

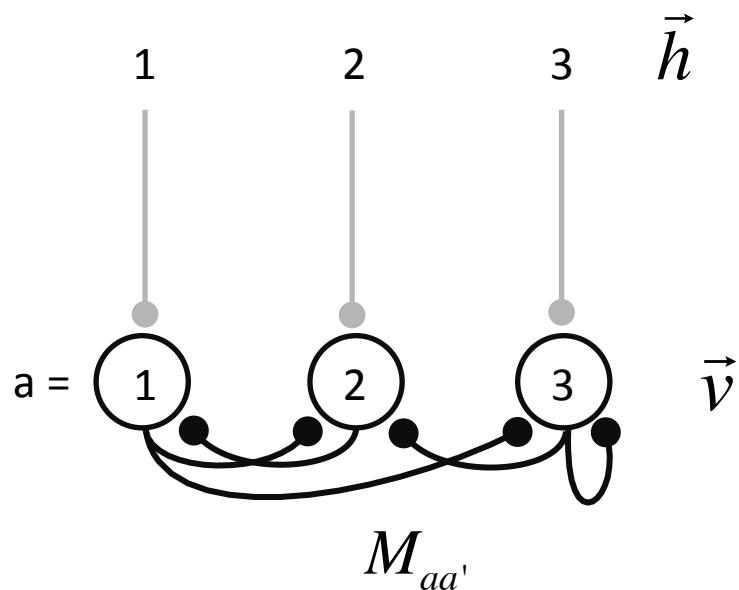
INTRODUCTION TO NEURAL COMPUTATION

Prof. Michale Fee
MIT BCS 9.40 — 2018

Lecture 20
Hopfield Networks

RECURRENT NETWORKS

- We have been considering the case where there are connections between different neurons in the output layer
- Develop an intuition for how recurrent networks respond to their inputs
- Examine computations performed by recurrent networks
- Use all the powerful linear algebra tools we have developed

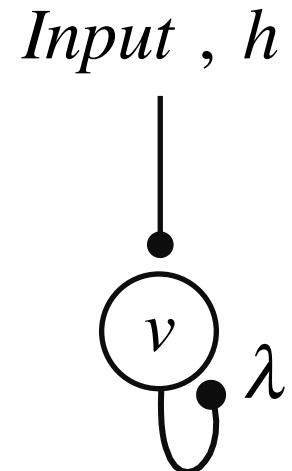


RECURRENT NETWORKS

- Simplest recurrent network – a neuron with an autapse

$$\tau_n \frac{dv}{dt} = -v + h + \lambda v$$

$$\tau_n \frac{dv}{dt} = -(1-\lambda)v + h$$

- We examined three cases:

$$\lambda < 1$$

$$\lambda = 1$$

$$\lambda > 1$$

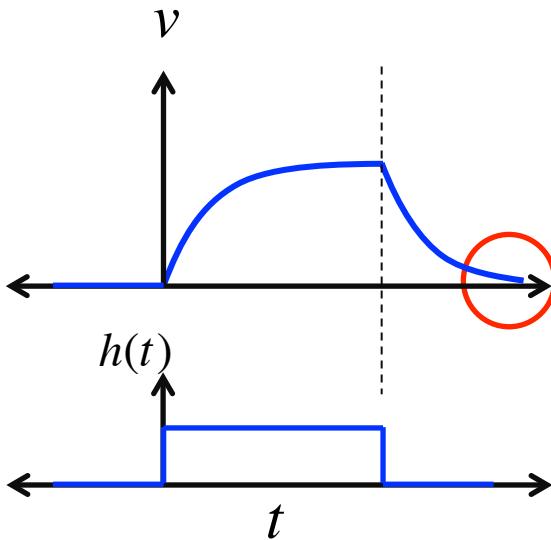
RECURRENT NETWORKS

$$\tau_n \frac{dv}{dt} = -(1-\lambda)v + h$$

- The behavior of the network depends critically on λ

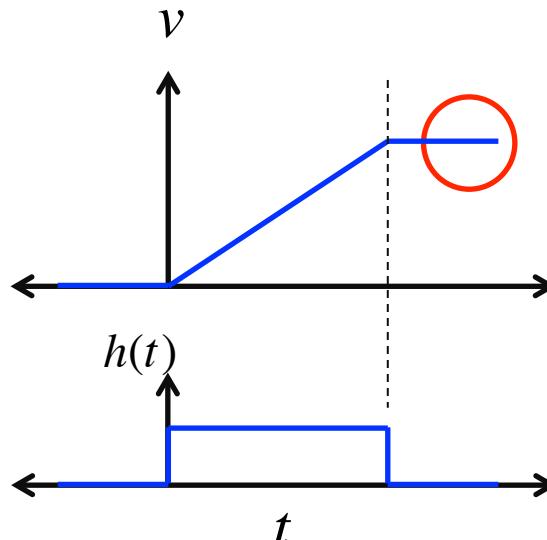
$$\lambda < 1$$

Exponential relaxation



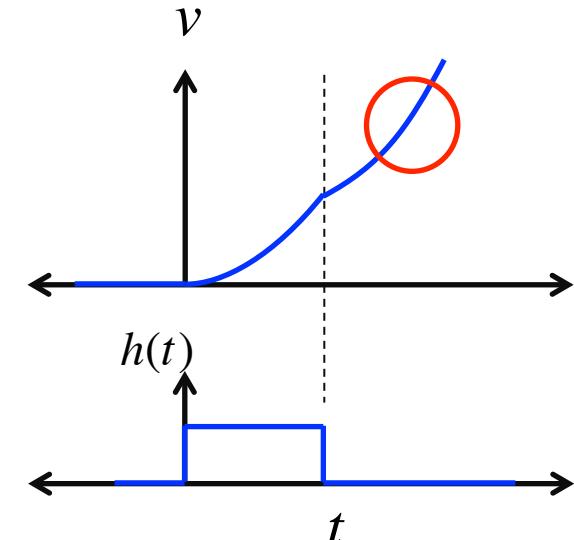
$$\lambda = 1$$

Integration



$$\lambda > 1$$

Exponential growth



With zero input...
relaxation back to zero

With zero input...
persistent activity!

MEMORY!

LEARNING OBJECTIVES FOR LECTURE 20

- Recurrent networks with lambda greater than one
 - Attractors
- Winner-take-all networks
- Attractor networks for long-term memory (Hopfield model)
- Energy landscape
- Hopfield network capacity

LEARNING OBJECTIVES FOR LECTURE 20

- Recurrent networks with lambda greater than one
 - Attractors
- Winner-take-all networks
- Attractor networks for long-term memory (Hopfield model)
- Energy landscape
- Hopfield network capacity

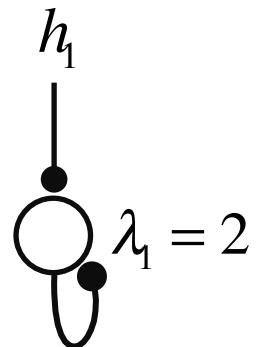
RECURRENT NETWORKS

- Networks with $\lambda \geq 1$ have memory!

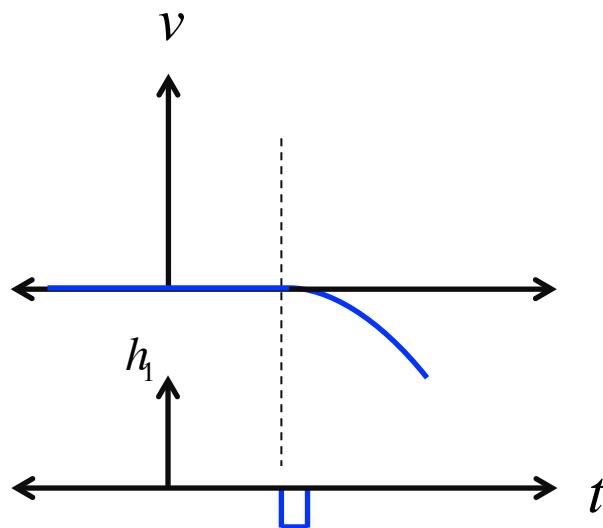
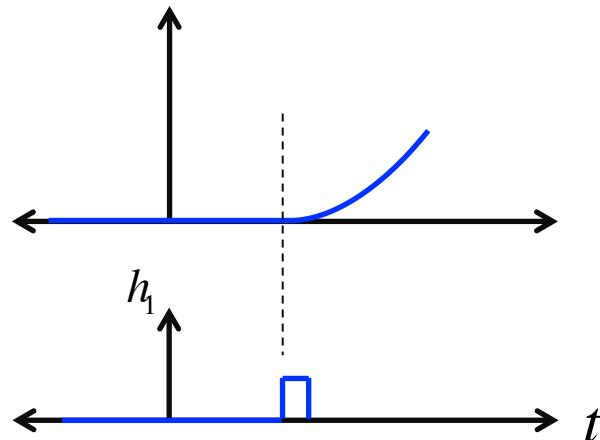
$$\tau_n \frac{dv}{dt} = (\lambda_1 - 1)v + h(t)$$

