

CS772: Deep Learning for Natural Language Processing (DL-NLP)

BPTT, Hopfield-net, LSTM

Pushpak Bhattacharyya

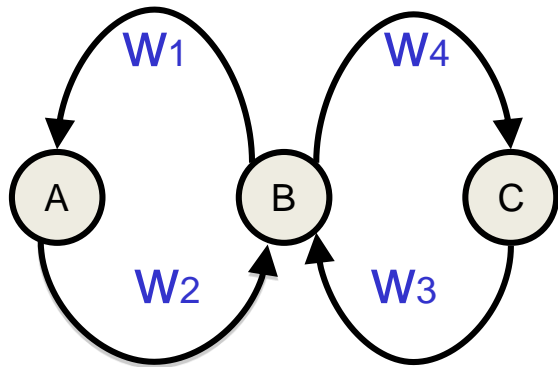
Computer Science and Engineering
Department

IIT Bombay

Week 7 of 14th Feb, 2022

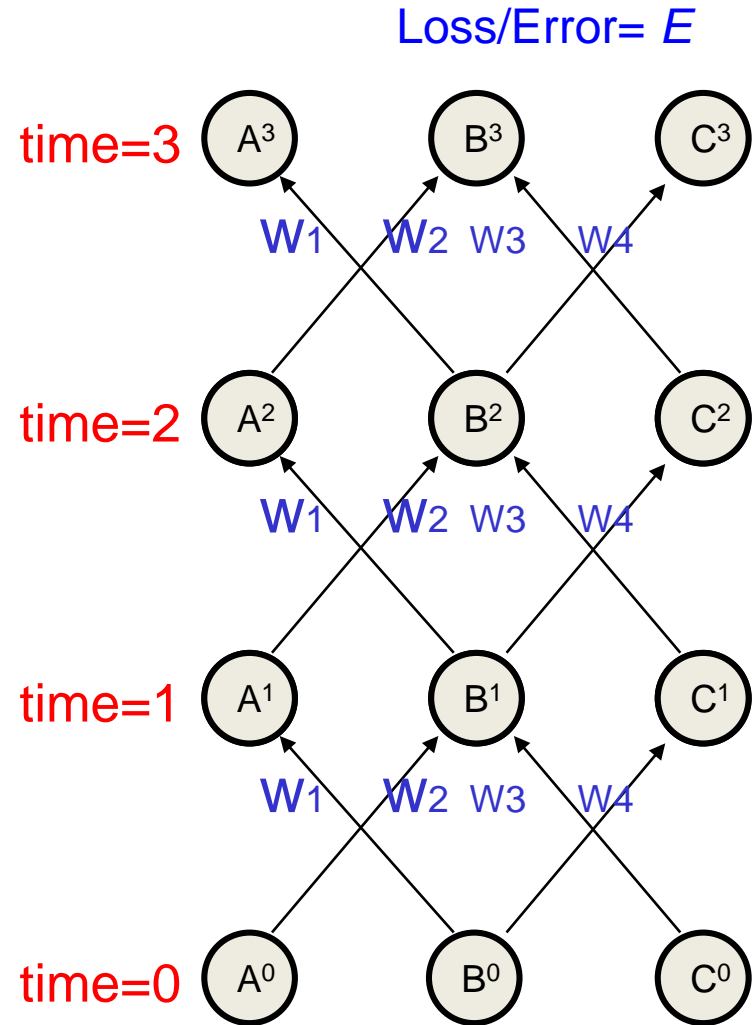
RNN-BPTT

The equivalence between feedforward nets and recurrent nets



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

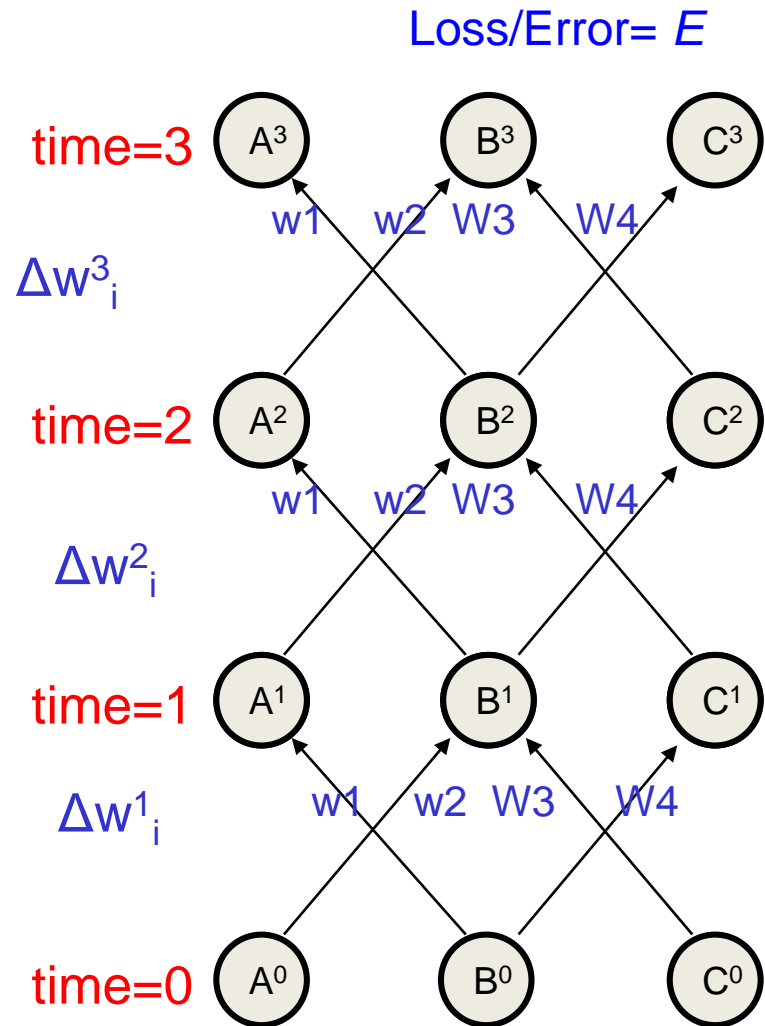
The recurrent net is just a layered net that keeps reusing the same weights.



