

# INTRODUCTION TO NEURAL NETWORKS

---

E0 270: Machine Learning

*by*

Ambedkar Dukkipati

Department of Computer Science and Automation

Indian Institute of Science, Bangalore

March 1, 2020

- ▶ There might be typos so it is important to refer to the actual lecture
- ▶ This is meant only for E0 270: Machine Learning course and do not distribute (due to the first bullet)
- ▶ Some images are taken from various resources and acknowledgments are assumed.

## INTRODUCTION

---

## Features

- ▶ Go beyond the curve fitting...
- ▶ Amazing results with raw data...
- ▶ Pay and get the tagged data...

## Popular Models

- ▶ Feed Forward Neural Networks and Convolutional Neural Networks
- ▶ Recurrent Neural Networks and Long Short Term Memory Networks
- ▶ Restricted Boltzmann Machines and Deep Boltzmann Machines
- ▶ Autoencoders
- ▶ Generative Adversarial Networks, Variational Autoencoders
- ▶ Memory networks and Neural Turing machines
- ▶ a lot more....

### Tools

- ▶ PyTorch
- ▶ Theano
- ▶ Caffe
- ▶ TensorFlow

### Consequences

- ▶ Brought “AI” back into computer science \*\*\*forey\*\*\*
- ▶ If it works, we accept....we do not mind waiting for “why”

Until recently most machine learning and signal processing techniques had exploited “Shallow-Structured Architectures”

- ▶ Shallow: typically contain at most one or two layers of nonlinear function transformations
  - ▶ Gaussian mixture models
  - ▶ Conditional random fields
  - ▶ Linear or nonlinear dynamical systems
  - ▶ Maximum entropy models
  - ▶ Support vector machines
  - ▶ Logistic regression
  - ▶ Kernel regression
  - ▶ Multilayer perceptron with single hidden layers.
- ▶ Deep: More nonlinear “hidden” layers



- ▶ Shallow architectures have been effective in solving many simple or well constrained problems
- ▶ Shallow architectures have limited modeling and “representation power” can cause difficulties when dealing with more complicated real world applications involving natural signals:
  - ▶ human speech,
  - ▶ natural language,
  - ▶ natural images and scenes
- ▶ These shallow architectures work well given very good “hand crafted features” (may require signal processing techniques)
- ▶ Advantage is that training is easy and may end up mostly with a “convex optimization problem”.

Human information processing mechanism (eg. vision and audio) suggest the need of deep architectures for extracting complex structures and building internal representations from rich sensory inputs.

- ▶ Human speech production and perception systems are equipped with “layered hierarchical structures” in transforming the information from the waveform level to the linguistic level.
- ▶ Human visual systems on the perception side is hierarchical but also “generative.”

- ▶ The concept of deep learning is obtained from neural networks
- ▶ Feedforward neural networks or MLPs with many hidden layers, which are often referred to as deep neural networks (DNNs)
- ▶ Back propagation is popularized in 1980 and has been well known algorithm for learning parameters.

- ▶ Unfortunately BP alone did not work because nonconvex nature of resulting optimization problems
- ▶ The bigger problem: **Vanishing gradient problem**
- ▶ This has steered away most of the ML researchers from natural networks to shallow models that have convex loss functions like
  - ▶ support vector machines (SVM),
  - ▶ conditional random fields (CRF) and
  - ▶ maximum entropy models (MaxEnt)

for which global optimum can be efficiently obtained at the cost of reduced modeling power.

The optimization difficulties associated with the deep models was empirically alleviated when a “reasonably efficient” unsupervised learning algorithms were introduced by Hinton (2006)

- ▶ DBN is composed of a stack of restricted Boltzmann machines (RBM)
- ▶ A greedy, layer by layer algorithm optimizes DBM weights at the at the time complexity linear to the size and depth of the networks.
- ▶ DBMs can be used to initialize the training of Deep Neural Networks (DNN)
- ▶ Advantages of DBMs:
  - ▶ Supply of good initialization for DNNs
  - ▶ Learning algorithm makes effective use of unlabeled data

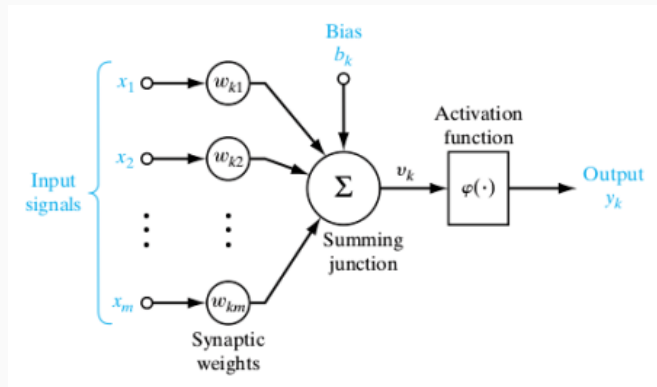
**Most importantly GPUs and Tools like Torch and TensorFlow**

## History:

- ▶ McCulloch and Pitts (1943) introduced idea of neural networks as computing machines
- ▶ Hebb (1949) postulated the first rule for self organizing maps
- ▶ Rosenblatt (1958) invented perceptron that algorithmically described neural networks



# Perceptron: History



- ▶ Nonlinearity
- ▶ Input-Output mapping
- ▶ Adaptivity
- ▶ Fault Tolerance
- ▶ VLSI implementability

- ▶ Perceptron
- ▶ Feed forward deep networks and Back propagation algorithm
- ▶ later CNNs, LSTMs, RBMs (if time permits)

## HYPERPLANE BASED CLASSIFIERS AND PERCEPTRON

---

## Supervised Learning Problem

- ▶ Given data  $\{(x_n, y_n)\}_{n=1}^N$  find  $f : \mathcal{X} \rightarrow \mathcal{Y}$  that best approximates the relation between  $X$  and  $Y$ .
- ▶ Determine  $f$  such a way that loss  $l(y, f(x))$  is minimum.
- ▶  $f$  and  $l$  are specific to the problem and the method that we choose.

- ▶ Data:  $\{(x_n, y_n)\}_{n=1}^N$ 
  - ▶  $x_n \in \mathbb{R}^D$  is a  $D$  dimensional input
  - ▶  $y_n \in \mathbb{R}$  is the output

Aim is to find a **hyperplane** that fits **best** these points.

- ▶ Here hyperplane is a model of choice i.e.,

$$f(x) = \sum_{j=1}^D x_j w_j + b = w^T x + b$$

- ▶ Here  $w_1, \dots, w_D$  and  $b$  are model parameters
- ▶ Best is determined by some loss function

$$Loss(w) = \sum_{n=1}^N [y_n - f(x_n)]^2$$

- ▶ **Aim** : Determine the model parameters that minimize the loss.

## Problem Set-Up

- ▶ Two class classification
- ▶ Instead of the exact labels estimate the probabilities of the labels i.e.

$$\text{Predict} \quad P(y_n = 1|x_n, w) = \mu_n$$

$$P(y_n = 0|x_n, w) = 1 - \mu_n$$

- ▶ Here  $(x_n, y_n)$  is the input output pair.

### Problem

Find a function  $f$  such that,

$$\mu = f(x_n)$$

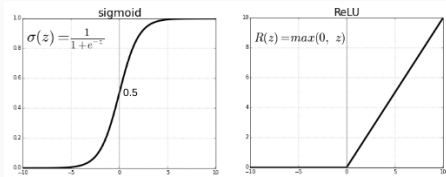
### Model

$$\begin{aligned}\mu_n = f(x_n) = \sigma(w^T x_n) &= \frac{1}{1 + \exp(-w^T x_n)} \\ &= \frac{\exp(w^T x_n)}{1 + \exp(w^T x_n)}\end{aligned}$$



## Sigmoid Function

- Here  $\sigma(\cdot)$  is the sigmoid function.



- The model first computes a real valued score  $w^T x = \sum_{i=1}^D w_i x_i$  and then nonlinearly “squashes” it between (0,1) to turn into a probability.