$$\tau_n \frac{dv}{dt} = v$$

$$v(t) = 0$$

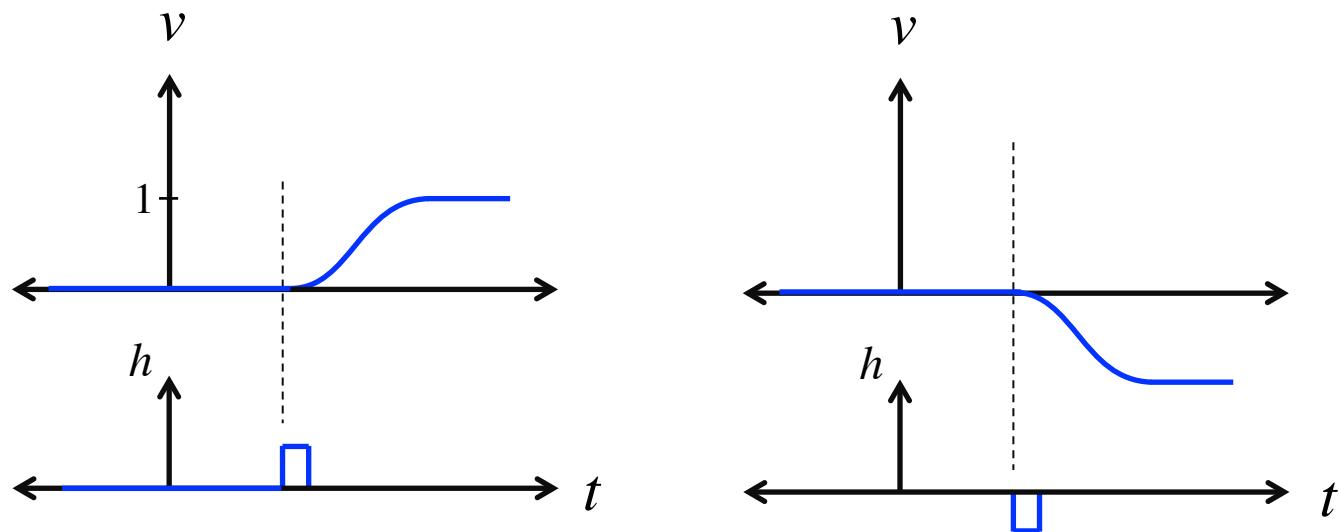
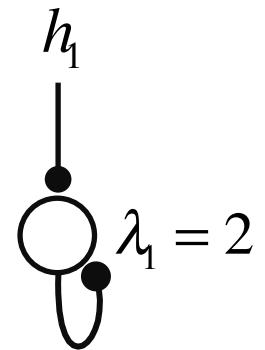
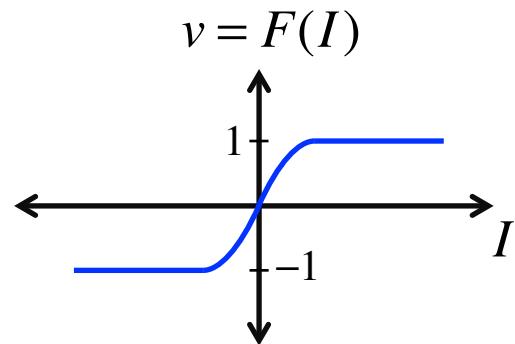


- With zero input, zero is an 'unstable fixed point' of the network v



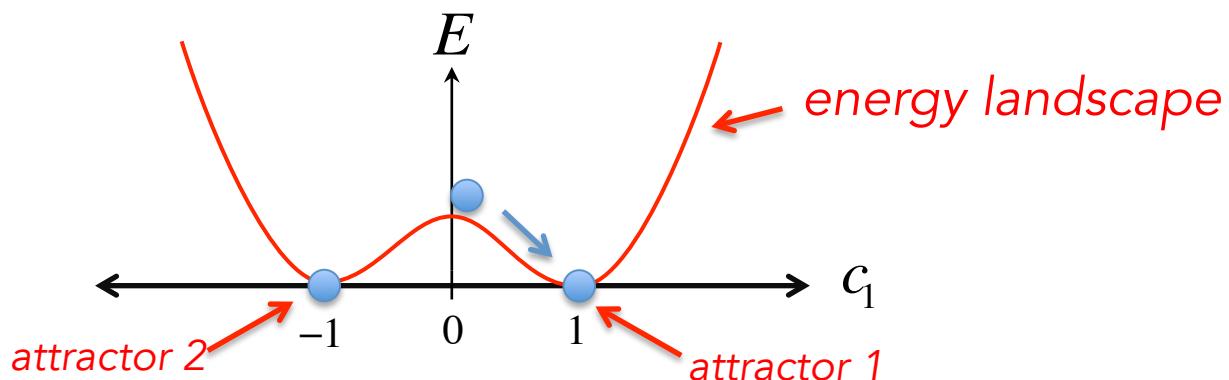
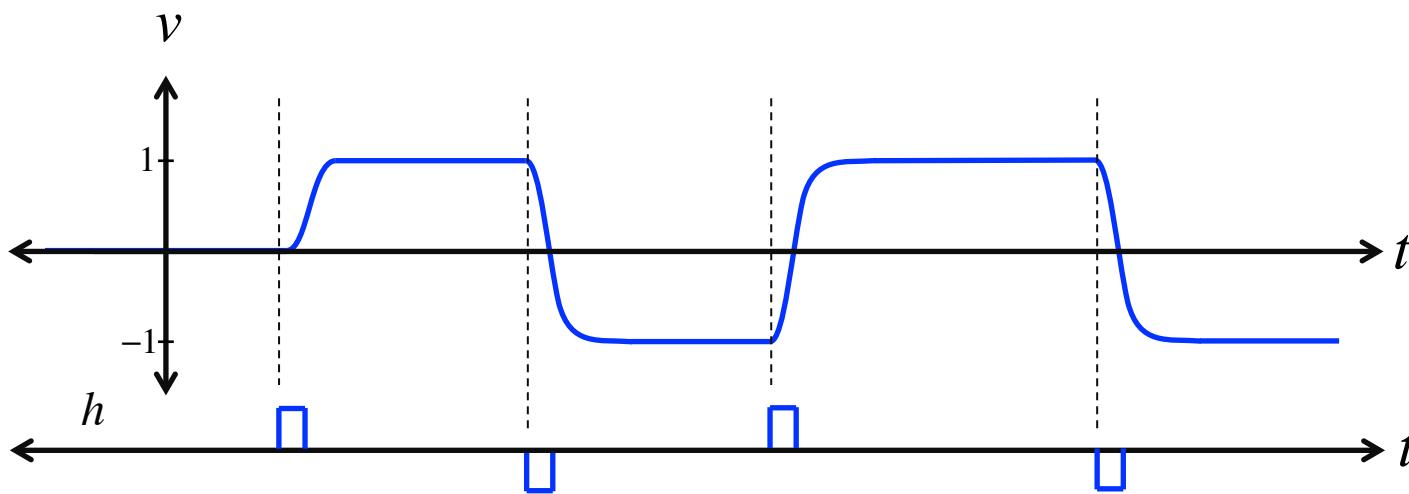
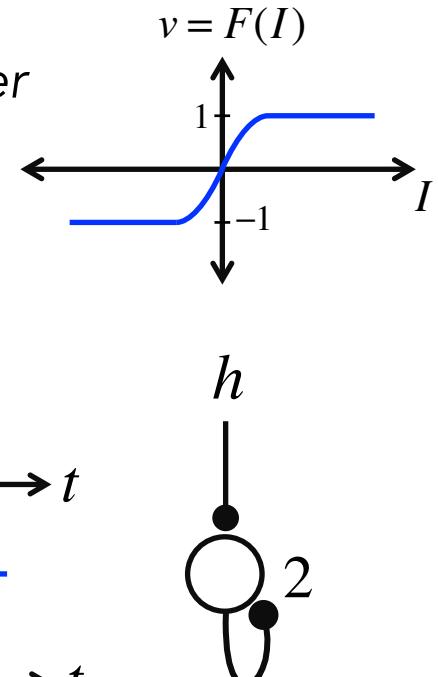
RECURRENT NETWORKS

