

# **SIMULATING RNNs ON GPUS**

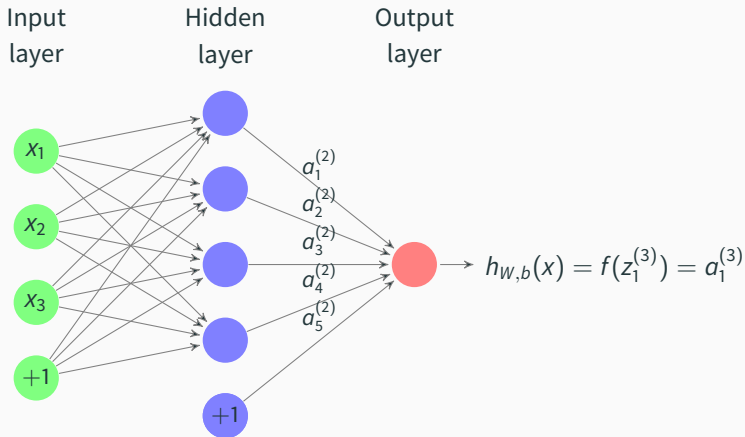
WORKING TITLE

---

Zhangsheng Lai

September 11, 2017

# MULT-LAYER PERCEPTRON

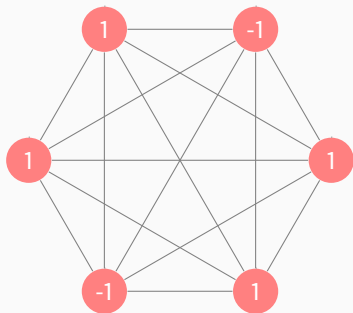


**Figure 1:** Multi-layer perceptron

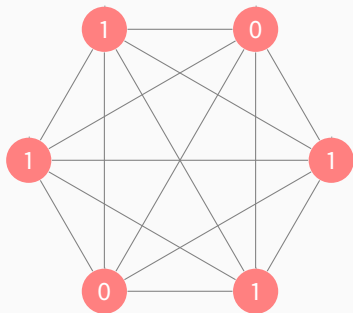
# **HOPFIELD NETWORKS AND BOLTZMANN MACHINES**

---

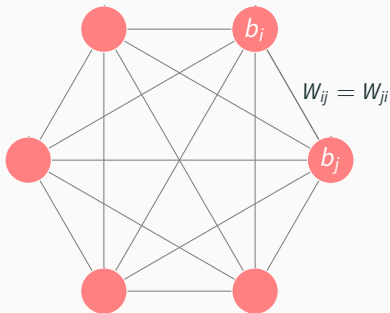
# HOPFIELD NETWORKS



# HOPFIELD NETWORKS



# HOPFIELD NETWORKS



$$\text{Energy configuration, } E = - \sum_{i < j} W_{ij} x_i x_j - \sum_i b_i x_i$$

$$\text{Energy gap, } \Delta E_i = E(x_i = 0) - E(x_i = 1) = \sum_j W_{ij} x_j + b_i$$

$$\text{Update rule, } x_i := \begin{cases} +1 & \sum_j W_{ij} x_j + b_i \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

# Simulating RNNs on GPUs

## Hopfield Networks



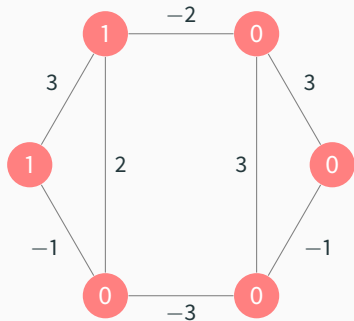
$$\text{Energy configuration, } E = - \sum_{i < j} W_{ij} x_i x_j - \sum_i b_i x_i$$

$$\text{Energy gap, } \Delta E = E(x_i = 0) - E(x_i = 1) = \sum_j W_{ij} x_j + b_i$$

$$\text{Update rule, } x_i = \begin{cases} +1 & \sum_j W_{ij} x_j + b_i \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

