

Synaptic Field Theory

Jaeok Yi

Department of Physics, KAIST

August 5, 2025

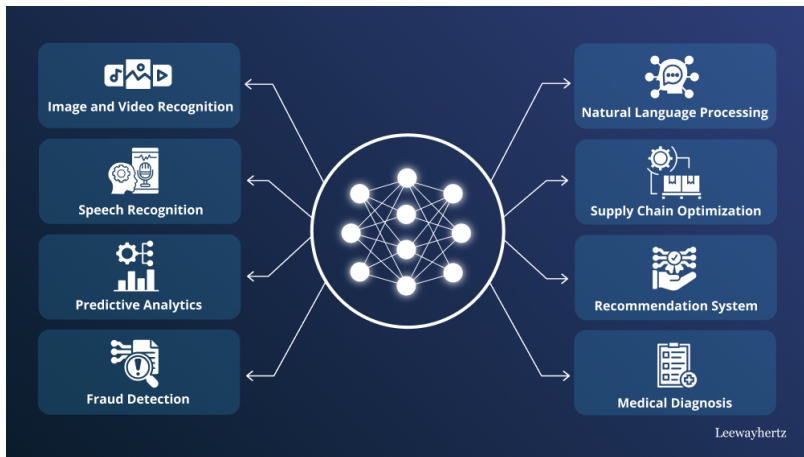
based on arXiv:2503.08827
with Donghee Lee and Hye-Sung Lee

Overview

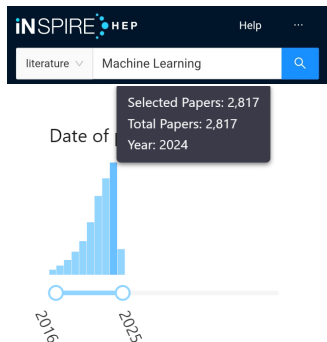
- 1 Introduction
- 2 Machine Learning 101
- 3 Synaptic Field Theory
- 4 Realization
- 5 Summary

I. Introduction

Machine Learning



Machine Learning and High Energy Physics



- Machine learning is also a topic of great interest in high energy physics.
 - Parton Distribution Function
 - Jet Classification
 - Constraining Effective Field Theories
 - Anomaly Detections
 - ...

Machine Learning is Still a Mystery

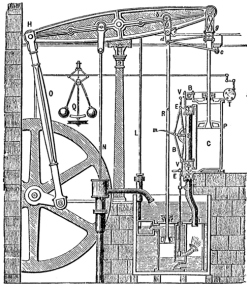


We have a rough idea of what it's doing,
but when it gets complicated,
we don't know what's going on,
similar to our understanding of the brain.

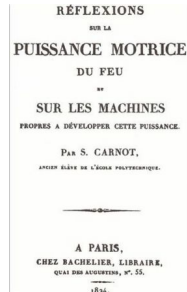
Geoffrey Hinton
(2024 Nobel Laureate in Physics)

Technology and Physics

- Technological development sometimes comes before full theoretical understanding.
 - Steam Engine & Thermodynamics



[James Watt, 1774]



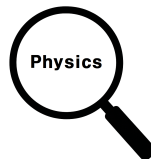
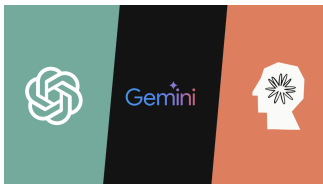
[Sadi Carnot, 1824]

Technology and Physics

- Once the physics is clear, progress tends to accelerate.
 - Steam Engine & Thermodynamics:
Invention (18C) → Thermodynamics (19C)
⇒ Steam locomotive and First industrial revolution.
 - Electromagnetic Phenomena & Maxwell's Theory:
Static Electricity, Compass (Ancient) → Maxwell's Theory (19C)
⇒ Powerplant, Telephone and Second Industrial Revolution.
 - Transistor & Semiconductor Physics:
Invention (1947) → Semiconductor Theory (1950s)
⇒ Computer, Internet and Third Industrial Revolution.