- Add a saturating activation function $F(x)$



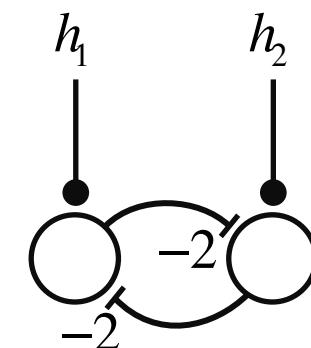
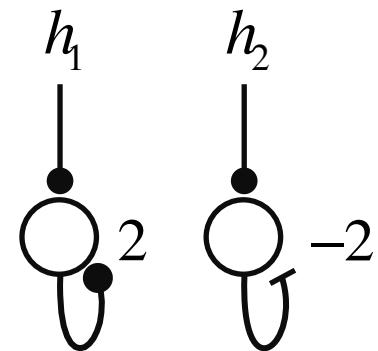
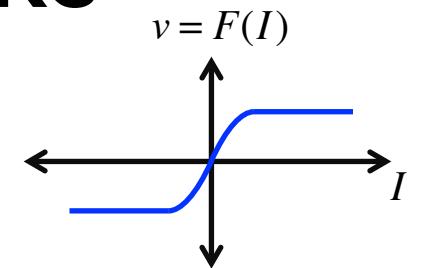
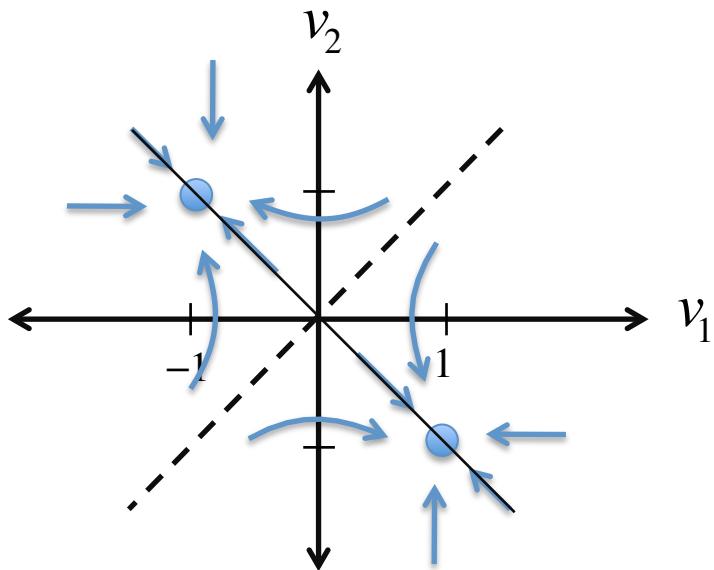
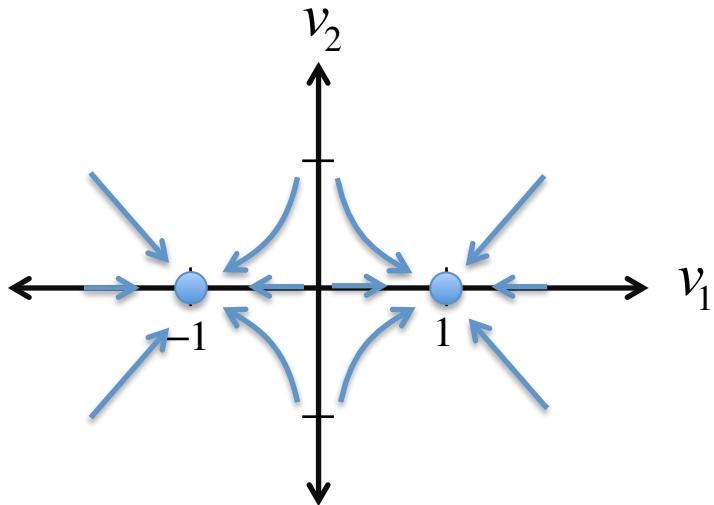
RECURRENT NETWORKS

- Saturating activation function plus eigenvalues greater than 1 lead to stable states other than zero!



RECURRENT NETWORKS

- Two-neuron network that has two attractors



LEARNING OBJECTIVES FOR LECTURE 20

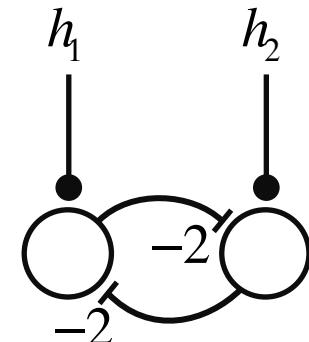
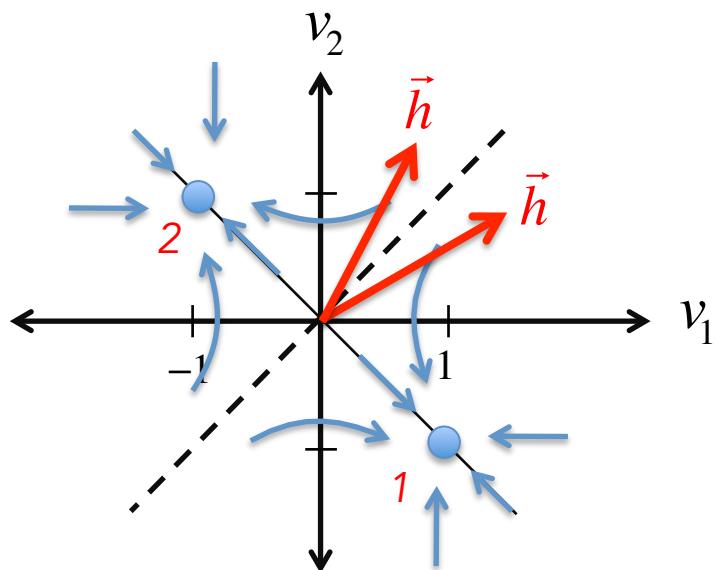
- Recurrent networks with lambda greater than one
 - Attractors
- Winner-take-all networks
- Attractor networks for long-term memory (Hopfield model)
- Energy landscape
- Hopfield network capacity

WINNER-TAKE-ALL NETWORK

- Implements decision making

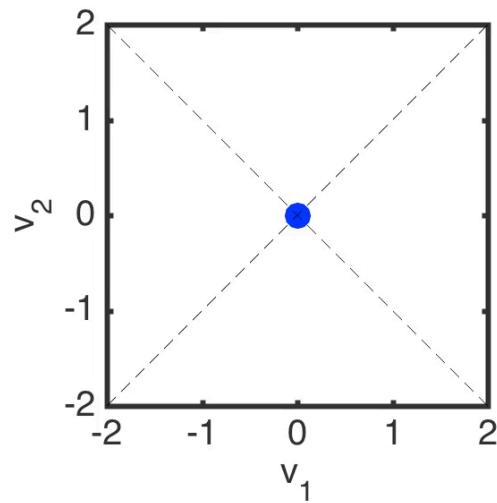
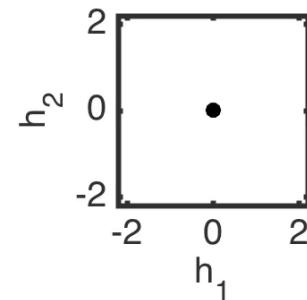
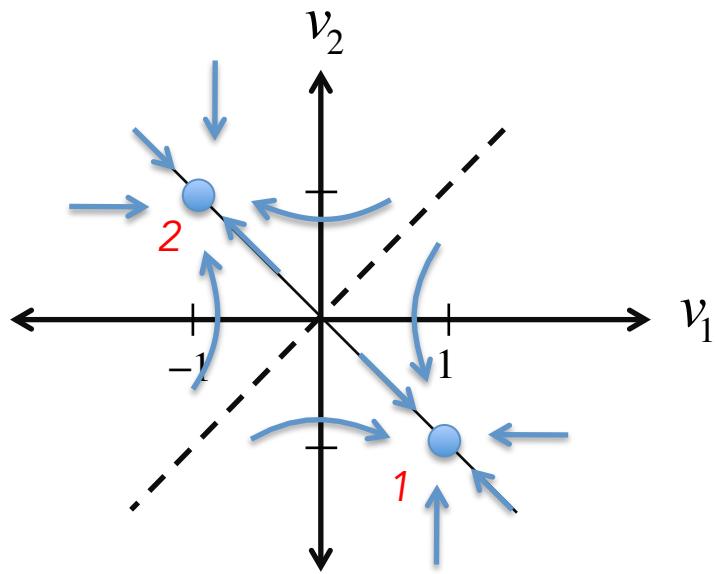
Network state will move to attractor 1 if $h_1 > h_2$

Network state will move to attractor 2 if $h_2 > h_1$



WINNER-TAKE-ALL NETWORK

- *Implements decision making*

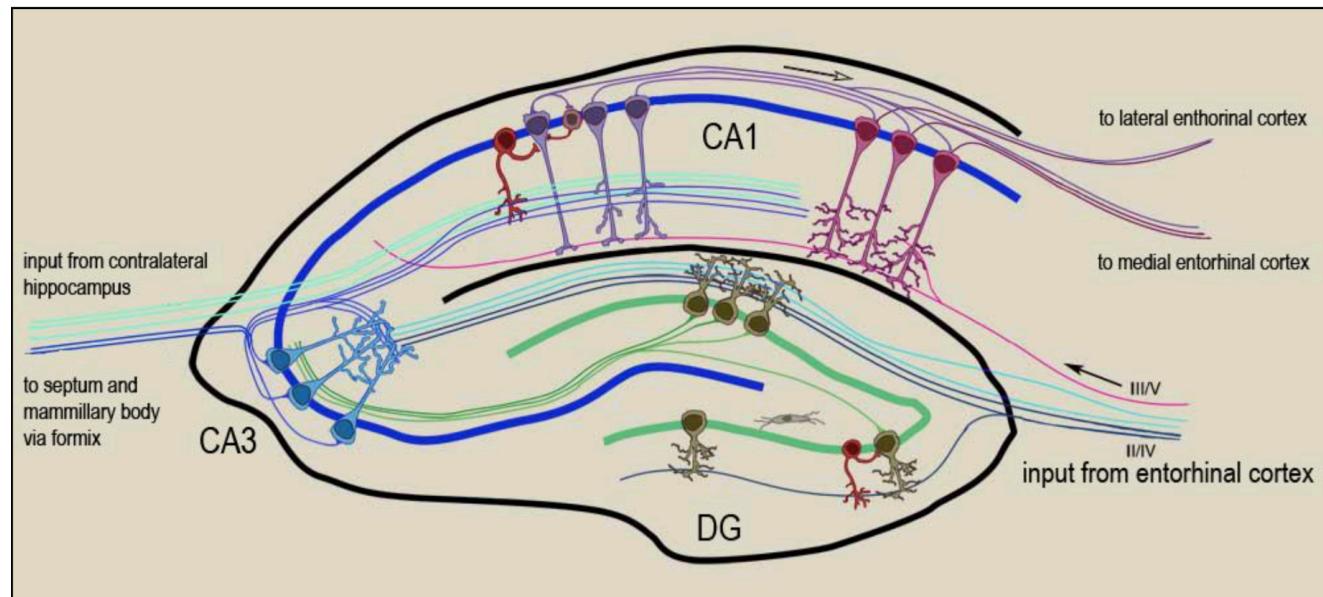


LEARNING OBJECTIVES FOR LECTURE 20

- Recurrent networks with lambda greater than one
 - Attractors
- Winner-take-all networks
- Attractor networks for long-term memory (Hopfield model)
- Energy landscape
- Hopfield network capacity