## Loss Function

Here we use cross entropy loss instead of squared loss.



$$\begin{aligned} L(y_n, f(x_n)) &= \begin{cases} -\log(\mu_n) & \text{when } y_n = 1 \\ -\log(1 - \mu_n) & \text{when } y_n = 0 \end{cases} \\ &= -y_n \log(\mu_n) - (1 - y_n) \log(1 - \mu_n) \\ &= -y_n \log(\sigma(w^\top x_n)) - (1 - y_n) \log(1 - \sigma(w^\top x_n)) \end{aligned}$$

►  $L(w) = -\sum_{n=1}^N [y_n w^\top x_n - \log(1 + \exp(w^\top x_n))]$

Here  $L(w)$  is the loss over all the data.

By taking the derivative w.r.t  $w$

$$\frac{\partial L}{\partial w} = \sum_{n=1}^N (\mu_n - y_n) x_n$$

► Here the Gradient Descent Algorithm is

- 1 Initialize  $w^{(1)} \in \mathbb{R}^D$  randomly
- 2 Iterate until the convergence

$$\underbrace{w^{(t+1)}}_{\text{New learned parameter or weights}} = \underbrace{w^{(t)}}_{\text{previous value}} - \underbrace{\eta}_{\text{Learning rate}} \sum_{n=1}^N \underbrace{(\mu_n^{(t)} - y_n) x_n}_{\text{Gradient at previous value}}$$

- Note: Here  $\mu^{(t)} = \sigma(w^{(t)\top} x_n)$
- **Problem:** Calculating gradient in each iteration requires all the data. When  $N$  is large this may not be feasible.

- Note: Gradient decent requires all the data to calculate the gradients.
- Strategy: Approximate gradient using randomly chosen data point  $(x_n, y_n)$

$$w^{(t+1)} = w^{(t)} - \eta_t(\mu_n^{(t)} - y_n)x_n$$

- .
- Also: Replace predicted label probability  $\mu_n^{(t)}$  by predicted binary label  $\hat{y}_n^{(t)}$ , where

$$\hat{y}_n^{(t)} = \begin{cases} 1 & \text{if } \mu_n^{(t)} \geq 0.5 \text{ or } w^{(t)\top} x_n \geq 0 \\ 0 & \text{if } \mu_n^{(t)} < 0.5 \text{ or } w^{(t)\top} x_n < 0 \end{cases}$$

- Hence: Update rule becomes

$$w^{(t+1)} = w^{(t)} - \eta_t(\hat{y}_n^{(t)} - y_n)x_n$$

- This is mistake driven update rule
- $w^{(t)}$  gets updated only when there is a misclassification i.e.  $\hat{y}_n^{(t)} \neq y_n$
- Change: the class labels to  $\{-1, +1\}$

$$\Rightarrow \hat{y}_n^{(t)} - y_n = \begin{cases} -2y_n & \text{if } \hat{y}_n^{(t)} \neq y_n^{(t)} \\ 0 & \text{if } \hat{y}_n^{(t)} = y_n^{(t)} \end{cases}$$

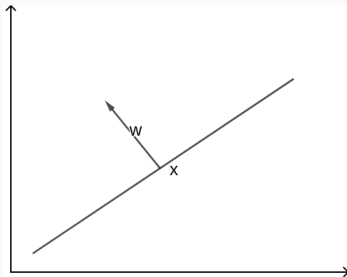
- Hence: Whenever there is a misclassification.

$$w^{(t+1)} = w^{(t)} - 2\eta_{(t)}y_nx_n$$

$\Rightarrow$  This is a perceptron learning algorithm which is a hyperplane based learning algorithm.

- Separates a d-dimensional space into two half spaces(positive and negative)
- Equation of the hyperplane is

$$w^T x = 0$$



- By adding bias  $b \in \mathbb{R}$

$$w^T x + b = 0$$

$b > 0$  moving the hyperplane parallelly along  $w$

$b < 0$  opposite direction

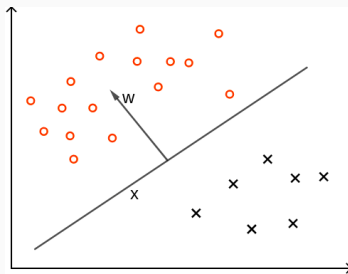
- Classification rule

$$y = \text{sign}(w^T x + b)$$

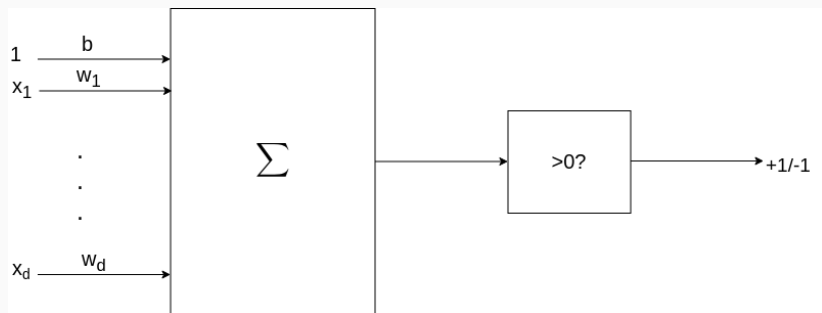


$$w^T x + b > 0 \implies y = +1$$

$$w^T x + b < 0 \implies y = -1$$



## Hyperplane based classification





- ▶ Aim is to learn a linear hyperplane to separate two classes.
- ▶ Mistake drives online learning algorithm
- ▶ Guaranteed to find a separating hyperplane if data is linearly separable.
- ▶ If data is not linearly separable
  - ▶ Make linearly separable using kernel methods.
  - ▶ (Or) Use multilayer perceptron.

- ▶ Given training data  $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- ▶ Initialize  $w_{old} = [0, \dots, 0]$ ,  $b_{old} = 0$
- ▶ Repeat until convergence.
  - ▶ For a random  $(x_n, y_n) \in \mathcal{D}$
  - ▶ If  $y_n(w^T x_n + b) \leq 0$   
[Or  $\text{sign}(w^T x + b) \neq y_n$  i.e mistake mode]  
 $w_{new} = w_{old} + y_n x_n$   
 $b_{new} = b_{old} + y_n$

"Roughly" : If the data is linearly separable perceptron algorithm converges.

## FEED FORWARD NEURAL NETWORKS

---

- ▶ Network will have hidden layers.
- ▶ Since perceptron works only for linearly separable data, each neuron has a non-linear activation function: Activation function is differentiable.
- ▶ Network exhibit a high degree of connectivity.

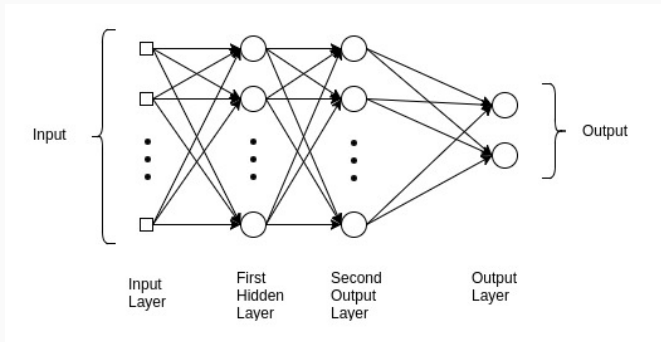
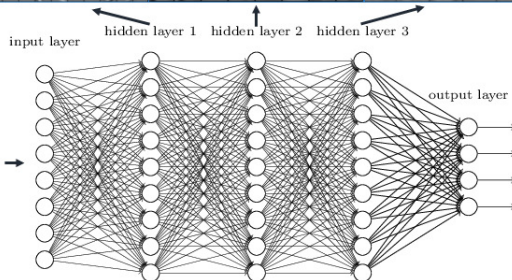
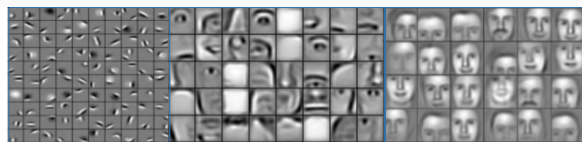


Figure: Fully connected multilayer perceptron with two hidden layers

# Why hidden Layers?

Deep neural networks learn hierarchical feature representations



- ▶ Hidden layers can automatically learn features from data
- ▶ The bottom-most hidden layer captures very low level features (e.g., edges). Subsequent hidden layers learn progressively more high-level features (e.g., parts of objects) that are composed of previous layers features

# Two important steps in training neural network

## 1 Forward step:

- ▶ Input is fed to the first layer.
- ▶ Input signal is propagated through the network layer by layer.
- ▶ Synaptic weights of the network are fixed i.e. no learning happens in this step.
- ▶ Error is calculated at the output layer by comparing the observed output with "desired output" (Ground truth)

## 2 Backward step:

- ▶ The observed error at the output layer is propagated "backwards", layer by layer. (How?)
- ▶ Error is propagated "backwards", layer by layer.
- ▶ In this step, successive adjustments are made to the synaptic weights.