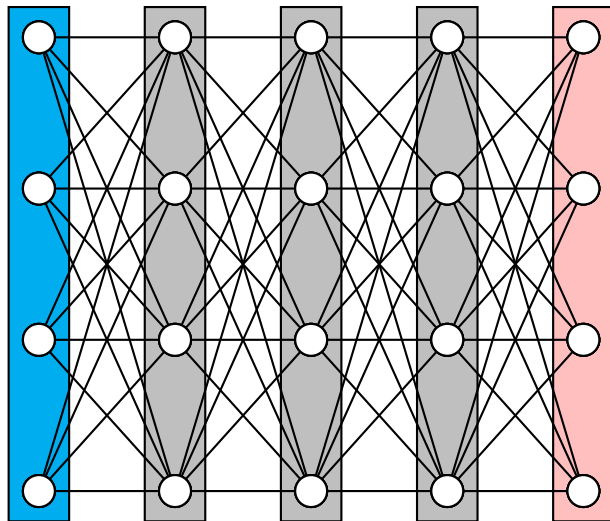
Technology and Physics

- Understanding the physics behind machine learning could drive its future progress.

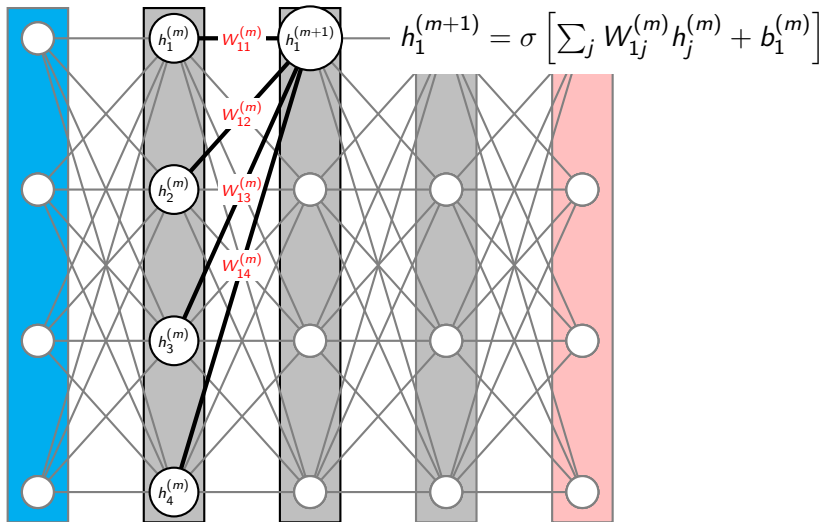


II. Machine Learning 101

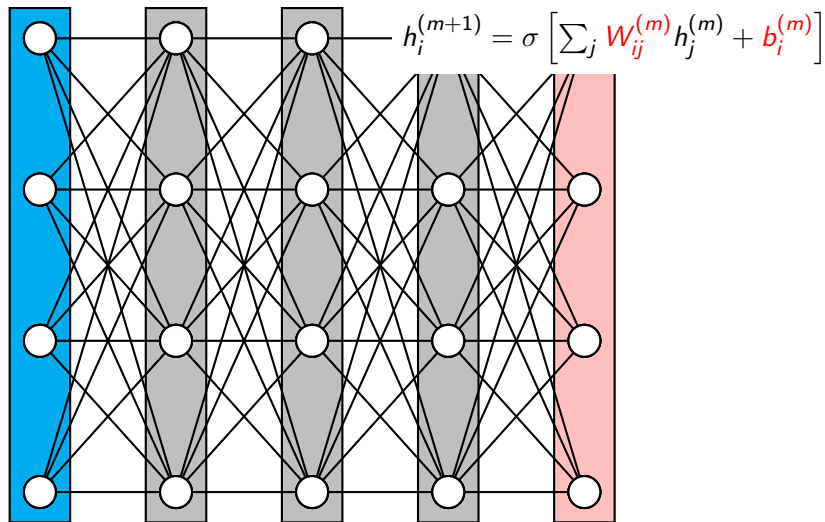
Neural Networks



Neural Networks



Neural Networks



Universal Approximation Theorem

- For any arbitrary continuous function, there exists a set of synaptic weights such that a neural network can approximate it.
- Infinite width cases: proved

