

# Pathway for Energy Based Models

Germán González, gonzalezv.germanh@javeriana.edu.co

**Abstract—**

**Index Terms—**Gabor filters, Modern Hopfield networks, Convolutional Networks, Associative memory.

## I. INTRODUCTION

One of the main field of research on AI is becoming the analysuis of Energy Based Models. The article "Attention is All You need" has had a tremendous impact in the advance of artificial intelligence. And the recent "Hopfield Networks is all you need" has shown that the Energy based models is the way to go in terms of modeling. However, the theory that backs them is not clearly understood, eventhough it promises that there is a resemblance between biology and the models that path is difficult to follow up. In order to cover that gap, this article will approach the models from two main stating points: The linear associative models that back the idea of attention mechanisms and energy models starting from the Ising model. In general article will be structured in the following mode:

- 1) **Previous work**
- 2) **Associative memories**
- 3) **Differential equations, eigenvalues and eigenvectors**
- 4) **Hamiltonian dynamics**
- 5) **Clasical Hopfield network**
- 6) **Modern Hopfield networks**

## II. PREVIOUS WORK

Energy Based Models can be traced back to the clasical Hopfield networks [8]. The article "Hopfield Networks is all you need" treats an extension of the classical Hopfield networks allowing the storage of exponentially more records that the ones possible in the classical networks. This has led then to focus on the investigation about energy Based Models as depicted in the recent article by Yann LeCun "A path towards autonomous machile learning" [5].

## III. ASSOCIATIVE MEMORIES

It is know that the human memory is associative. That is to say, that the human needs to have a cue of information in order to correalte the rest of a recall. The memory is built up day after day from a colection of information that somehow gets store in the brain (Hebbian Learning).

The intuition of association can be explained as in the following figure. A person might have two recalls that are totally uncorrelated. However, the presence of a third recall might bring a link between the two so that getting to remeber one induces the recall of another one. When the process only leads to getting to get the full recovery of the recall, then we just have a completion process while, in the case of linking

to another recall might lead to and induction process.

However, in this article, the process of induction will not be analyzed.

To understand the process of associative learning it is necessary to go back to the theory of eigenvalues and eigenvectors. The reason why the theory was developed was help resolving systems of diferential equations of first order.

$$\dot{x} = \mathbf{A}x \quad (1)$$

The objective is to transform the system in a system where the  $\mathbf{A}$  matrix is replaced with a diagonal  $\mathbf{D}$  matrix that is easy to solved.

$$\dot{z} = \mathbf{D}z \quad (2)$$

The eingenvector and eigenvalues comes from a change of coordinates as follows:

$$x = Tz \quad (3)$$

Deriving it:

$$\dot{x} = T\dot{z} \quad (4)$$

From equation 1:

$$T\dot{z} = ATz \quad (5)$$

$$T^{-1}T\dot{z} = T^{-1}ATz \quad (6)$$

$$\dot{z} = T^{-1}ATz \quad (7)$$

In this last equation  $T^{-1}AT$  can by identified as:

$$T^{-1}AT = \Lambda \quad (8)$$

So, in the equation 2 we can identify  $\Lambda$  as the  $\mathbf{D}$  matrix which is a diagonal matrix of the eigenvalues of  $\mathbf{A}$ .

For example, consider the following system of non linear differential equations:

$$\begin{aligned} \dot{x} &= -0.4x + 0.02xy \\ \dot{y} &= 0.8y - 0.01y^2x - 0.1xy \end{aligned} \quad (9)$$

The equation has a complex eigenvalue 0.8-2.8i near the point (6,20) when the Jacobian of the equation 9 is evaluated in that point, meaning that any starting point in the right quadrant will converge to it.

Now that is clear that a system of differetial equation with complex eigenvalues can represent memory let's analyze the propertis of a matrix to store them. Consider a set of orthogonal vectors  $v^{(j)}$ ,  $j = 1, 2, 3, \dots, n$  where  $v^{(i)T}v^{(j)} = 0$

$$\mathbf{W} = \sum_{j=1}^r v^{(j)T}v^{(j)} = \mathbf{V}\mathbf{V}^T \quad (10)$$

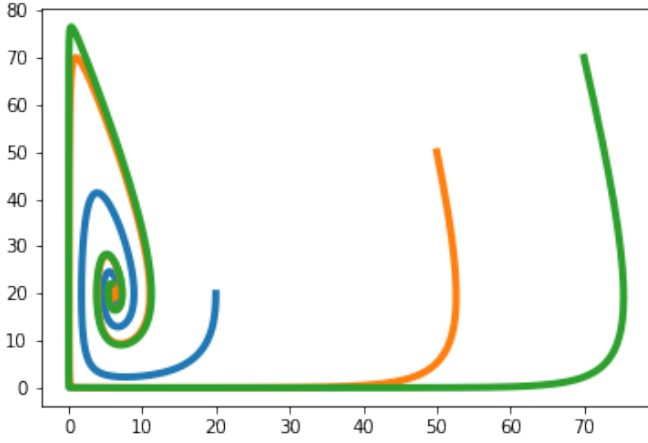


Fig. 1: Phase diagram of the system non linear differential equations 9. The lines blue, orange and green indicate the trayectories of three starting points converging to the sink point 6,20. This can be seen as an associative memory that once given certain inofrmation is able to get to the stored memory by itself.

