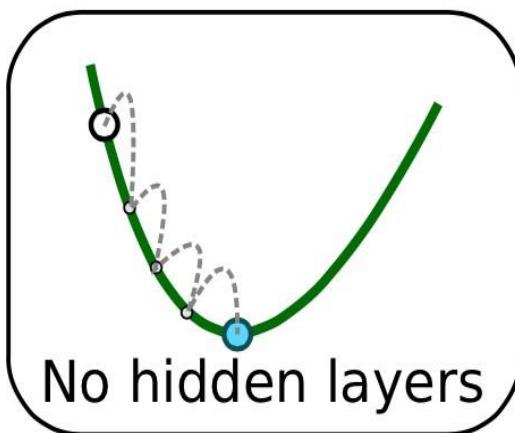


- In backpropagation, the gradient becomes **zero** in the layers nearest to the input.
- **Intuition** : layers near to the output have larger influence in the error rate.
- Training deep layers successfully is terribly slow.
- Until deep layers learn some sensible weights, subsequent layers are trying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

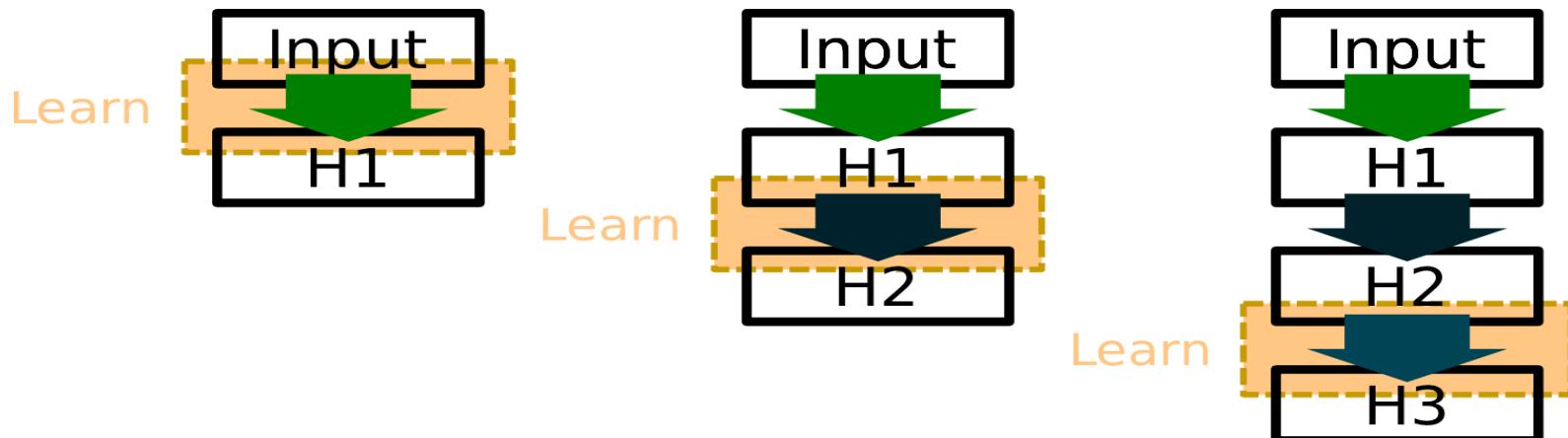
- Smart weight initialization
- Use RLand Maxout units
- Train for a longer time .
- Use restarts .

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 \simeq 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

