

Natural Language Processing

Pushpak Bhattacharyya
CSE Dept,
IIT Patna and Bombay

Introduction to DL-NLP

Motivation for DL-NLP

Text classification using Deep NN

- Text classification at the heart of many NLP tasks
- **Soft** decisions needed
- A piece of text may belong to multiple classes
 - “Recent curb on H1-B in the US visas- promised in Trump’s election speeches- are causing IT industries to re-orient their business plan”
- Belongs to BOTH economics and Politics- *more* to former

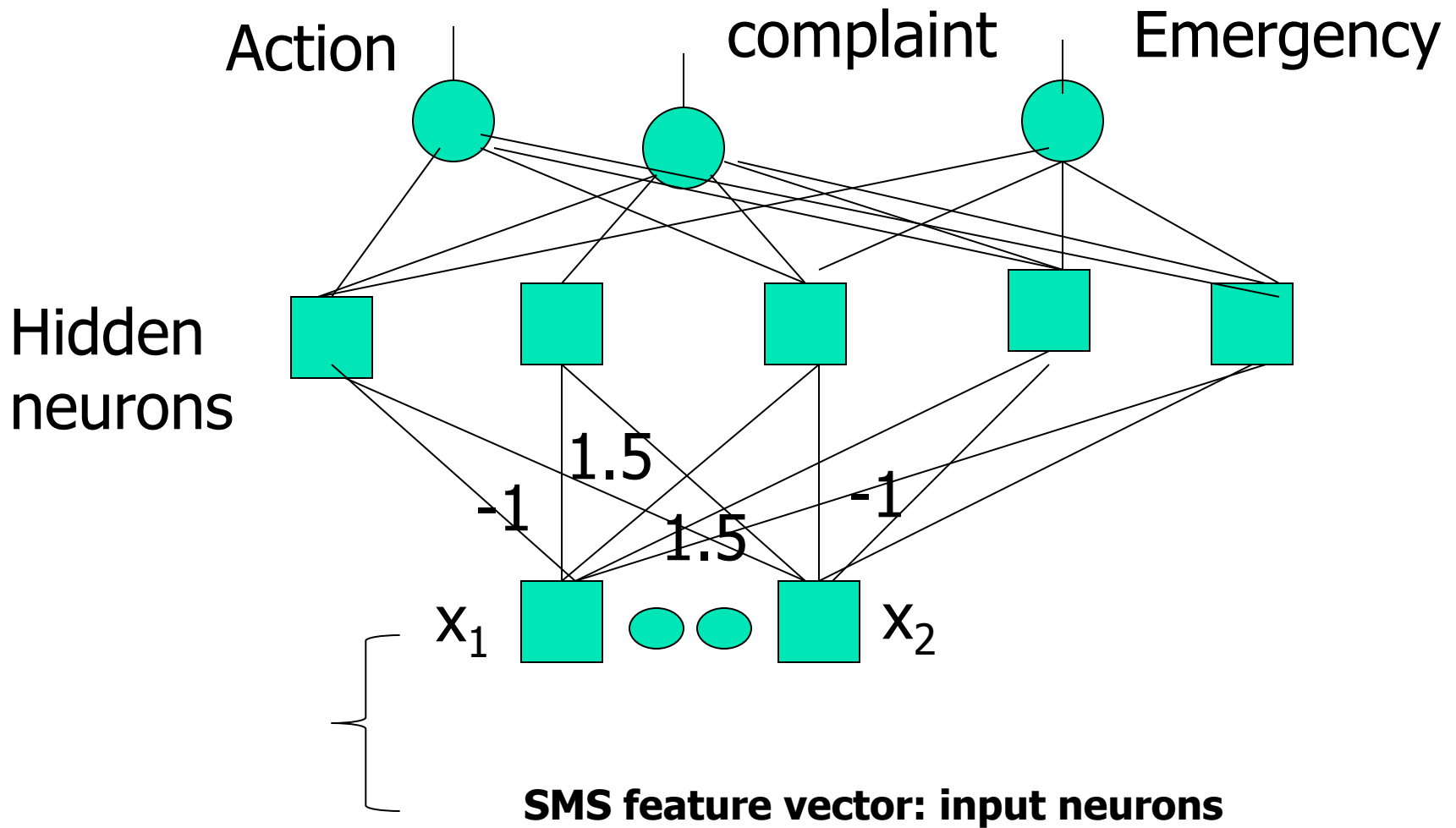
An example SMS complaint

- I have purchased a 80 litre Videocon fridge about 4 months ago when the freeze go to sleep that time compressor give a sound (khat khat khat khat) what is possible fault over it is normal I can't understand please help me give me a suitable answer.

Significant words (in red): after stop word removal

- I have purchased a 80 litre Videocon fridge about 4 months ago when the freeze go to sleep that time compressor give a sound (khat khat khat khat) what is possible fault over it is normal I can't understand please help me give me a suitable answer.

SMS classification

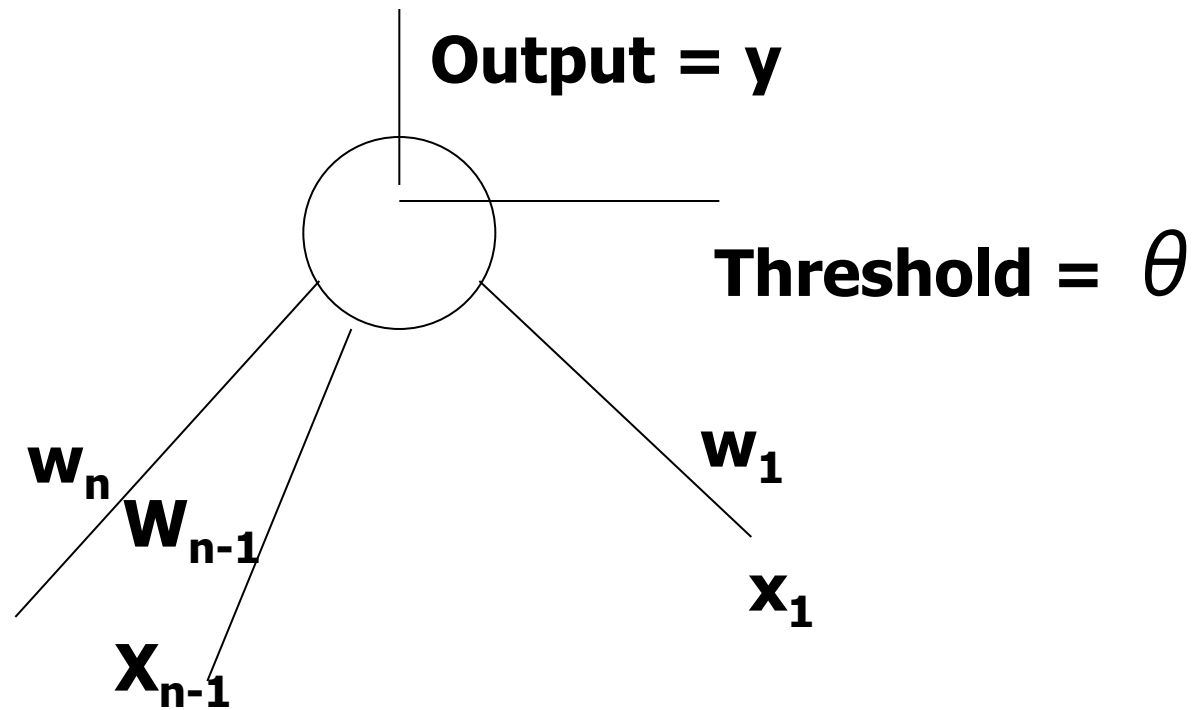


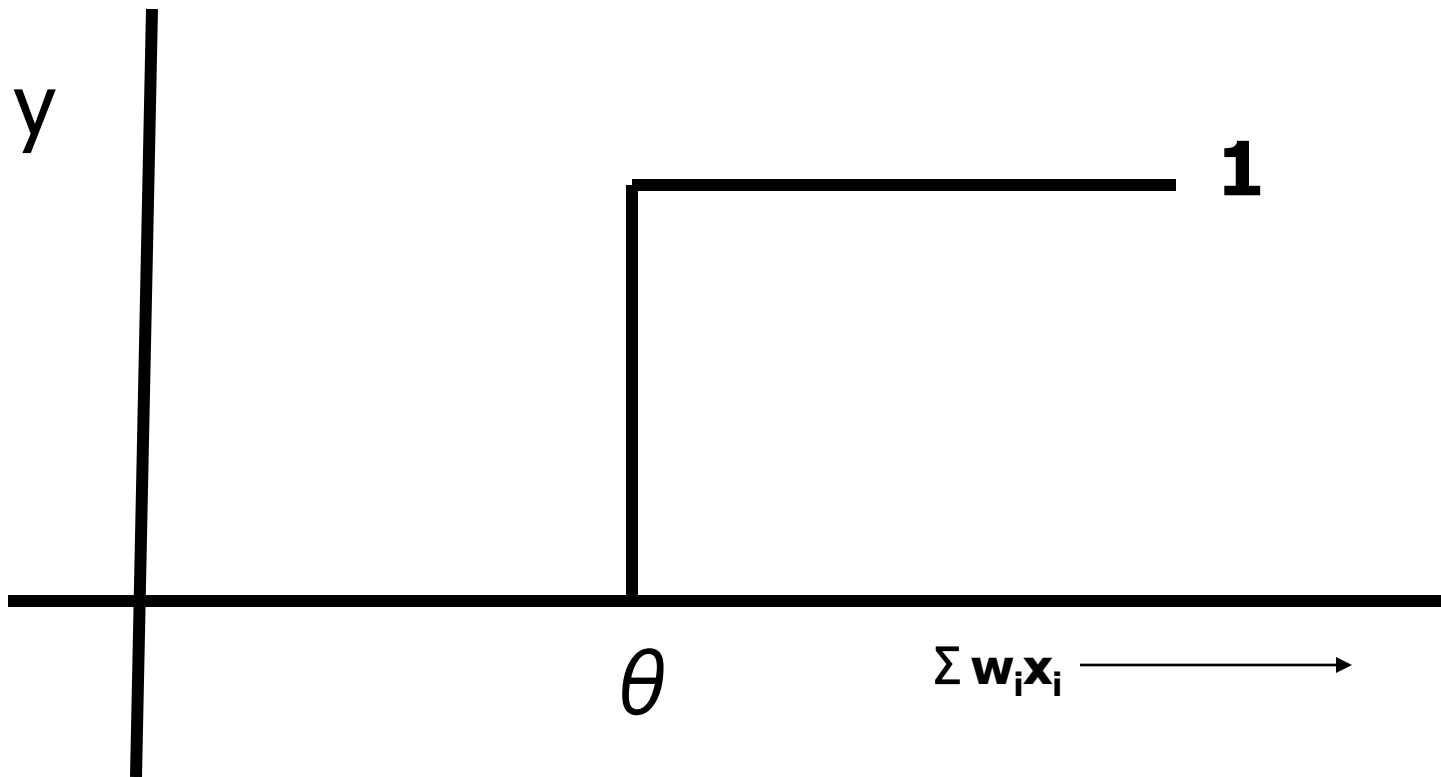
Basic Neural Models

- Precursor to Deep Learning
- Models
 - Perceptron and PTA
 - Feedforward Network and Backpropagation
 - Boltzmann Machine
 - Self Organization and Kohonen's Map
 - Neocognitron

Perceptron

The Perceptron Model





Step function / Threshold function

$$y = \begin{cases} 1 & \text{for } \sum w_i x_i \geq \theta \\ 0 & \text{otherwise} \end{cases}$$

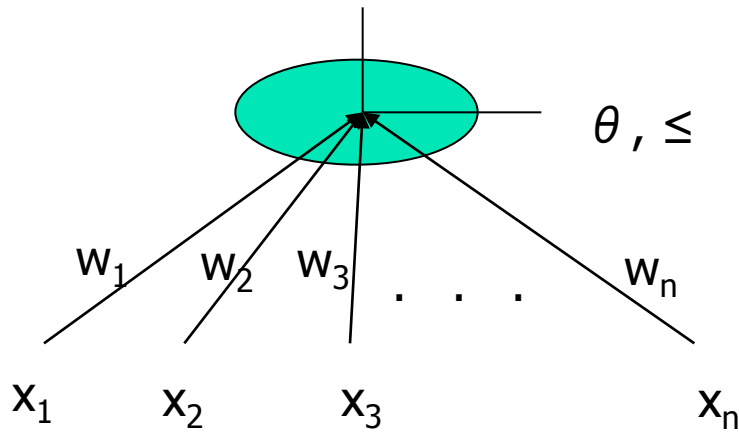
Perceptron Training Algorithm (PTA)

Preprocessing:

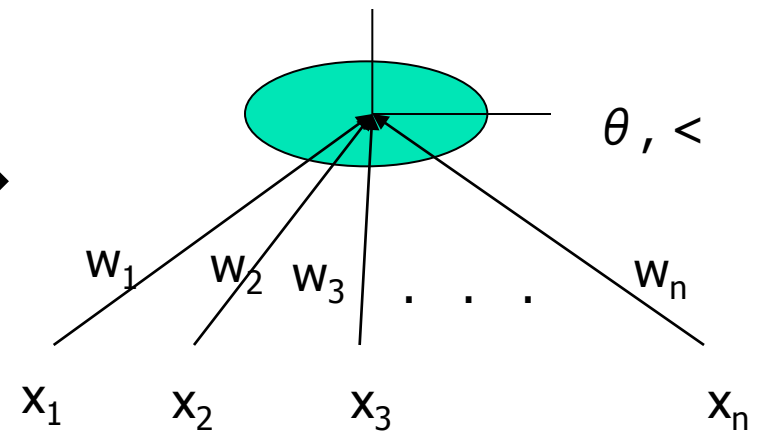
1. The computation law is modified to

$$y = 1 \text{ if } \sum w_i x_i > \theta$$

$$y = 0 \text{ if } \sum w_i x_i < \theta$$



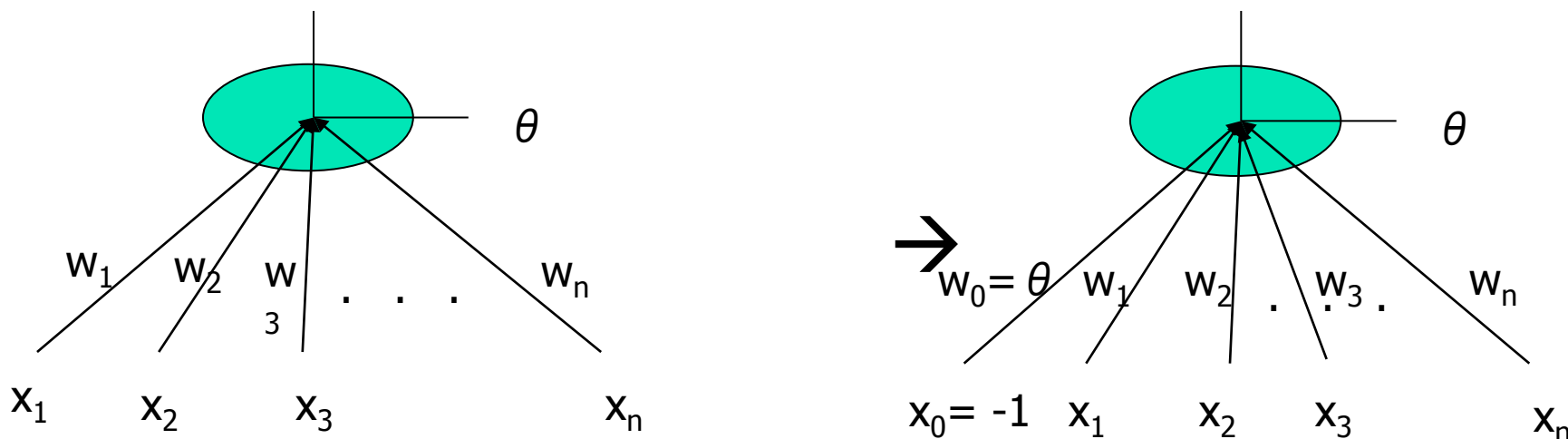
12 June, 2017



LG:nlp:DL:pushpak

PTA – preprocessing cont...

2. Absorb θ as a weight



3. Negate all the zero-class examples

Example to demonstrate preprocessing

■ OR perceptron

1-class $\langle 1,1 \rangle$, $\langle 1,0 \rangle$, $\langle 0,1 \rangle$

0-class $\langle 0,0 \rangle$

Augmented x vectors:-

1-class $\langle -1,1,1 \rangle$, $\langle -1,1,0 \rangle$, $\langle -1,0,1 \rangle$

0-class $\langle -1,0,0 \rangle$

Negate 0-class:- $\langle 1,0,0 \rangle$

Example to demonstrate preprocessing cont..

Now the vectors are

	X_0	X_1	X_2
X_1	-1	0	1
X_2	-1	1	0
X_3	-1	1	1
X_4	1	0	0

Perceptron Training Algorithm

1. Start with a random value of w
ex: $\langle 0, 0, 0 \dots \rangle$
2. Test for $w x_i > 0$
If the test succeeds for $i=1, 2, \dots, n$
then return w
3. Modify w , $w_{\text{next}} = w_{\text{prev}} + X_{\text{fail}}$

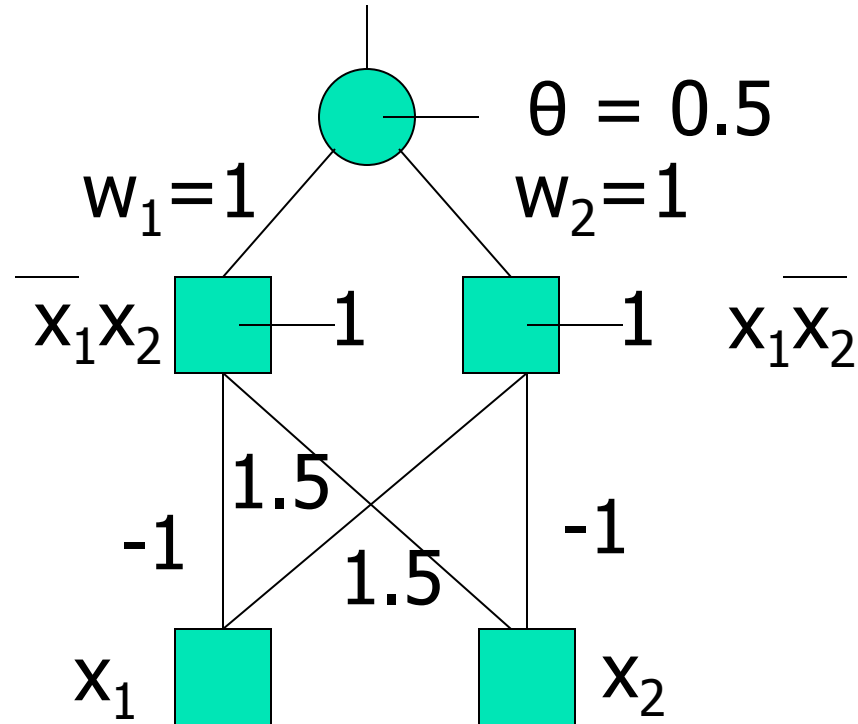
Convergence of PTA

- Statement:

Whatever be the initial choice of weights and whatever be the vector chosen for testing, PTA converges if the vectors are from a linearly separable function.

Feedforward Network and Backpropagation

Example - XOR



Gradient Descent Technique

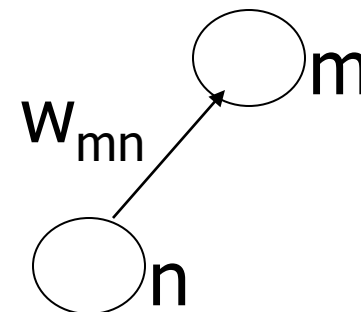
- Let E be the error at the output layer

$$E = \frac{1}{2} \sum_{j=1}^p \sum_{i=1}^n (t_i - o_i)_j^2$$

- t_i = target output; o_i = observed output
- i is the index going over n neurons in the outermost layer
- j is the index going over the p patterns (1 to p)
- Ex: XOR:— $p=4$ and $n=1$

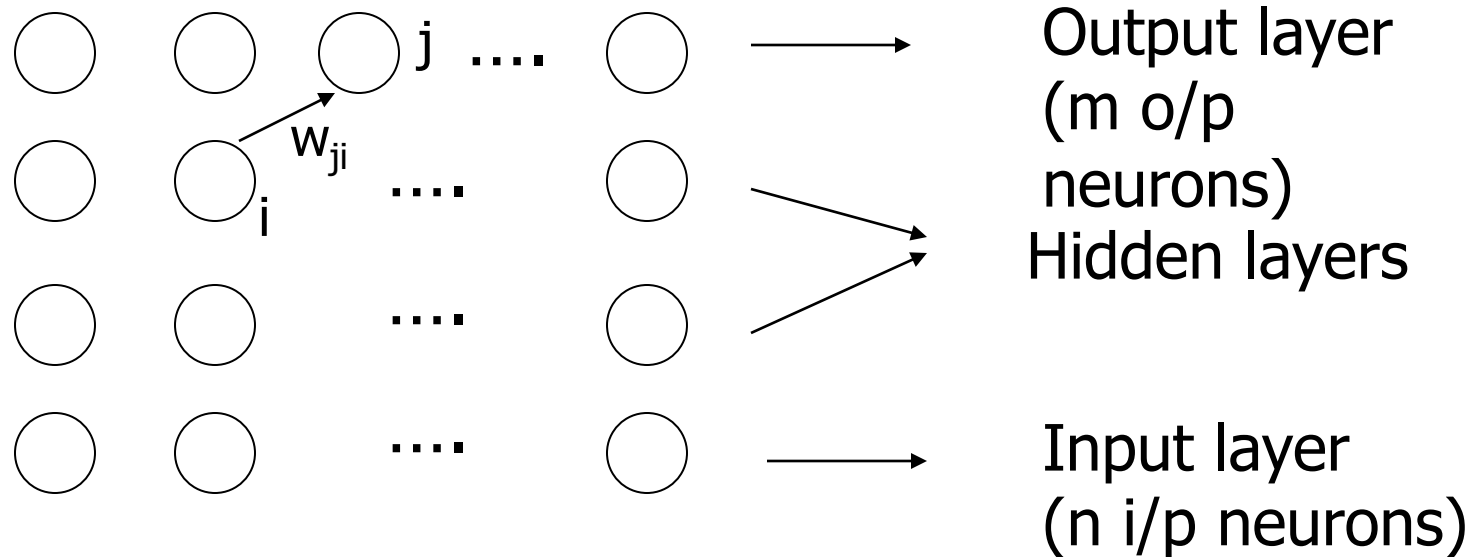
Weights in a FF NN

- w_{mn} is the weight of the connection from the n^{th} neuron to the m^{th} neuron
- E vs \bar{W} surface is a complex surface in the space defined by the weights w_{ij}
- $-\frac{\delta E}{\delta w_{mn}}$ gives the direction in which a movement of the operating point in the w_{mn} coordinate space will result in maximum decrease in error



$$\Delta w_{mn} \propto -\frac{\delta E}{\delta w_{mn}}$$

Backpropagation algorithm



- Fully connected feed forward network
- Pure FF network (no jumping of connections over layers)

Gradient Descent Equations

$$\Delta w_{ji} = -\eta \frac{\delta E}{\delta w_{ji}} \quad (\eta = \text{learning rate, } 0 \leq \eta \leq 1)$$

$$\frac{\delta E}{\delta w_{ji}} = \frac{\delta E}{\delta net_j} \times \frac{\delta net_j}{\delta w_{ji}} \quad (net_j = \text{input at the } j^{th} \text{ layer})$$

$$\frac{\delta E}{\delta net_j} = -\delta_j$$

$$\Delta w_{ji} = \eta \delta_j \frac{\delta net_j}{\delta w_{ji}} = \eta \delta_j o_i$$

Backpropagation – for outermost layer

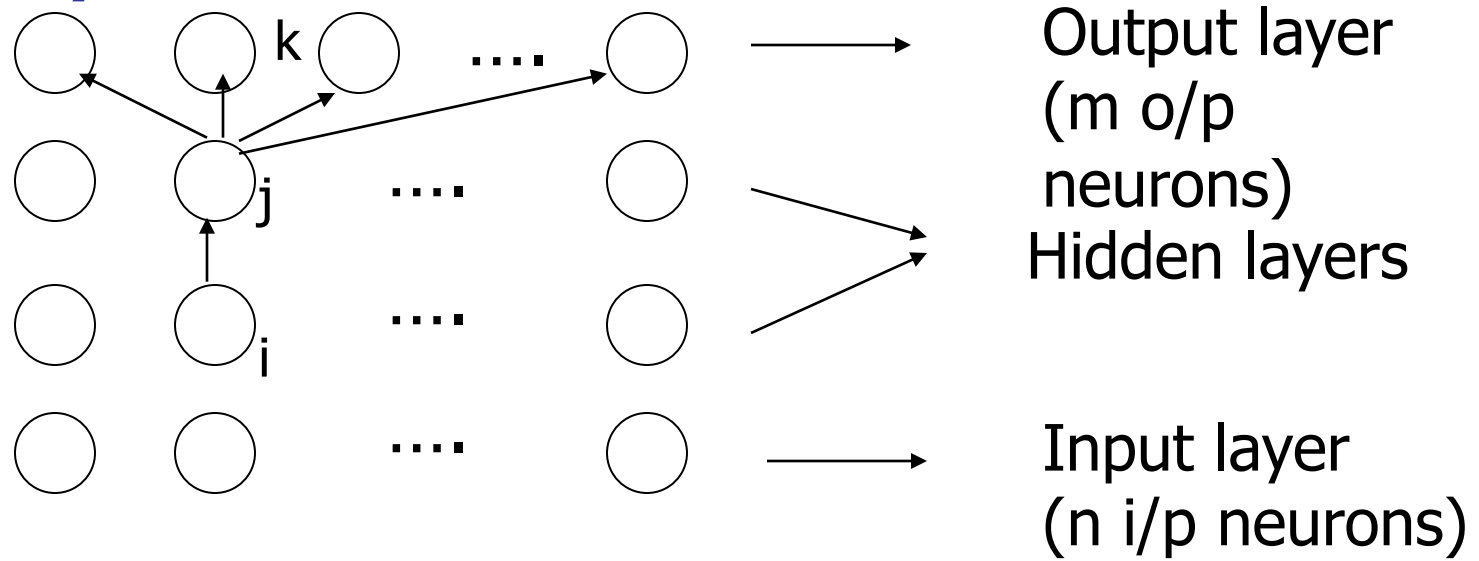
$$\delta_j = -\frac{\delta E}{\delta net_j} = -\frac{\delta E}{\delta o_j} \times \frac{\delta o_j}{\delta net_j} \quad (net_j = \text{input at the } j^{th} \text{ layer})$$

$$E = \frac{1}{2} \sum_{p=1}^m (t_p - o_p)^2$$

$$\text{Hence, } \delta_j = -(-(t_j - o_j)o_j(1 - o_j))$$

$$\Delta w_{ji} = \eta(t_j - o_j)o_j(1 - o_j)o_i$$

Backpropagation for hidden layers



δ_k is propagated backwards to find value of δ_j

Backpropagation – for hidden layers

$$\Delta w_{ji} = \eta \delta_j o_i$$

$$\delta_j = -\frac{\delta E}{\delta net_j} = -\frac{\delta E}{\delta o_j} \times \frac{\delta o_j}{\delta net_j}$$

$$= -\frac{\delta E}{\delta o_j} \times o_j (1 - o_j)$$

Cause of **Vanishing Gradient** problem!