Universal approximation theorem—Let $C(X, \mathbb{R}^m)$ denote the set of [continuous functions](#) from a subset X of a Euclidean \mathbb{R}^n space to a Euclidean space \mathbb{R}^m . Let $\sigma \in C(\mathbb{R}, \mathbb{R})$. Note that $(\sigma \circ x)_i = \sigma(x_i)$, so $\sigma \circ x$ denotes σ applied to each component of x .

Then σ is not [polynomial if and only if](#) for every $n \in \mathbb{N}$, $m \in \mathbb{N}$, [compact](#) $K \subseteq \mathbb{R}^n$, $f \in C(K, \mathbb{R}^m)$, $\varepsilon > 0$ there exist $k \in \mathbb{N}$, $A \in \mathbb{R}^{k \times n}$, $b \in \mathbb{R}^k$, $C \in \mathbb{R}^{m \times k}$ such that

$$\sup_{x \in K} \|f(x) - g(x)\| < \varepsilon$$

where $g(x) = C \cdot (\sigma \circ (A \cdot x + b))$

- Infinite depth or bounded depth and width cases: partially proved
- The universal approximation theorem guarantees the existence of a solution.

Training of Neural Networks

- However, the universal approximation theorem does not provide a method for finding the solution.
 - “We are not guaranteed, however, that the training algorithm will be able to learn that function.”
[Goodfellow, I., Bengio, Y., & Courville, A. (2018). Deep learning. MITP.]
- The commonly used training algorithms are **gradient descent** and its variants.
- It is still unknown whether training algorithms actually find the solutions guaranteed by the universal approximation theorem.

Gradient Descent

- Prepare the training set $(X_i^{[l]}, Y_i^{[l]})$ and then define the cost function:

$$C = \sum_{i,l} (Y_i^{[l]} - Z_i^{[l]})^2$$

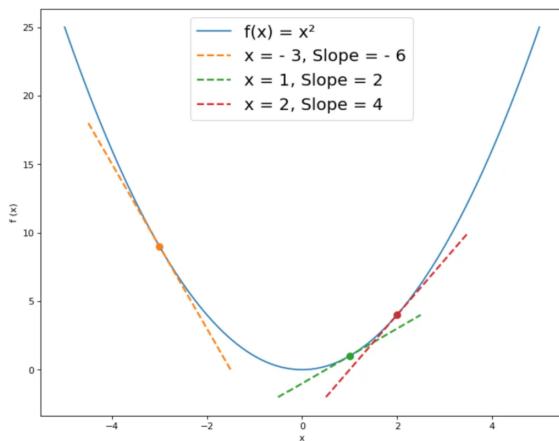
where $Z_i^{[l]}$ is the result of the neural network for $X_i^{[l]}$.

- Update the synaptic weights and biases using gradient descent:

$$\Delta W_{ij}^{(m)} = -\eta \frac{\partial C}{\partial W_{ij}^{(m)}}, \quad \Delta b_i^{(m)} = -\eta \frac{\partial C}{\partial b_i^{(m)}},$$

with the step size η .

Minimizing Cost Function



$$\Delta W = -\eta \frac{\partial C}{\partial W}$$

Issues on Gradient Descent

- Gradient descent easily gets stuck in local minima.
- Even if it reaches a global minimum, it's just the minimum of the given cost function—not necessarily the solution guaranteed by the universal approximation theorem.

- Training Dataset -

$$7 + 2 = 9$$

$$5 + 3 = 8$$

$$4 + 2 = 6$$

$$3 + 1 = 4$$

- Test Artificial intelligence -

$$5 + 4 = ?$$

Importance of Gradient Descent

- Nonetheless, almost all training algorithms are based on gradient descent.
- Nearly all of deep learning is powered by one very important algorithm: stochastic gradient descent. Stochastic gradient descent is an extension of the gradient descent algorithm.