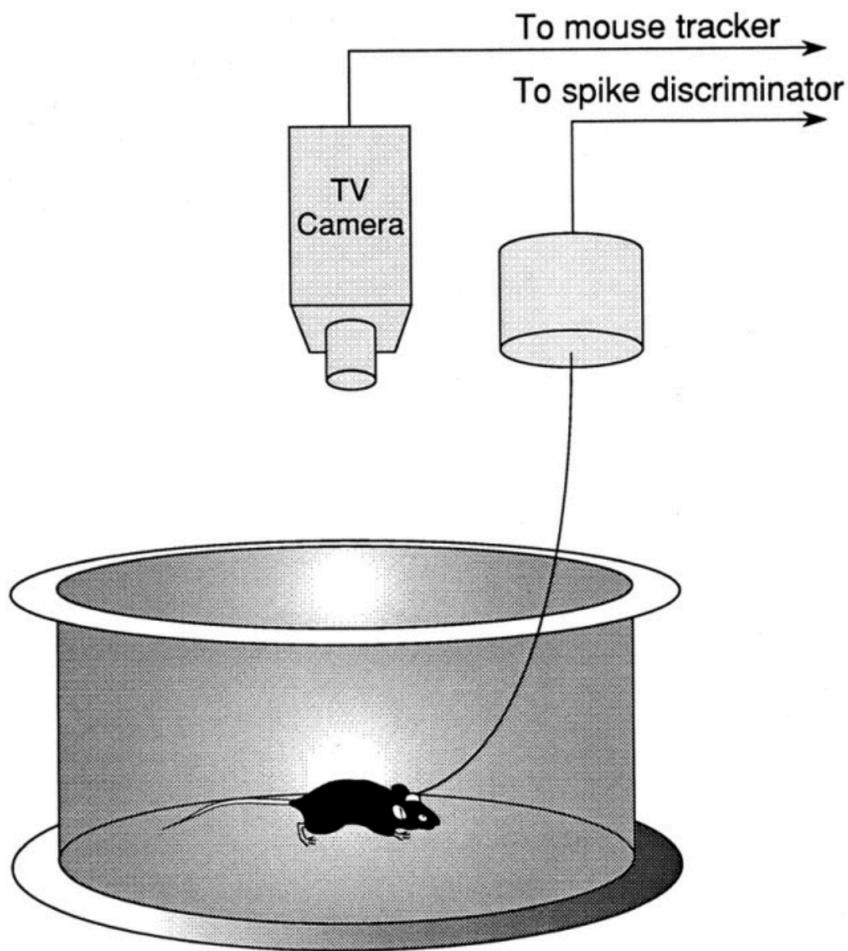
HIPPOCAMPUS

- *The region CA3 of the hippocampus gets sensory inputs and forms a dense recurrently connected network*



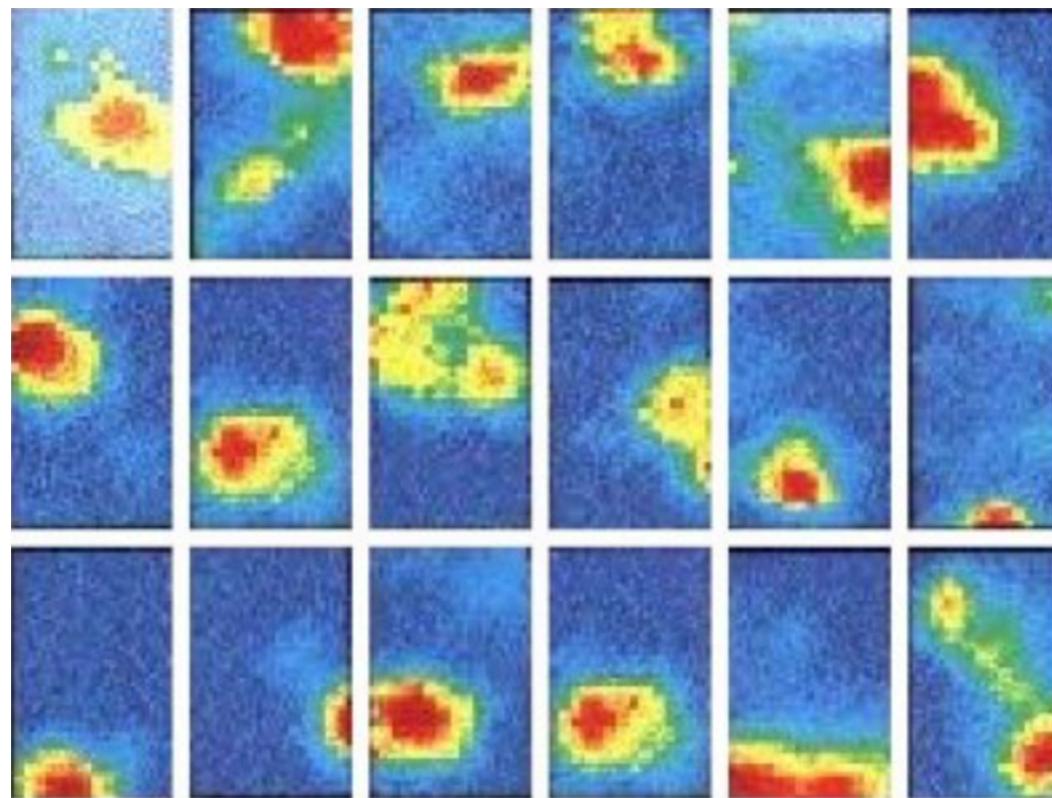
HIPPOCAMPUS

- Some neurons in CA3 represent 'memories' of locations in space (place cells).



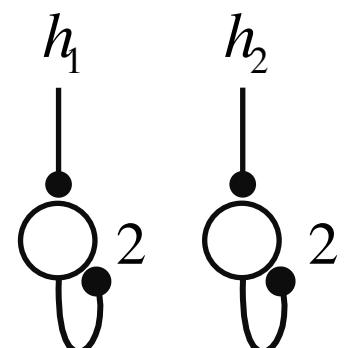
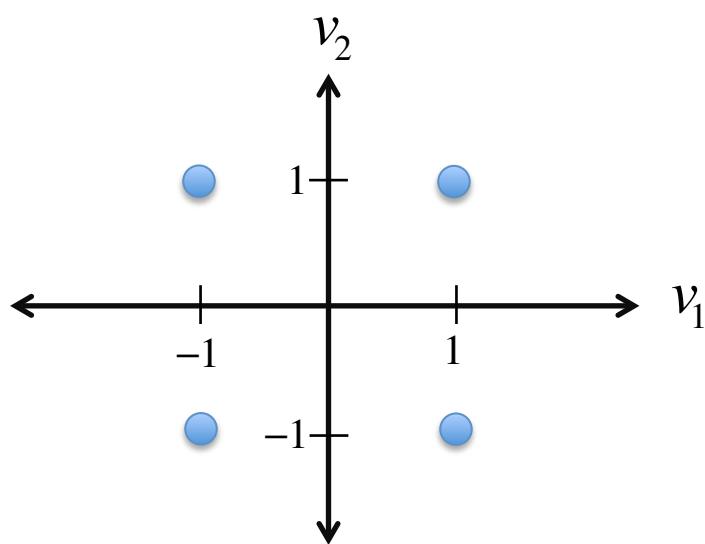
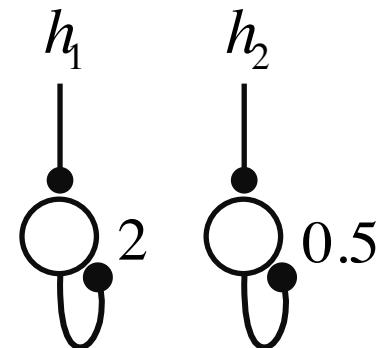
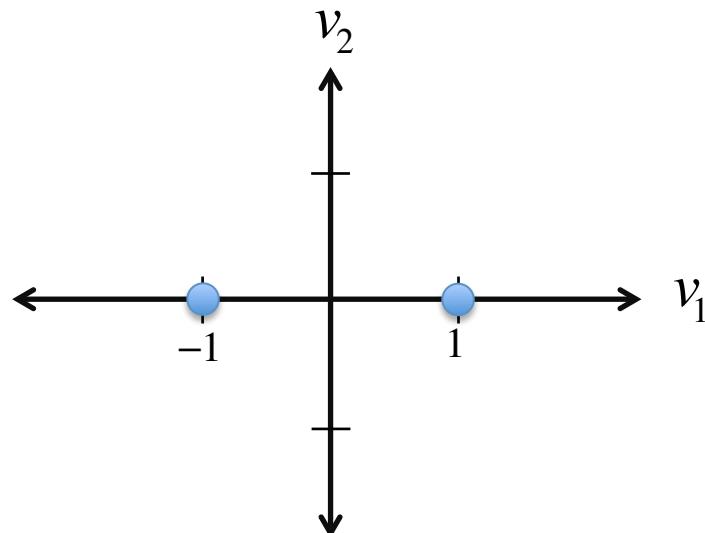
HIPPOCAMPUS