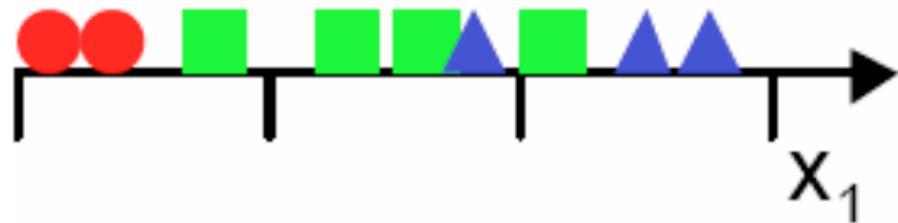


UNSUPERVISED LEARNING

Curse of Dimensionality

Single feature

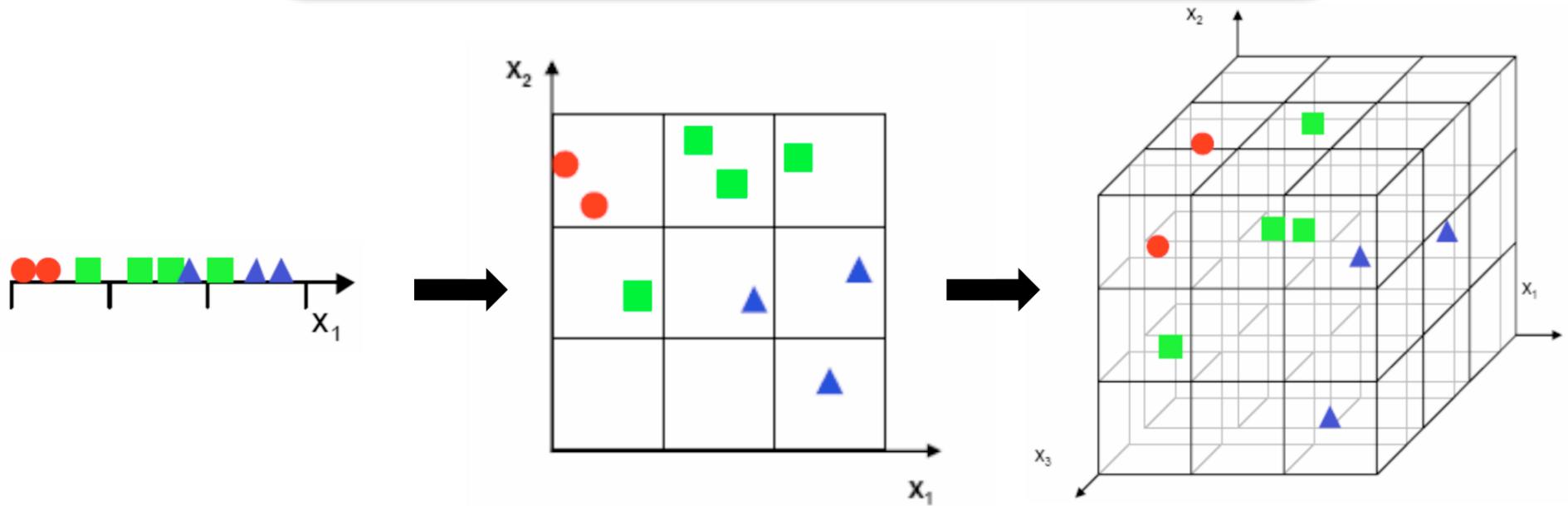
Divide feature space into 3 simple bins



Too simplistic –
Too much overlap between classes.

Curse of Dimensionality

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



Feature Selection

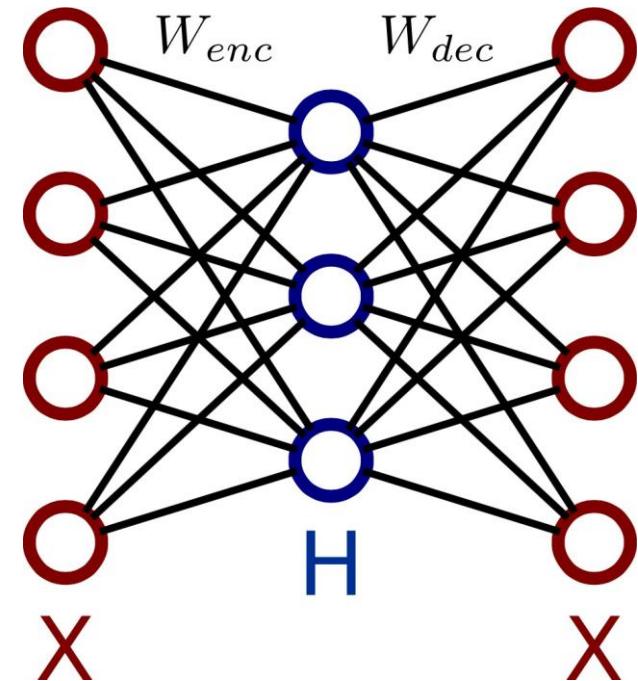
Need: Automated feature selection

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).