# Propagation of information in neural network

Two kinds of signals:

- 1 Function signal (leads to observed error)
- 2 Error signal (leads to updation of weights or parameters)

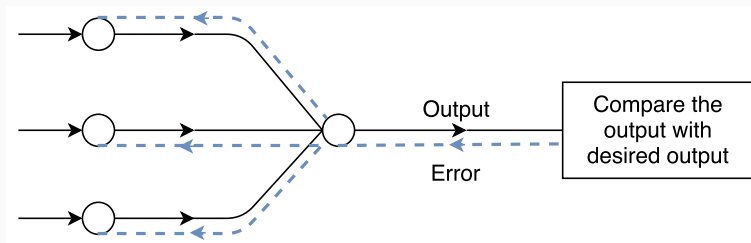
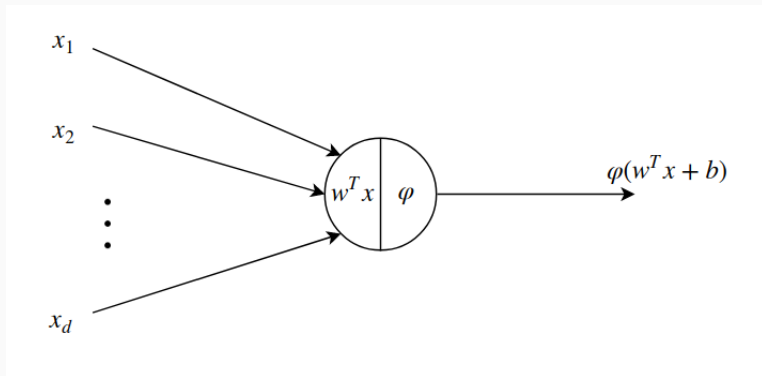


Figure: Propagation of information in neural network

## 1 Computation of function signal (in the forward step)



## 2 Computation of gradient

- gradients of the "error surface" w.r.t weights (we will see this later how)

- ▶  $\mathcal{D} = \{(x(n), z(n))\}_{n=1}^N$  be a training sample where  $x(n)$  is an input and  $z(n)$  is the desired output.
  - ▶  $x(n) \in \mathbb{R}^D$ , we write  $x(n) = (x_1(n), \dots, x_D(n))$ .
  - ▶  $z(n) \in \mathbb{R}^M$ , we write  $z(n) = (z_1(n), \dots, z_M(n))$ .
- ▶ Suppose the output of the network is  $y(n) = (y_1(n), \dots, y_M(n))$  when  $x(n)$  is the input.
- ▶ Error at the  $j^{th}$  output neuron is

$$e_j(n) = |z_j(n) - y_j(n)|, \quad j = 1, 2, \dots, M$$

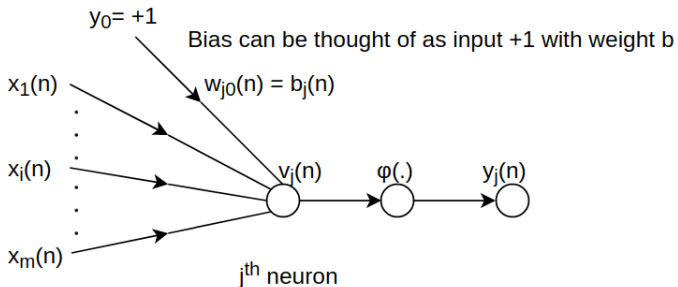
- ▶ The total error per sample is

$$\mathcal{E}(n) = \frac{1}{2} \sum_{j=1}^M (z_j(n) - y_j(n))^2$$

- ▶ Average error for the training data or empirical risk

$$\bar{\mathcal{E}} = \frac{1}{N} \sum_{n=1}^N \mathcal{E}(n) = \frac{1}{2N} \sum_{n=1}^N \sum_{j=1}^M (z_j(n) - y_j(n))^2$$

# The Backpropagation Algorithm



- ▶  $x_1(n), \dots, x_i(n), \dots, x_m(n)$ : Function signals that are produced by the previous layer which is an input to the  $j^{th}$  neuron.
- ▶  $v_j(n) = \sum_{i=0}^m w_{ji}(n)x_i(n)$ 
  - ▶  $v_j(n)$  is the induced local field.
  - ▶  $m$  is the size of the input (i.e. in the previous layer there are  $m$  neurons)
- ▶  $y_j(n) = \varphi(v_j(n))$ 
  - ▶ Function signal appearing at the output of neuron  $j$ .

- ▶ BPA applies a correction  $\Delta w_{ji}(n)$  to the synaptic weight proportional to  $\frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)}$ ,  $i = 1, 2, \dots, m$ .
  - ▶ Note that we are trying to update  $j^{th}$  neuron out of all  $M$  neurons.
  - ▶ For  $n^{th}$  data point, the error is

$$\mathcal{E}(n) = \frac{1}{2} \sum_{j=1}^M (z_j(n) - y_j(n))^2 = \frac{1}{2} \sum_{j=1}^M e_j^2$$

We compute the derivative of

$$\mathcal{E}(n) = \frac{1}{2} \sum_{j=1}^M (z_j(n) - y_j(n))^2$$

w.r.t  $w_{ji}(x)$  (apply chain rule)



The derivative

$$\frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)} = \frac{\partial \mathcal{E}(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \frac{\partial v_j(n)}{\partial w_{ji}(n)}$$

Since

- ▶  $\mathcal{E}$  is a function of  $e_j$
- ▶  $e_j$  is a function of  $y_j$  ( $y_j$  is the output)
- ▶  $y_j$  is a function of  $v_j$  ( $v_j$  is the local field)
- ▶  $v_j$  is a function of  $w_{ji}$

- The derivative is

$$\frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)} = \frac{\partial \mathcal{E}(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \frac{\partial v_j(n)}{\partial w_{ji}(n)}$$

- $\mathcal{E}(n) = \frac{1}{2} \sum_{j=1}^M e_j^2(n) \implies \frac{\partial \mathcal{E}(n)}{\partial e_j(n)} = e_j(n)$
- $e_j(n) = z_j(n) - y_j(n) \implies \frac{\partial e_j(n)}{\partial y_j(n)} = -1$
- $y_j(n) = \varphi(v_j(n)) \implies \frac{\partial y_j(n)}{\partial v_j(n)} = \varphi'_j(v_j(n))$
- $v_j(n) = \sum_{i=0}^m w_{ji}(n)x_i(n) \implies \frac{\partial v_j(n)}{\partial w_{ji}(n)} = x_i(n)$
- $\implies$

$$\frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)} = -e_j(n)\varphi'_j(v_j(n))x_i(n)$$

## Update rule for $j$ th output neuron

- We have

$$\frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)} = -e_j(n) \varphi_j'(v_j(n)) x_i(n)$$

- Hence, the update rule is

$$\begin{aligned} w_{ji}(n+1) &= w_{ji}(n) - \eta \frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)} \\ &= w_{ji}(n) + \eta e_j(n) \varphi_j'(v_j(n)) x_i(n) \end{aligned}$$

## Local Gradient

- Define local gradient  $\delta_j(n)$  for  $j^{\text{th}}$  neuron as

$$\begin{aligned}\delta_j(n) &= -\frac{\partial \mathcal{E}(n)}{\partial v_j(n)} \\ &= -\frac{\partial \mathcal{E}(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} = e_j(n) \varphi'_j(v_j(n))\end{aligned}$$