$$= -\sum_{k \in \text{next layer}} \left(\frac{\delta E}{\delta net_k} \times \frac{\delta net_k}{\delta o_j} \right) \times o_j (1 - o_j)$$

$$\text{Hence, } \delta_j = -\sum_{k \in \text{next layer}} (-\delta_k \times w_{kj}) \times o_j (1 - o_j)$$

$$= \sum_{k \in \text{next layer}} (w_{kj} \delta_k) o_j (1 - o_j)$$

General Backpropagation Rule

- General weight updating rule:

$$\Delta w_{ji} = \eta \delta_j o_i$$

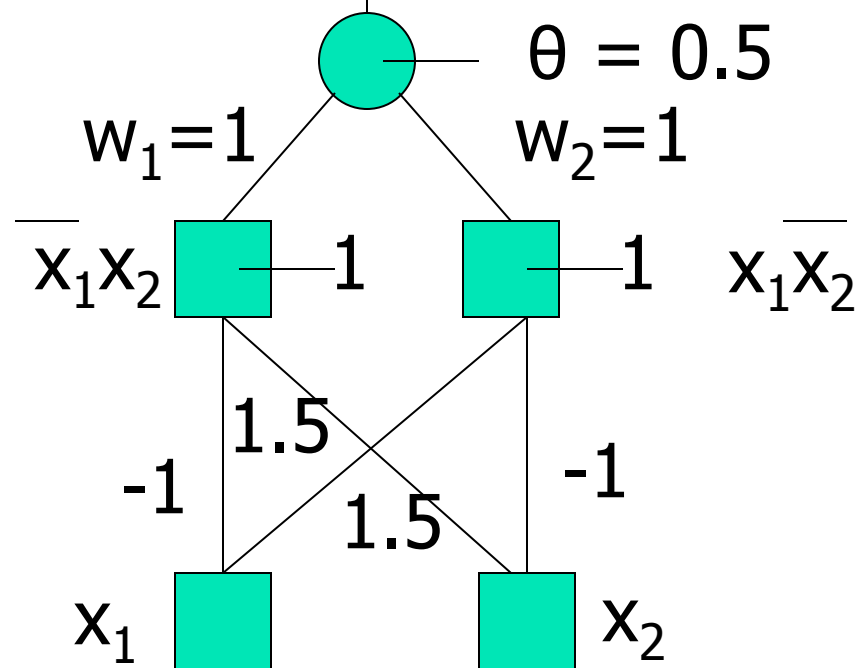
- Where

$$\delta_j = (t_j - o_j) o_j (1 - o_j) \quad \text{for outermost layer}$$

$$= \sum_{k \in \text{next layer}} (w_{kj} \delta_k) o_j (1 - o_j) o_i \quad \text{for hidden layers}$$

How does it work?

- Input propagation forward and error propagation backward (e.g. XOR)



Recurrent Neural Network

Acknowledgement:

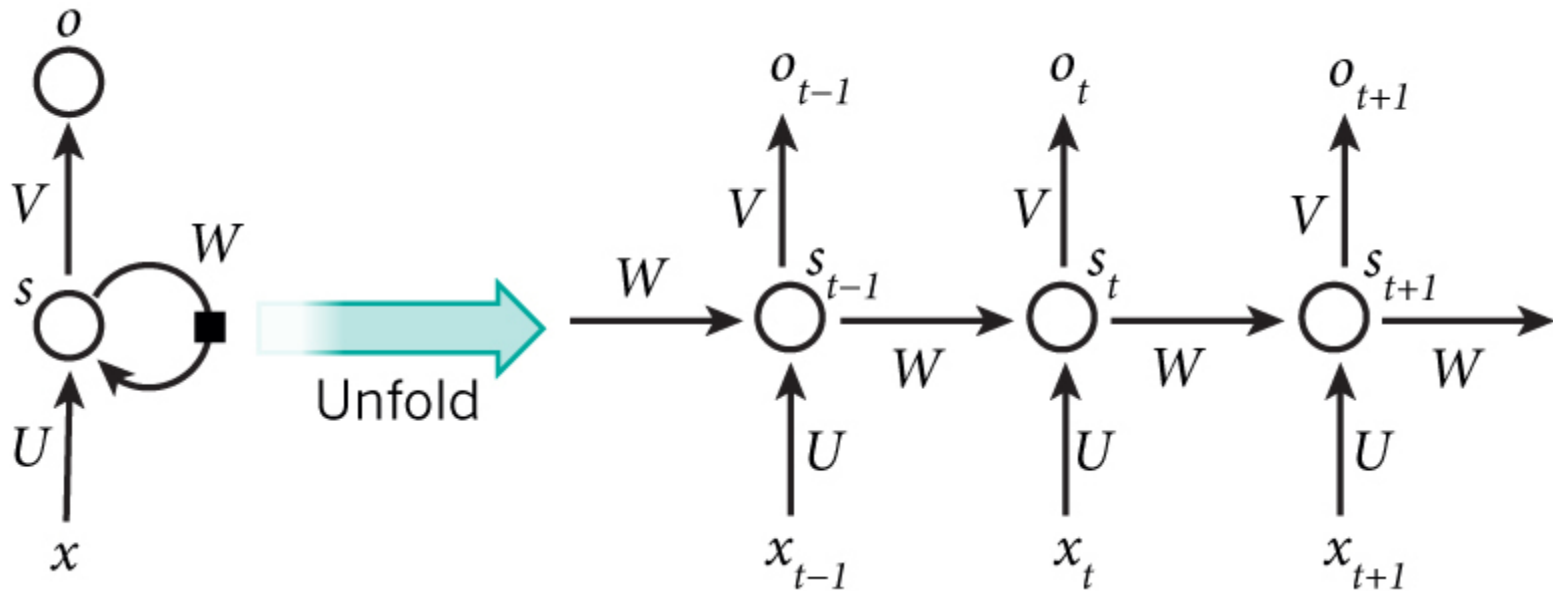
1. <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>

By Denny Britz

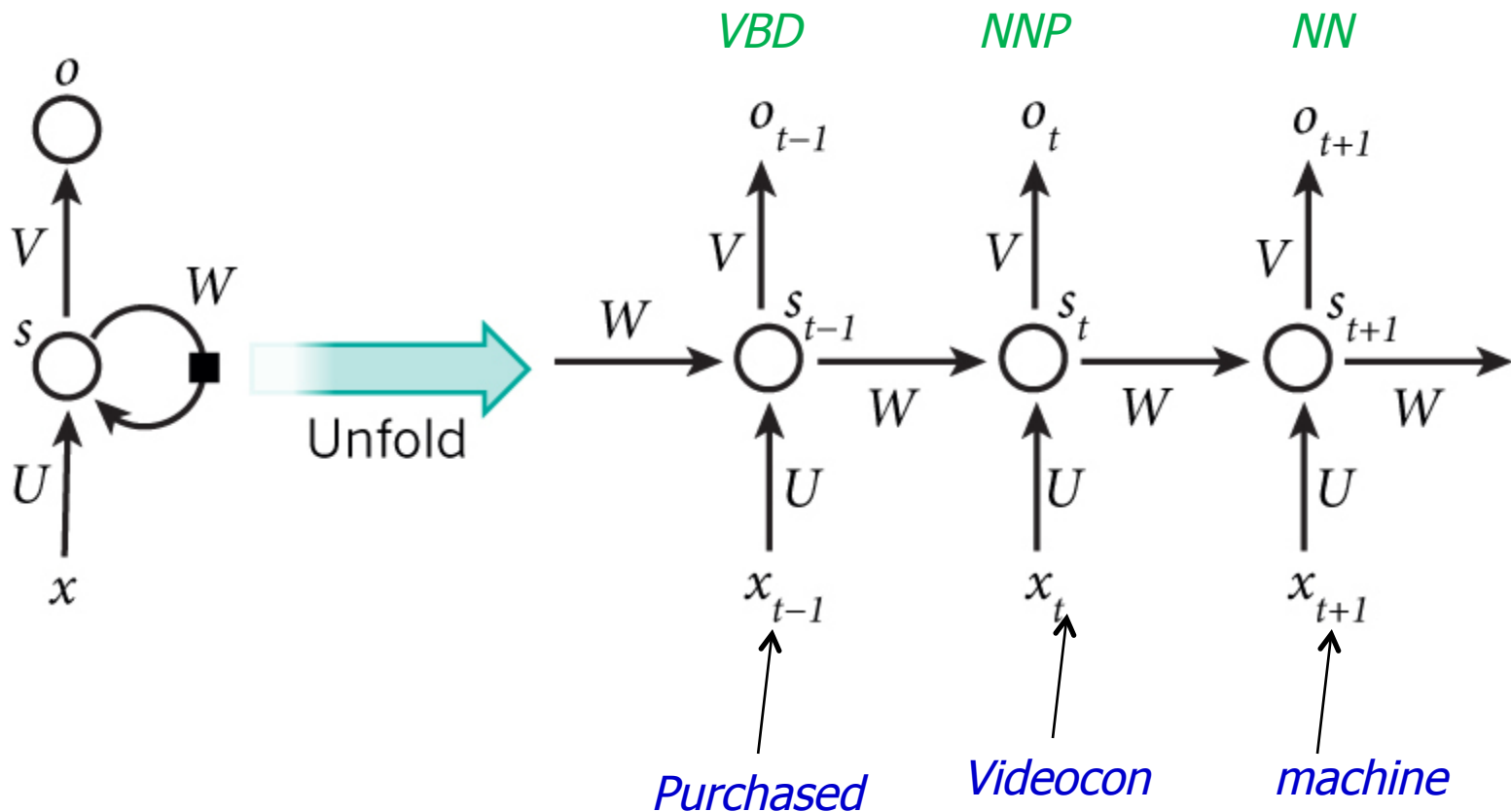
2. Introduction to RNN by Jeffrey Hinton