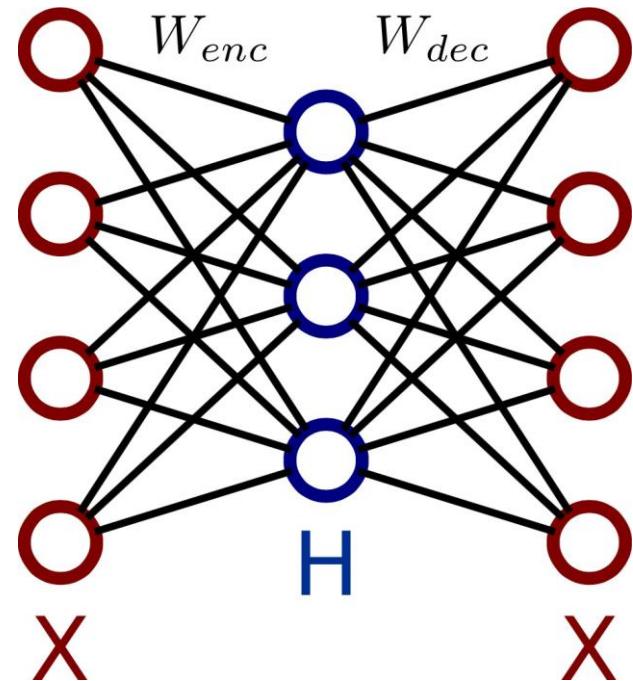


Autoencoders - Output

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*).

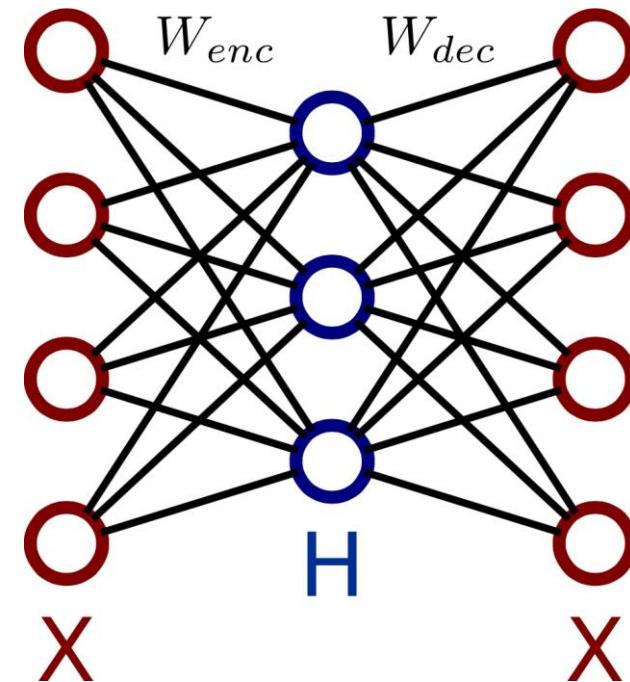


Autoencoders – To make them effective

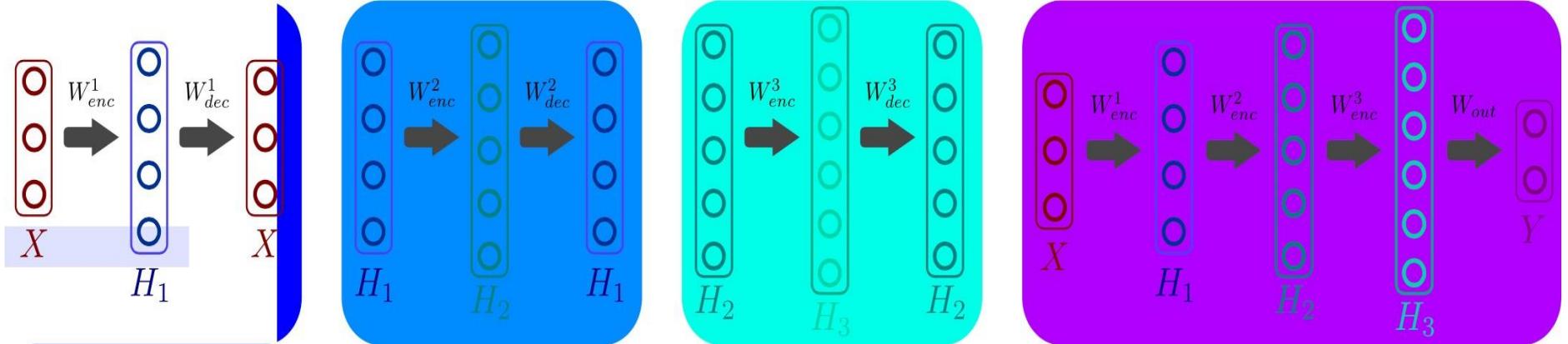
Enforce sparsity .