- $w_{ji}(n+1) = w_{ji}(n) + \underbrace{\eta \delta_j(n)}_{\text{Local Gradient}} x_i(n)$

- Output neuron has an “easy” access to the error

$$e_j(n) = d_j(n) - y_j(n)$$

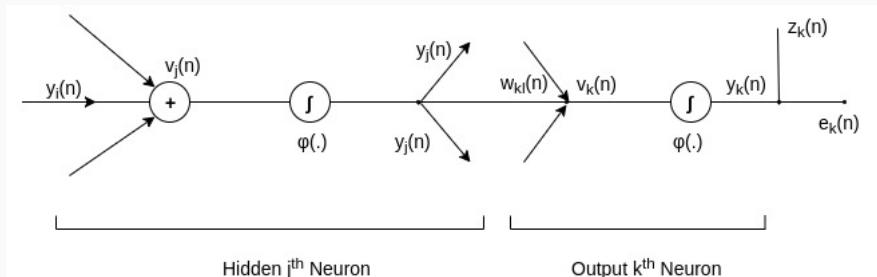
$$\mathcal{E}(n) = \frac{1}{2} \sum_{j=1}^M e_j^2(n)$$

- Update rule

$$\begin{aligned} w_{ji}(n+1) &= w_{ji}(n) - \eta \frac{\partial \mathcal{E}(n)}{\partial w_{ji}(n)} \\ &= w_{ji}(n) - \eta \underbrace{e_j(n) \varphi'(v_j(n))}_{\text{Local Gradient}} x_i(n) \\ &= w_{ji}(n) - \eta \delta_j(n) x_i(n) \end{aligned}$$

- ▶ Unlike in the case of output neuron, hidden neuron does not have a direct access to the "error".
- ▶ **TRICK**
  - ▶ Error signal for a hidden neuron will be determined recursively.
  - ▶ They expect the next hidden neuron (or output neuron), they are connected to, to share "some" error.
  - ▶ Error propagates by working backwards.

## BPA: Case 2: Neuron j is a hidden node (contd...)



### Strategy

- First compute the local gradient  $\delta_j(n)$  for  $j^{\text{th}}$  hidden neuron.
  - How we will see ...
- Then use the update that is similar to output neuron

$$\begin{aligned}\Delta w &= \text{Learning rate} \times \text{Local gradient} \times \text{Input} \\ &= \eta \delta_j(n) x_i(n)\end{aligned}$$

Local field of  $j^{th}$  hidden neuron:

$$\begin{aligned}\delta_j(n) &= -\frac{\partial \mathcal{E}(n)}{\partial v_j(n)} = -\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \\ &= -\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} \varphi'_j(v_j(n)) \quad \because y_j(n) = \varphi_j(v_j(n))\end{aligned}$$

**Note:** If this had been the output neuron, we would have had

$$\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} = \frac{\partial \mathcal{E}(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} = -e_j(n)$$

Since j is hidden, it does not have access to the error.



- We are trying to compute local gradient

$$\delta_j(n) = -\frac{\partial \mathcal{E}(n)}{\partial v_j(n)} = -\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} \varphi'_j(v_j(n))$$

- Let us compute  $\frac{\partial \mathcal{E}(n)}{\partial y_j(n)}$

- We have  $\mathcal{E}(n) = \frac{1}{2} \underbrace{\sum_{k \in C} e_k^2(n)}$

Summation over all the output neurons

- Then

$$\begin{aligned} \frac{\partial \mathcal{E}(n)}{\partial y_j(n)} &= \sum_{k \in C} e_k(n) \frac{\partial e_k(n)}{\partial y_j(n)} \\ &= \sum_{k \in C} e_k(n) \frac{\partial e_k(n)}{\partial v_k(n)} \frac{\partial v_k(n)}{\partial y_j(n)} \end{aligned}$$

We are computing  $\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} = \sum_{k \in C} e_k(n) \frac{\partial e_k(n)}{\partial v_k(n)} \frac{\partial v_k(n)}{\partial y_j(n)}$

► We have

$$\begin{aligned} e_k(n) &= z_k(n) - y_k(n) \\ &= z_k(n) - \varphi_k(v_k(n)) \end{aligned}$$

►  $\frac{\partial e_k(n)}{\partial v_k(n)} = -\varphi'(v_k(n))$

► We have  $v_k(n) = \sum_{l=1}^m w_{kl}(n) y_l(n)$

Note that  $j \in \{1, 2, \dots, m\}$  and  $j^{th}$  neuron output along with the other neurons in that layer are fed to the  $k^{th}$  output neuron.

$$\implies \frac{\partial v_k(n)}{\partial y_j(n)} = w_{kj}(n)$$

We are computing  $\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} = \sum_{k \in C} e_k(n) \frac{\partial e_k(n)}{\partial v_k(n)} \frac{\partial v_k(n)}{\partial y_j(n)}$

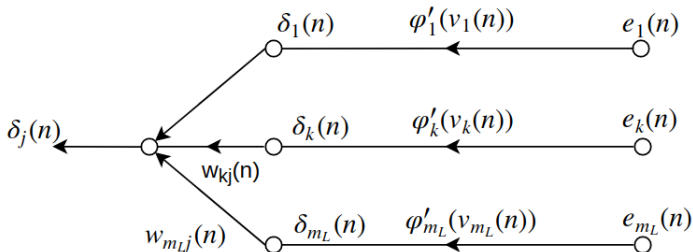
$$\begin{aligned} \frac{\partial \mathcal{E}(n)}{\partial y_j(n)} &= - \sum_{k \in C} e_k(n) \varphi_k(v_k(n)) w_{kj}(n) \\ &= - \sum_{k \in C} \delta_k(n) w_{kj}(n) \end{aligned}$$

where  $\delta_k(n) = e_k(n) \varphi_k(v_k(n))$  is the local gradient of the  $k^{th}$  neuron.

We have  $\frac{\partial \mathcal{E}(n)}{\partial y_j(n)} = - \sum_k \delta_k(n) w_{kj}(n)$

and  $\frac{\partial y_j(n)}{\partial v_j(n)} = \varphi'_j(v_j(n))$

Hence,  $\delta_j(n) = \varphi'_j(v_j(n)) \sum_k \delta_k(n) w_{kj}(n)$



- Now we have local gradient at the  $j^{th}$  hidden node  
i.e.  $\delta_j(n) = \varphi'_j(v_j(n)) \sum_k \delta_k(n) w_{kj}(n)$   
where  $\delta_k(n) = e_k(n) \varphi'_k(v_k(n))$  is the local field at  $k^{th}$  output layer.
- Hence,

$$\begin{aligned} w_{ji}(n+1) &= w_{ji}(n) + \eta \delta_j(n) x_i(n) \\ &= w_{ji}(n) + \eta \left[ \varphi'_j(v_j(n)) \sum_{k \in C} \delta_k(n) w_{kj}(n) \right] x_i(n) \end{aligned}$$

**Case 1:**  $j^{th}$  neuron is an output neuron

$$\Delta w_{ji}(n) = \eta \underbrace{e_j(n) \varphi'_j(v_j(n))}_{\text{Local gradient at the } j^{th} \text{ neuron}} x_i(n)$$

**Case 2:**  $j^{th}$  neuron is a hidden neuron

$$\Delta w_{ji}(n) = \eta \underbrace{\left( \sum_{k \in C} \delta_k(n) w_{kj}(n) \right) \varphi'_j(v_j(n))}_{\text{Local gradient at the } j^{th} \text{ neuron}} x_i(n)$$

**Case 1:**  $j^{th}$  neuron is an output neuron

$$\Delta w_{ji}(n) = \eta \underbrace{e_j(n) \varphi'_j(v_j(n))}_{\text{Local gradient at the } j^{th} \text{ neuron}} x_i(n)$$

**Case 2:**  $j^{th}$  neuron is a hidden neuron

$$\Delta w_{ji}(n) = \eta \underbrace{\left( \sum_{k \in C} \delta_k(n) w_{kj}(n) \right) \varphi'_j(v_j(n))}_{\text{Local gradient at the } j^{th} \text{ neuron}} x_i(n)$$

## Batch Learning

- ▶ Each adjustment to the weights is performed after the presentation of all the  $N$  examples in the training samples are presented.
- ▶ That is, cost function is average error or empirical risk.

$$\begin{aligned}\bar{\mathcal{E}} &= \frac{1}{N} \sum_{n=1}^N \mathcal{E}(n) = \frac{1}{2N} \sum_{n=1}^N \sum_{j=1}^M e_j^2(n) \\ &= \frac{1}{2N} \sum_{n=1}^N \sum_{j=1}^M (z_j(n) - y_j(n))^2\end{aligned}$$



### Batch Learning

- ▶ This constitutes one epoch of training.
- ▶ In each epoch of training, samples are randomly shuffled.
- ▶ The learning curve in this case is  $\bar{\mathcal{E}}$  vs epoch number.
- ▶ Advantage: It can be easily parallelized.
- ▶ Disadvantage: Memory requirements are very high.

