

RECURRENT NEURAL NETWORKS

— HOPFIELD NETWORK.

Multilayer (feedforward) NN \rightarrow relations are Algebraic.

$$o = \Gamma(Wx)$$

Recurrent (Feedback) NN \rightarrow relations are dynamic.

$$\text{ODE} \rightarrow \dot{x} = f(x)$$

$$\text{Difference} \rightarrow x(k+1) = f(x(k))$$

Initial condition \rightarrow input to the system.

How can this be useful? \rightarrow asymptotic behaviour.

$$\begin{array}{l} x(t) \rightarrow \\ x(k) \rightarrow \end{array} \text{ a bounded set } S \quad \left(\begin{array}{l} t \rightarrow \infty \\ k \rightarrow \infty \end{array} \right)$$

$S \rightarrow$ an attractor. (attracts solutions)

$S \rightarrow$ a point, a periodic solution, a set, chaotic attractor ...

\rightarrow information is stored in attractor.

$S \rightarrow$ "point" $\{x^*\}$

$$x(t) \rightarrow x^* \Rightarrow \dot{x} \rightarrow 0, \quad \dot{x} = f(x)$$

$$\Rightarrow \boxed{f(x^*) = 0}$$

(Equilibrium point of f .
= DC operating point)

$$x(k) \rightarrow x^* \Rightarrow x(k+1) \Rightarrow x^*, \quad x(k+1) = f(x(k))$$

$$\Rightarrow \boxed{f(x^*) = x^*}$$

Fixed point of f .

Example:

PROBLEM:

min $E(x)$

\Rightarrow gradient descent
(back propagation)

$$\boxed{x(k+1) = x(k) - h \frac{dE(x)}{dx}}$$

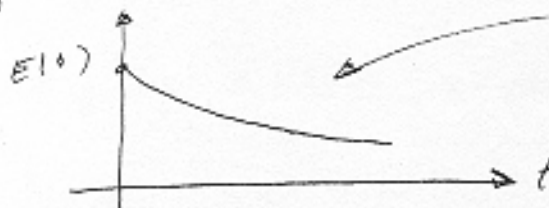
$$\Rightarrow \frac{x(k+1) - x(k)}{h} = -\frac{dE}{dx} \Rightarrow$$

$$\boxed{\dot{x} = -\frac{dE}{dx}}$$

claim: solution $x(t) \rightarrow S$ since $E(x)$ is locally minimum
(Gradient network)

$$\frac{dE}{dt} = \frac{dE}{dx} \cdot \frac{dx}{dt} = \frac{dE}{dx} \cdot \left(-\frac{dE}{dx}\right) = -\left(\frac{dE}{dx}\right)^2 \leq 0$$

\Rightarrow $E(t)$



\Rightarrow Along the solutions, $E(t)$ decreases

\Rightarrow under certain conditions

$x(t) \rightarrow S$ since $E(x)$ is
locally minimum.

Example.

$$E(x) = \frac{1}{2} (x - x^*)^2 \Rightarrow \min E(x) \Rightarrow \arg \min E(x) = x^*$$

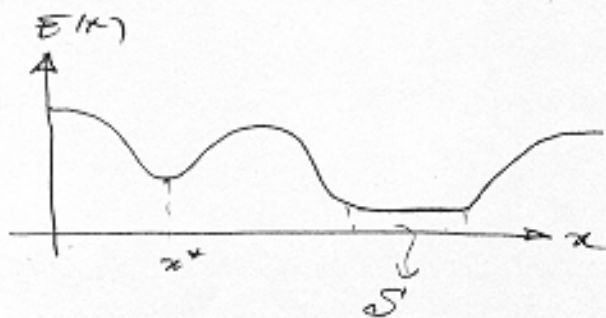
$$\frac{dE}{dx} = (x - x^*) \Rightarrow \boxed{\dot{x} = -(x - x^*)} \quad \text{Gradient eqn.}$$

$$y = x - x^* \Rightarrow \dot{y} = \dot{x} = -y \Rightarrow y(t) = e^{-t} y(0) \rightarrow 0$$

$$\Rightarrow \boxed{x(t) \rightarrow x^*}$$

$$\text{Equilibrium pt.} \Rightarrow f(x) = -(x - x^*) \quad \boxed{f(x^*) = 0}$$