(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.

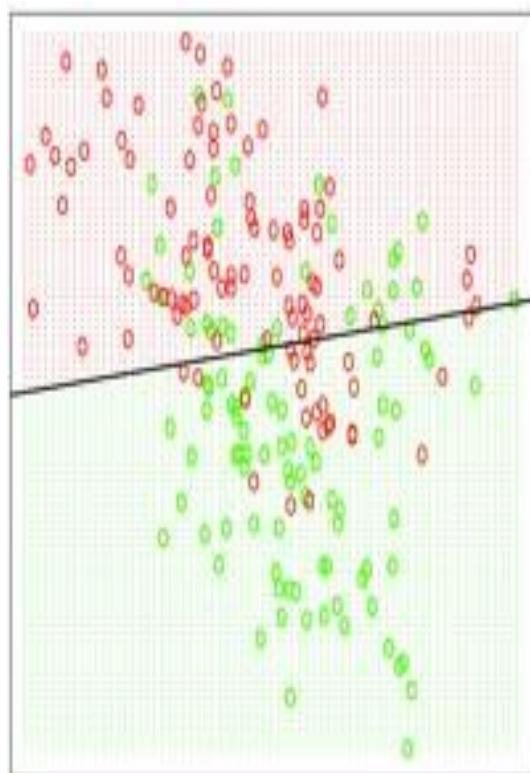
The background of the slide is a photograph of a vast, calm lake or sea. The water is a deep blue at the bottom, transitioning to a lighter, more turquoise hue towards the horizon. The sky above is a clear, pale blue with a few wispy, white clouds near the top. In the far distance, a dark, horizontal line of trees or hills marks the horizon.