[Goodfellow, I., Bengio, Y., & Courville, A. (2018). Deep learning. MITP.]

Various Examples of Neural Networks

- There are many ways to develop neural networks.

Various Examples of Neural Networks

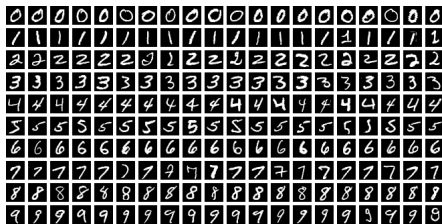
- There are many ways to develop neural networks.

(M)NIST

Input: Image of Numbers

Output: Numbers

NIST(1990), MNIST(1994)



Various Examples of Neural Networks

- There are many ways to develop neural networks.

`http://thispersondoesnotexist.com`
(2019)



Input: Noised Image, Keywords
Output: Original Image

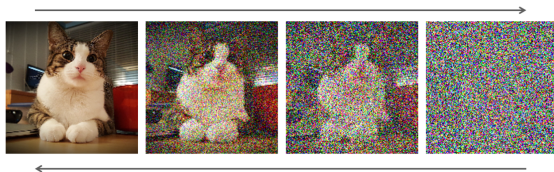
Various Examples of Neural Networks

- There are many ways to develop neural networks.

`http://thispersondoesnotexist.com`
(2019)

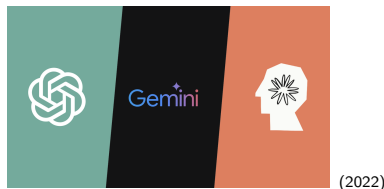


Input: Noised Image, Keywords
Output: Original Image



Various Examples of Neural Networks

- There are many ways to develop neural networks.



Input: Previous texts
Output: Next word

Various Examples of Neural Networks

- There are many ways to develop neural networks.

Text: Second Law of Robotics: A robot must obey the orders given it by human beings



Generated training examples

Example #	Input (features)	Correct output (labels)
1	Second law of robotics :	a
2	Second law of robotics : a	robot
3	Second law of robotics : a robot	must
...		

Various Examples of Neural Networks

- There are many ways to develop neural networks.

List of datasets for machine-learning research

[🌐 3 languages](#) ▾[Article](#) [Talk](#)[Read](#) [Edit](#) [View history](#) [Tools](#) ▾

From Wikipedia, the free encyclopedia

This article contains [dynamic lists](#) that may never be able to satisfy particular standards for completeness. You can help by [adding missing items with reliable sources](#).

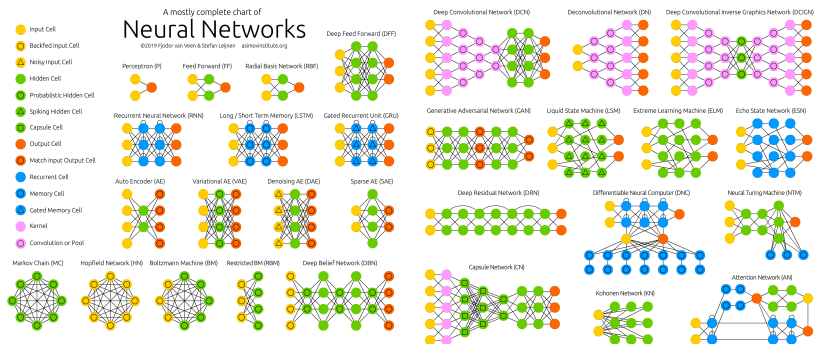
These [datasets](#) are used in [machine learning](#) (ML) research and have been cited in [peer-reviewed academic journals](#). Datasets are an integral part of the field of machine learning. Major advances in this field can

Part of a series on

**Machine learning
and data mining**

Various Examples of Neural Networks

- There are many ways to develop neural networks.