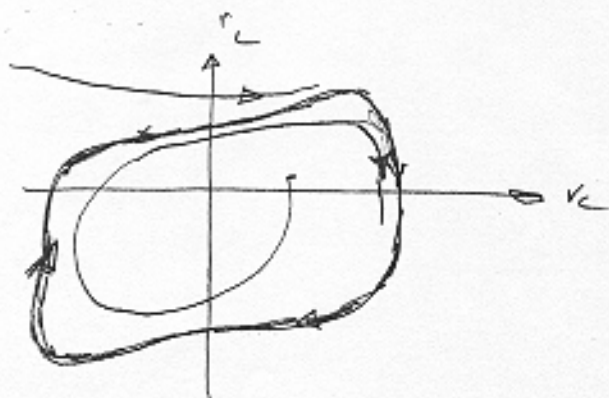
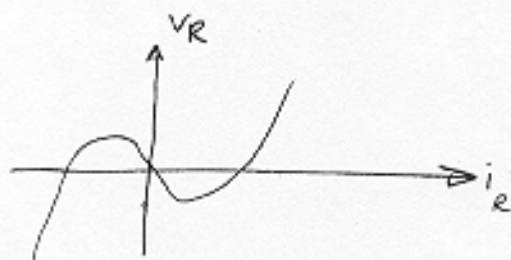
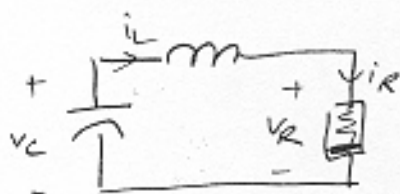
In general:



$$\boxed{\dot{x} = -\frac{dE}{dx}}$$

$$x(t) \rightarrow S \quad \text{where} \quad \frac{dE}{dx} = 0$$

Attractor can be a periodic solution: (oscillators)



IN DISCRETE TIME (1985)

$$\begin{aligned}x(k+1) &= x(k) - h \cdot \frac{dE}{dx} = x(k) - h \cdot (x(k) - a) \\ &= (1-h)x(k) + ah\end{aligned}$$

For convergence $1-h = \rho$ should satisfy $|\rho| < 1$

Check convergence.

$$\begin{aligned}y &= x - a \\ \Rightarrow y(k+1) &= x(k+1) - a = (1-h)x(k) + ah - a = (1-h)x(k) - (1-h)a \\ &= (1-h)[x(k) - a] = (1-h)y(k) = \rho y(k)\end{aligned}$$

$$y(k) = \rho^k y(0) \rightarrow 0 \quad \text{if } |\rho| < 1$$

$$x(k) \rightarrow 0$$

Check equilibrium.

$$(1-h)x^* + ah = x^* \Rightarrow -hx^* + ah = 0 \Rightarrow \boxed{x^* = a}$$

Stability of attractor: if $x(0)$ is sufficiently close to the attractor ($d(x(0), S) < \epsilon$), do the solution $x(t)$ $\rightarrow S$? If yes, S is stable.

If no \Rightarrow not stable.

In stable case $\Rightarrow S$ has a basin of attraction

$B = \{ x(0) \mid x(t) \rightarrow S \}$. (All initial conditions such that $x(t) \rightarrow S$)

Pattern storage and recognition

idea: Design a network such that patterns become stable attractors

Recognition: $x(0) \neq S \Rightarrow$ view $x(0)$ as distorted pattern, stability \Rightarrow if $x(0)$ is sufficiently close to S $x(t) \rightarrow S$.

Example



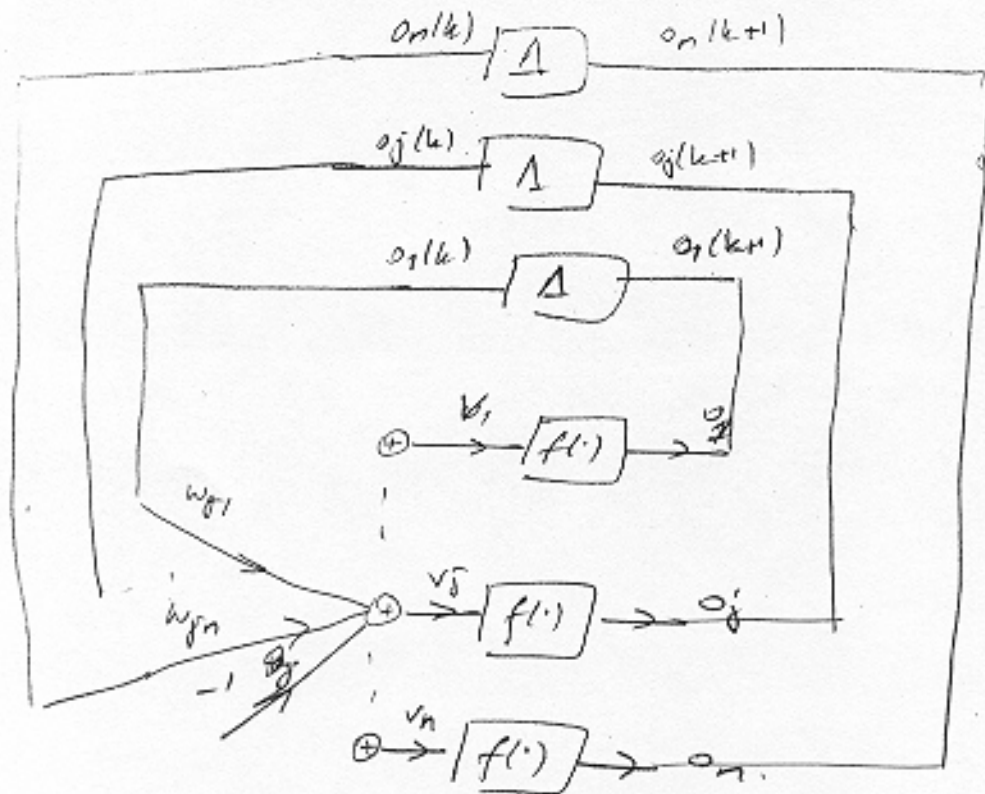
$$\Rightarrow x^* = \begin{pmatrix} 0 \\ \vdots \\ 1 \end{pmatrix}$$