## Online Learning

- ▶ Each adjustment to the weights is performed example by example in the training data.
- ▶ The cost function is error obtained in each sample.

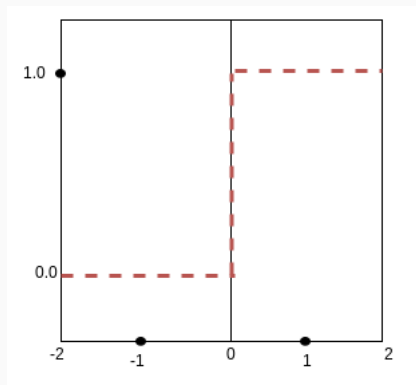
$$\mathcal{E} = \frac{1}{2} \sum_{j=1}^M e_j^2(n) = \frac{1}{2} \sum_{j=1}^M (z_j(n) - y_j(n))^2$$

- ▶ The learning curve in this case is  $\mathcal{E}(n)$  vs epoch.
- ▶ Learning curve is significantly different from that of batch learning.
- ▶ Online learning take advantage of redundant data (multiple copies of data).
- ▶ Online learning is simple to implement.

**Heaviside step function:**

$$\varphi(x) = 0 \quad \text{if } x < 0$$

$$\varphi(x) = 1 \quad \text{if } x \geq 0$$



- This is useful in the case of perceptron which works only when the data is

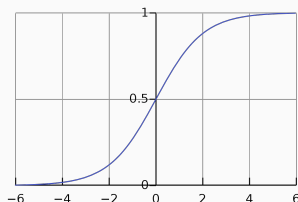
## Heaviside step function (contd.)

The reasons why we cannot use Heaviside step function in feedforward neural networks:

- ▶ We train neural network using backpropagation algorithm which requires differential activation function. For Heaviside step function, it is not differentiable at  $x=0$  and it has 0 derivation everywhere else.  
     $\implies$  The gradient descent will not be able to make progress in updating weights.
- ▶ We want our neural network weights to be modified continuously so that predictions can be as close as real values. Having a function that can only generate either 0 or 1 will not help to achieve this objective.

### Sigmoid Function

- ▶ Sigmoid function is also known as the logistic function.
- ▶ It non-linearly squashes a number to a value between 0 and 1.
- ▶  $Sigmoid(x) = \frac{1}{1+e^{-x}}$
- ▶ Activations are bounded between 0 and 1.



### Sigmoid Function (contd...)

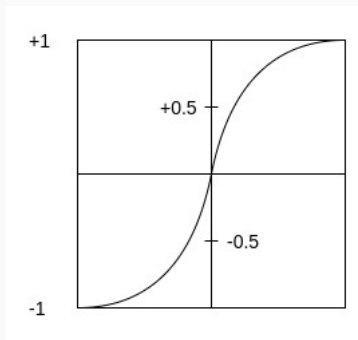
#### Disadvantages:

- 1 When the input is too small (towards  $-\infty$ ) the gradient is zero.
  - ▶ Hence while executing the backpropagation algorithm weights will not get updated i.e. there is no learning.
  - ▶ Vanishing gradient problem.
- 2 Though computing activation functions are less computationally expensive than matrix multiplication or convolutions, still computing exponential is expensive.

### Tanh Function (or Hyperbolic tangent)

- It is similar to sigmoid function but squashes the values, non-linearly between -1 and 1.

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

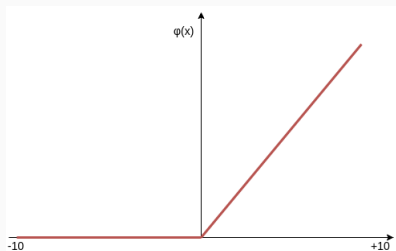


- Shares some disadvantages of sigmoid.

### Rectified Linear Unit (ReLU)

- Given an input if it is negative or zero, it outputs zero. Otherwise it outputs the number same as input.

$$\text{ReLU}(x) = \max(0, x)$$





### Rectified Linear Unit (ReLU)

- Is it non-linear? Yes.
  - A linear function should satisfy the property that

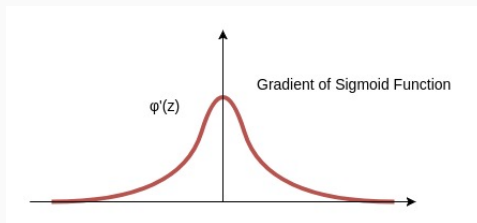
$$f(x + y) = f(x) + f(y)$$

$$\text{But } \varphi(-1) + \varphi(+1) \neq \varphi(0)$$

- But it is piece-wise linear.

## ReLU (contd...)

- ▶ It is a non-bounded function.
- ▶ Change from sigmoid to ReLU as an activation function in hidden layer is possible - hidden neurons need not have bounded values.
- ▶ The issue with sigmoid function is that it has very small values (near zero) everywhere except near 0.



## ReLU (contd...)

- At the  $j^{th}$  neuron which is a hidden layer

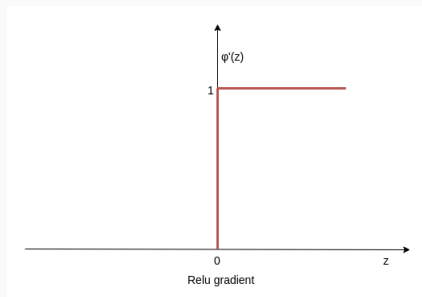
$$\Delta w_{ji}(n) = \eta \underbrace{\left( \sum_{k \in C} \delta_k(n) w_{kj}(n) \right)}_{\text{Local gradient from above}} \underbrace{\varphi'_j(v_j(n))}_{\text{Derivation of sigmoid function}}$$

We multiply the gradient from the layer above with partial derivative of the sigmoid function.

- Since in the case of sigmoid, it has very small values everywhere except for when the input has values close to 0.  
 $\implies$  Lower layers will likely have smaller gradients in terms of magnitude compared to higher layers.

## ReLU (contd...)

- The reason is in the case of sigmoid  $\varphi'(\cdot)$  is always less than 1 with most values being 0.  
 $\implies$  This imbalance in gradient magnitude makes it difficult to change the parameters of the neural networks with stochastic descent.



### ReLU (contd...)

- ▶ This problem can be adjusted by the use of rectified linear activation function because derivative of ReLU can have many non-zero values.
  - ⇒ Which in turn means that magnitude of the gradient is more balanced throughout the network.
- ▶ Dying ReLU: A neuron in the network is permanently dead due to inability to fire in forward pass.
  - ⇒ When activation is zero in the forward pass all the weights will get zero gradient.
  - ⇒ In backpropagation, the weights of neurons never get updated.
- ▶ Using ReLU in RNNs can blow up the computations to infinity as activations are not bounded.

"Backpropagation algorithm provides an 'approximation' to the trajectory in weight space computed by the method of steepest gradient."