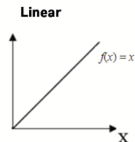
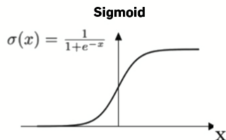
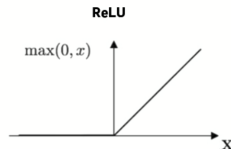
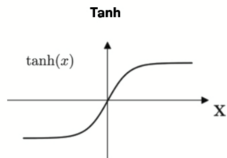
Various Examples of Neural Networks

- There are many ways to develop neural networks.

Task	Error type	Loss function	Note
Regression	Mean-squared error	$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$	Easy to learn but sensitive to outliers (MSE, L2 loss)
	Mean absolute error	$\frac{1}{n} \sum_{i=1}^n y_i - \hat{y}_i $	Robust to outliers but not differentiable (MAE, L1 loss)
Classification	Cross entropy = Log loss	$-\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] =$	Quantify the difference between two probability

Various Examples of Neural Networks

- There are many ways to develop neural networks.



Anyway, Gradient Descent

- As shown before, various neural networks can be developed by exploiting various architectures, activations, cost functions and training datasets.
- However, once those are chosen, training or optimizing the cost function is proceed according to the gradient descent or its variants.
- The gradient descent is a good beginning point to study machine learning.

Continuum Limit of Gradient Descent

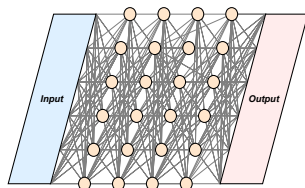
- In the continuum limit, the equation for gradient descent becomes

$$\dot{W} = -\eta \frac{\partial C}{\partial W}.$$

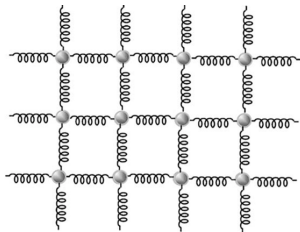
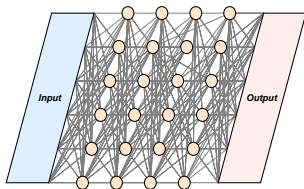
- From the perspective of physics, this equation is the equation of motion that determines the dynamics of synaptic weights and biases.
- Also, it reminds us that the fundamental degrees of freedom of neural network are synaptic weights and biases.

III. Synaptic Field Theory

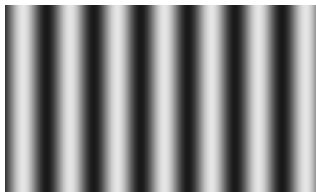
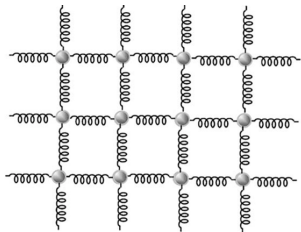
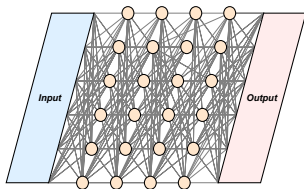
Motivation



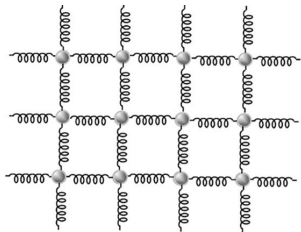
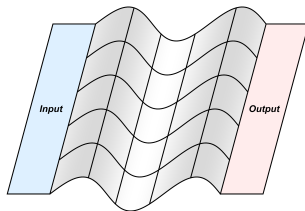
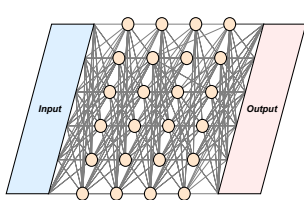
Motivation



Motivation



Motivation



Field Theoretic Approach to Neural Networks

- There exist some previous works trying to apply field theory to neural networks.
- Krippendorf and Spannowsky attempted to develop an effective theory of outputs of neural network and proposed a relationship between neural networks and cosmology.
[\[S. Krippendorf and M. Spannowsky Mach.Learn.Sci.Tech. 3 \(2022\) 3, 035011\]](#)
- To do so, they considered the limit where the effect from synaptic weights and biases becomes a constant.