- *Different neurons represent different remembered locations.*



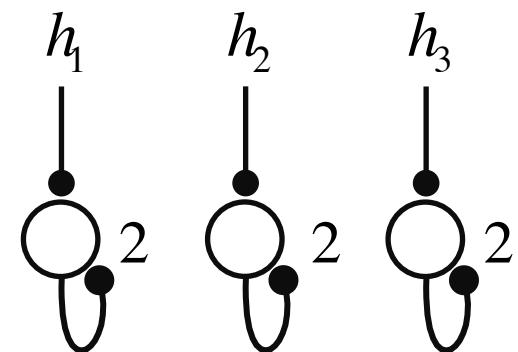
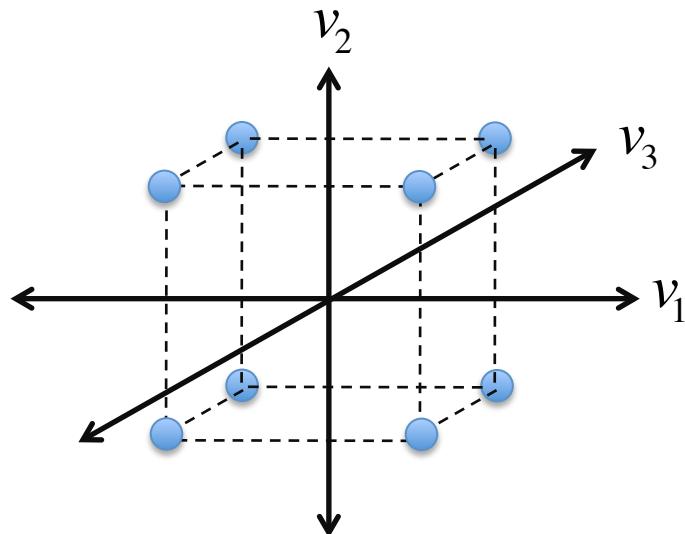
RECURRENT NETWORKS

- Networks with many attractors...



HOPFIELD NETWORKS

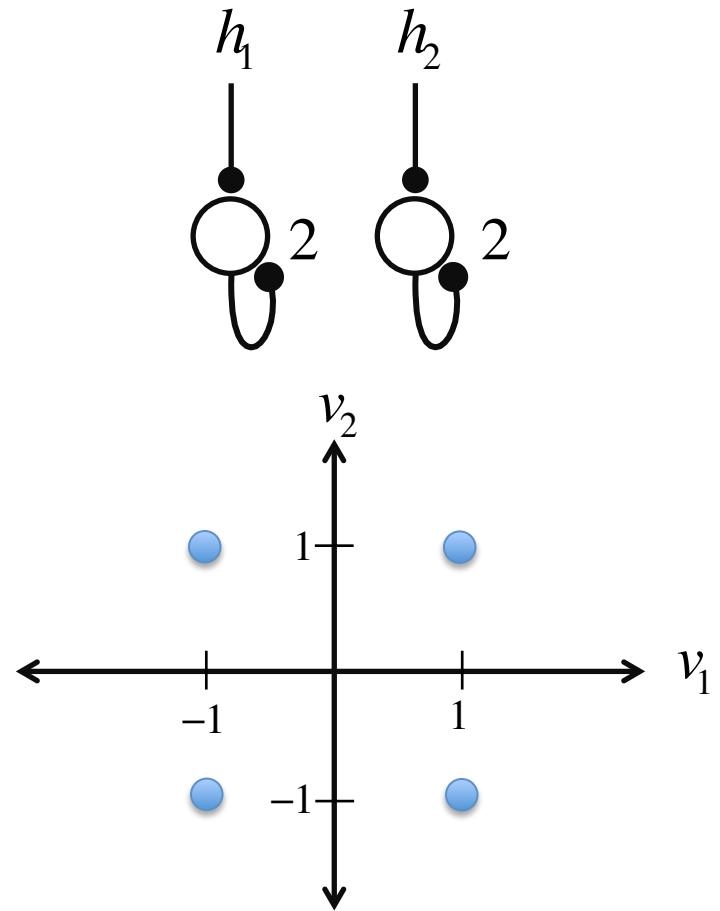
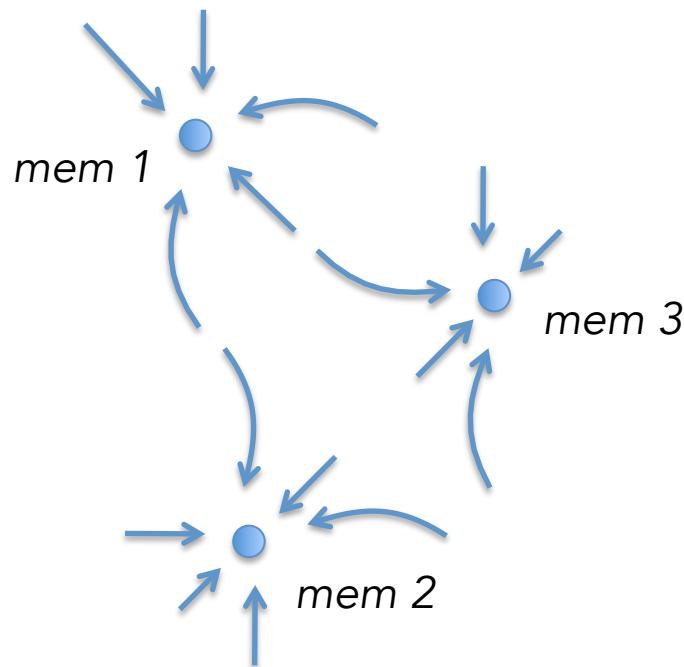
- Networks with many attractors...



2^n possible states !

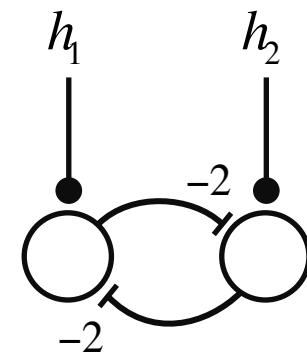
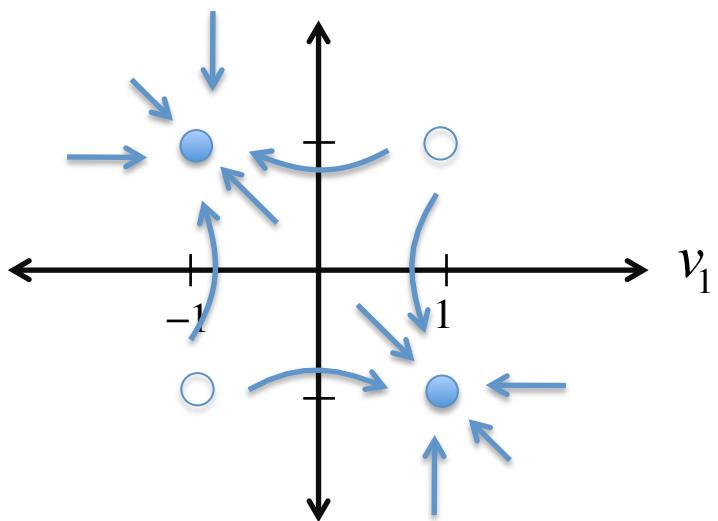
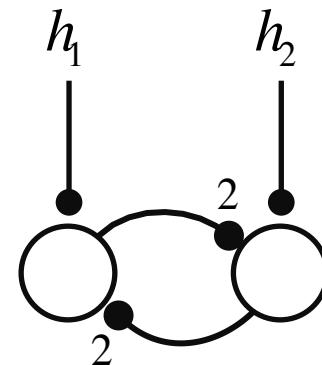
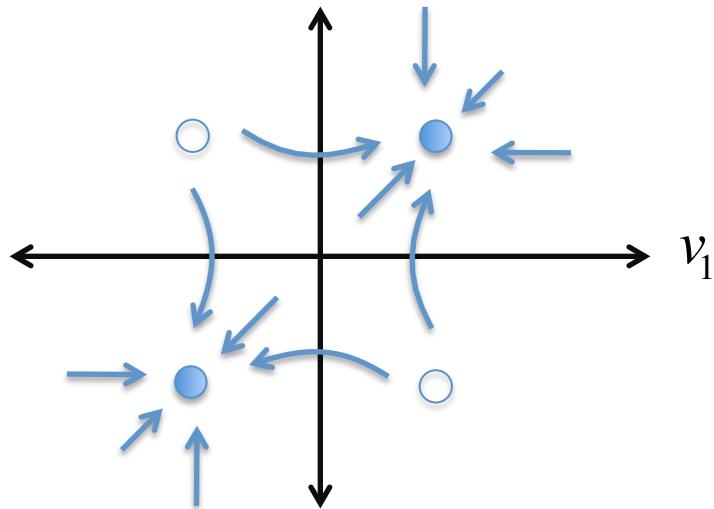
HOPFIELD NETWORKS

- We only want some of the possible states to be stable



HOPFIELD NETWORKS

- Networks with many attractors...