**Case 1:-** Smaller the learning rate  $\Rightarrow$  smaller the change to the synoptic weights in the network



Smoother the learning will be

**Case 2:-** Larger learning rate  $\Rightarrow$  The network may become unstable or oscillatory

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \eta \delta_j(n) y_i(n)$$

where  $\alpha$  is the momentum constant

$\alpha = 0$  gives us original delta rule

"In general backpropagation algorithm cannot be shown to converge"



Hence this is not a well defined criteria for stopping its operation.

**Some good Criteria:**

1. Euclidean norm of  $\frac{\partial \mathcal{E}}{\partial w}$  reaches a sufficiently small gradient threshold.

∴ The necessary condition for  $w^*$  to be global maximum or local minimum is

$$\left. \frac{\partial \mathcal{E}}{\partial w} \right|_{w^*} = 0$$



2. Absolute rate of change is the average squared error per epoch is sufficiently small.
3. Stop when "generalizing performance" is adequate or it apparent that generalization performance is peaked.

$$w_{ji}(n+1) = w_{ji}(n) + \eta \delta_j(n) x_i(n)$$

- ▶  $\eta$  : learning rate
- ▶  $\delta_j(n)$  : local gradient
- ▶  $x_i(n)$  : Input to the  $j^{th}$  neuron

**Local gradient:**

$$\begin{aligned} \delta_j(n) &= e_j(n) \varphi'_j(v_j(n)) && \text{if } j \text{ is a output neuron} \\ &= \left[ \sum_k \underbrace{\delta_k(n)}_{\substack{\text{local gradient} \\ \text{of } k^{th} \text{ neuron} \\ \text{at the output}}} w_{kj}(n) \right] \varphi'_j(v_j(n)) && \text{if } j \text{ is a hidden neuron} \end{aligned}$$

- Rosenblatt's single layer perceptron has no hidden layer hence it cannot classify input pattern that are NOT linearly separable.

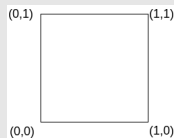
**Consider XOR Problem :-**

$$0 \oplus 0 = 0$$

$$1 \oplus 0 = 0$$

$$0 \oplus 1 = 1$$

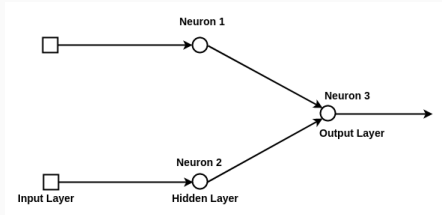
$$1 \oplus 1 = 1$$



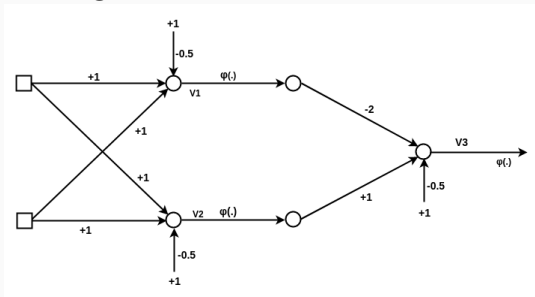
- $\{(0,0), (1,1)\}$  and  $\{(0,1), (1,0)\}$  are not linearly separable. Hence single layer neural network cannot solve this problem.

# XOR Problem (contd. . .)

- We use a hidden layer



- Consider the following neural network



- The function of output neuron : construct a linear combination of decision boundaries formed by the two hidden neurons.

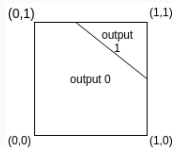
**For various inputs :**

- For input (1,1):  $v_1 = (1)(+1) + (1)(+1) + (+1)(-1.5)$   
 $= 1 + 1 - 1.5 = 0.5 \Rightarrow \varphi(v_1) = 1$   
 $v_2 = (1)(+1) + 1(+1) + (+1)(-0.5)$   
 $= 1 + 1 - 0.5 = 1.5 \Rightarrow \varphi(v_2) = 1$   
 $v_3 = 1(-2) + 1(+1) + (+1)(-0.5)$   
 $= -2 + 1 - 0.5 = -1.5 \Rightarrow \varphi(v_3) = 0$

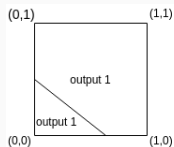
### For various inputs :

- ▶ For input (0,0):  
 $v_1 = (+1)(-1.5) = -1.5 \Rightarrow \varphi(v_1) = 0$   
 $v_2 = (+1)(-0.5) = -0.5 \Rightarrow \varphi(v_2) = 0$   
 $v_3 = (+1)(-0.5) = -0.5 \Rightarrow \varphi(v_3) = 0$
- ▶ For input (1,0):  
 $v_1 = (1)(+1) + (+1)(-0.5)$   
 $\quad = 1 - 1.5 = -0.5 \Rightarrow \varphi(v_1) = 0$   
 $v_2 = 1(+1) + (+1)(-0.5)$   
 $\quad = 1 - 0.5 = 0.5 \Rightarrow \varphi(v_2) = 1$   
 $v_3 = 1(+1) + (+1)(-0.5)$   
 $\quad = 1 - 0.5 = 0.5 \Rightarrow \varphi(v_3) = 1$
- ▶ For input (0,1):  
 $v_1 = (1)(+1) + (+1)(-1.5)$   
 $\quad = 1 - 1.5 = -0.5 \Rightarrow \varphi(v_1) = 0$   
 $v_2 = 1(+1) + (+1)(-0.5)$   
 $\quad = 1 - 0.5 = 0.5 \Rightarrow \varphi(v_2) = 1$   
 $v_3 = 1(+1) + (+1)(-0.5)$   
 $\quad = 1 - 0.5 = 0.5 \Rightarrow \varphi(v_3) = 1$

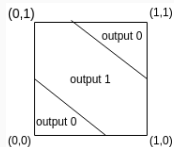
- Decision Boundary of neuron 1



- Decision Boundary of neuron 2



- Decision Boundary of neuron 3



Let  $\varphi(\cdot)$  be a non-constant, bounded and monotonic-increasing function. Let  $I_{m_0}$  denotes the  $m_0$  -dimensional unit hypercube  $[0, 1]^{m_0}$ . Let the space of continuous functions on  $I_{m_0}$  is denoted by  $C(I_{m_0})$ . Given any function  $f \in C(I_{m_0})$  and  $\epsilon > 0$ , there exists an integer  $m_1$  and sets of real numbers  $\alpha_i$ ,  $b_i$  and  $w_{ij}$ , where  $i = 1, 2, \dots, m_1$  and  $j = 1, 2, \dots, m_0$ .

Such that

$$F(x_1, \dots, x_{m_0}) = \sum_{i=1}^{m_1} \alpha_i \varphi\left(\sum_{j=1}^{m_0} w_{ij} x_j + b_i\right)$$

and  $F$  arbitrarily approximates  $f(\cdot)$ . That is

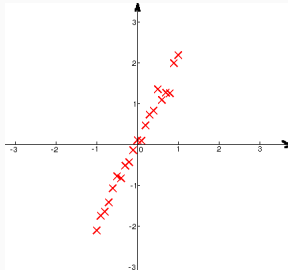
$$|F(x_1, \dots, x_{m_0}) - f(x_1, \dots, x_{m_0})| < \epsilon_0 \quad \forall x_1, \dots, x_{m_0} \in I_{m_0}$$



## AUTOENCODERS

---

- ▶ The origin of deep learning (post neural networks) since early 2000 was the use of Deep Belief Networks to “pretrain” deep networks.
- ▶ This approach is based on the observation that random initialization is not a good idea, and that pretraining each layer with an unsupervised learning algorithm can allow for better initial weights.
- ▶ Examples such unsupervised algorithms are
  - ▶ Deep Belief Networks based on Restricted Boltzmann Machines
  - ▶ Deep autoencoders



- ▶ Aim is to transmit this data: that is we have to send both the first and second dimension
- ▶ If we observe carefully, value at the second dimension is just twice the first dimension
- ▶ Hence we can just transmit first dimension (can be thought of as encoding of the data) and compute the value of the second dimension (can be thought of as decoding the data)

The process...