And getting back to a linear system as in [12]:

$$\mathbf{y} = \mathbf{W}\mathbf{x} \quad (11)$$

Where  $\mathbf{x} = \mathbf{v}^k$ , then:

$$\mathbf{y} = \mathbf{W}\mathbf{v}^{(k)} = \sum_{j=1}^r v^{(j)T} v^j v^{(k)} \quad (12)$$

$$= v^{(k)T} v^k v^{(k)} = c v^{(k)} \quad (13)$$

So, the output can be taken as a linear combination of the inputs. If the vector  $\mathbf{v}^{(k)} = \mathbf{v}^{(k1)} + \mathbf{v}^{(k2)}$ , where  $\mathbf{v}^{(k1)}$  and  $\mathbf{v}^{(k2)}$  are orthogonal, presenting partial information in the form of  $\mathbf{v}^{(k1)}$ :

$$\mathbf{y} = \mathbf{W}\mathbf{v}^{(k1)} \quad (14)$$

Produces:

$$\mathbf{W} = \sum_{j=1}^r v^{(j)T} v^j v^{k1} \quad (15)$$

$$= v^k v^{(k)T} v^{k1} \quad (16)$$

$$= (v^{k1} + v^{(k2)})(v^{k1} + v^{(k2)})^T v^{k1} \quad (17)$$

$$= (v^{k1} + v^{(k2)})v^{(k1)T} v^{k1} \quad (18)$$

$$= d v^{k1} \quad (19)$$

The equation (19) is what is called a Content addressable memory.

#### IV. ATTENTION MECHANISM

As mentioned at the begining in [24], the attention mechanism is derived from:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (20)$$

The equation above means that attention is a linear combination as expressed down below:

$$\mathbf{y} = \sum_{j=1}^N \alpha_j \mathbf{v}_j \quad (21)$$

Where:

$$\sum_{j=1}^N \alpha_j = 1 \quad (22)$$

To avoid scaling up the input. In equation 20 we see that there is a factor  $\sqrt{(d)}$  that is similar to the factor  $d$  in equation 19, so it leads to consider also that the product  $\mathbf{QK}^T$  can be replaced by  $\mathbf{V}^T \mathbf{V}$  and finally multiplying by  $\mathbf{V}$  and adding over all  $N$  we get to an equation similar to 12 where the softmax accounts for the non-orthogonality of the  $v^k$  vectors.

#### V. SOFTMAX INTERPRETATION

The intepretation of the softmax can be deduced from the entropy function [4]. Assuming the multiplicity  $W$  of a combinatorial event:

$$W = \frac{N!}{n_1! n_2! \dots n_t!} \quad (23)$$

Using the Stirling approximation  $x! \approx (x/e)^x$ :

$$W = \frac{(N/e)^N}{(n_1/e)^{n_1} \dots (n_t/e)^{n_t}} \quad (24)$$

$$= \frac{(N)^N}{(n_1)^{n_1} \dots (n_t)^{n_t}} = \frac{1}{p_{n_1}^{n_1} \dots p_{n_t}^{n_t}} \quad (25)$$

Then:

$$\ln W = - \sum_{i=1}^t n_i \ln p_i = \frac{S}{k} \quad (26)$$

The equality (27) is the entropy divided by the Boltzmann constant. In the mean time wil will no worry about it. What is important it to to find a probability distribution that obeys that constrains of an energy system and the probabily sum and also maximizes the entropy. That is to say:

$$\sum_{i=1}^t p_i = 1 \quad (27)$$

And the average energy:

$$\epsilon = \frac{E}{N} = \sum_{i=1}^t \epsilon_i p_i \quad (28)$$

The solution to that is:

$$p_i^* = \frac{e^{\beta \epsilon_i}}{\sum_{i=1}^t n_i e^{\beta \epsilon_i}} \quad (29)$$

Equation (29) shows that the softmax function is indeed the probability of an event that maximizes the entropy. In other words is telling you what is the most probable macrostate possible (in terms of thermodynamics) taking certain input.