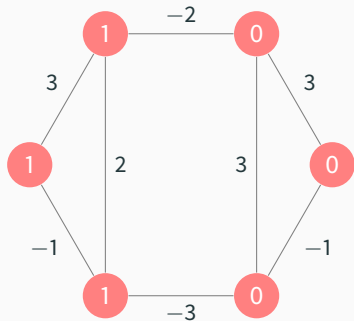
- composed of primitive computing elements called units
- units has two states, on or off, represented by  $\{1, -1\}$  or  $\{1, 0\}$
- connected to each other by bi-directional links
- adopts these states as a function of the states of its neighbouring units and weights of its links to them, it is a probabilistic function for a Boltzmann machine.
- weights can take on any real value
- a unit being on or off is taken to mean that the system currently accepts or rejects some elemental hypothesis of the domain
- weight on a link represents a weak pairwise constrain between two hypothesis
- positive (negative) weights indicate that two hypothesis support (contradict) one another with other things being equal
- link weights are symmetric, having the same strength in both directions

# HOPFIELD NETWORKS

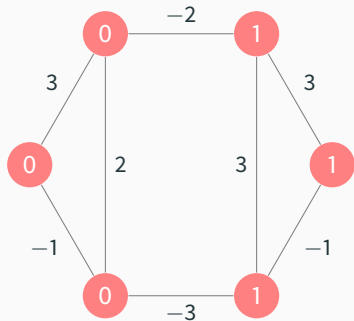




# HOPFIELD NETWORKS



# HOPFIELD NETWORKS



# Simulating RNNs on GPUs

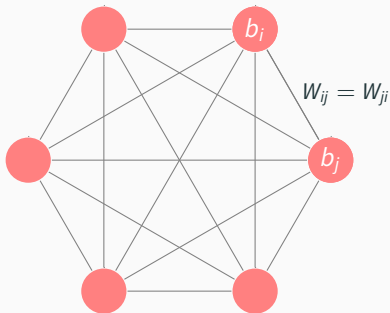
2017-09-11

## Hopfield Networks



- Updating of Hopfield networks is done sequentially usually in a randomized order. Parallel updating might increase the energy instead.
- Hopfield networks always make decisions to reduce the energy and makes it impossible to escape from local minima.
- random noise can help us escape from poor minima, by starting with lots of noise so its easy to cross energy barriers and gradually decrease the noise so the system ends in a deep minimum. This is called simulated annealing.

# BOLTZMANN MACHINES



$$E = - \sum_{i < j} w_{ij} x_i x_j - \sum_i b_i x_i$$

$$\Delta E_i = E(x_i = 0) - E(x_i = 1) = \sum_j w_{ij} x_j + b_i$$

$$\mathbb{P}(x_i = 1) = \frac{1}{1 + e^{-\Delta E_i / \tau}}$$

# Simulating RNNs on GPUs

## Boltzmann Machines



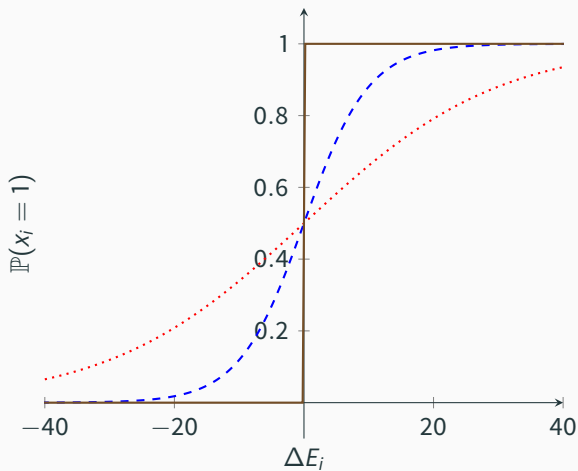
$$E = - \sum_{i,j} w_{ij} x_i x_j - \sum_i b_i x_i$$

$$\Delta E_i = E(x_i = 0) - E(x_i = 1) = \sum_j w_{ij} x_j + b_i$$

$$P(x_i = 1) = \frac{1}{1 + e^{-\Delta E_i / \tau}}$$

- replace the binary threshold units by binary stochastic units that make biased random decisions
- temperature variable controls the amount of noise
- when  $\tau \rightarrow 0$  we get back the Hopfield network
- for  $\tau_1 > \tau_2$ , we are less likely to go to a lower energy state compared to in  $\tau_1$  compared to  $\tau_2$ , i.e. more likely to go to a higher energy state when the temperature is higher. This allows us to escape from local minimum and arrive at the global minimum