BPTT illustration

$$\Delta w_i = \Delta w^3_i + \Delta w^2_i + \Delta w^1_i$$

Vanishing/Exploding
Gradient can strike!!!



An NLP example (1/2)

- The RNN in the previous slide can be used for POS tagging of trigrams.
- “Purchased Videocon Machine” → VBD
NNP NN
- We would know the loss value after this sequence
- Can apply gradient descent
- Obtain all Δw s
- Post them

An NLP example (2/2)

- Similarly for “The Blue Sky” → DT JJ NN
- This way was obtain Δ_{ws} for each input-output pattern
- At the end the RNN will learn Tri-gram POS tagging

BPTT important points

- 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.

Long word sequences

- The famous book by Charles Dickens “A Tale of Two Cities” starts the book with the famous sentence “This was the best of times, this was the worst of times....”
- The sentence has 119 words
- “It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair, we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way--in short, the period was so far like the present period that some of its noisiest authorities insisted on its being received, for good or for evil, in the superlative degree of comparison only.

The “best of times...” sentence

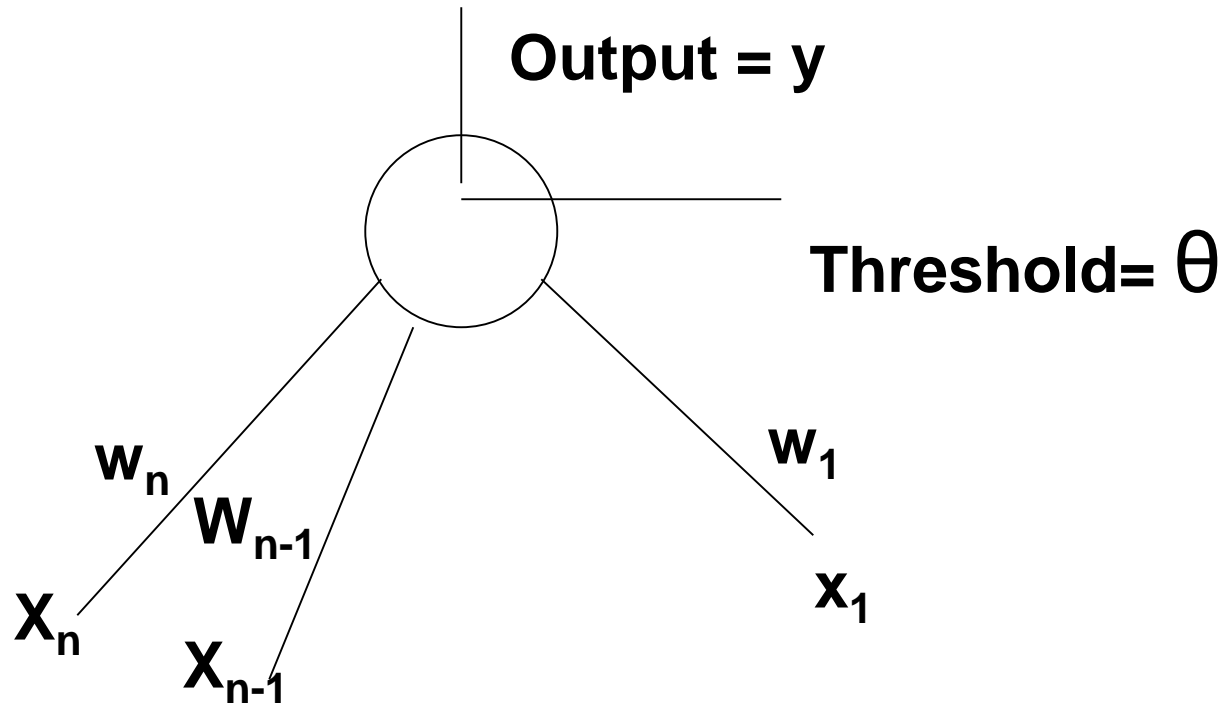
- Vanishing gradient will surely strike!!
- Exercise: give an example from NLP, where exploding gradient will strike!!