## VI. LAGRANGIAN AND HAMILTONIAN DYNAMICS

In order to study certain Energy based models it is necessary to understand the foundation of classical mechanics. The idea is that knowing the differential equations of a neural network it is then possible to derive an energy model that can help to model a system.

### A. Lagrangian mechanics

The root of Lagrangian mechanics come from the minimization of the action. Assuming that a system is described in generalized coordinates. The action is expressed as in equation (30)

$$S = \int_{x_1}^{x_2} L(\mathbf{q}, \dot{\mathbf{q}}, t) dt \quad (30)$$

The value of  $S$  is maximized or minimized when the Euler-Lagrange (31) holds:

$$\frac{\partial L}{\partial q} - \frac{d}{dx} \frac{\partial L}{\partial \dot{q}} = 0 \quad (31)$$

### B. Hamiltonian mechanics

As known from classical mechanics, the Hamiltonian is the total energy of a system as in equation (32).

$$H = T + V \quad (32)$$

And the Lagrangian is:

$$L = T - V \quad (33)$$

The Legendre transformation allows us to move from the Lagrangian to the Hamiltonian as follows:

$$E = \left( \sum_n^{i=1} \dot{q}_i \frac{\partial L}{\partial \dot{q}_i} \right) - L \quad (34)$$

## VII. BIOLOGICAL MODEL

The article from Hodgkin-Huxley describes the dynamic equations of the action potential [7]. The figure 2 shows the diagram that resembles the model. The equations that govern the model are the following:

$$I = C_m \frac{dV}{dt} + I_{input} \quad (35)$$

The total input current is:

$$I = I_{Na} + I_K + I_l \quad (36)$$

The Currents as in equations (35) are the Sodium (Na), Potassium (K) and a Leakage current (l). Those current have

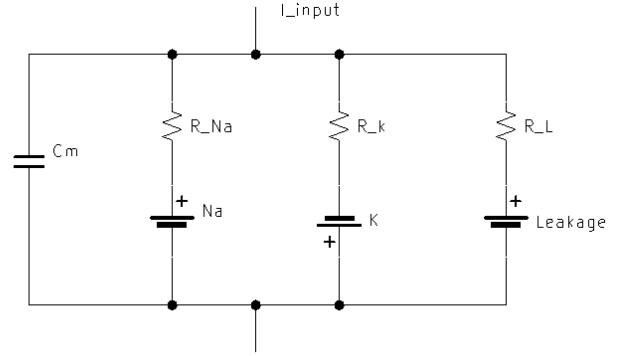


Fig. 2: The electrical equivalent circuit of the Hodgkin Huxley Model [?].

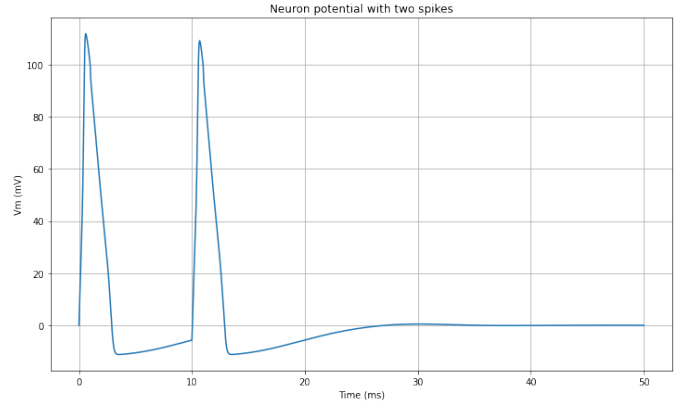


Fig. 3: The electrical response of a Hodgkin Huxley Model.

associated conductances  $g_{Na}$ ,  $g_K$  and  $g_l$ . However, the conductances are non-linear functions as they obey diffusion laws. Therefore:

$$I = C_m \frac{dV}{dt} + \bar{g}_K n^4 (V - V_K) + \bar{g}_{Na} m^3 (V - V_{Na}) + \bar{g}_l (V - V_l) \quad (37)$$

Where:

$$\frac{dn}{dt} = \alpha_n (1 - n) - \beta_n n \quad (38)$$

$$\frac{dm}{dt} = \alpha_m (1 - n) - \beta_m m, \quad (39)$$

$$\frac{dh}{dt} = \alpha_h (1 - n) - \beta_h h \quad (40)$$

In the above equations, the fact that they vary according to the voltage, allows the spikes to present transient forms different from the ones present in an electric RC circuit as in the figure 3.