CHALLENGES IN DEEP LEARNING

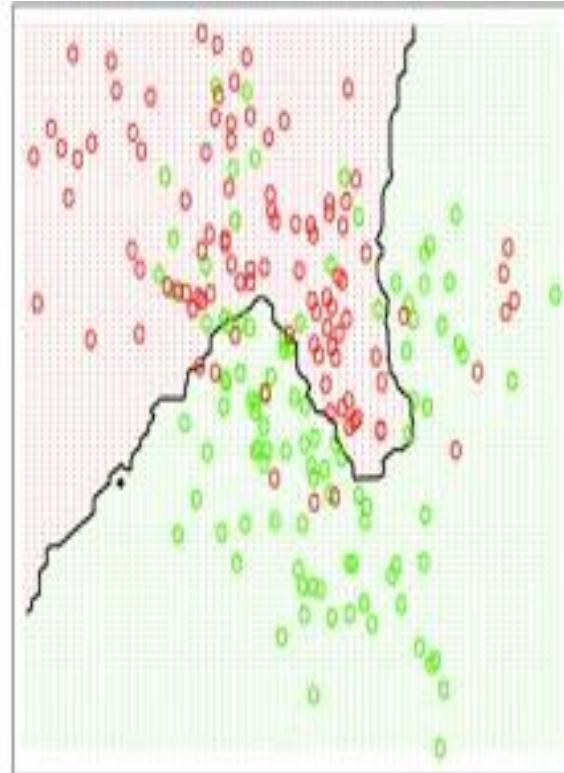
- Continued

Overfitting

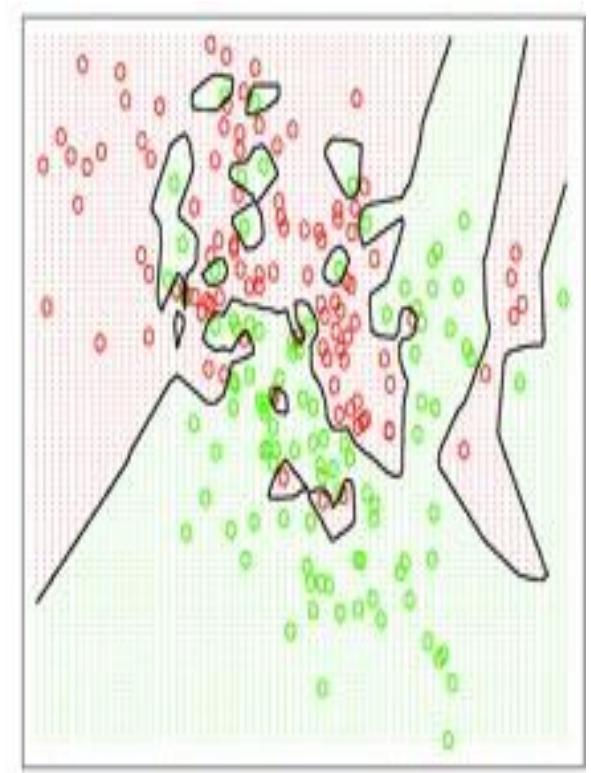
Underfitted
Model



Good
Model



Overfitted
Model

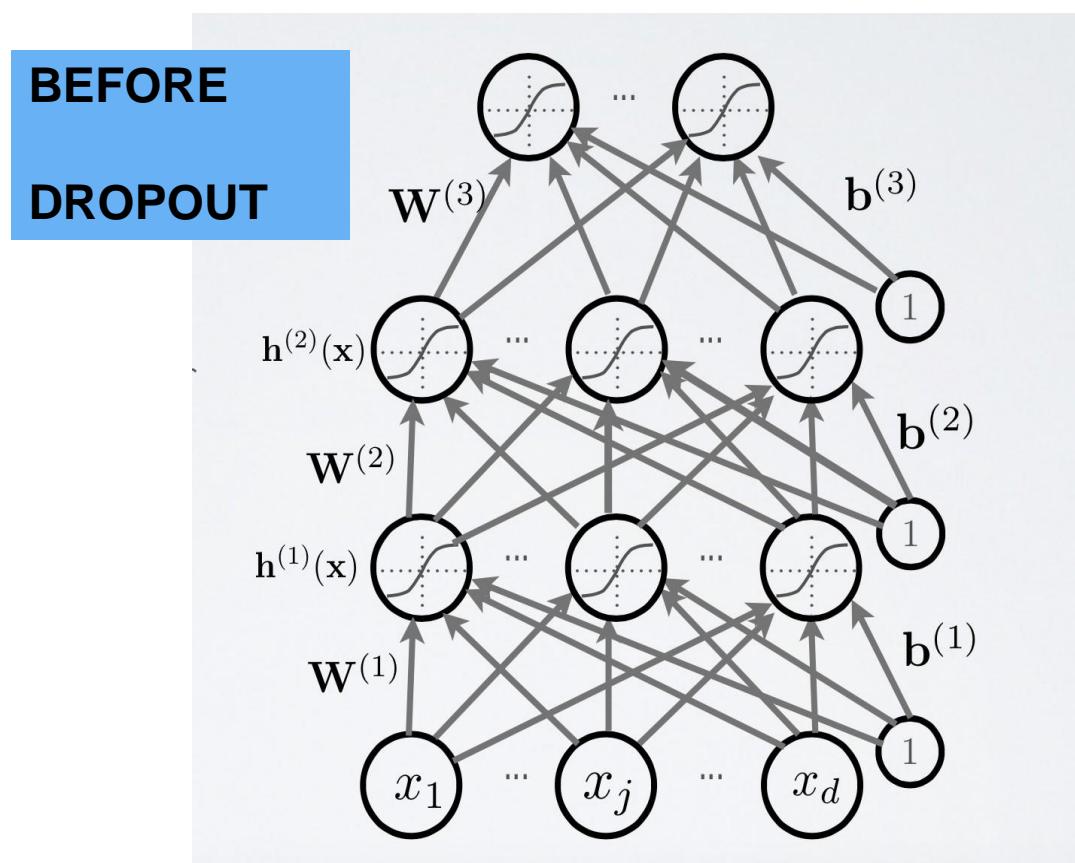


Some ways to address Overfitting