A Classical RNN: Hopfield Net

Hopfield net

- Inspired by associative memory which means memory retrieval is not by address, but by part of the data.
- Consists of
 - N neurons fully connected with symmetric weight strength $w_{ij} = w_{ji}$
- No self connection. So the weight matrix is 0-diagonal and symmetric.
- Each computing element or neuron is a linear threshold element with threshold = 0
- **Binary neurons with activations as +1 or -1**

Computation



Example

$$W_{12} = W_{21} = 5$$

$$W_{13} = W_{31} = 3$$

$$W_{23} = W_{32} = 2$$

$\theta=0$ for each neuron

At time $t=0$

$$s_1(t) = 1$$

$$s_2(t) = -1$$

$$s_3(t) = 1$$

Unstable state: Neuron 1 will flip.

A stable pattern is called an
attractor for the net.

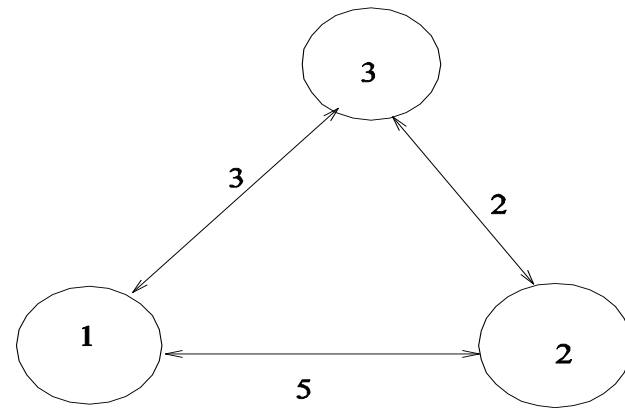


Figure: An example Hopfield Net

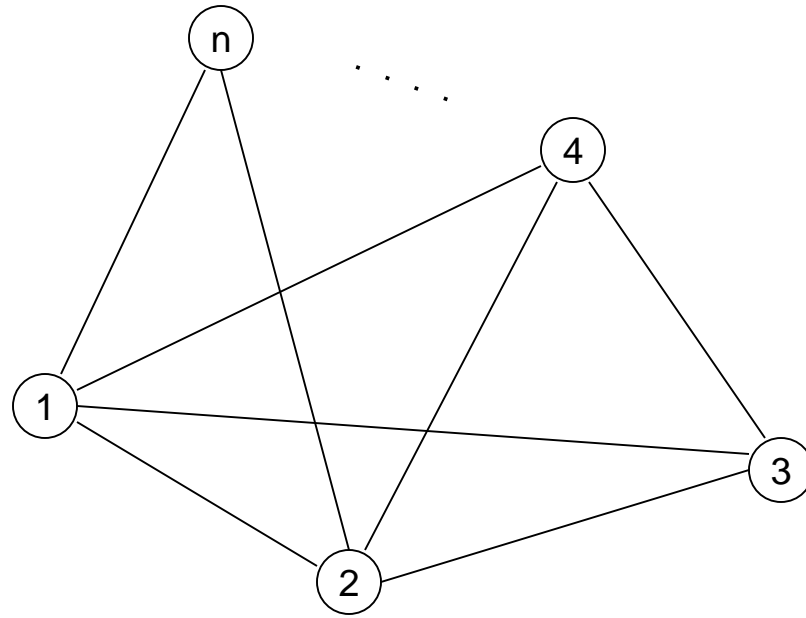
Stability

- Asynchronous mode of operation: at any instant a randomly selected neuron compares the net input with the threshold.
- In the *synchronous* mode of operation all neurons update themselves simultaneously at any instant of time.
- Since there are feedback connections in the Hopfield Net the question of *stability* arises. At every time instant the network evolves and finally settles into a stable state.
- How does the Hopfield Net function as *associative* memory ?
- One needs to store or stabilize a vector which is the memory element.

Energy consideration

- Stable patterns correspond to minimum energy states.
- Energy at state $\langle x_1, x_2, x_3, \dots, x_n \rangle$
- $$E = -1/2 \sum_j \sum_{j \neq i} w_{ji} x_i x_j$$
- Change in energy always comes out to be negative in the asynchronous mode of operation. Energy *always* decreases.
- Stability ensured.