# BOLTZMANN MACHINES

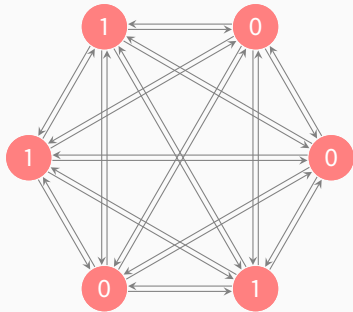


**Figure 2:**  $\tau = 0$  (solid),  $\tau = 5$  (dashed),  $\tau = 15$  (dotted)

# **McCULLOH-PITTS MACHINES**

---

# McCULLOH-PITTS MACHINES





# Simulating RNNs on GPUs

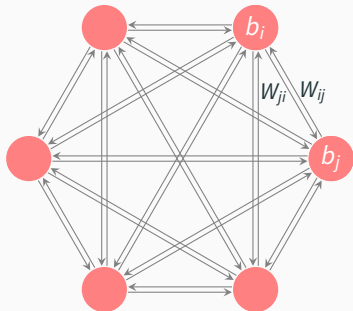
2017-09-11

## └─ McCulloch-Pitts Machines



- state 1 is the refractory state, the neuron just fired and is unable to fire till it recovers
- state 0 is the armed state, the neuron just recovered and is waiting to fire
- here we model the units with the Nossenson-Messer neuron model, which explains biological firing rates in response to external stimuli

# McCULLOH-PITTS MACHINES



$$\text{Transition Energy, } E(y, x|\theta) = - \sum_{ji \in E} w_{ji} y_j x_i - \sum_{j \in V} b_j s_j - \sum_{i \in V} b_i s_i$$

$$\Gamma_{yx} = \exp \left( -\frac{1}{2\tau} E(y, x|\theta) + \frac{1}{2\tau} E(x, x|\theta) \right)$$

# Simulating RNNs on GPUs

2017-09-11

## └ McCulloch-Pitts Machines

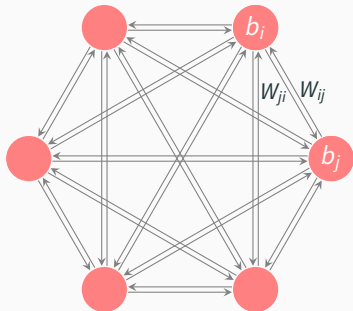


$$\text{Transition Energy, } E(y, x| \theta) = - \sum_{j \in \mathcal{I}} w_{ij} y_i x_j - \sum_{j \in \mathcal{I}} b_j y_j - \sum_{j \in \mathcal{I}} b_j x_j$$

$$\Gamma_{yx} = \exp \left( - \frac{1}{2\sigma^2} E(y, x| \theta) + \frac{1}{2\sigma^2} E(x, x| \theta) \right)$$

- allow transitions where  $y$  and  $x$  differ by only one bit
- here the  $W$  matrix need not be symmetrical with zero diagonals like what we had in the Hopfield network and Boltzmann machine models
- for each  $y \neq x$ , start a Poisson process with rate  $\Gamma_{yx}$

# McCULLOH-PITTS MACHINES



$$\text{Transition Energy, } E(y, x | \theta) = - \sum_{ji \in E} W_{ji} y_j x_i - \sum_{j \in V} b_j s_j - \sum_{i \in V} b_i s_i$$

$$\Gamma_{yx} := \exp \left( \frac{1}{2\tau} s_j z_j \right)$$

where  $s_j = 1 - 2x_j$ ,  $z_j = \sum_i W_{ji} x_i + b_j$  and  $x, y$  differ by the  $j$ th unit.

# Simulating RNNs on GPUs

## └ McCulloch-Pitts Machines



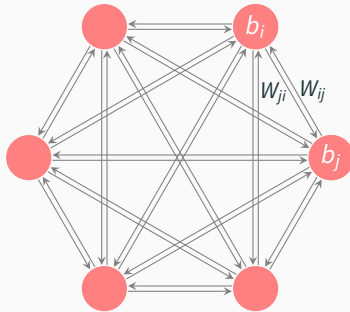
$$\text{Transition Energy, } E(y, x(i)) = - \sum_{j \in J} W_{ij} y_j x_i - \sum_{j \in J} b_j y_j - \sum_{i \in I} b_i x_i$$

$$r_{ij} := \exp\left(\frac{1}{\sum_{k \in J} W_{ik} y_k}\right)$$

where  $y_i = 1 - 2x_i$ ,  $x_i = \sum_{j \in J} W_{ij} y_j + b_i$  and  $x, y$  differ by the  $j$ th unit.