vector corresponding to pattern

$$\Rightarrow x(k+1) = f(x(k)) \quad (x^0 = f(x))$$

such that x^* is a stable equilibrium (fixed point)

DISCRETE TIME HOPFIELD NETWORK



Componentwise:

$$v_j^i(k) = w_{j1} \cdot o_1(k) + w_{j2} \cdot o_2(k) + \dots + w_{jn} \cdot o_n(k) - \theta_j$$

$$o_j^i(k+1) = f(v_j^i(k)) \quad \left| \quad \begin{matrix} j=1, 2, \dots, n \\ i=1, 2, \dots, n \end{matrix} \right.$$

$$k = 0, 1, 2, \dots, \infty$$

Matrix form:

$$\begin{cases} v(k) = W o(k) - \theta \\ o(k+1) = \Gamma(v(k)) \end{cases}$$

$$v = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix}, \quad o = \begin{pmatrix} o_1 \\ o_2 \\ \vdots \\ o_n \end{pmatrix}, \quad W = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{j1} & w_{j2} & \dots & w_{jn} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nn} \end{pmatrix}, \quad \theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{pmatrix}$$

for $f(\cdot) \Rightarrow \text{sgn} \cdot \text{fn}$

$$f(v) = \begin{cases} +1 & v \geq 0 \\ -1 & v < 0 \end{cases}$$

1) In Hopfield's original work $\Rightarrow w = w^T$
(symmetric). Not essential for storing patterns,
but important in stability.

2) In Hopfield's original work $\Rightarrow w_{ii} = 0$
(no self feedback in i th neuron). Important for
asynchronous update stability.

Synchronous update: All neuron outputs are updated
at the same time.

Asynchronous updates: Only one neuron output is updated
at a time. When this update takes place, most recent
neuron outputs are used.

\rightarrow Suppose that we have $o_1(k), \dots, o_n(k)$.

\rightarrow update o_1 :
$$o_1(k+1) = \text{sgn}(w_{11} o_1(k) + w_{12} o_2(k) + \dots + w_{1n} o_n(k) - \theta_1)$$

\rightarrow update o_2 :
$$o_2(k+1) = \text{sgn}(w_{21} \underline{o_1(k+1)} + w_{22} o_2(k) + \dots + w_{2n} o_n(k) - \theta_2)$$

\rightarrow update o_3 :
$$o_3(k+1) = \text{sgn}(w_{31} \underline{o_1(k+1)} + w_{32} \underline{o_2(k+1)} + \dots + w_{3n} o_n(k) - \theta_3)$$

\rightarrow update o_n :
$$o_n(k+1) = \text{sgn}(w_{n1} o_1(k+1) + w_{n2} o_2(k+1) + \dots + w_{n,n-1} o_{n-1}(k+1) + w_{nn} o_n(k) - \theta_n)$$

\rightarrow Go to the beginning of a new update

Example-

$$W = \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix} \quad \theta = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$o(0) = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

Synch. update

$$o(1) = \tau (W o(0) - \theta) = \tau \begin{pmatrix} -1 \\ 1 \end{pmatrix} = \underline{\underline{\begin{pmatrix} -1 \\ 1 \end{pmatrix}}}$$

asynch. update

$$o_1(1) = \text{sgn}(1 \cdot o_1(0) + 2 \cdot o_2(0)) = -1$$
$$o_2(1) = \text{sgn}(2 \cdot o_1(1) + 1 \cdot o_2(0)) = -1$$

ie.