Hopfield Net is a fully connected network



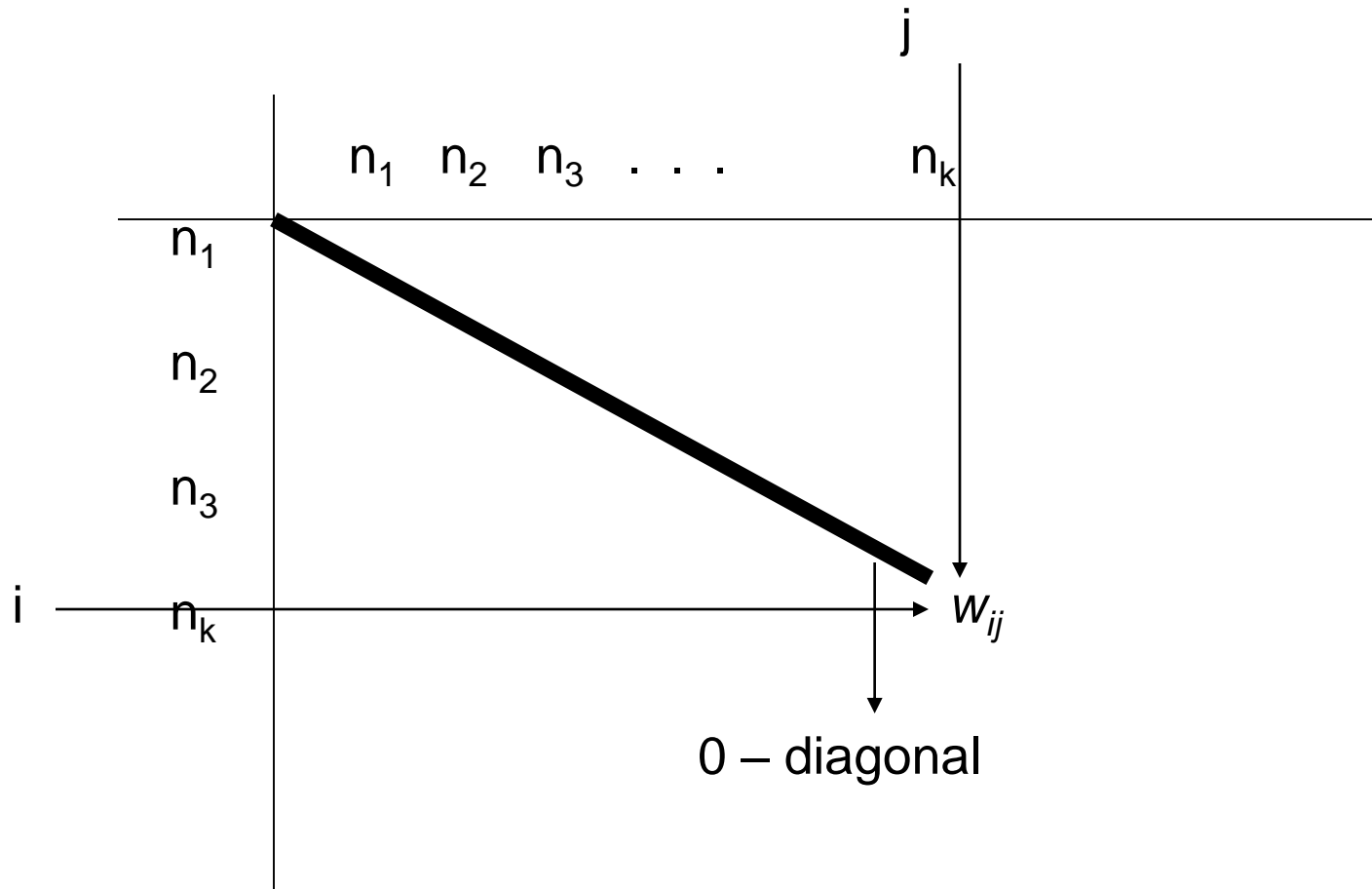
- i^{th} neuron is connected to $(n-1)$ neurons

Concept of Energy

- Energy at state s is given by the equation:

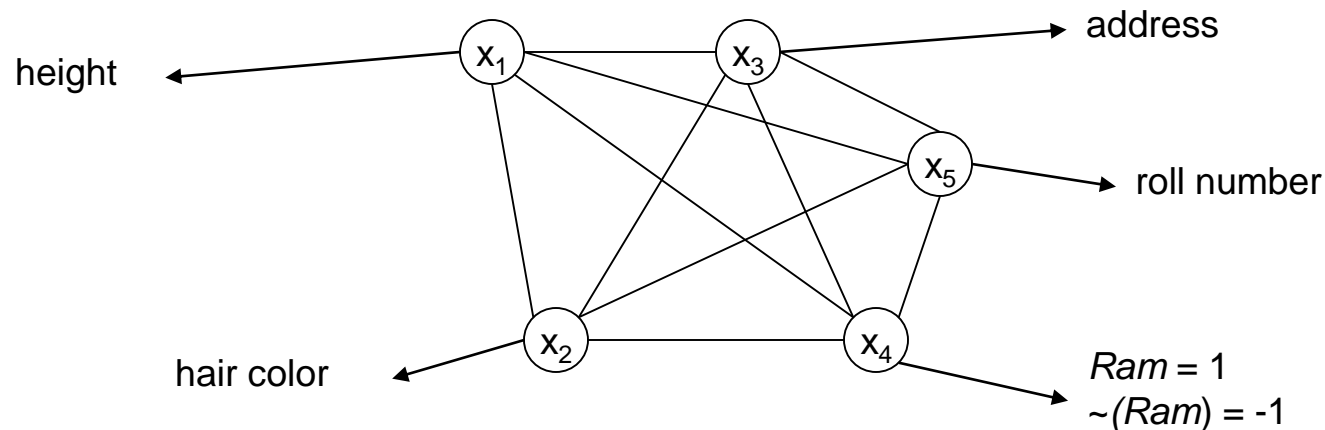
$$E(s) = -\left[w_{12}x_1x_2 + w_{13}x_1x_3 + \dots + w_{1n}x_1x_n \right. \\ \left. + w_{23}x_2x_3 + \dots + w_{2n}x_2x_n + \right. \\ \left. \vdots \right. \\ \left. + w_{(n-1)n}x_{(n-1)}x_n \right]$$