- ▶ Encoding: Map the data  $x_n$  by means of some method to compressed data  $z_n$
- ▶ Transmit
- ▶ Decoding: Map from compressed data  $z_n$  to  $\tilde{x}_n$

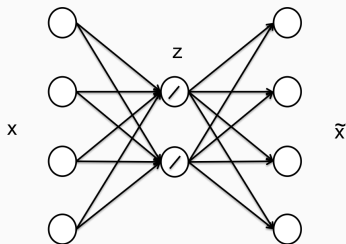
A linear encoding and decoding

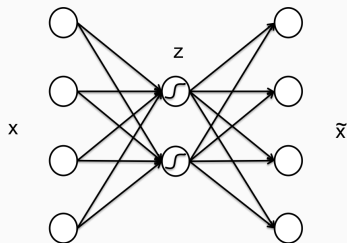
► Encoding:  $z_n = W_1 x_n + b_1$

► Decoding:  $\tilde{x}_n = W_2 z_n + b_2$

Objective function:

$$J(W_1, b_1, W_2, b_2) = \sum_{i=1}^N (\tilde{x}_n - x_n)^2$$





- ▶ If the data lie on a nonlinear surface, we use nonlinear activation functions.
- ▶ If the data is highly nonlinear, one could add more hidden layers to the networks to have a deep encoder.
- ▶ Note that this is an unsupervised learning.

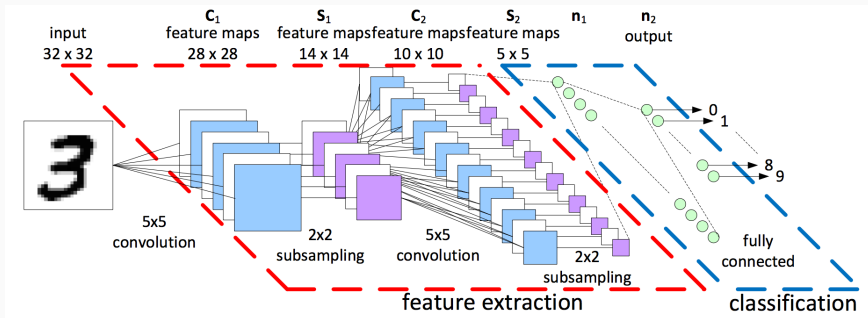
# CONVOLUTIONAL NEURAL NETWORKS

---

- ▶ Convolutional Neural Network (CNN) came into limelight in 2012
  - ▶ Alex Krizhevsky used CNN to win 2012 Imagenet competition.
  - ▶ The classification error has been improved from 20% to 15%.
  - ▶ **Paper:** Krizhevsky, Sutskever and Hinton : Imagenet classification with Deep Convolutioinal Neural Network, NIPS 2012.
- ▶ CNN were first proposed in the paper by Lecun, Bottou, Bengio, Haffner : Gradient based Learning Applied to Document Recognition, 1998 (*Proceedings of IEEE*)

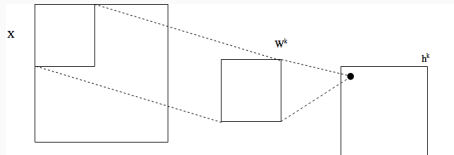


- ▶ Experiment by Hubel and Wiesel (1962)
- ▶ Some individual neuronal cells in the brain fire only in the presence of edges of certain Orientation.
- ▶ For example, Some neurons fired when exposed to vertical and some fired when exposed to horizontal edges.
- ▶ Hubel and Wiesel found that all these neurons were organized in a column architecture and that together they were able to produce visual perception.



- ▶ **Convolution:** Extract "Local" properties of the signal, using "filters" that have to be learned.
- ▶ **Pooling:** Down Samples the output to reduce the size of representation.
- ▶ **Nonlinearity:** Non-linearity is used after the convolution layer.

- ▶ This operation extracts local spacial properties of input



Figure

- ▶ The operation is defined as

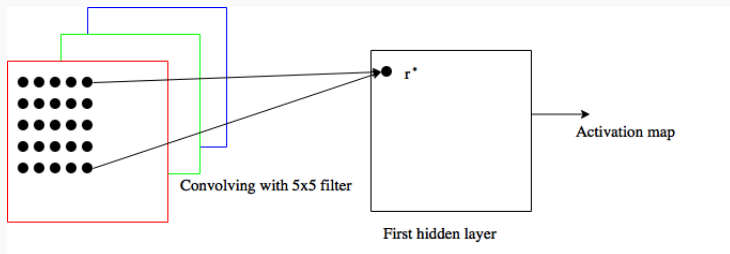
$$h_{ij}^k = f((W^k \odot X)_{ij} + b_k)$$

where  $W^k$  is a filter,  $\odot$  is the convolution operation and  $f$  is a nonlinear function.

- ▶ Second filter  $W^k$ ,  $k = 1, 2, 3, \dots$  are applied which need to be learned. Size of filter also need to be specified.

# Convolution Layer

- ▶ A small portion of the image that we look at the image from a "lens".
- ▶ Suppose size of this lens is 5x5.



$$r^* = \sum_{i=1}^5 \sum_{j=1}^5 a_{ij} b_{ij}$$

► **Example:**

Input	Filter	Output
32x32x3	5x5x3	28x28
32x32x32	5x5x3 (2 nos.)	28x28x2

- Each filter can be thought of as a "Feature identifier".
- **Intuition:** In the input image, if there is a shape that generally resembles the curve that the particular filter is representing thus all the multiplications summed together will result in a large Value.

**Stride:** Stride is size of the shift of the filter across the image (previously we kept stride as 1).

Ex.

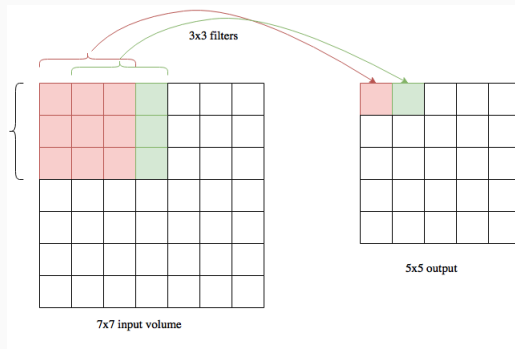


Figure: 3x3 convolution with a stride of 1

Stride (contd...) Example :

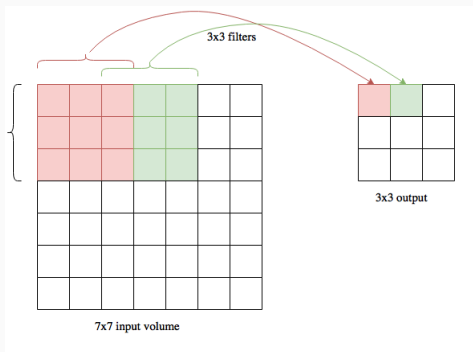


Figure: Convolution with stride of 2

►

$$\text{Size of output} = \frac{\text{Size of input} - \text{Size of filter}}{\text{Stride}} + 1$$



**Padding:** If we want output to be same size as input then we pad the output with zeros.

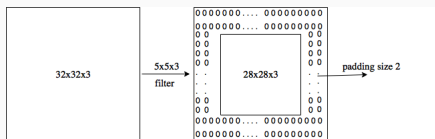


Figure: Padding of two to the output

- To enforce size of input and output to be same we need the padding size to be

$$\text{size of padding} = \frac{\text{size of filter} - 1}{2}$$

In general,

$$\text{size of output} = \frac{\text{size of input} - \text{size of filter} + 2 * \text{size of padding}}{\text{size of stride}} + 1$$

- A recent advance (not very recent) : Use ReLU,  $y(z) = \max(0, z)$  as the activation function instead of traditional sigmoid function.

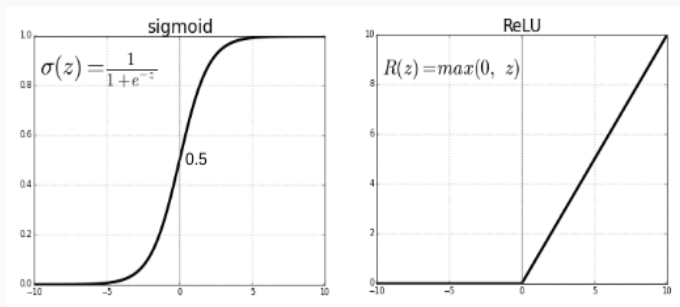


Figure: Activation functions

- ReLU improves performance of many networks.