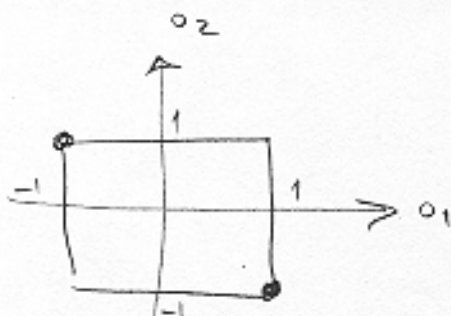
$$\begin{pmatrix} 1 \\ -1 \end{pmatrix} \rightarrow \begin{pmatrix} -1 \\ 1 \end{pmatrix} \quad \text{Synch. update.}$$

$$\begin{pmatrix} 1 \\ -1 \end{pmatrix} \rightarrow \begin{pmatrix} -1 \\ -1 \end{pmatrix} \rightarrow \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

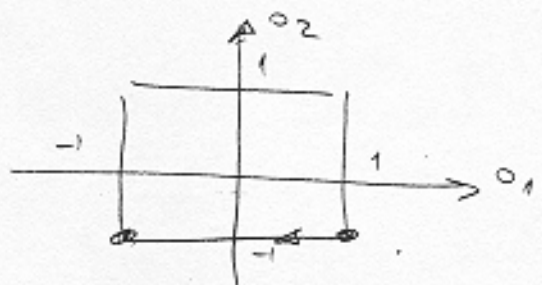
first component update second component update

Synch. update : $o(k+1)$ and $o(k)$ may differ by more than 1 Hamming distance

Asynch. update - At any intermediate step, the iterations differ by one Hamming distance.



Synch.



asynch.

- Relation with Hopfield Network:

$$\begin{cases} v(k) = W o(k) - \theta \\ o(k+1) = \Gamma(v(k)) \end{cases}$$

* Computational Energy of Hopfield Network:

$$E = -\frac{1}{2} o^T W o + \theta^T o$$

quadratic form.

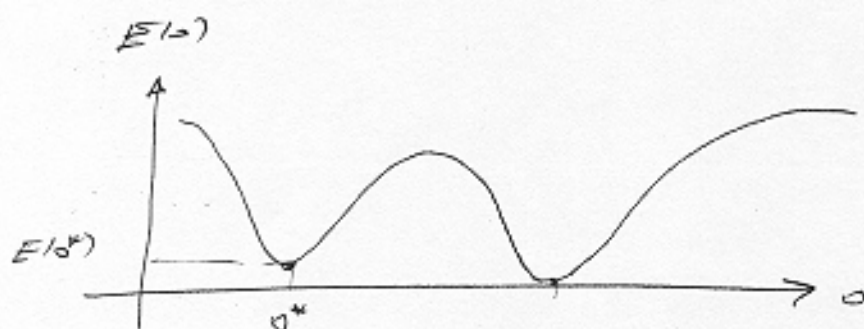
$$E = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_{ji} o_i o_j + \sum_{i=1}^n \theta_i o_i$$

claim: Hopfield Network solves the following problem:

$$\arg \min_o E(o)$$

i.e. start with $o(0)$, iterate; then $o(k) \rightarrow o^*$

where $E(o^*)$ is a local minimum of $E(o)$.



Some Linear Algebra

1) Let $A = -A^T$ (skew symmetric matrix).

$$\Rightarrow x^T A x = 0, \quad \forall x.$$

Proof: $(x^T A x)^T = x^T A^T x = -x^T A x$

$$\Rightarrow x^T A x = 0$$

$$2) \quad A = \underbrace{\frac{1}{2}(A+A^T)}_V + \underbrace{\frac{1}{2}(A-A^T)}_Q = V + Q$$

$$V = V^T \quad ((A+A^T)^T = A^T + A)$$

$$Q = -Q^T \quad ((A-A^T)^T = A^T - A = -(A-A^T))$$

i.e. Any matrix is the sum of a symmetric and a skew symmetric matrices. (This decomposition is unique).

$$3) \quad x^T A x = x^T (V+Q)x = x^T V x \quad (x^T Q x = 0).$$

\Rightarrow For quadratic forms, only the symmetric part of the matrix is important.

4) Gradient of a function:

$$f(x): \mathbb{R}^n \rightarrow \mathbb{R}.$$

$$f(x + \Delta x) = f(x) + (\Delta x)^T \cdot \nabla f(x) + h.o.t.$$

$$\nabla f(x): \text{Gradient} = \frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$