Connection matrix of the network, 0-diagonal and symmetric

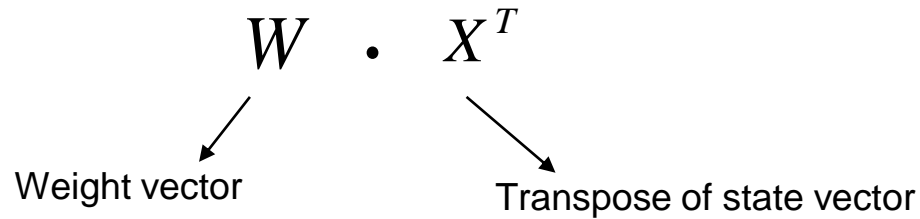


State Vector

- Binary valued vector: value is either 1 or -1, $X = \langle x_n \quad x_{n-1} \quad \dots \quad x_3 \quad x_2 \quad x_1 \rangle$
- e.g. Various attributes of a student can be represented by a state vector



Relation between weight vector W and state vector X



For example, in figure 1,
At time $t = 0$, state of the neural network is:
 $s(0) = \langle 1, -1, 1 \rangle$ and corresponding vectors are as shown.

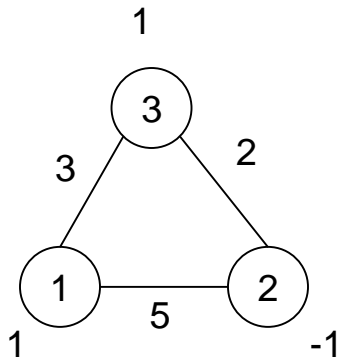


Fig. 1

$$W = \begin{bmatrix} 0 & 5 & 3 \\ 5 & 0 & 2 \\ 3 & 2 & 0 \end{bmatrix} \quad X^T = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$

$$W \cdot X^T = \begin{bmatrix} 0 & 5 & 3 \\ 5 & 0 & 2 \\ 3 & 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$

$W.X^T$ gives the inputs to the neurons at the next time instant

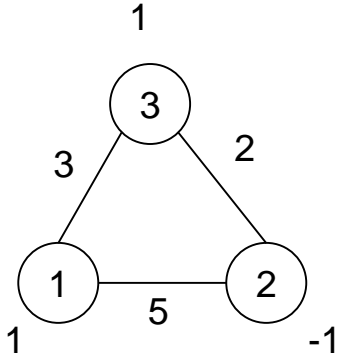
$$W \cdot X^T = \begin{bmatrix} 0 & 5 & 3 \\ 5 & 0 & 2 \\ 3 & 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} -2 \\ 7 \\ 1 \end{bmatrix}$$

$$\text{sgn}(W.X^T) = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}$$

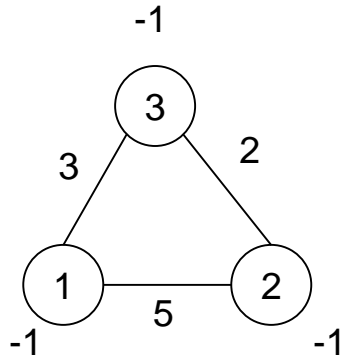
This shows that the n/w will change state

Energy Consideration



At time $t = 0$, state of the neural network is:
 $s(0) = \langle 1, -1, 1 \rangle$

- $E(0) = -[(5 * 1 * -1) + (3 * 1 * 1) + (2 * -1 * 1)] = 4$



The state of the neural network under stability is $\langle -1, -1, -1 \rangle$

$$E(\text{stable state}) = - - [(5 * -1 * -1) + (3 * -1 * -1) + (2 * -1 * -1)] = -10$$

State Change

- $s(0) = \langle 1, -1, 1 \dots \rangle$
- $s(1) = \text{compute by comparing and summing}$
- $x_1(t=1) = \text{sgn}[\sum_{j=2}^n w_{1j} x_j]$
 $= 1 \text{ if } \sum_{j=2}^n w_{1j} x_j > 0$
 $= -1 \text{ otherwise}$

Theorem

- In the asynchronous mode of operation, the energy of the Hopfield net always decreases.
- Proof:

$$E(t_1) = -\left[w_{12}x_1(t_1)x_2(t_1) + w_{13}x_1(t_1)x_3(t_1) + \dots + w_{1n}x_1(t_1)x_n(t_1) \right. \\ \left. + w_{23}x_2(t_1)x_3(t_1) + \dots + w_{2n}x_2(t_1)x_n(t_1) + \right. \\ \left. \vdots \right. \\ \left. + w_{(n-1)n}x_{(n-1)}(t_1)x_n(t_1) \right]$$