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

Dropout

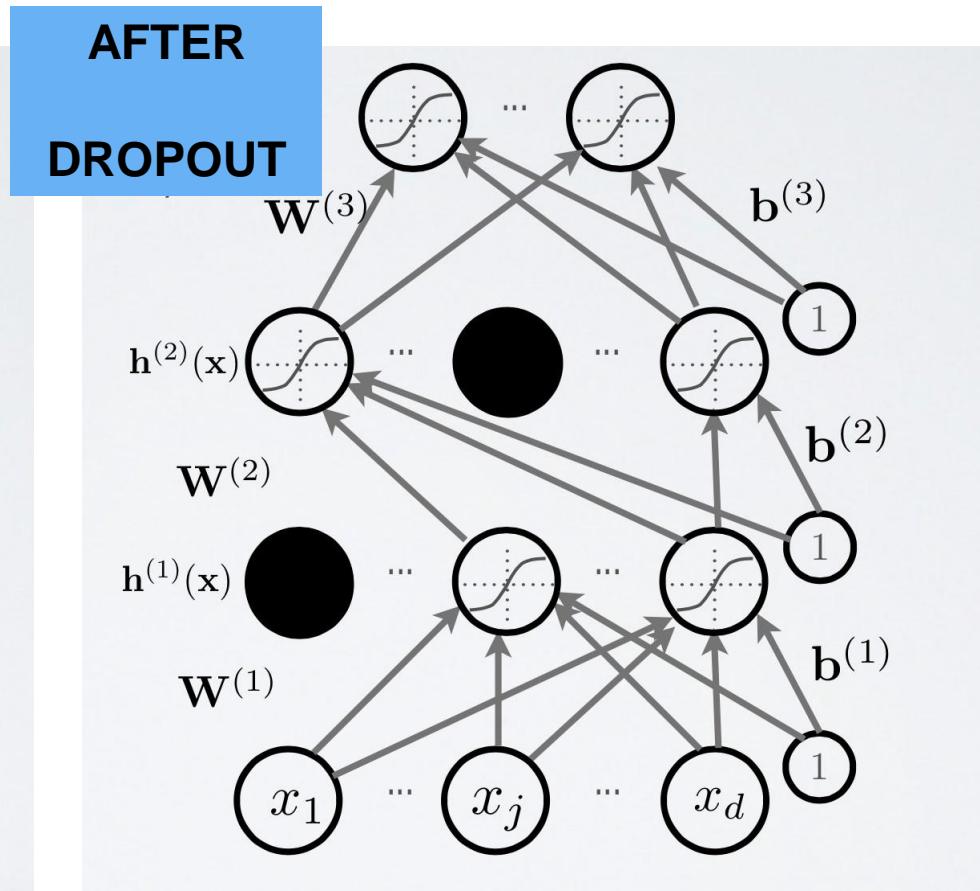
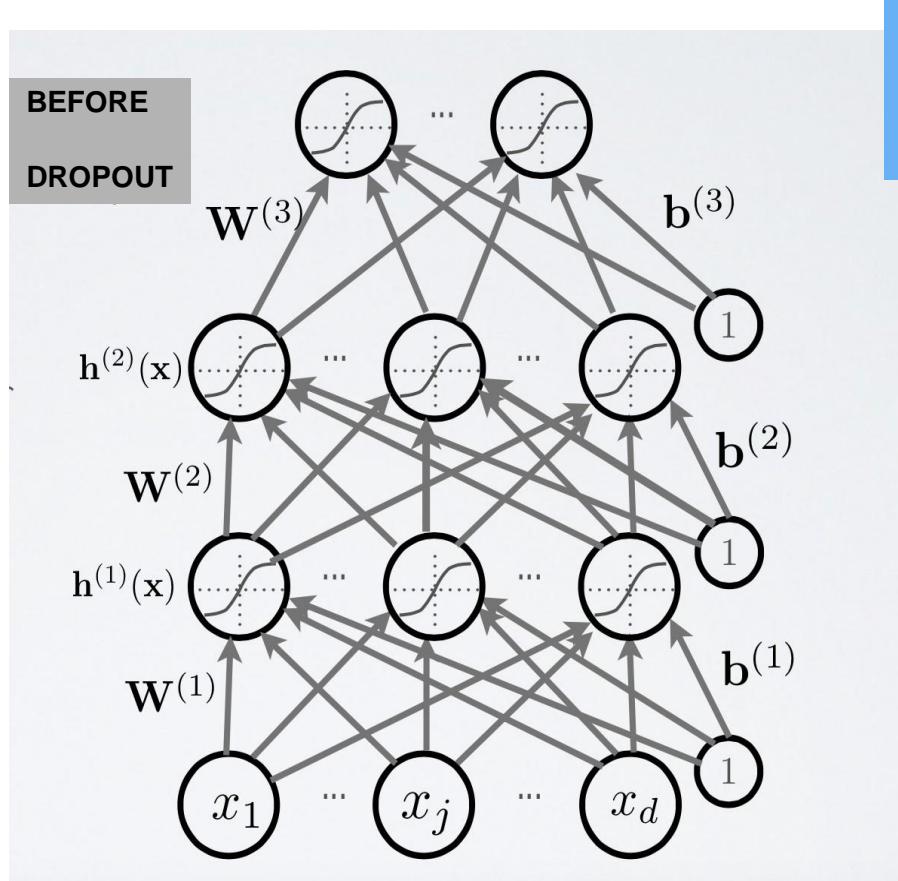
- Cripple the network by removing hidden units stochastically
- In practice, probability of 0.5 works well.