<http://www.cs.toronto.edu/~hinton/csc2535/lectures.html>

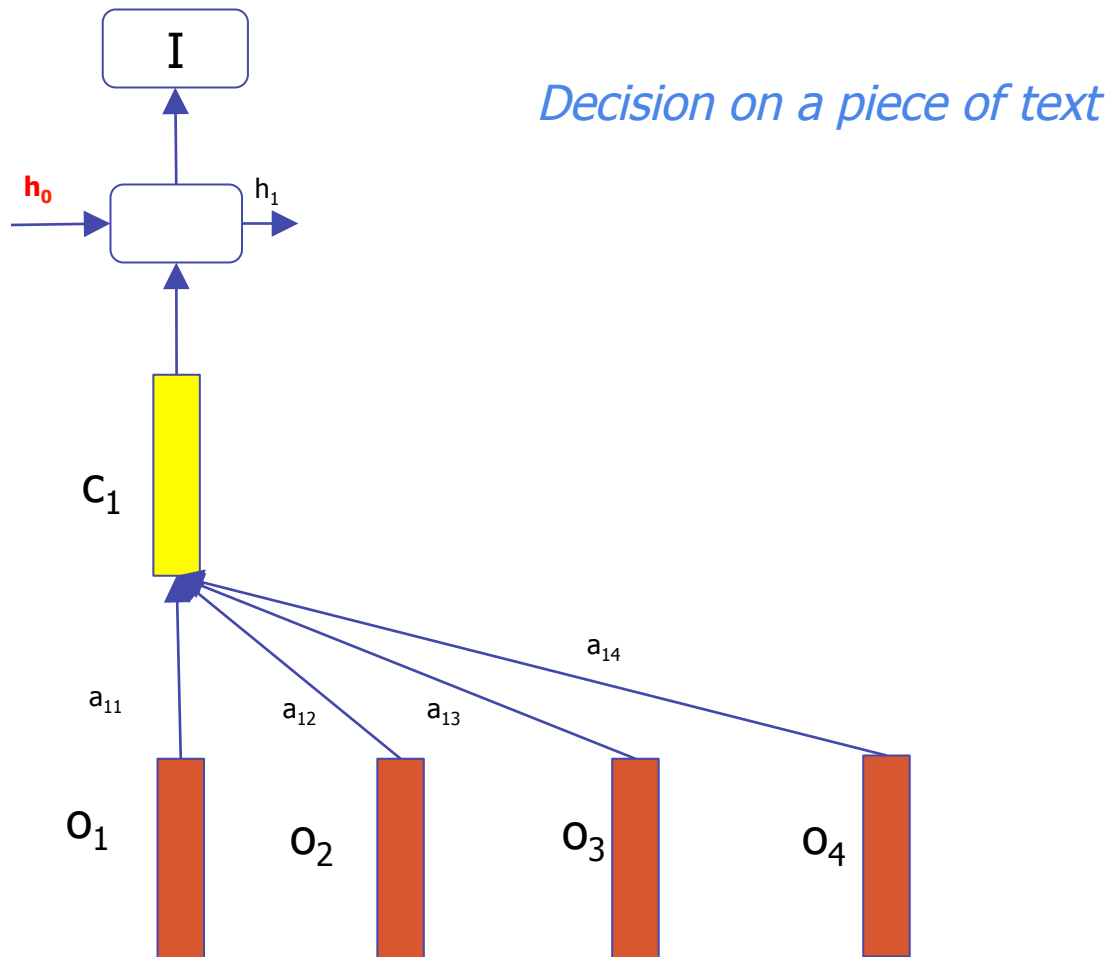
Sequence processing m/c

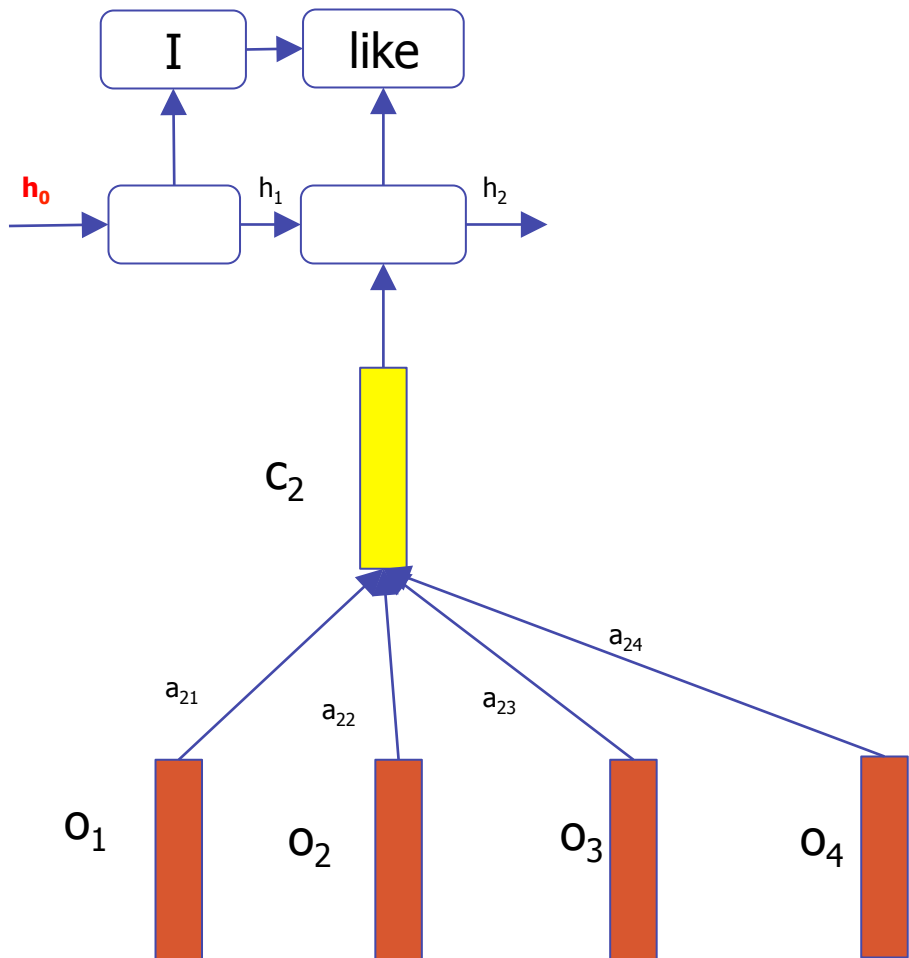


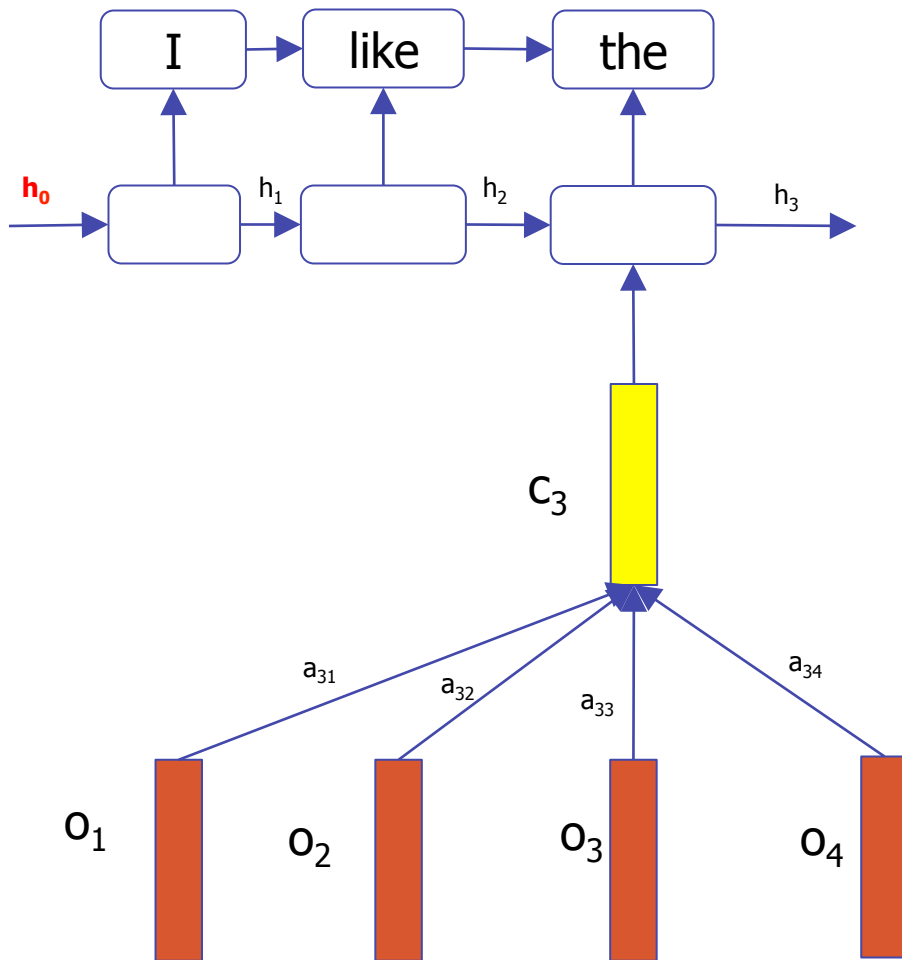
E.g. POS Tagging

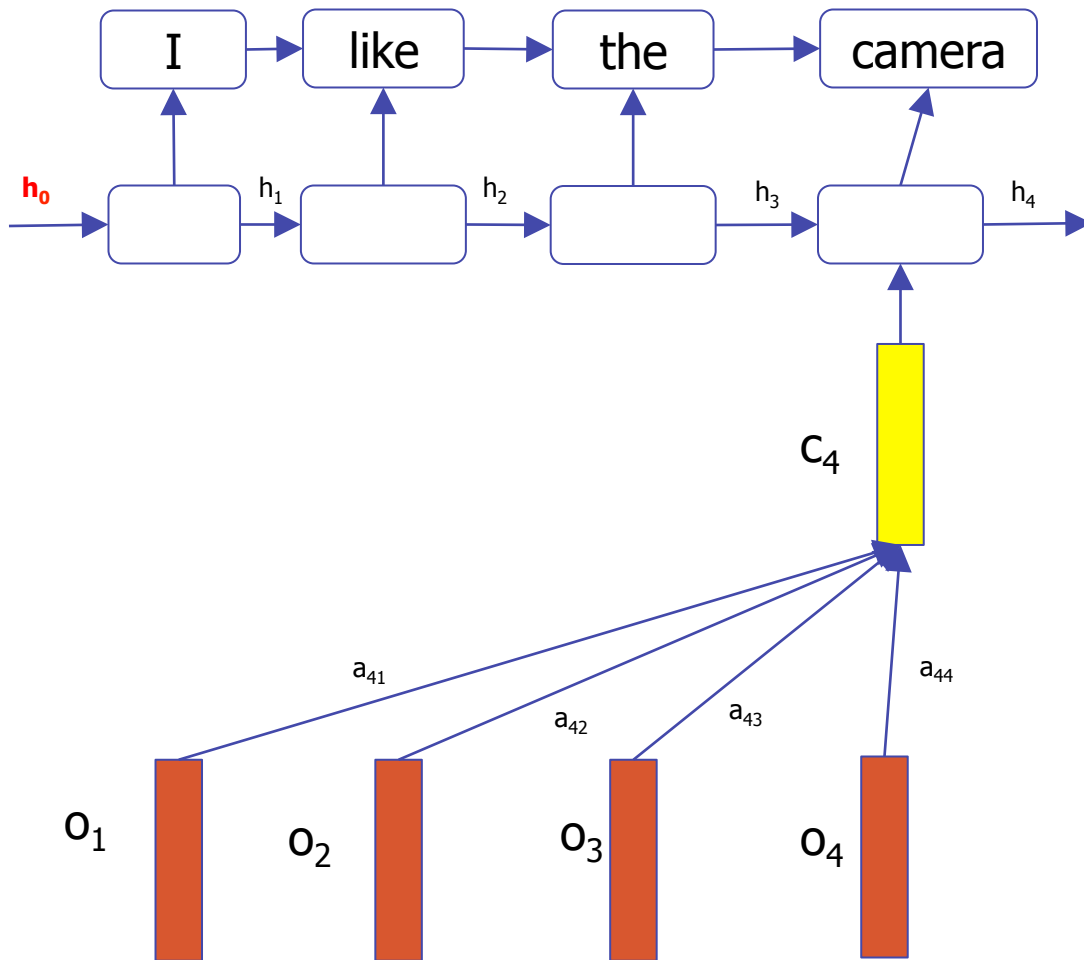


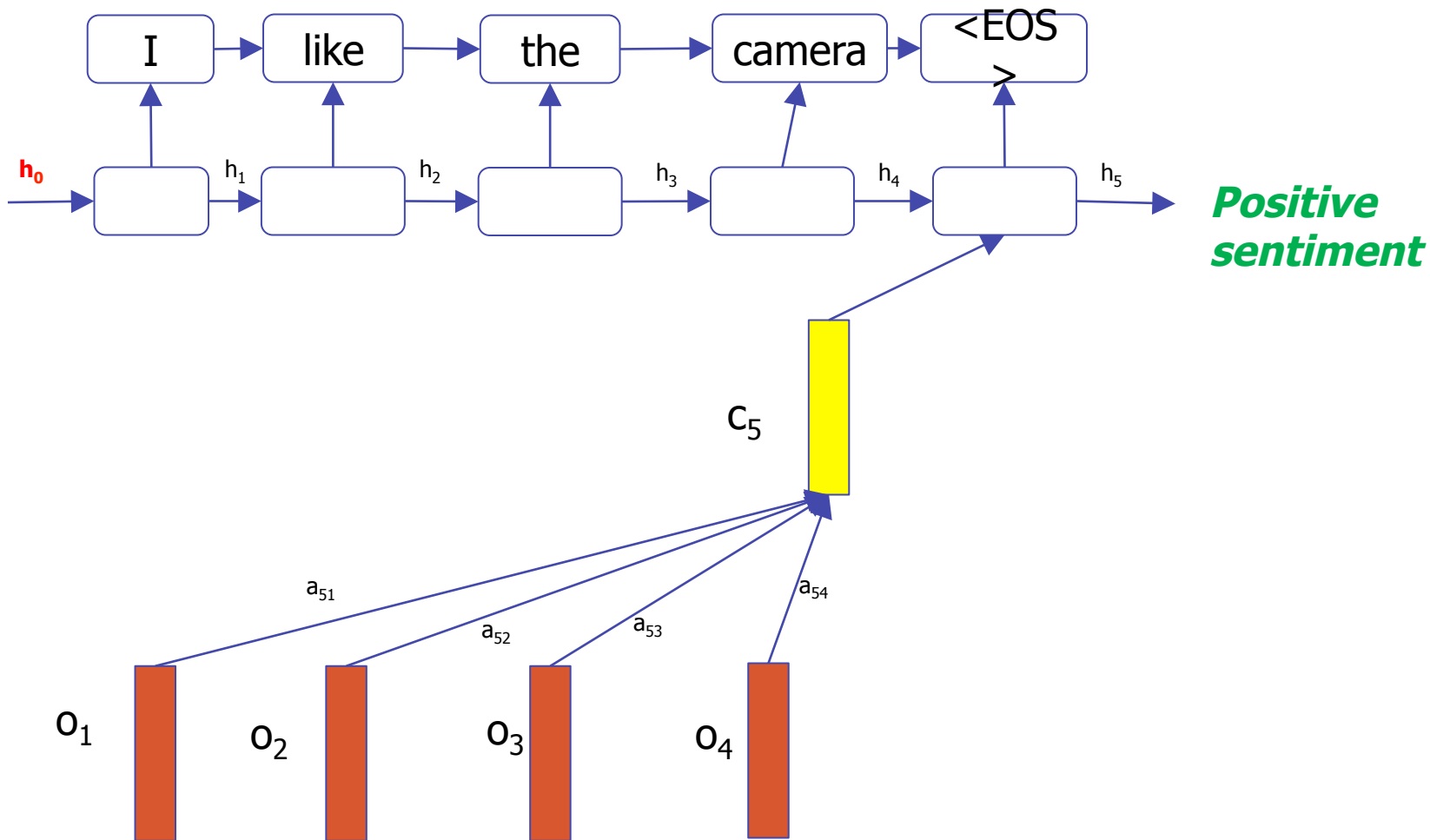
E.g. Sentiment Analysis



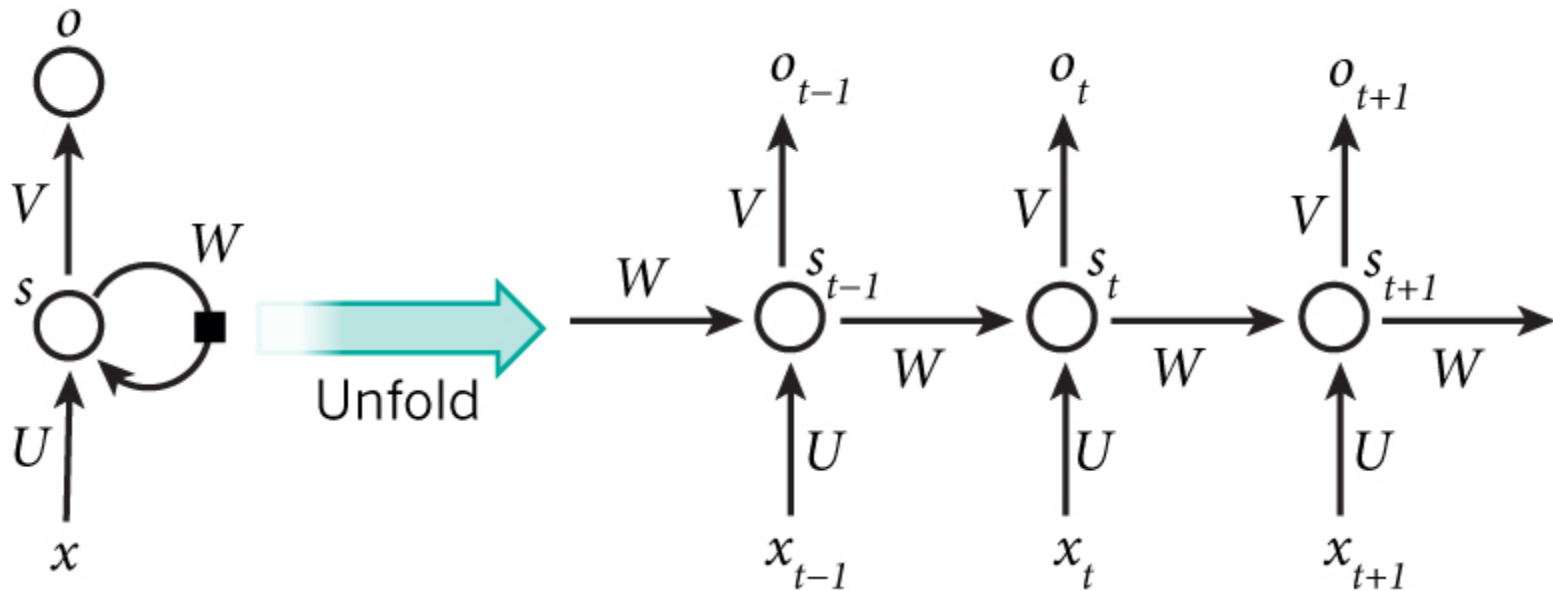








Back to RNN model

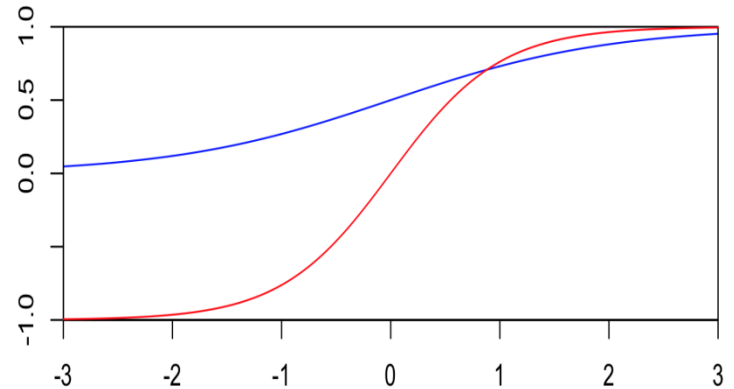


Notation: input and state

- x_t is the input at time step t . For example, could be a one-hot vector corresponding to the second word of a sentence.
- s_t is the hidden state at time step t . It is the “memory” of the network.
- $s_t = f(U \cdot x_t + W s_{t-1})$ **U and W matrices are learnt**
- f is a function of the input and the previous state
- Usually \tanh or $ReLU$ (approximated by *softplus*)

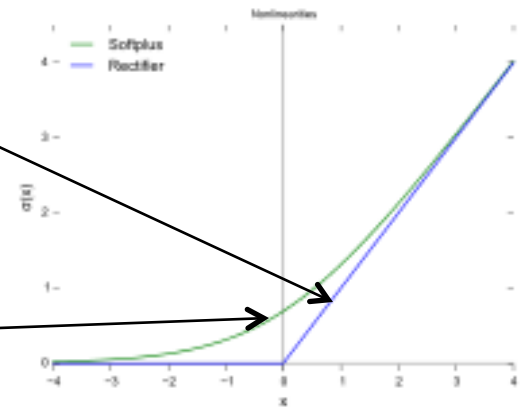
Tanh, ReLU (rectifier linear unit) and Softplus

$$\tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



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

$$g(x) = \ln(1 + e^x)$$



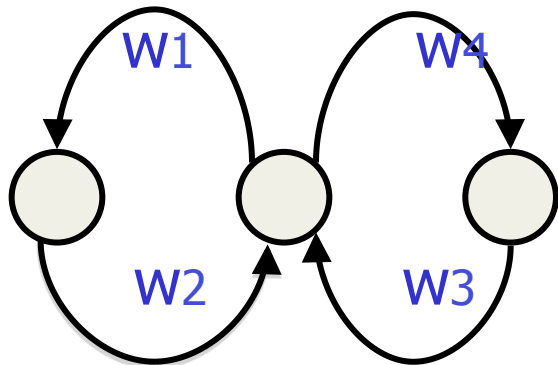
Notation: output

- o_t is the output at step t
- For example, if we wanted to predict the next word in a sentence it would be a vector of probabilities across our vocabulary
- $o_t = \text{softmax}(V \cdot s_t)$

Operation of RNN