Proof: *note that only neuron 1 changes state*

$$\begin{aligned}\Delta E &= E(t_2) - E(t_1) \\ &= -\left\{ \left[w_{12}x_1(t_2)x_2(t_2) + w_{13}x_1(t_2)x_3(t_2) + \dots + w_{1n}x_1(t_2)x_n(t_2) \right] \right. \\ &\quad \left. - \left[w_{12}x_1(t_1)x_2(t_1) + w_{13}x_1(t_1)x_3(t_1) + \dots + w_{1n}x_1(t_1)x_n(t_1) \right] \right\} \\ &= -\sum_{j=2}^n w_{1j} \left[x_1(t_2) \cdot x_j(t_2) - x_1(t_1) \cdot x_j(t_1) \right]\end{aligned}$$

Since only neuron 1 changes state, $x_j(t_1)=x_j(t_2)$, $j=2, 3, 4, \dots, n$, and hence

$$= \sum_{j=2}^n \left[w_{1j} \cdot x_j(t_1) \right] \left[x_1(t_1) - x_1(t_2) \right]$$

The Hopfield net has to “converge” in the asynchronous mode of operation

- As the energy E goes on decreasing, it has to hit the bottom, since the weight and the state vector have finite values.
- That is, the Hopfield Net has to converge to an energy minimum.
- Hence the Hopfield Net reaches stability.

Probabilistic state change

- If the state change happens with a probability, we have what is called a Boltzmann Machine
- State change is controlled by Boltzmann Distribution
- Randomness controlled by “Temperature”

Boltzmann Machine

Comparative Remarks

Feed forward n/w with BP	Hopfield net	Boltzmann m/c
<p>Mapping device: (i/p pattern --> o/p pattern), <i>i.e.</i> Classification</p>	<p>Associative Memory + Optimization device</p>	<p>Constraint satisfaction. (Mapping + Optimization device)</p>
<p>Minimizes total sum square error, typically</p>	<p>Energy</p>	<p>Entropy (<u>Kullback–Leibler divergence</u>)</p>

Comparative Remarks (contd.)

Feed forward n/w with BP	Hopfield net	Boltzmann m/c
Deterministic neurons	Deterministic neurons	Probabilistic neurons
Learning to associate i/p with o/p <i>i.e.</i> equivalent to a function	Pattern	Probability Distribution

Boltzmann Machine

- Hopfield net
- Probabilistic neurons
- Energy expression = $-\sum_i \sum_{j>i} w_{ij} x_i x_j$
where x_i = activation of i^{th} neuron
- Used for optimization
- Central concern is to ensure global minimum
- Based on simulated annealing

Meaning of “Learning” a Probability Distribution

- If we have standard distributions like normal, binomial, poisson etc., learning the distribution means estimating the parameters like mean, standard deviation etc.
- A feedforward NN with softmax, also learns a probability distribution, though there is no closed form expression per se for the learnt distribution

For example

Say, a FFNN with softmax learns POS tags of trigrams

Then the softmax values at the outer layer give

$$P(T|W)$$

Where W is the word sequence, T is the tag sequence

This learning is ensured by fixing the weight values of the FNNN

So $P(\text{VBD NNP NN}|\text{purchased Videocon machine})$ will be high after training, and $P(\text{JJ NN VB}|\text{-same-})$ will be low

Comparative Remarks (contd.)

Feed forward n/w with BP	Hopfield net	Boltzmann m/c
Can get stuck in local minimum (Greedy approach)	Local minimum possible	Can come out of local minimum
Credit/blame assignment (consistent with Hebbian rule)	Activation product (consistent with Hebbian rule)	Probability and activation product (consistent with Hebbian rule)

Explaining Hebb's Rule

- Law of nature
- Comes from neurophysiology
- States, “if two connected neurons are both in a state of excitation, the change in weight will be so as to maintain the state of excitation”
- This means two connected neurons in states +1 each, the change in weight will be positive