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 (for Images)

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

A silhouette of a person in mid-air, performing a dynamic pose, set against a backdrop of a mountain range at sunset. The sky is filled with dramatic, colorful clouds transitioning from blue to orange and yellow. The mountains are dark silhouettes against the bright horizon.

HANDS-ON

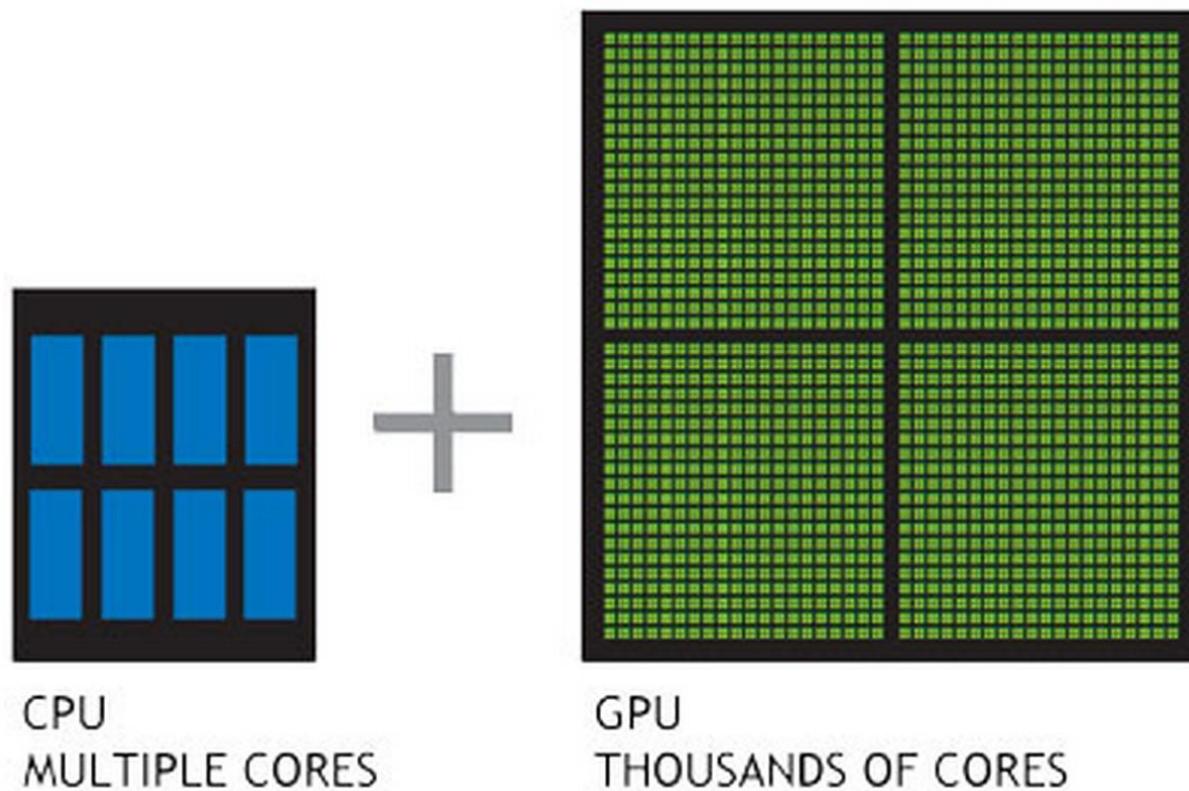


Graphical Processing Unit

Impact of GPU

Compared to CPUs 20x speedups are typical

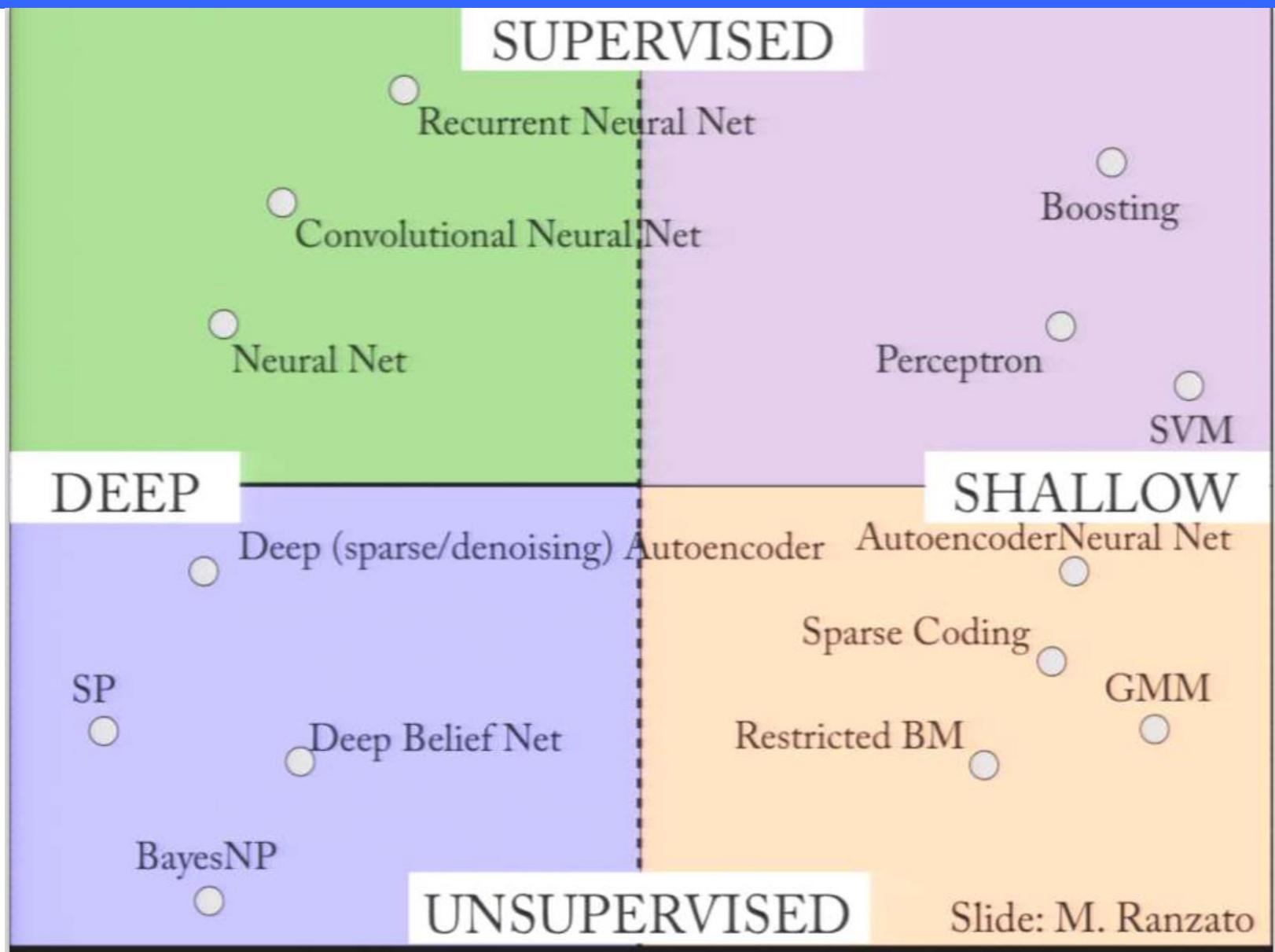
GPUs have thousands of cores to process parallel workloads efficiently



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)
- CNN – 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



Thank you !