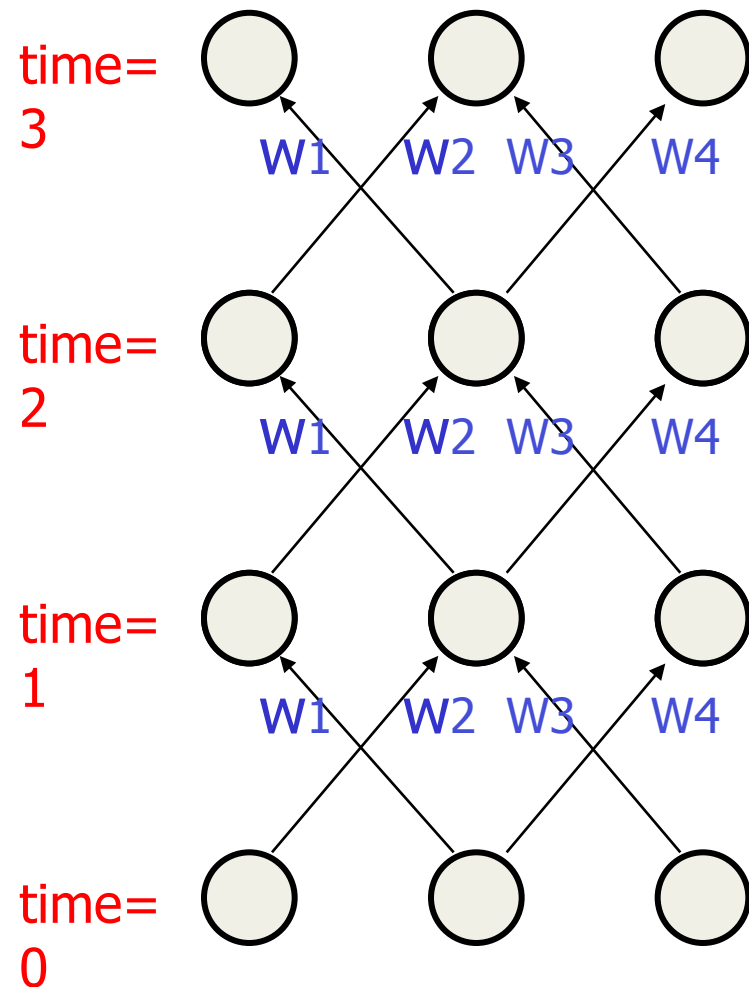
- RNN shares the same parameters (U , V , W) across all steps
- Only the input changes
- Sometimes the output at each time step is not needed: e.g., in sentiment analysis
- Main point: the **hidden states** !!

The equivalence between feedforward nets and recurrent nets



Assume that there is a time delay of 1 in using each connection.

The recurrent net is just a layered net that keeps reusing the same weights.



Reminder: Backpropagation with weight constraints

Example

- **Linear constraints between the weights.**
- **Compute the gradients as usual**
- **Then modify the gradients so that they satisfy the constraints.**
- **So if the weights started off satisfying the constraints, they will continue to satisfy them.**

To constrain: $w_1 = w_2$

we need: $\Delta w_1 = \Delta w_2$

compute: $\frac{\partial E}{\partial w_1}$ and $\frac{\partial E}{\partial w_2}$

use $\frac{\partial E}{\partial w_1} + \frac{\partial E}{\partial w_2}$ for w_1 and w_2

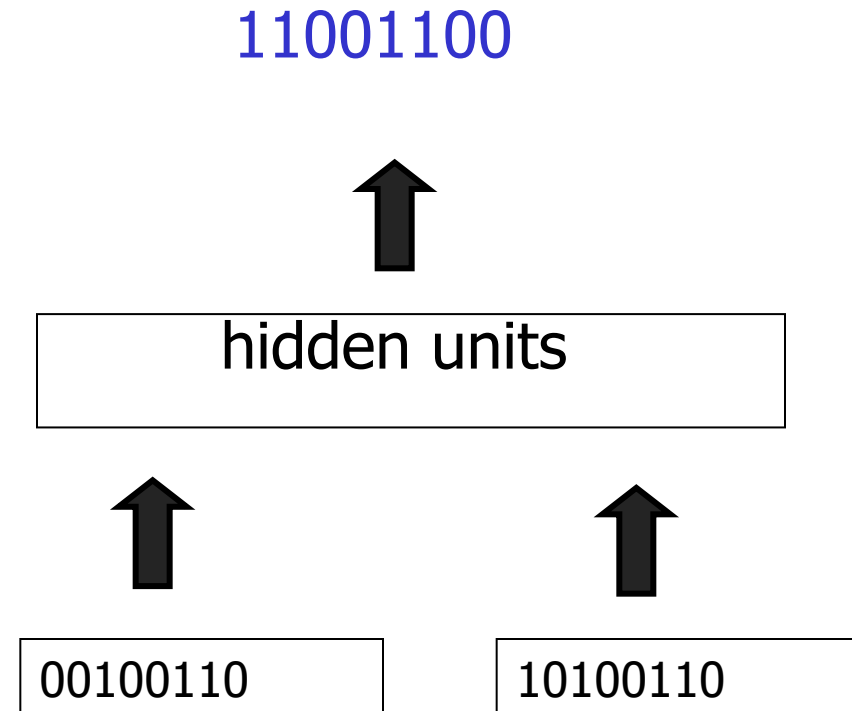
Backpropagation through time (BPTT algorithm)

- The forward pass at each time step.
-
- The backward pass computes the error derivatives at each time step.

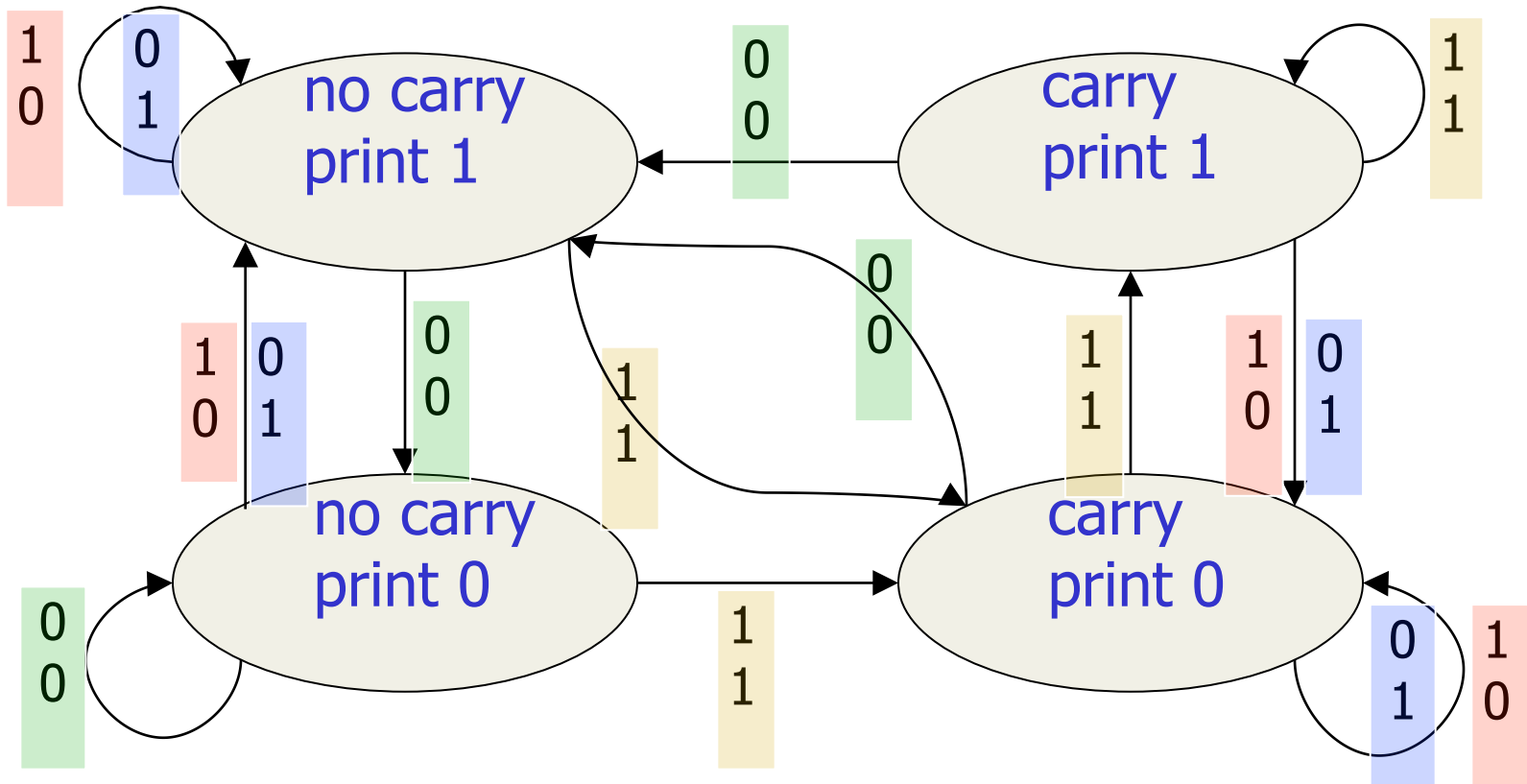
- After the backward pass we add together the derivatives at all the different times for each weight.