HOPFIELD NETWORKS

- We started with this dynamical equation

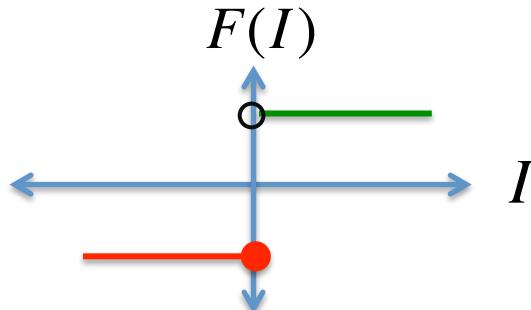
$$\tau_n \frac{d\vec{v}}{dt} = -\vec{v} + F[\vec{h} + M \vec{v}]$$

- We are going to simplify this as follows:

$$\vec{v}(t+1) = F[M \vec{v}(t)]$$

$$v_i(t+1) = F\left[\sum_{j=1}^N M_{ij} v_j(t)\right]$$

where the neuronal activation function is



binary threshold neuron

$$F(x) = \text{sgn}(x) = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x \leq 0 \end{cases}$$

HOPFIELD NETWORKS

- Our goal is to make a network that evolves so that it approaches any desired pattern $\vec{\xi}$

Where ξ_i is the activity (1 or -1) of the i^{th} neuron.

- The condition for $\vec{\xi}$ to be a stable pattern is

$$v_i(t+1) = \operatorname{sgn} \left[\sum_{j=1}^N M_{ij} \xi_j \right] = \xi_i$$

HOPFIELD NETWORKS

- Let's try the weight matrix $M_{ij} = \alpha \xi_i \xi_j$ where $\alpha > 0$

$$v_i(t+1) = \operatorname{sgn} \left[\sum_{j=1}^N M_{ij} \xi_j \right]$$

NOTE: This is a symmetric matrix!

$$= \operatorname{sgn} \left[\sum_{j=1}^N (\alpha \xi_i \xi_j) \xi_j \right]$$

$$= \operatorname{sgn} \left[\alpha \xi_i \sum_{j=1}^N \xi_j \xi_j \right]$$

$$= \operatorname{sgn} [\alpha N \xi_i]$$

$$M_{ij} = \frac{1}{N} \xi_i \xi_j !$$

$$M = \frac{1}{N} \vec{\xi} \vec{\xi}^T$$

$$v_i(t+1) = \xi_i$$

HOPFIELD NETWORKS

- Let's take an example: Design a network of 3 neurons that remembers a pattern $(1,1,-1)$.

$$\vec{\xi} = \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}$$

$$M_{ij} = \frac{1}{N} \vec{\xi} \vec{\xi}^T = \frac{1}{3} \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix} \begin{pmatrix} 1 & 1 & -1 \end{pmatrix} = \frac{1}{3} \begin{pmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \\ -1 & -1 & 1 \end{pmatrix}$$

- Is the pattern $(1,1,-1)$ a stable state?

$$v(t+1) = \text{sgn} \left[\frac{1}{3} \begin{pmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \\ -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix} \right] = \text{sgn} \left[\frac{1}{3} \begin{pmatrix} 3 \\ 3 \\ -3 \end{pmatrix} \right] = \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}$$

Yes!

HOPFIELD NETWORKS

- Does our network have an 'attractor' at the pattern $(1,1,-1)$?
 - Let's start the network at a different state and see what happens...

$$\vec{v}_0 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$v(t+1) = \text{sgn} \left[\frac{1}{3} \begin{pmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \\ -1 & -1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \right] = \text{sgn} \left[\frac{1}{3} \begin{pmatrix} 3 \\ 3 \\ -3 \end{pmatrix} \right] = \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}$$

- The network evolves toward the 'attractor' state $(1,1,-1)$!

HOPFIELD NETWORKS

- Let's prove that $\vec{\xi}$ is an attractor of the network $M_{ij} = \frac{1}{N} \vec{\xi} \vec{\xi}^T$

Calculate total input onto the i^{th} neuron starting in arbitrary state \vec{v}

$$k_i = \sum_{j=1}^N M_{ij} v_j \quad \text{where } v_j \text{ is the firing rate of the } j^{th} \text{ neuron} \quad \vec{k} = M \vec{v}$$

$$\begin{aligned} k_i &= \sum_{j=1}^N \left(\frac{1}{N} \xi_i \xi_j \right) v_j &= \frac{1}{N} \xi_i \sum_{j=1}^N \xi_j v_j \\ && \sum_{j=correct} \xi_j v_j + \sum_{j=incorrect} \xi_j v_j \\ k_i &= \frac{1}{N} \xi_i (N_{correct} - N_{incorrect}) & v_i(t+1) = \text{sgn}[k_i] = \xi_i \end{aligned}$$

The total input has the correct sign if the majority of the neurons have the correct state!

LEARNING OBJECTIVES FOR LECTURE 20

- Recurrent networks with lambda greater than one
 - Attractors
- Winner-take-all networks
- Attractor networks for long-term memory (Hopfield model)
- Energy landscape
- Hopfield network capacity

THE ENERGY FUNCTION

$$k_i = \sum_{j=1}^N M_{ij} v_j \quad \vec{k} = M \vec{v}$$

- Each possible state of the network has an energy given by:

$$H = -\frac{1}{2} \vec{v} \cdot \vec{k} \quad H = -\frac{1}{2} \vec{v}^T M \vec{v}$$

This is just the overlap of the current state of the network with the pattern of inputs to all the neurons!

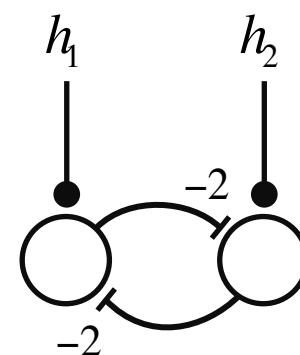
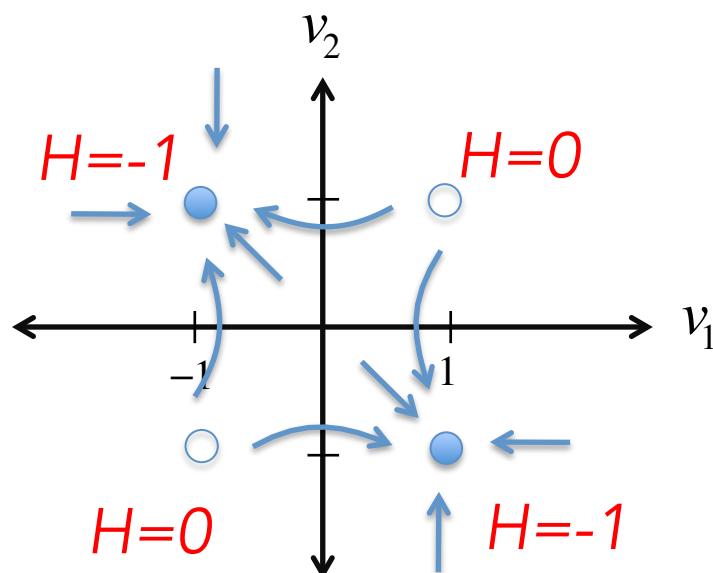
The energy is lowest when current state has high overlap with the synaptic drive to the next state ---> in an attractor