Divide layer into partitions and get a max or average of each partition

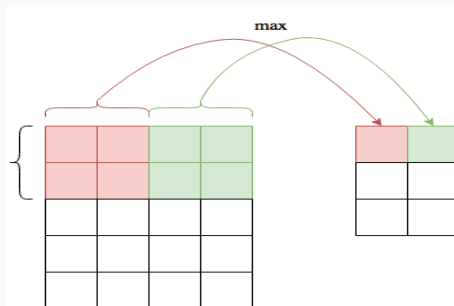


Figure: Maxpooling

- ▶ Max Pooling
- ▶ Average Pooling
- ▶  $L_2$  norm Pooling

- ▶ Advantages
  - ▶ Reduces the dimension of representation
  - ▶ Controls overfitting

Lookout : If you have 99% to 100% accuracy on training set and only 40% to 50% of test accuracy it is a cause of concern.

- ▶ This layer drops random set of activation in that layer by setting them to zero.
- ▶ Helps as a Regularizer.

## Alex Net (2012)

- ▶ Trained on 15 million images.
- ▶ Achieved test error 15.4% (The next best was only 26%).
- ▶ 5 convolution layer, max-pooling later, dropout layer and 3 fully connected layer.
- ▶ Used ReLU for activation.
- ▶ Used data augmentation techniques that consists of image translation, horizontal reflection and patch extraction.
- ▶ Implemented dropout layer in order to control overfitting to the training data.

### Alex Net (2012) (contd...)

- ▶ Trained the model using batch stochastic gradient descent with specific values of momentum and decay.
- ▶ Trained on two GTx580 GPU's for five to six days.

### ZF Net (2013) Zieler and Fergus

- ▶ Error of 11.2%.
- ▶ More of a fine tuning of Alex Net.
- ▶ Provided visualizations which provided better intuitions.
- ▶ ZF trained using 1.3 million images.



- ▶ Used 7x7 filters instead of 11x11 filters (as in Alex Net) also with decreased stride value.
  - ▶ Smaller filters in convolution layer help retain a lot of original pixel information in input image.
- ▶ ReLU for activation, cross entropy loss, training using batch stochastic descent.
- ▶ Trained on a GTx580 GPU for 12 days
- ▶ Deconvolutional network helps to visualize the feature maps.

### VGG Net (2014) Simonyan and Zisserman

- ▶ Error 7.3%.
- ▶ 19 layers of convolutional layer, 3x3 filters, padding of 2, max pooling with stride of 2.
- ▶ Trained for two to three weeks.

### Google Net (2015)

- ▶ Error 6.7%
- ▶ 22 layer CNN.
- ▶ User Inception modules.

### Microsoft ResNet (2015)

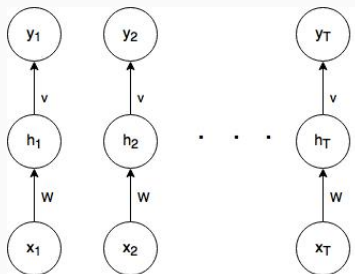
- ▶ Error 3.6%, (Human accuracy is 5 – 10%).
- ▶ Residual blocks.
- ▶ 152 layers.
- ▶ Trained on 8 GPU machines for 2 to 3 weeks.

## RECURRENT NEURAL NETWORKS

---

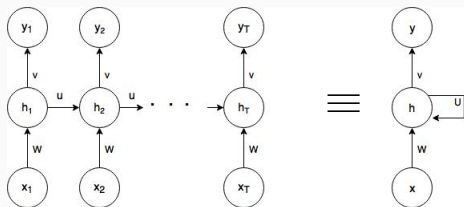
## Feed Forward Neural Networks(Recall)

- Consists of input,hidden and output layers
- Given sequential data, FF networks does not take sequential structure in the data
- Given a sequence of observations,  $x_1, \dots, x_T$ , then corresponding hidden units are  $h_1, \dots, h_T$  are assumed independent of each other (i.i.d data?)

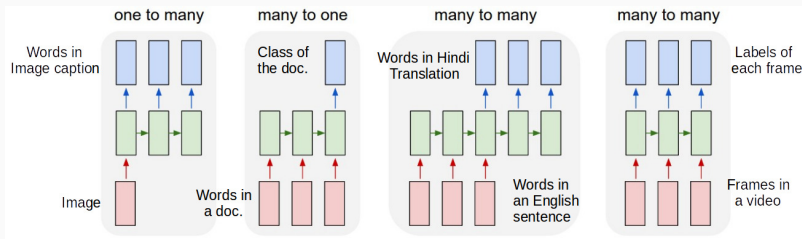


- Examples of sequential data:, text,audio,video.

- Since we have sequential data, hidden state at each step depends on the hidden states of the previous

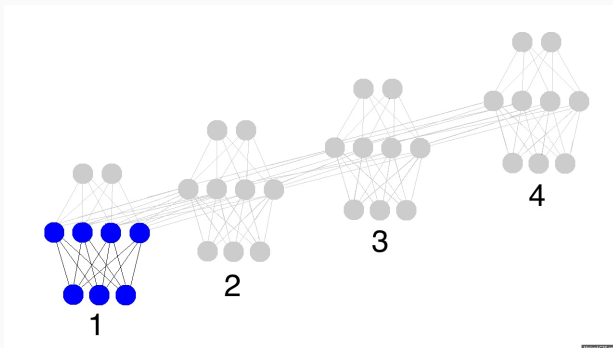


- Hence,  $h_t = \varphi(Wx_t + Uh_{t-1})$  where  $U$  acts as a transition matrix and  $\varphi$  is a nonlinear activation function
- $h_t$  acts as a memory



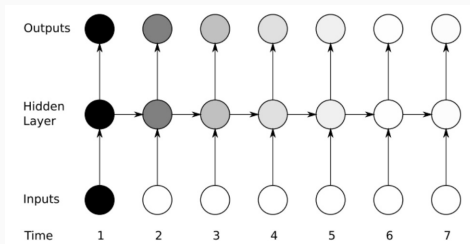
- ▶ RNN have many applications in modeling the sequential data
  - ▶ Input, Output or both can be sequences (possibly of different lengths)
  - ▶ Different inputs and different outputs need not be of the same length
- ▶ Regardless of the length of the input, RNN will learn fixed sized embedding for the input

- ▶ Trained using Backpropagation Through Time (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- ▶ Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets



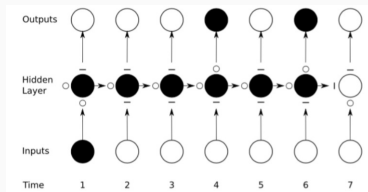


# Vanishing Gradient Problem



- Learnability of hidden states and outputs become weaker as we move away from them along the sequence  $\Rightarrow$  Weak Memory
- New inputs "overwrite" the activations of the previous hidden states
- Repeated multiplications can cause the gradients to vanish or explode (with ReLU)

# Capturing the Long range Dependencies



- ▶ Augment hidden states with gates
  - ▶ The gates involves some parameters which needs to be learned
- ▶ These gates will help the model to remember and target the information selectively
- ▶ The hidden states has three types of gates
  - ▶ Input(bottom), Forget(left) and Output(top)
- ▶ Open 'o', close '-'

- ▶ Proposed by Hochreiter and Schmidhuber
- ▶ Essentially RNN, except that the hidden states are computed differently
- ▶ In RNN, hidden states as  $h_t = \varphi(Wx_t + Uh_{t-1})$
- ▶ For RNN, state space is multiplicative (Weak Memory)

## CONCLUSION

---

What we have learned so far

- ▶ Perceptron
- ▶ Feed forward neural networks
- ▶ back propagation algorithm
- ▶ Autoencoders
- ▶ CNNs
- ▶ RNNs