Binary addition using recurrent network (Jeffrey Hinton's lecture)

- Feed forward n/w
- But problem of variable length input



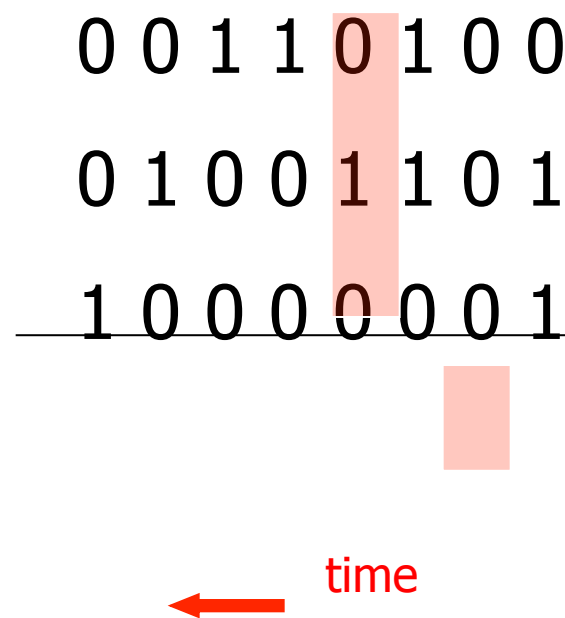
The algorithm for binary addition



This is a finite state automaton. It decides what transition to make by looking at the next column. It prints after making the transition. It moves from right to left over the two input numbers.

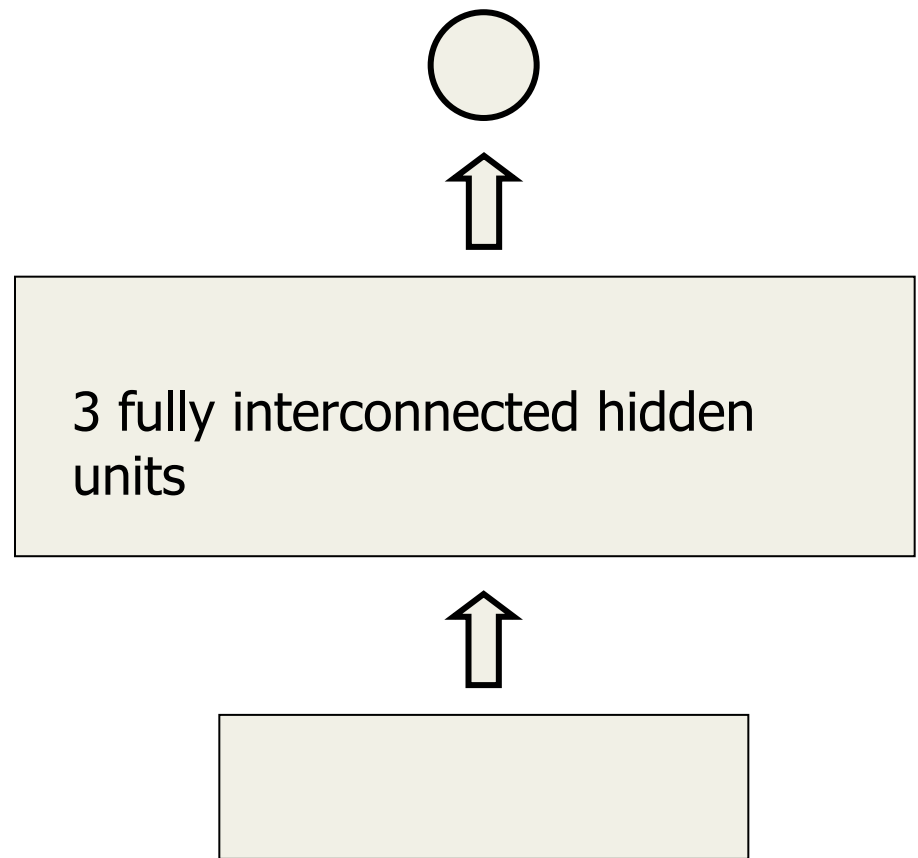
A recurrent net for binary addition

- Two input units and one output unit.
- Given two input digits at each time step.
- The desired output at each time step is the output for the column that was provided as input two time steps ago.
 - It takes one time step to update the hidden units based on the two input digits.
 - It takes another time step for the hidden units to cause the output.



The connectivity of the network

- The input units have feed forward connections
- Allow them to vote for the next hidden activity pattern.

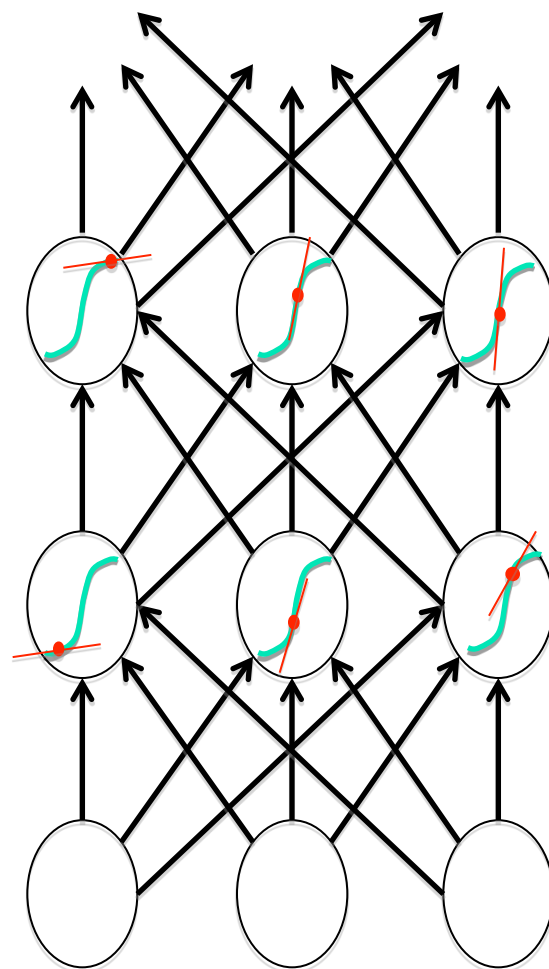


What the network learns

- Learns four distinct patterns of activity for the 3 hidden units.
- **Patterns** correspond to the nodes in the finite state automaton
- Nodes in FSM are like activity vectors
- The automaton is restricted to be in exactly one **state** at each time
- The hidden units are restricted to have exactly one **vector** of activity at each time.

The backward pass is linear

- The backward pass, is completely **linear**. If you double the error derivatives at the final layer, all the error derivatives will double.
- The forward pass determines the slope of the **linear** function used for backpropagating through each neuron.



Recall: Backpropagation Rule

- General weight updating rule:

$$\Delta w_{ji} = \eta \delta_j o_i$$

- Where

$$\delta_j = (t_j - o_j) o_j (1 - o_j) \quad \text{for outermost layer}$$

$$= \sum_{k \in \text{next layer}} (w_{kj} \delta_k) o_j (1 - o_j) o_i \quad \text{for hidden layers}$$

The problem of exploding or vanishing gradients (1/2)

- If the weights are small, the gradients shrink exponentially
- If the weights are big the gradients grow exponentially.
- Typical feed-forward neural nets can cope with these exponential effects because they only have a few hidden layers.

The problem of exploding or vanishing gradients (2/2)

- In an RNN trained on long sequences (*e.g.* sentence with 20 words) the gradients can easily explode or vanish.
 - We can avoid this by initializing the weights very carefully.
- Even with good initial weights, its very hard to detect that the current target output depends on an input from many time-steps ago.
 - So RNNs have difficulty dealing with long-range dependencies.

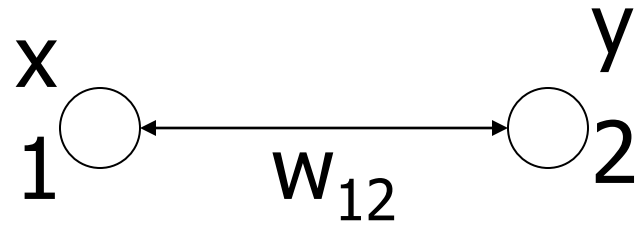
Vanishing/Exploding gradient: solution

- LSTM
- Error becomes “trapped” in the memory portion of the block
- This is referred to as an “error carousel”
- Continuously feeds error back to each of the gates until they become trained to cut off the value
- (to be expanded)

Boltzmann Machine

Illustration of the basic idea of Boltzmann Machine

- Illustrative task: To learn the identity function
- The setting is probabilistic, $x = 1$ or $x = -1$, with uniform probability, *i.e.*,
 - $P(x=1) = 0.5, P(x=-1) = 0.5$
- For, $x=1, y=1$ with $P=0.9$
- For, $x=-1, y=-1$ with $P=0.9$



x	y
1	1
-1	-1

Illustration of the basic idea of Boltzmann Machine (contd.)

- Let α = output neuron states

β = input neuron states

$P_{\alpha|\beta}$ = observed probability distribution

$Q_{\alpha|\beta}$ = desired probability distribution

Q_{β} = probability distribution on input states β

Illustration of the basic idea of Boltzmann Machine (contd.)

- The divergence D is given as:

$$D = \sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} \ln Q_{\alpha|\beta} / P_{\alpha|\beta}$$

called KL divergence formula

$$D = \sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} \ln Q_{\alpha|\beta} / P_{\alpha|\beta}$$

$$\geq \sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} (1 - P_{\alpha|\beta} / Q_{\alpha|\beta})$$

$$\geq \sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} - \sum_{\alpha} \sum_{\beta} P_{\alpha|\beta} Q_{\beta}$$

$$\geq \sum_{\alpha} \sum_{\beta} Q_{\alpha\beta} - \sum_{\alpha} \sum_{\beta} P_{\alpha\beta}$$

{ $Q_{\alpha\beta}$ and $P_{\alpha\beta}$ are joint distributions}

$$\geq 1 - 1 = 0$$

Precursor to
Softmax
Layer

Gradient descent for finding the weight change rule

$$P(S_\alpha) \propto \exp(-E(S_\alpha)/T)$$

$$P(S_\alpha) = (\exp(-E(S_\alpha)/T)) / (\sum_{\beta \in \text{all states}} \exp(-E(S_\beta)/T))$$

$$\ln(P(S_\alpha)) = (-E(S_\alpha)/T) - \ln Z$$

$$D = \sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} \ln (Q_{\alpha|\beta} / P_{\alpha|\beta})$$

$$\Delta w_{ij} = \eta (\delta D / \delta w_{ij}); \text{ gradient descent}$$

Calculating gradient: 1/2

$$\begin{aligned}\delta D / \delta w_{ij} &= \delta / \delta w_{ij} [\sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} \ln (Q_{\alpha|\beta} / P_{\alpha|\beta})] \\ &= \delta / \delta w_{ij} [\sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} \ln Q_{\alpha|\beta} \\ &\quad - \sum_{\alpha} \sum_{\beta} Q_{\alpha|\beta} Q_{\beta} \ln P_{\alpha|\beta}]\end{aligned}$$

$$\delta(\ln P_{\alpha|\beta}) / \delta w_{ij} = \delta / \delta w_{ij} [-E(S_{\alpha})/T - \ln Z]$$

Constant
With respect
To w_{ij}

$$Z = \sum_{\beta} \exp(-E(S_{\beta})) / T$$

Calculating gradient: 2/2

$$\begin{aligned} \delta \left[\frac{-E(S_\alpha)}{T} \right] / \delta w_{ij} &= (-1/T) \delta / \delta w_{ij} \left[- \sum_i \sum_{j>i} w_{ij} s_i s_j \right] \\ &= (-1/T) [-s_i s_j]_\alpha \\ &= (1/T) [s_i s_j]_\alpha \end{aligned}$$

$$\delta (\ln Z) / \delta w_{ij} = (1/Z) (\delta Z / \delta w_{ij})$$

$$Z = \sum_\beta \exp(-E(S_\beta)/T)$$

$$\begin{aligned} \delta Z / \delta w_{ij} &= \sum_\beta \left[\exp(-E(S_\beta)/T) (\delta(-E(S_\beta)/T) / \delta w_{ij}) \right] \\ &= (1/T) \sum_\beta \exp(-E(S_\beta)/T) \cdot s_i s_j |_\beta \end{aligned}$$

Final formula for Δw_{ij}

$$\Delta w_{ij} = [1/T] [s_i s_j |_{\alpha} - \underbrace{(1/Z) \sum_{\beta} \exp(-E(S_{\beta})) / T}_{\substack{\text{Expectation of } i^{\text{th}} \text{ and } j^{\text{th}} \\ \text{Neurons being on together}}} \cdot s_i s_j |_{\beta}]$$

Issue of Hidden Neurons

- Boltzmann machines
 - can come with hidden neurons
 - are equivalent to a Markov Random field
 - with hidden neurons are like a Hidden Markov Machines
- Training a Boltzmann machine is equivalent to running the Expectation Maximization Algorithm