Theory of Boltzmann m/c

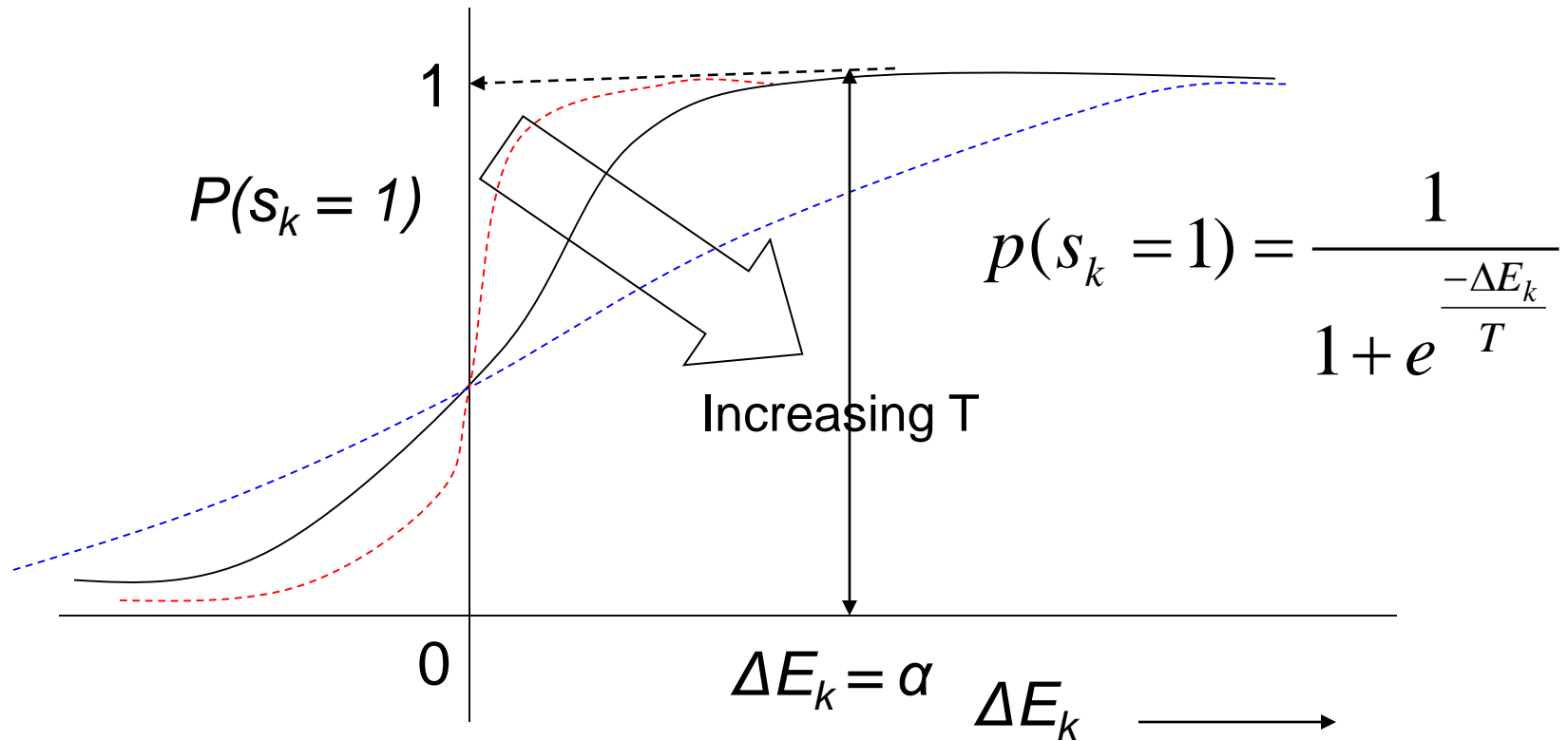
For the m/c the computation means the following: at any time instant, make the state of the k^{th} neuron (s_k) equal to 1, with probability:

$$P(s_k = 1) = \frac{1}{1 + e^{\frac{-\Delta E_k}{T}}}$$

ΔE_k = change in energy of the m/c when the k^{th} neuron changes state

T = temperature which is a parameter of the m/c

Theory of Boltzmann m/c (contd.)



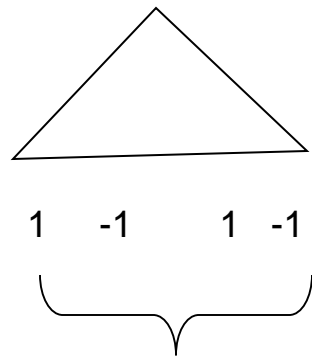
Theory of Boltzmann m/c (contd.)

$$\begin{aligned}\Delta E_k &= (E_f - E_i) |_k \\ &= (S_i - S_k) |_k \cdot \sum_{j \neq k} w_{kj} S_j\end{aligned}$$

We observe:

1. The higher the temperature, lower is $P(S_k=1)$
2. at $T = \textit{infinity}$, $P(S_k=1) = P(S_k=0) = 0.5$, equal chance of being in state 0 or 1. Completely random behavior
3. If $T \rightarrow 0$, then $P(S_k=1) \rightarrow 1$
4. The derivative is proportional $P(S_k=1) \cdot (1 - P(S_k=1))$