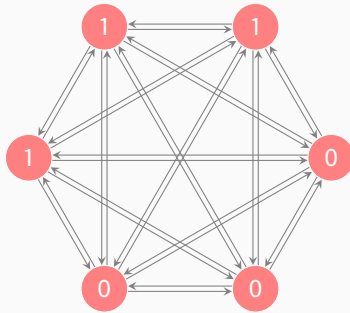
- when doing the updates we can just update the linear responses  $z_j$  and apply softmax on the  $\lambda_j$ 's to get the probability distribution of the transitions.

# McCULLOH-PITTS MACHINES

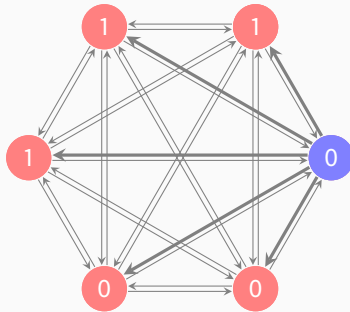


Transition probability from  $x$  to  $y$ ,  $p_{yx} = \frac{\lambda_j}{\sum_{j'} \lambda_{j'}}$

# McCULLOH-PITTS MACHINES

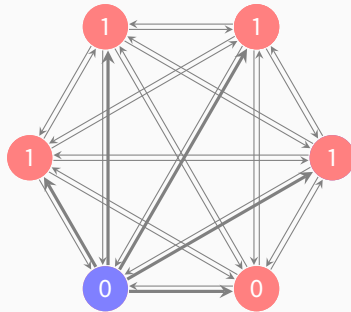


# McCULLOH-PITTS MACHINES

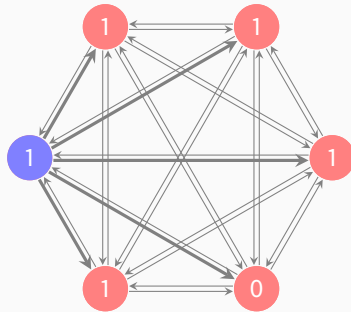




# McCULLOH-PITTS MACHINES



# McCULLOH-PITTS MACHINES



# Simulating RNNs on GPUs

2017-09-11

## └─ McCulloch-Pitts Machines



- digraph,  $G = (V, E)$ , weights  $W : V \rightarrow \mathbb{R}$ , biases  $b : E \rightarrow \mathbb{R}$ , binary states  $\mathbb{B} = \{0, 1\}$  with an initial distribution  $\mathbb{B}^{|V|} \rightarrow \delta$  and a temperature  $\tau$
- allow transitions where  $y$  and  $x$  differ by only one bit
- for each  $y \neq x$ , start a Poisson process with rate  $\Gamma_{yx}$
- if the Poisson process for state  $y'$  is the first to fire, or the interarrival timing from  $x$  to  $y'$  is the shortest,



$$E(x^{(n+1)}, x^{(n)} | \theta) = - \sum_{ji \in E} W_{ji} x_j^{(n+1)} x_i^{(n)} - \sum_{j \in V} b_j s_j - \sum_{i \in V} b_i s_i$$

$$\Gamma_{x^{(n+1)} x^{(n)}} = \exp \left( -\frac{1}{2\tau} E(x^{(n+1)}, x^{(n)} | \theta) + \frac{1}{2\tau} E(x^{(n)}, x^{(n)} | \theta) \right)$$



# REFERENCES



J. Macor.

***A Brief Introduction to Type Theory and the Univalence Axiom***

<http://math.uchicago.edu/~may/REU2015/REUPapers/Macor.pdf>



The Univalent Foundations Program

***Homotopy Type Theory: Univalent Foundations of Mathematics.***

<https://homotopytypetheory.org/book>



The n-Category Café

***From Set Theory to Type Theory***

[https://golem.ph.utexas.edu/category/2013/01/from\\_set\\_theory\\_to\\_type\\_theory.html](https://golem.ph.utexas.edu/category/2013/01/from_set_theory_to_type_theory.html)



The nLab

***Function Type***

<https://ncatlab.org/nlab/show/function+type>