Use of Boltzmann machine

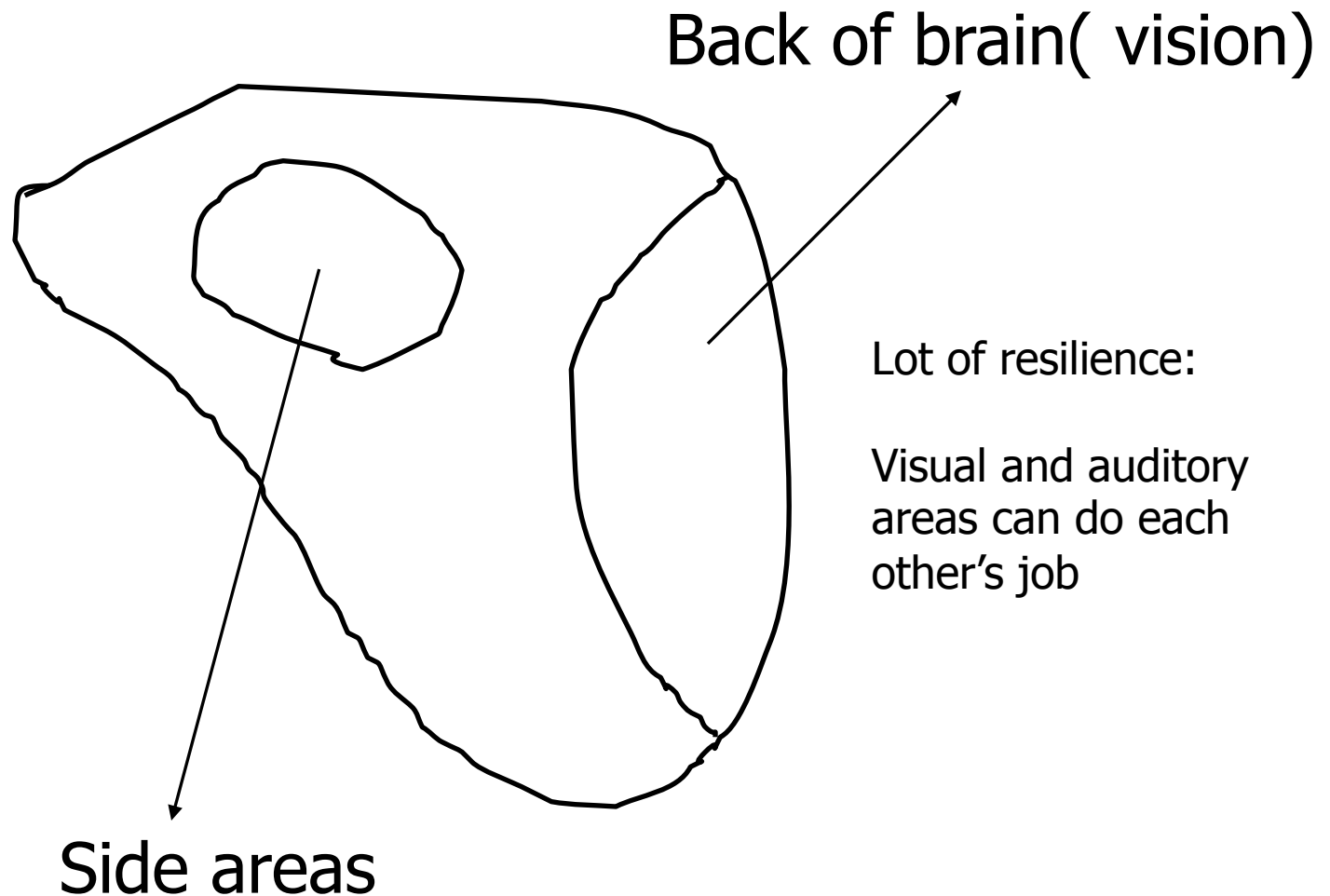
- Computer Vision
 - Understanding scene involves what is called “Relaxation Search” which gradually minimizes a cost function with progressive relaxation on constraints
- Boltzmann machine has been found to be slow in the training
 - Boltzmann training is NP-hard.

Questions

- Does the Boltzmann machine reach the global minimum? What ensures it?
- Why is simulated annealing applied to Boltzmann machine?
 - local minimum \rightarrow increase $T \rightarrow$ n/w runs \rightarrow gradually reduce $T \rightarrow$ reach global minimum.
- Understand the effect of varying T
 - Higher $T \rightarrow$ small difference in energy states ignored, convergence to local minimum fast.

Self Organization and Kohonen Net

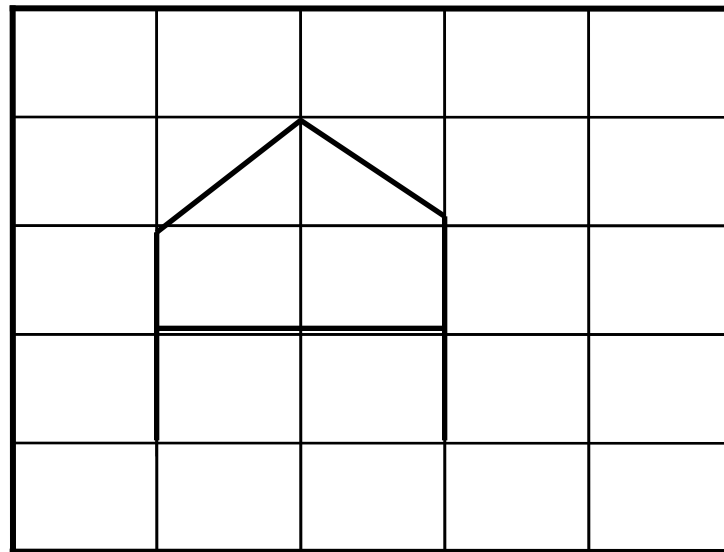
Mapping of Brain



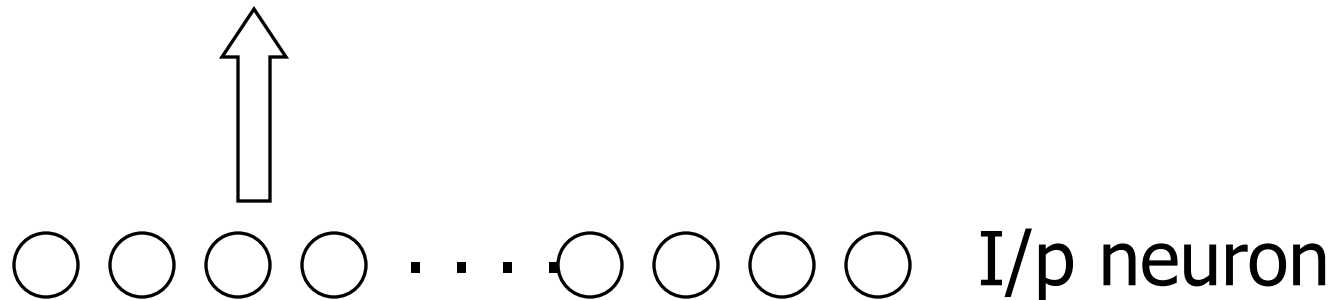
For auditory information processing

Character Recognition:

A, **A**, A, A, A

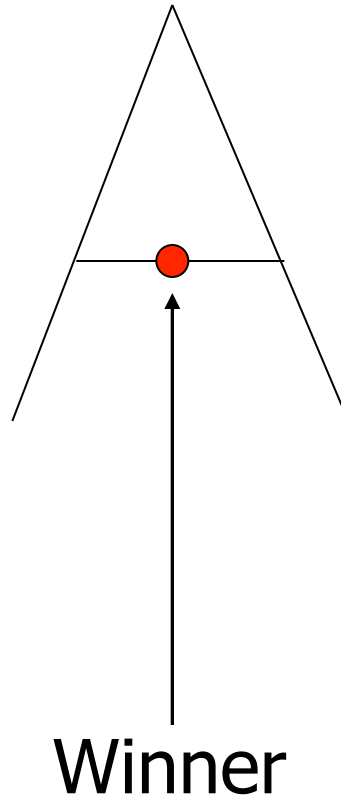


' O/p grid



Kohonen Net

- Self Organization or Kohonen network fires a group of neurons instead of a single one.
- The group “some how” produces a “picture” of the cluster.
- Fundamentally SOM is competitive learning.
- But weight changes are incorporated on a neighborhood.
- Find the winner neuron, apply weight change for the winner and its “neighbors”.



Neurons on the contour are the
"neighborhood" neurons.

Weight change rule for SOM

$$W^{(n+1)}_{P+\delta(n)} = W^{(n)}_{P+\delta(n)} + \eta^{(n)}_{P+\delta(n)} (I^{(n)}_{P+\delta(n)} - W^{(n)}_{P+\delta(n)})$$

Neighborhood: function of n

Learning rate: function of n

$\delta(n)$ is a decreasing function of n

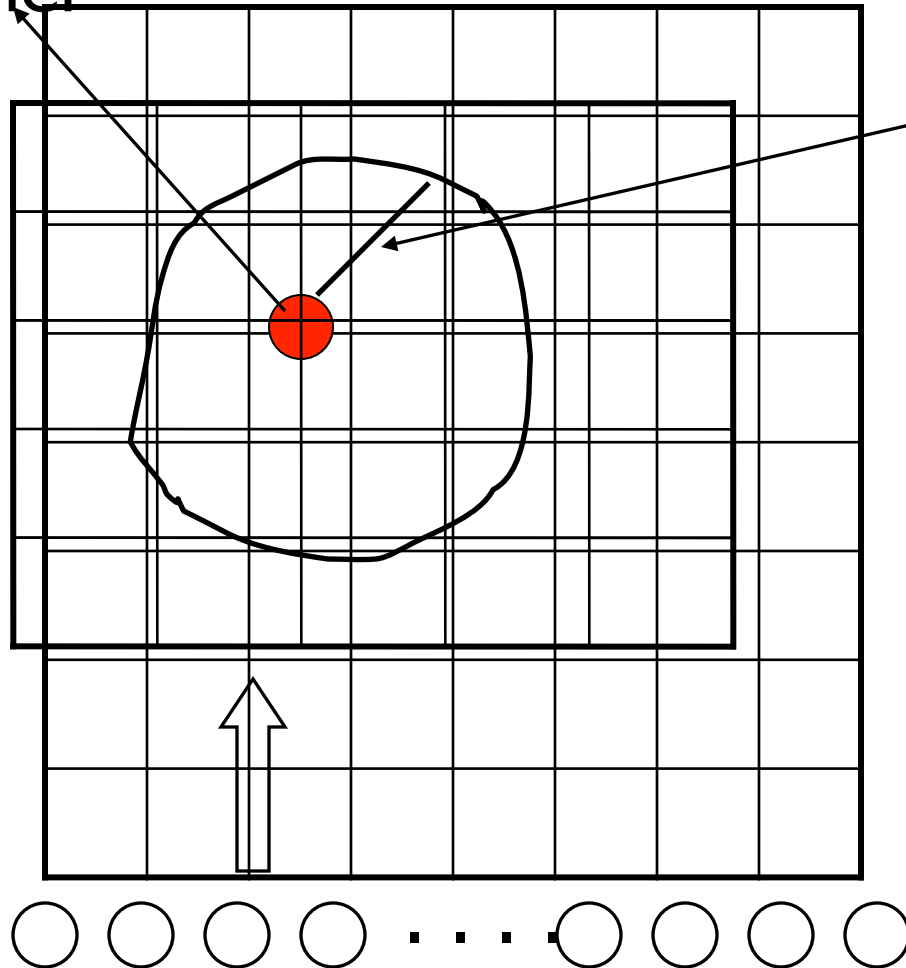
$\eta(n)$ learning rate is also a decreasing function of

n

$$0 < \eta(n) < \eta(n-1) \leq 1$$

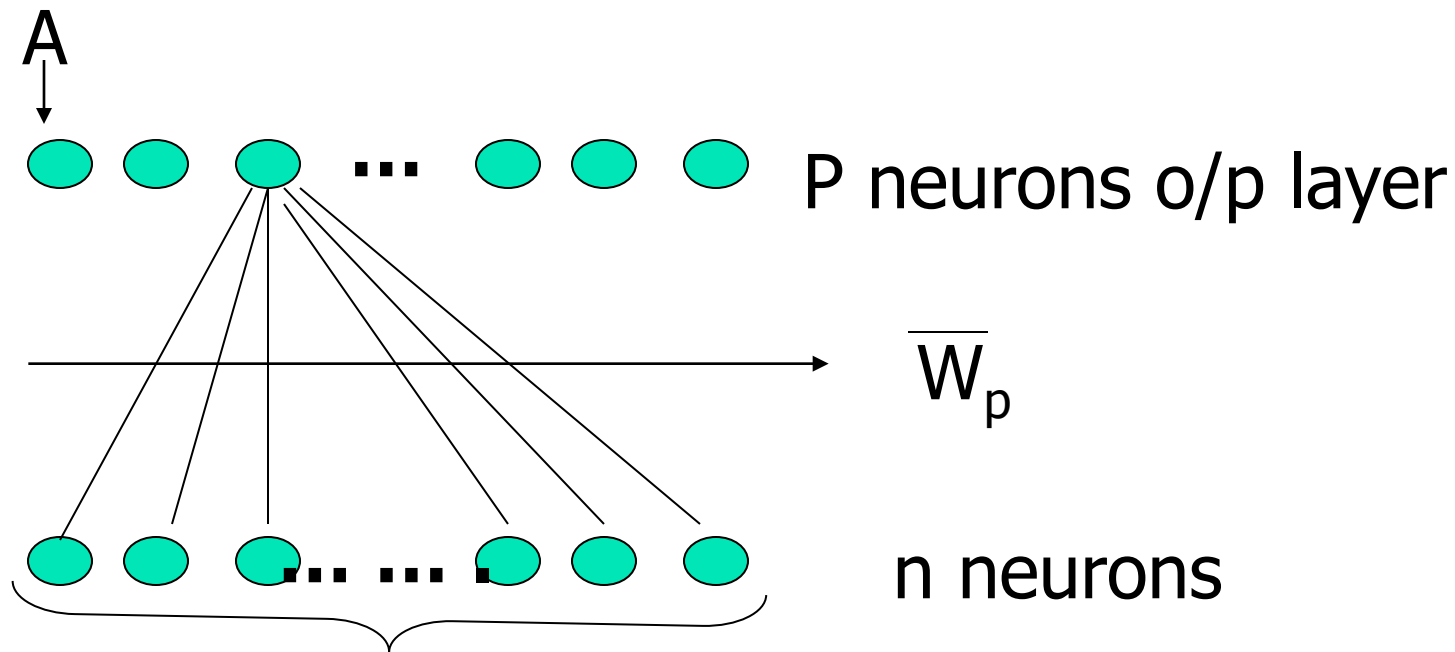
Pictorially

Winner



$\delta(n)$

Convergence for kohonen not proved except for uni-dimension



Clusters:

A : **A** A A

B : :