THE ENERGY FUNCTION

- Each possible state of the network has an energy given by:

$$H = -\frac{1}{2} \vec{v}^T M \vec{v}$$

$$M = \begin{pmatrix} 0 & -2 \\ -2 & 0 \end{pmatrix}$$

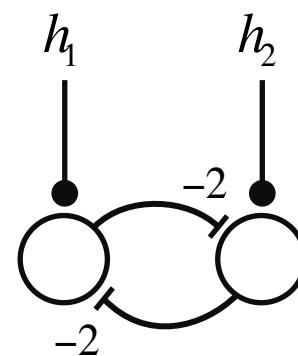
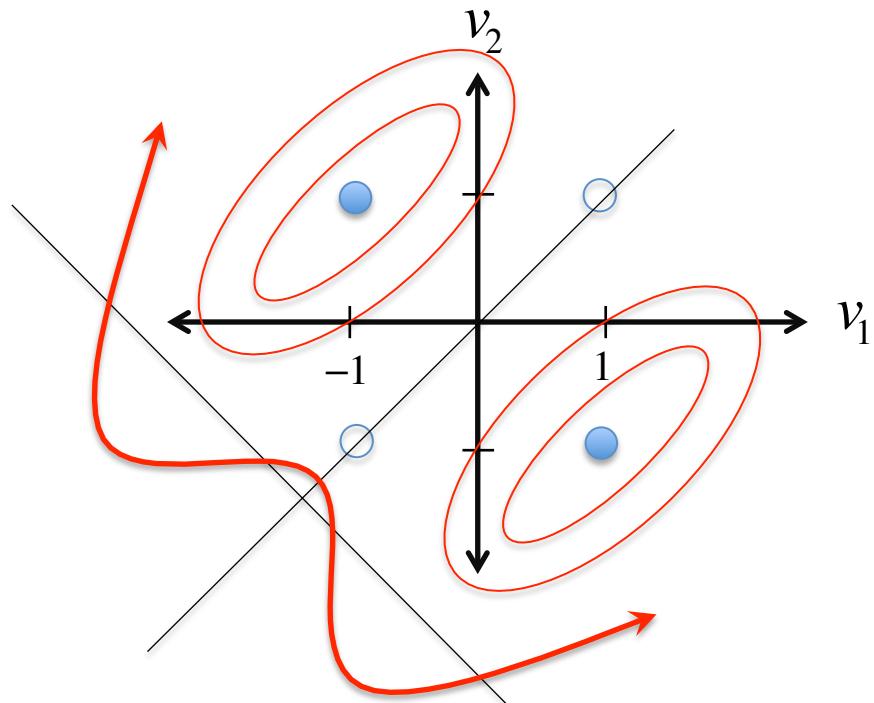


THE ENERGY FUNCTION

- Each possible state of the network has an energy given by:

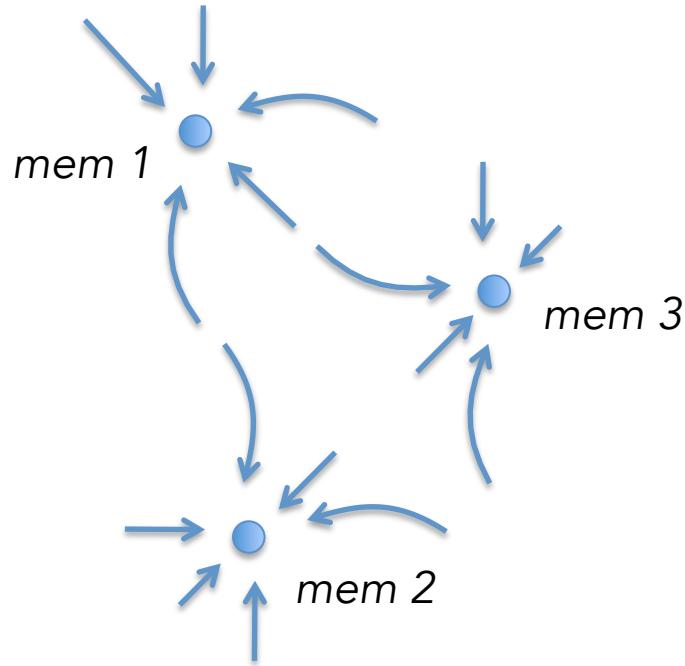
$$H = -\frac{1}{2} \vec{v}^T M \vec{v}$$

$$M = \begin{pmatrix} 0 & -2 \\ -2 & 0 \end{pmatrix}$$



HOPFIELD NETWORKS

- A Hopfield network can 'reconstruct' a memory or pattern from a partial pattern
- It will evolve 'downhill' toward whichever memory is closest to the input pattern (content addressable memory)



HOPFIELD NETWORKS CAN BE USED TO STORE A MEMORY OF AN IMAGE

- ‘Content addressable memory’



LEARNING OBJECTIVES FOR LECTURE 20

- Recurrent networks with lambda greater than one
 - Attractors
- Winner-take-all networks
- Attractor networks for long-term memory (Hopfield model)
- Energy landscape
- Hopfield network capacity

MULTIPLE MEMORIES

- We want our network to remember P different patterns $\vec{\xi}^\mu$
 $\mu = 0, 1, 2, \dots, P-1$
- We compute the contribution to the weight matrix from each pattern...

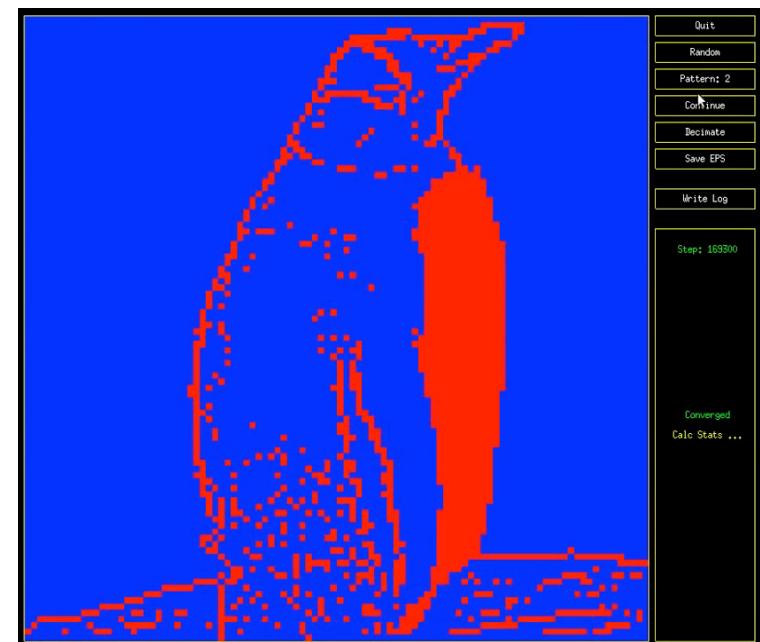
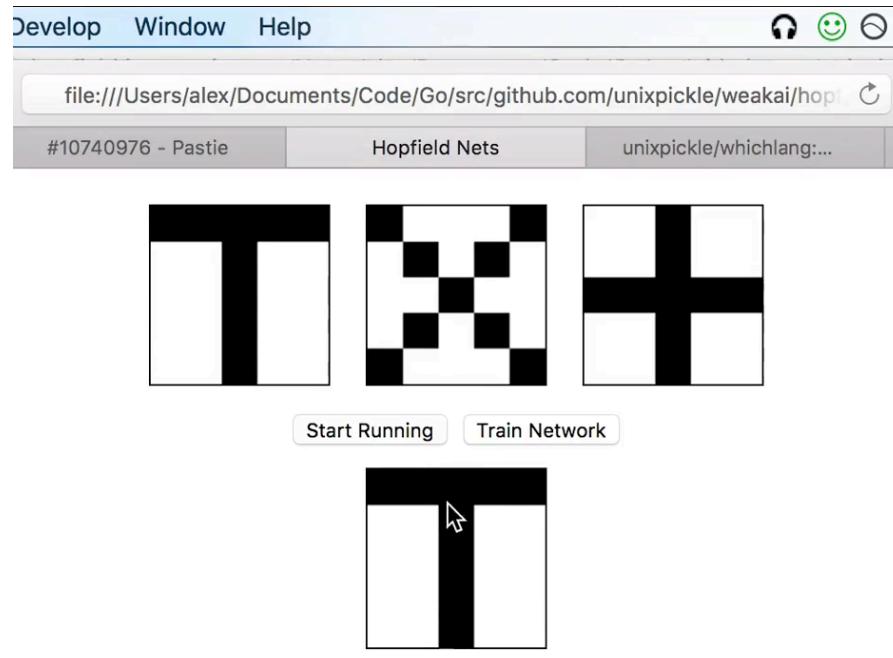
$$M_{ij}^\mu = \frac{1}{N} \xi_i^\mu \xi_j^\mu$$

and add them up!

$$M_{ij} = \frac{1}{N} \sum_{\mu=0}^{P-1} \xi_i^\mu \xi_j^\mu$$

- The network state will evolve to the attractor that is 'closest' to the initial state (that with the biggest overlap).

MULTIPLE MEMORIES



WHAT IS THE CAPACITY?

- Do the same stability analysis we did before...