Consequence of the form of $P(S_k=1)$



N - bits

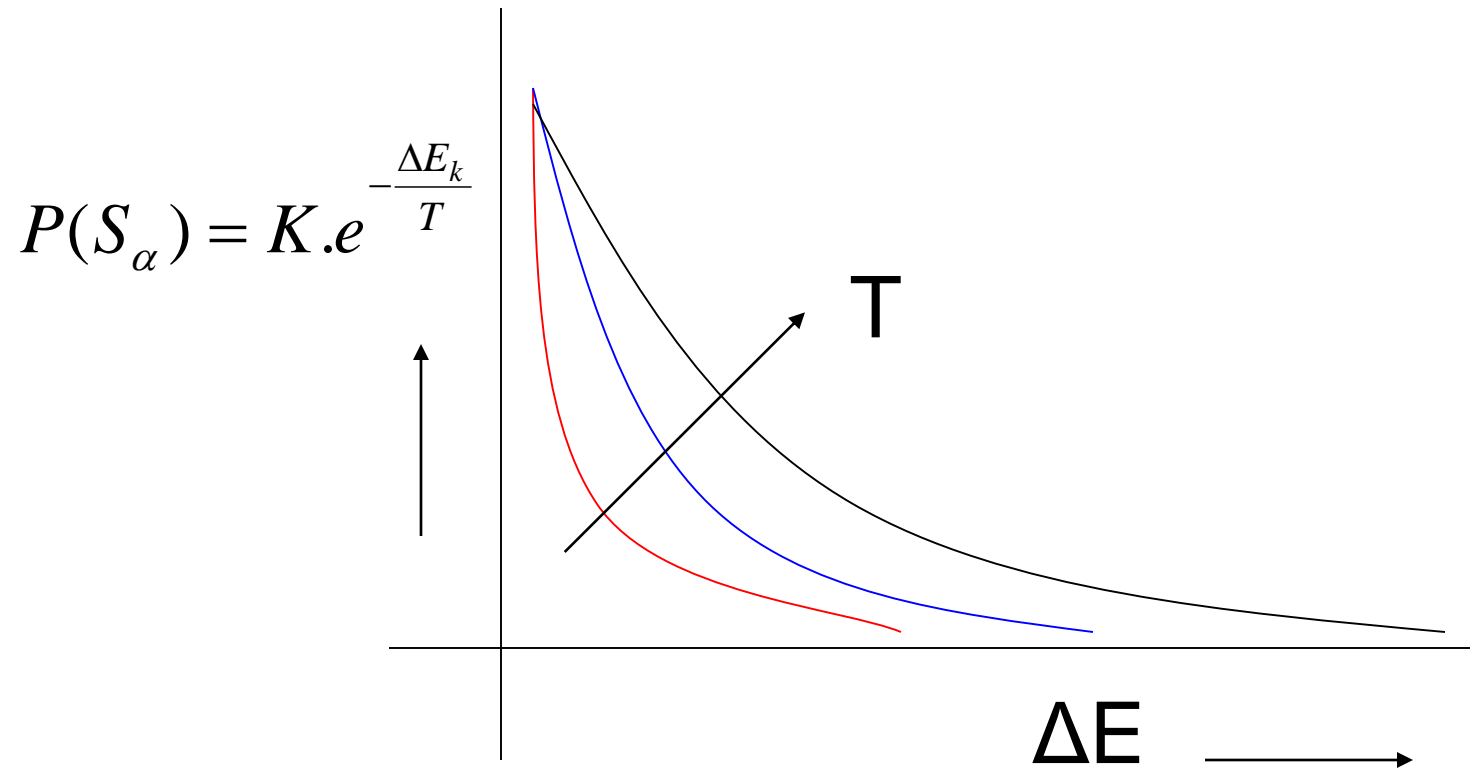
$$P(S_{\alpha}) = K \cdot e^{-\frac{E_{\alpha}}{T}}$$

Probability Distribution called as Boltzmann Distribution

$P(S_{\alpha})$ is the probability of the state S_{α}

Local “sigmoid” probabilistic behavior leads to global Boltzmann Distribution behaviour of the n/w

(Ex: try to prove the above)



Ratio of state probabilities

Normalizing (α and β are states)

$$\begin{aligned} P(S_\alpha) &= K \cdot e^{-\frac{E(S_\alpha)}{T}} \\ &= \frac{e^{-\frac{E(S_\alpha)}{T}}}{\sum_{\beta} e^{-\frac{E(S_\beta)}{T}}} \end{aligned}$$

Learning a probability distribution

- Digression: Estimation of a probability distribution Q by another distribution P
- $D = \text{deviation} = \sum_{\text{sample space}} Q \ln Q/P$
- $D \geq 0$, which is a required property (just like sum square error ≥ 0)

To prove $D \geq 0$

Lemma: $\ln(1/x) \geq (1-x)$

Let, $x = 1 / (1 - y)$

$\ln(1 - y) = -[y + y^2/2 + y^3/3 + y^4/4 + \dots]$

Proof (contd.)

$$\begin{aligned}(1 - x) &= 1 - 1/(1 - y) \\ &= -y(1 - y)^{-1} \\ &= -y(1 + y + y^2 + y^3 + \dots) \\ &= -[y + y^2 + y^3 + y^4 + \dots]\end{aligned}$$

But, $y + y^2/2 + y^3/3 + \dots \leq y + y^2 + y^3 + \dots$

Lemma proved.

Proof (contd.)

$$D = \sum Q \ln Q / P$$

$$\geq \sum_{\text{over the sample space}} Q(1 - P/Q)$$

$$= \sum (Q - P)$$

$$= \sum Q - \sum P$$

$$= 1 - 1 = 0$$