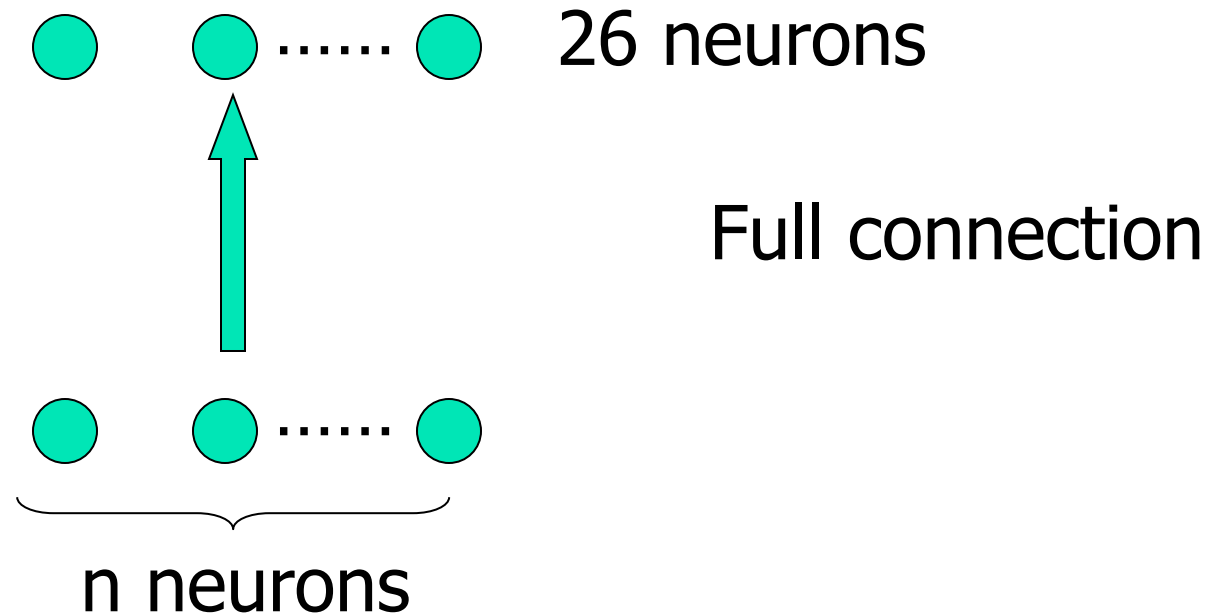
C : :

Clustering Algos

1. Competitive learning
2. K – means clustering
3. Counter Propagation

K – means clustering

K o/p neurons are required from the knowledge of K clusters being present.

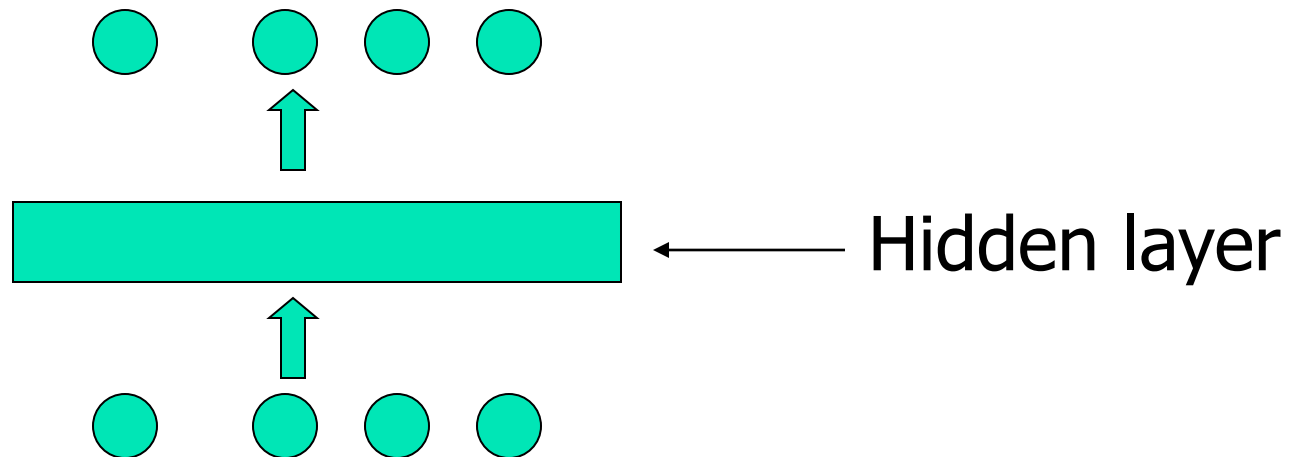


Steps

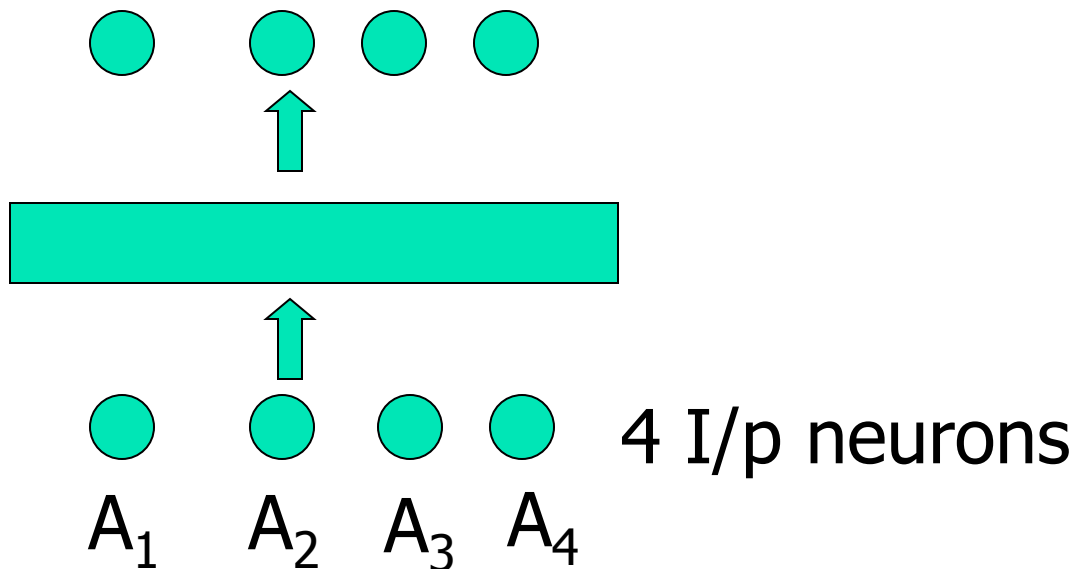
1. Initialize the weights randomly.
2. I^k is the vector presented at k^{th} iteration.
3. Find w^* such that
$$|w^* - I^k| < |w_j - I^k| \text{ for all } j$$
4. make $W^{*(\text{new})} = W^{*(\text{old})} + \eta (I^k - w^*)$.
- 5 $k \leftarrow k + 1$.
6. Go to 2 until the error is below a threshold.

Two part assignment

Supervised



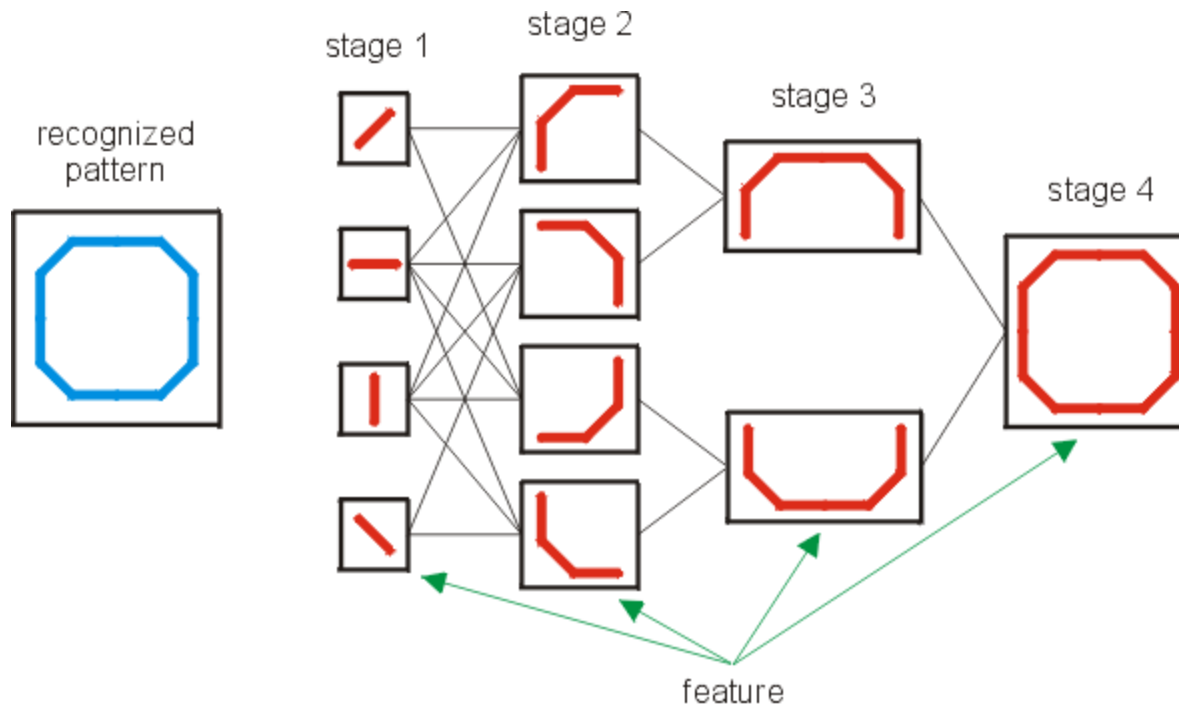
Cluster Discovery By SOM/Kohonen Net

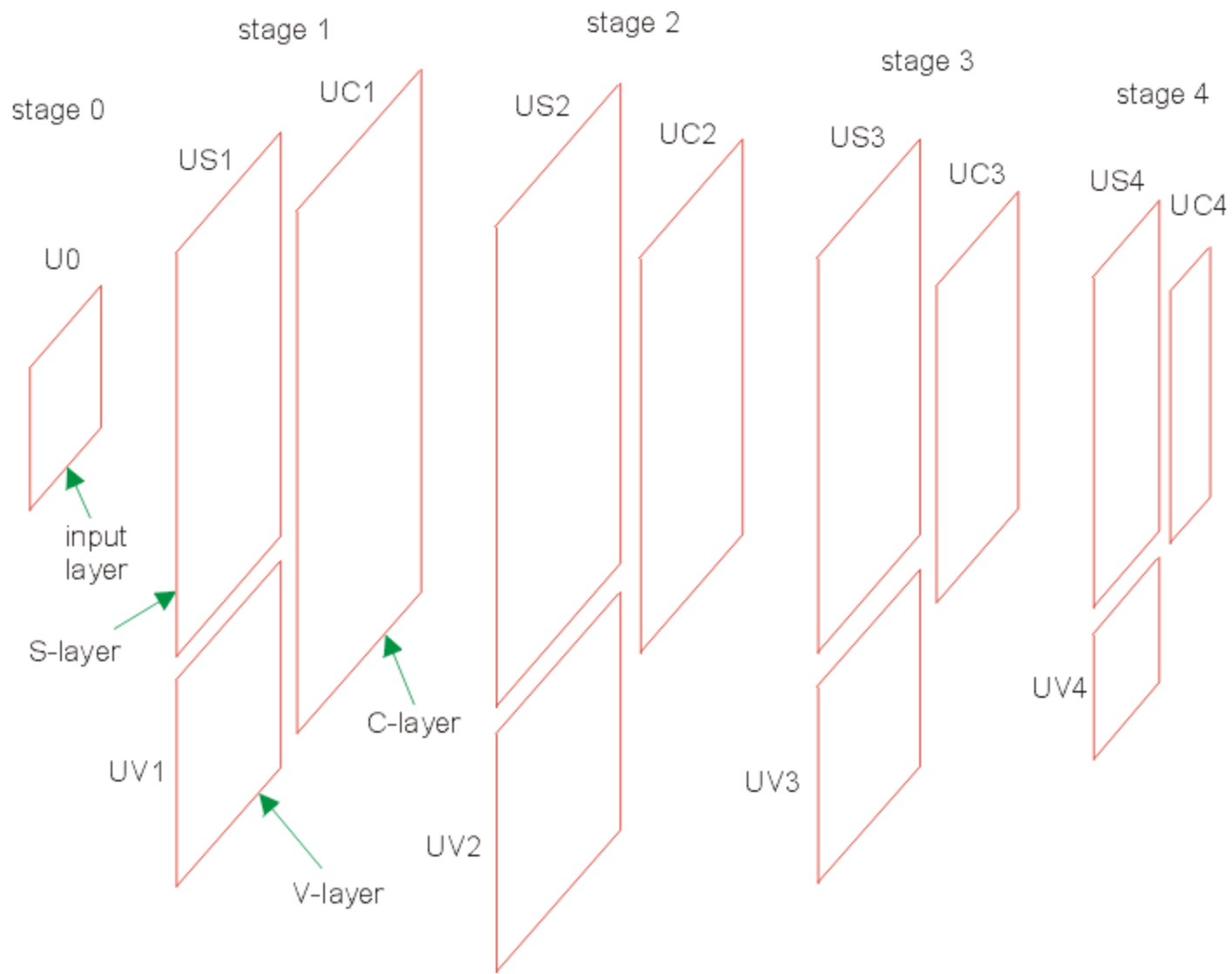


NeoCognitron

(Fukusima et. al., 1980)

Hierarchical feature extraction based





S-Layer

- Each **S-layer** in the neocognitron is intended for **extraction of features** from corresponding stage of hierarchy.
- Particular S-layers are formed by distinct number of S-planes. Number of these S-planes depends on the number of extracted features.

V-Layer

- Each **V-layer** in the neocognitron is intended for **obtaining of informations about average activity** in previous C-layer (or input layer).
- Particular V-layers are always formed by only one V-plane.

C-Layer

- Each **C-layer** in the neocognitron is intended for **ensuring of tolerance of shifts of features** extracted in previous S-layer.
- Particular **C-layers are formed by distinct number of C-planes**. Their number depends on the number of features extracted in the previous S-layer.

Hands on for afternoon

- Implement the binary adder on RNN
- Create a FF-BP POS tagger