Field Theoretic Approach to Neural Networks

- Since weights and biases are fundamental building blocks, their effects should not be neglected.
- Although the effective field theory is promising framework, it is still worth studying the fundamental theory.
- The theory dealing with fields developed by the continuum limit of weights and biases is worth studying.

Lagrangian Approach to Gradient Descent

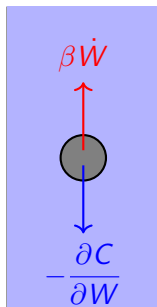
- The equation for gradient descent is

$$\dot{W} = -\eta \frac{\partial C}{\partial W}$$

- It can be considered as the high-viscosity limit ($\gamma = \eta^{-1} \gg 1$) of

$$\ddot{W} + \gamma \dot{W} + \frac{\partial C}{\partial W} = 0.$$

High-viscosity Limit



- High Viscosity Medium
- Large Drag Force
- Terminal Velocity
- $\ddot{W} = 0$

Lagrangian Approach to Gradient Descent

- This equation can be derived from the action given as

$$S = \int dt \, e^{\gamma t} \left[\frac{1}{2} \dot{W}^2 - C \right]$$

- Since the terms in the bracket have the form of a kinetic term minus a potential term, let us introduce the shorthand notation.

$$S = \int dt \, \sqrt{-g} L_W$$

with $\sqrt{-g} = e^{\gamma t}$ and $L_W = \frac{1}{2} \dot{W}^2 - C$.

Gradient Descent as the de Sitter Dynamics

- Assume that L_W admits a continuum limit, meaning it can be expressed as an integral of a Lagrangian density composed of fields:

$$L_W = \int d^d \mathbf{x} \mathcal{L}_w[w(t, \mathbf{x})].$$

- The action has the form of the action of fields in the curved spacetime:

$$S = \int d^{d+1}x \sqrt{-g} \mathcal{L}_w.$$

- In particular, $\sqrt{-g} = e^{\gamma t}$ matches that of a universe dominated by a positive cosmological constant, a typical example of de Sitter space.

Synaptic Field Theory

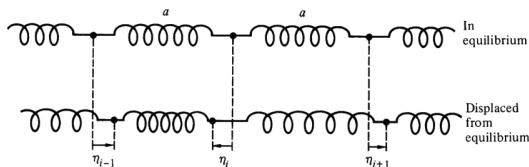
- L_W includes a sum over the indices of synaptic weights and biases.
- By taking the continuum limit of this summation to a spatial integral, we can develop a field theory in de Sitter spacetime.
- The training dataset behaves as the external sources $J(\mathbf{x})$, $K(\mathbf{x})$ in the synaptic field theory.
- The resulting **synaptic field theory** would be a familiar framework to those who study high energy physics or cosmology.

IV. Realization

Continuum Limit

- Here is a typical example of taking continuum limit.

[H. Goldstein, C. Poole, J. Safko (2002). Classical Mechanics, Pearson.]



$$L = \frac{1}{2} \sum_i [m\dot{\eta}_i^2 - k(\eta_{i+1} - \eta_i)^2] \quad \Rightarrow \quad L = \frac{1}{2} \int \left[\mu \dot{\eta}^2 - Y \left(\frac{d\eta}{dx} \right)^2 \right] dx$$

Continuum Limit

- In the continuum limit, we substitute the summation with the integral.
- The indices are promoted to the spatial coordinates.
- The difference between nearby degrees of freedom is promoted to the spatial derivative.

$$L = \frac{1}{2} \sum_i [u \dot{\eta}_i^2 - k(\eta_{i+1} - \eta_i)^2 - m \eta_i^2]$$

$$\Rightarrow L = \frac{1}{2} \int \left[U \dot{\eta}^2 - K \left(\frac{d\eta}{dx} \right)^2 - M \eta^2 \right] dx$$

Comments on Locality

- The previous example gives a local Lagrangian because every term involves only variables with the same index.
- Series expansion of cost function:

$$C = \sum J_{1 \ i_1 j_1}^{(m_1)} W_{i_1 j_1}^{(m_1)} + \sum J_{2 \ i_1 j_1 i_2 j_2}^{(m_1 m_2)} W_{i_1 j_1}^{(m_1)} W_{i_2 j_2}^{(m_2)} + \dots$$

The coefficients $J_{1 \ i_1 j_1}^{(m_1)}$ and $J_{2 \ i_1 j_1 i_2 j_2}^{(m_1 m_2)}$ depend on the data set.