Taking into account the equation (37). The model can be rewritten as:

$$I = C_m \frac{dV}{dt} + \sum_{i=1}^N M_{\mu i} g_i \quad (41)$$

Rewriting it:

$$\tau \frac{dV}{dt} = \sum_{i=1}^N \xi_{\mu i} g_i - I \quad (42)$$

## VIII. LARGE ASSOCIATIVE MEMORY PROBLEM

In the article [17], Krotov and Hopfield address the problem of creating a plausible energy model. The model is based on RBM<sup>1</sup> and deduce the equations for a feature and hidden layers. The equations are:

And get transformed into Hamiltonian energy models via a Legendre transformation.

## REFERENCES

- [1] Brain model. (2021, 11 25). *Brain Model*. [Notebook]. Available at [https://github.com/ghgv/Brain\\_model](https://github.com/ghgv/Brain_model)
- [2] H. K. Hartline, 'The response of single optic nerve fibers of the vertebrate eye to illumination of the retina', *Am J Physiol* vol 121, pp :400–415, 1938
- [3] D. H. Hubel, T. N. Wiesel, 'Receptive fields of single neurones in the cat's striate cortex' *J. Physiol.*, pp 574–591, 1959
- [4] Ken A. Dill, Sarina Bromberg, 'Molecular Driving Forces, Statistical Thermodynamics in Biology Chemistry, Physics, and Nanoscience, 2nd ed.' , pp 574–591, 1959
- [5] LeCun, Yann, 'A Path Towards Autonomous Machine Intelligence' , 2022.
- [6] McCulloch, Warren S and Pitts, Walter, 'A logical calculus of the ideas immanent in nervous activity' *The bulletin of mathematical biophysics.*, vol. 5, no. 4, pp. 115–133, 1943.
- [7] A. L. Hodgkin and A.F. Huxley, 'A quantitative description of membrane current and its application to conduction and excitation in nerve' *Physiological Laboratory, University of Cambridge.*, 1952
- [8] Hopfield, J J., 'Neural networks and physical systems with emergent collective computational abilities' *Proceedings of the National Academy of Sciences* , pp 2554–2558, 1982
- [9] Hubert Ramsauer, Bernhard Schaeffl, Johannes Lehner, Philipp Seidl , Michael Widrich , Lukas Gruber, Markus Holzleitner, Milena Pavlovic, Geir Kjetil Sandve, Victor Greiff, David P. Kreil, Michael Kopp, Genter Klambauer, Johannes Brandstetter, Sepp Hochreiter, 'Hopfield Networks is all you need' , 2021
- [10] Donald O. Hebb, 'The Organization of Behavior' , 1949
- [11] Fukushima, K., 'Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position' *Biological Cybernetics.*, pp. 193–202, 1980
- [12] Oh, Cheolhwan and Stanisław H. Żak. "The Generalized Brain-State-in-a-Box ( gBSB ) Neural Network : Model , Analysis , and Applications." (2005).
- [13] B. Olshausen, D. Field, 'Emergence of simple cell receptive field properties by learning a sparse code for natural images' *Nature*, 1996
- [14] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, ' Gradient-Based Learning Applied to Document Recognition' *Proc. of the IEEE*, 1998
- [15] Seppo Linnainmaa, 'Taylor expansion of the accumulated rounding error' *BIT , Numerical Mathematics*, vol 16, pages 146–160, 1976
- [16] Peter H. Schiller, 'Parallel information processing channels created in the retina' *Proceedings of the National Academy of Sciences* , 2010
- [17] D. Krotov, J.J. Hopfield, 'LARGE ASSOCIATIVE MEMORY PROBLEM IN NEUROBIOLOGY AND MACHINE LEARNING', *ICLR 2021*, 2021
- [18] N. V. Medathati, H. Neumann, G.S. Masson, P. Kornprobst, 'Bio-inspired computer vision: Towards a synergistic approach of artificial and biological vision', *Elsevier*, 2016
- [19] R. Rojas, 'Neural Network: A Systematic Introduction', *Springer*, 1996
- [20] A.L. Yuille, Anand Rangarajan, 'The Concave-Convex Procedure', *Neural Comput*, 2003
- [21] 'http : //vision.psych.umn.edu/users/kersten/kersten – lab/courses/Psy5036W2017/Lectures/17pythonForVision – /Demos/html/2b.Gabor.html'
- [22] M. Fee, 'https : //Introduction – to – neural – computation – spring – 2018/lecture – notes/MIT9.0S18\_Lec02.pdf'
- [23] Johannes Brandstetter, 'https : //ml – jku.github.io/hopfield – layers/'
- [24] Attention is all you need, A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, Kaiser, and I. Polosukhin. *Advances in Neural Information Processing Systems* , page 5998–6008. 2017

<sup>1</sup>Boltzmann Machines