For a pattern $\vec{\xi}^0$ to be an attractor, we want

$$v_i(t+1) = \operatorname{sgn}\left(\sum_j M_{ij} \xi_j^0\right) = \xi_i^0 \quad M_{ij} = \frac{1}{N} \sum_\mu \xi_i^\mu \xi_j^\mu$$

- Let's plug in our weight matrix and see what we get...

$$v_i(t+1) = \operatorname{sgn}\left(\sum_j \left[\frac{1}{N} \sum_\mu \xi_i^\mu \xi_j^\mu \right] \xi_j^0\right) = ??$$

WHAT IS THE CAPACITY?

$$v_i(t+1) = \operatorname{sgn} \left(\sum_j \left[\frac{1}{N} \sum_{\mu} \xi_i^{\mu} \xi_j^{\mu} \right] \xi_j^0 \right) = ??$$

- Rearrange...

$$= \operatorname{sgn} \left(\frac{1}{N} \sum_{\mu} \xi_i^{\mu} \sum_j \xi_j^{\mu} \xi_j^0 \right)$$

- Separate terms where $\mu = 0$ and $\mu \neq 0$

$$= \operatorname{sgn} \left(\frac{1}{N} \xi_i^0 \underbrace{\sum_j \xi_j^0 \xi_j^0}_N + \frac{1}{N} \sum_{\mu \neq 0} \xi_i^{\mu} \sum_j \xi_j^{\mu} \xi_j^0 \right)$$

WHAT IS THE CAPACITY?

$$\operatorname{sgn}\left(\sum_j M_{ij} \xi_j^0\right) = \xi_i^0 \quad ???$$

$$= \operatorname{sgn}\left(\xi_i^0 + \frac{1}{N} \sum_{\mu \neq 0} \xi_i^\mu \sum_j \xi_j^\mu \xi_j^0\right)$$
$$\underbrace{\xi_i^0}_{\text{Pattern}} + \underbrace{\frac{1}{N} \sum_{\mu \neq 0} \xi_i^\mu \sum_j \xi_j^\mu \xi_j^0}_{\text{Cross-talk}} \xrightarrow{\xi^\mu \cdot \xi^0}$$

Cross-talk between our pattern $\vec{\xi}^0$ and all the other memories
depends on how much overlap there is !

WHAT IS THE CAPACITY?

$$v_i(t+1) = \operatorname{sgn} \left(\xi_i^0 + \frac{1}{N} \sum_{\mu \neq 0} \xi_i^\mu \left(\vec{\xi}^\mu \cdot \vec{\xi}^0 \right) \right)$$

- You can see that if all the memories are orthogonal, then all are stable attractors.
- But if one of the memories (e.g. $\vec{\xi}^1$) is close to $\vec{\xi}^0$ then

$$\vec{\xi}^0 \cdot \vec{\xi}^1 \approx N$$

- And...

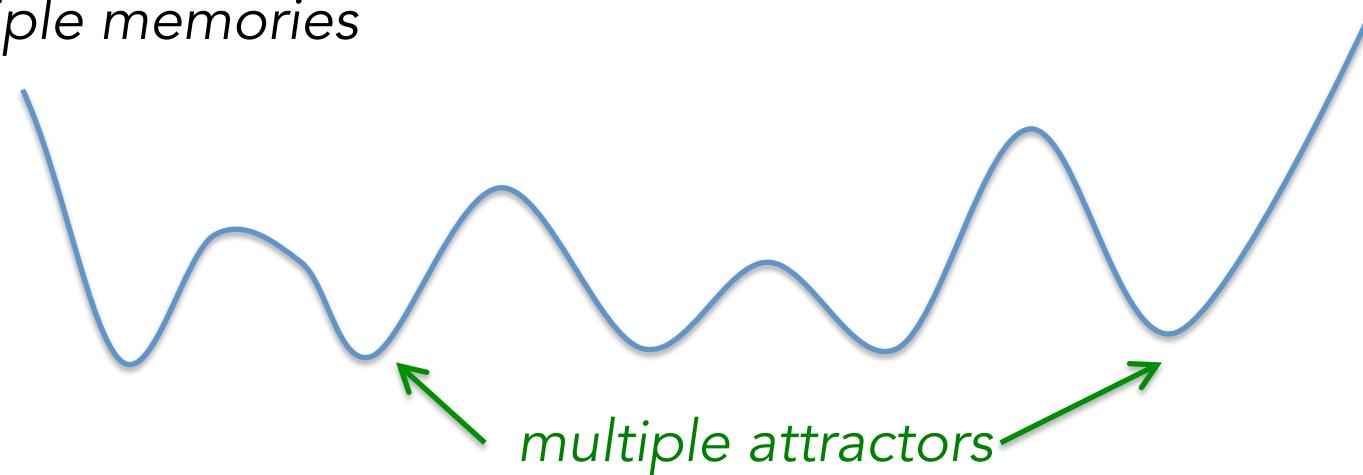
$$v_i(t+1) = \operatorname{sgn} \left(\xi_i^0 + \xi_i^1 \right) \neq \xi_i^0$$

WHAT IS THE CAPACITY?

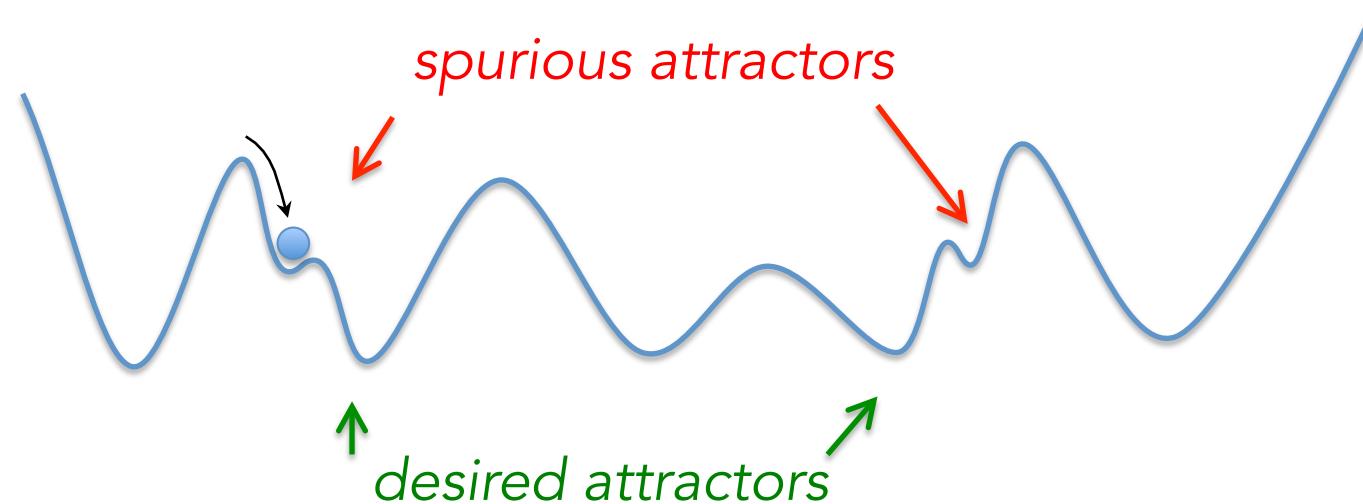
- Memories don't need to be orthogonal, as long as the cross-talk term is not large enough to change the sign of the inputs, the memories will not have any errors.
- For random values of ξ_i , a Hopfield network can store up to $0.15N$ memories and still have a very small probability ($p<0.01$) of having an error in one neuron.

WHAT IS THE CAPACITY?

- Multiple memories



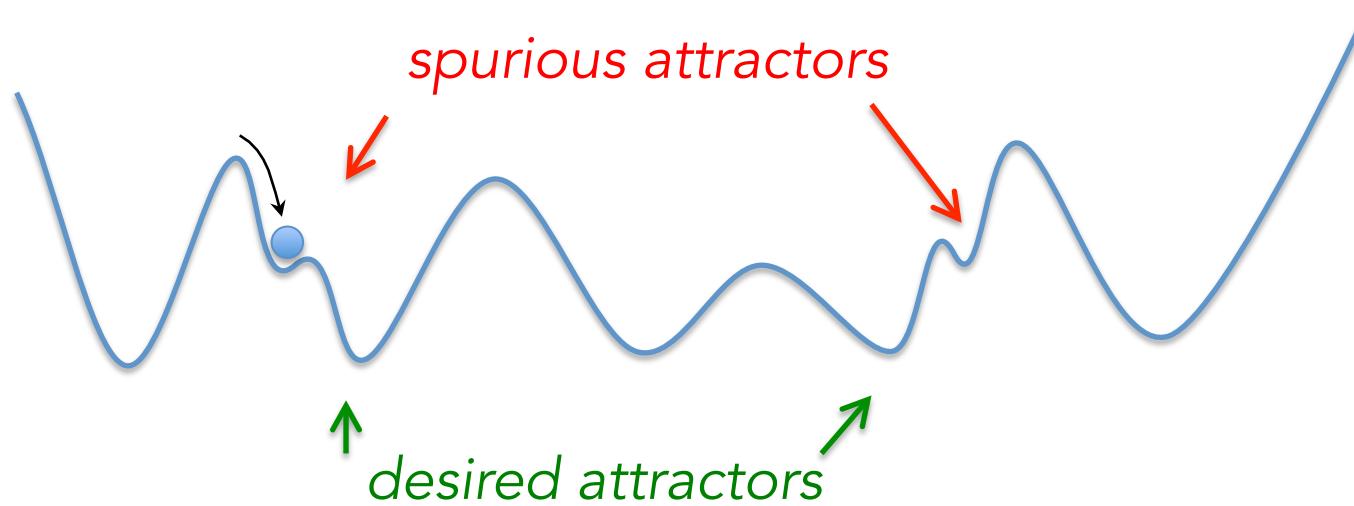
- Too many memories...



WHAT IS THE CAPACITY?



- Too many memories...



LEARNING OBJECTIVES FOR LECTURE 20

- Recurrent networks with lambda greater than one
 - Attractors
- Winner-take-all networks
- Attractor networks for long-term memory (Hopfield model)
- Energy landscape
- Hopfield network capacity