- Note that there are terms involving different indices.

Nonlocality of Neural Networks

- Taking the continuum limit,

$$L \supset \int d^3\mathbf{x} J_1(\mathbf{x})w(t, \mathbf{x}) + \int d^3\mathbf{x}d^3\mathbf{y} J_2(\mathbf{x}, \mathbf{y})w(t, \mathbf{x})w(t, \mathbf{y}) + \cdots$$

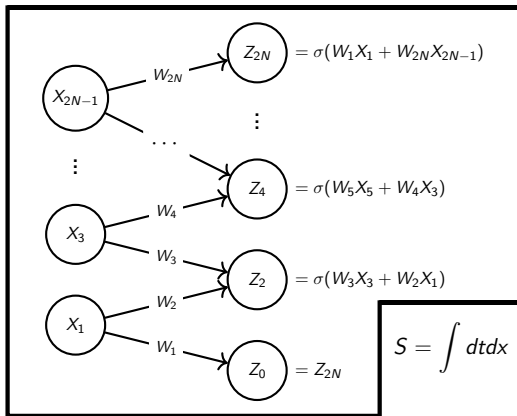
Here, J_1 and J_2 act as external sources determined by the training examples.

- Since the second term involves two spatial coordinates \mathbf{x} and \mathbf{y} , this Lagrangian is not local.
- In general, it is difficult to study a nonlocal theory.

Spacetime Geometry and Neural Network Architecture

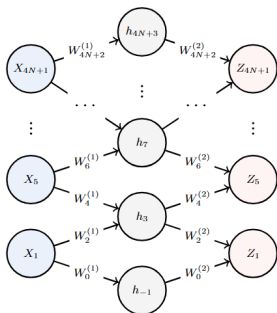
- This naive approach expects the nonlocal Lagrangian.
- The locality is related to the spacetime geometry.
- In the synaptic field theory, the spacetime is given as the continuum limit of the indices of parameters.
- The spatial geometry depends on how to construct the architecture and how to index the parameters.
- We may construct a neural network possessing locality.

Toy Neural Network



$$S = \int dt dx \sqrt{-g} \left[\frac{1}{2} [\partial_t w(t, x)]^2 - \frac{1}{2} K(x) [\partial_x w(t, x)]^2 - \frac{1}{2} J(x) w(t, x)^2 \right]$$

Toy Neural Network



$$S = \int dt dx \sqrt{-g} \left[\frac{1}{2} (\partial_t w_1)^2 + \frac{1}{2} (\partial_t w_2)^2 - \frac{1}{2} m^2 w_2^2 \right. \\ \left. - J_1 - J_2 w_2 - J_3 w_2 w_1 - K_1 w_2 \partial_x w_1 \right. \\ \left. - K_2 w_2 \partial_x^2 w_2 - K_3 w_2 \partial_x^2 w_1 - K_4 \partial_x w_2 \partial_x w_1 \right. \\ \left. - K_5 w_1 \partial_x w_2 - K_6 \partial_x^2 w_2 - K_7 w_1 \partial_x^2 w_2 + \dots \right].$$

$$m^2 = 8a^{-1} N_l r^2$$

$$J_1(x) = a^{-1} \sum_l (Y^{[l]})^2 \quad J_2(x) = 4ra^{-1} \sum_l Y^{[l]}$$

$$J_3(x) = 8qa^{-1} \sum_l (X^{[l]} + 4a^2 \partial_x^2 X^{[l]}) Y^{[l]}$$

$$K_1(x) = 48aq \sum_l Y^{[l]} \partial_x X^{[l]} \quad K_2 = 4aN_l r^2$$

$$K_3(x) = 20qa \sum_l X^{[l]} Y^{[l]} \quad K_4(x) = 16aq \sum_l X^{[l]} Y^{[l]}$$

$$K_5(x) = 16aq \sum_l Y^{[l]} \partial_x X^{[l]} \quad K_6(x) = 2ra \sum_l Y^{[l]}$$

$$K_7(x) = 4qa \sum_l X^{[l]} Y^{[l]}$$

Discussions on Toy Neural Network

- These examples may be too simple to behave as a practical artificial intelligence.
- However, it is an interesting example that shows the locality.
- This locality comes from the architecture of the neural network and the indexing convention.
- One may attempt to develop an architecture and indexing convention that enables locality to emerge for general neural networks.

Further Remarks

- External sources in the previous example only have the spatial dependence.
- For sources to have explicit time dependence, we may consider the training algorithm, such as the stochastic gradient descent, involving time dependence.
- By further pushing this possibility, we may consider the system with interesting symmetry, such as Lorentz symmetry.
- It would be interesting to embed cosmological dynamics into neural networks.

V. Summary

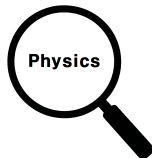
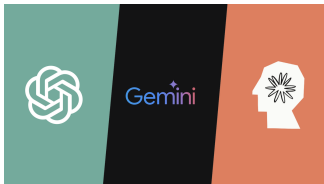
Summary Table

- The synaptic field theory suggests a friendly framework for physicists to study machine learning.

Neural Network	Synaptic Field Theory
Parameters $W_{ij}^{(m)}$	Field $w(t, x)$
Training examples (X, Y)	External sources K, J, \dots
Indices i, j, m	Space x
Training step T	Time t
Cost function C	Lagrangian L
Step size η	Hubble parameter H

Take-home Message

- Understanding the nature of deep learning is the mission of physicists so more physics is warranted.



Thank you for listening

Back-up Slides

Discussions

- The linear activation is used and results in the bilinear action.
 - If we use non-polynomial activations, then higher order interaction terms should be considered.
- The previous example is not Lorentz invariant.
 - With a time-dependent training algorithm such as stochastic gradient and carefully designed training datasets, a Lorentz-invariant theory may emerge.
- The previous example results in the local action due to the simple structure and the practical indexing convention.
 - For the finite depth neural networks, the naive approach of continuum limit will result in nonlocal terms like

$$\int d^d \mathbf{x} d^d \mathbf{y} A(\mathbf{x}, \mathbf{y}) w(t, \mathbf{x}) w(t, \mathbf{y})$$

- As shown in the example, the structure and the indexing convention will be related to the geometry and locality of the space of the theory.

High Viscosity Limit

- In the high-viscosity limit, $\eta = 1/\gamma$ is small, allowing a perturbative expansion:

$$W = \mathbb{W}^{(0)} + \eta \mathbb{W}^{(1)} + \mathcal{O}(\eta^2).$$

- The equation from the action becomes

$$\frac{1}{\eta} \dot{\mathbb{W}}^{(0)} + \left(\ddot{\mathbb{W}}^{(0)} + \dot{\mathbb{W}}^{(1)} + \left. \frac{\partial \mathcal{C}}{\partial W} \right|_{W=\mathbb{W}^{(0)}} \right) + \mathcal{O}(\eta) = 0.$$

- At $\mathcal{O}(\eta^{-1})$, we find $\dot{\mathbb{W}}^{(0)} = 0$ at any t , which implies $\ddot{\mathbb{W}}^{(0)} = 0$. Therefore, at $\mathcal{O}(\eta^0)$, we have

$$\dot{\mathbb{W}}^{(1)} + \left. \frac{\partial \mathcal{C}}{\partial W} \right|_{W=\mathbb{W}^{(0)}} = 0.$$

It is the equation of motion for the training of neural networks.