

# Interpreting Artificial Neural Networks in the Context of Theoretical Physics

Canada

NRC·CNRC

DERIMETER

NSTITUTE

**PIQUIL** 

Sebastian Johann Wetzel

# Success of Artifical Neural Networks

#### Image Classification (Convolutional Network)





Generative Modelling / Anomaly Detection (Autoencoders)

Similarity Detection (Siamese Network)



# (Supervised) Machine Learning with Neural Nets

"Machine learning is the subfield of computer science that gives computers the ability to learn without being explicitly programmed." - Wikipedia

#### **Training Data**



# (Supervised) Machine Learning with Neural Nets

"Machine learning is the subfield of computer science that gives computers the ability to learn without being explicitly programmed." - Wikipedia

**Training Data** 



 What does the Neural Network actually learn?
Consthis Kneudeders hale in Opiontifie

2) Can this Knowledge help in Scientific Discovery?

?

Cats

Dogs



Dog

# Overview

- X Artificial Neural Networks
- x Interpretation Techniques
- x Interpretation of Convolutional Neural Networks
- x Interpretation of Autoencoders
- x Interpretation of Siamese Networks

# **Artificial Neural Networks**

#### Feed forward neural network



Input: Data  $X = (\vec{x}_1, ..., \vec{x}_n)$ , Label  $Y = (y_1, ..., y_n)$ Output:  $Y_{pred} = F(X, w_{ij}^L, b_i^L)$ 

Goal: choose  $w_{ij}^L$  and  $b_i^L$  such that  $Y_{pred} \approx Y$ 



# Interpretation Techniques



# Bottleneck Interpretation +Correlation Probing Neural Network





# Looking at the weights

*×* No, works but only for the most simple problems.

# **Influence** Functions

#### Phase Detection with Neural Networks: Interpreting the Black Box

Anna Dawid,<sup>1,2</sup> Patrick Huembeli,<sup>2</sup> Michał Tomza,<sup>1</sup> Maciej Lewenstein,<sup>2,3</sup> and Alexandre Dauphin<sup>2</sup>

- \* Remove specific datapoints or features and measure the effect on the performance
- \* Largest change in performance indicates the most influential data point or feature

# **Dark Matter**

#### Discovering Symbolic Models from Deep Learning with Inductive Biases

Miles Cranmer <sup>1</sup>	Alvaro Sanchez-Gonzalez <sup>2</sup>	Peter Battaglia <sup>2</sup>	Rui Xu $^1$
Kyle Cranmer	<sup>3</sup> <b>David Spergel</b> <sup>4,1</sup>	Shirley Ho	4,3,1,5

#### × Simulate Dark Matter

\* Apply symbolic regression at the output of a graph neural network to recover force equation

Cranmer et al., Neurips 2020

# Condensed Matter+Correlator Network

#### Correlator Convolutional Neural Networks: An Interpretable Architecture for Image-like Quantum Matter Data

Cole Miles,<sup>1</sup> Annabelle Bohrdt,<sup>2, 3, 4</sup> Ruihan Wu,<sup>5</sup> Christie Chiu,<sup>2, 6, 7</sup> Muqing Xu,<sup>2</sup> Geoffrey Ji,<sup>2</sup> Markus Greiner,<sup>2</sup> Kilian Q. Weinberger,<sup>5</sup> Eugene Demler,<sup>2</sup> and Eun-Ah Kim<sup>1</sup>

#### \* Explicit feature engineering layer that probes for correlations

\* Dominant features correspond to dominant correlations in condensed matter system

Miles et al., Arxiv 2020

# **Physical Concepts**

#### Discovering physical concepts with neural networks

Raban Iten,<sup>\*</sup> Tony Metger,<sup>\*</sup> Henrik Wilming, Lídia del Rio, and Renato Renner *ETH Zürich, Wolfgang-Pauli-Str.* 27, 8093 Zürich, Switzerland. (Dated: January 24, 2020)

#### × Interpretation of autoencoder latent representation

× Ask physical questions to be extractable from latent space

Iten et al., PRL 2020



# Bottleneck Interpretation +Correlation Probing Neural Network



# **Bottleneck Interpretation**

Interpretation is often difficult since information is spread over several neurons and layers

If the neuron contains the information of <u>one</u> single — quantity/obervable Q(S)

- Idea: identify or enforce bottlenecks in the network
- Perform regression on the output of the bottleneck neuron

The output of the neuron can be mapped via a bijective function to the observable

F(S) = f(Q(S))



### Supervised Learning 2d Ising Model

- > Data: Monte Carlo samples
- Training at well known points in phase diagram
- Labels: Phase



- Testing in interval containing phase transition
- > Estimate within 1% of exact value  $T_c = \frac{2}{\ln(1+\sqrt{2})}$



# **Artificial Neural Networks**



Output:  $Y_{pred} = F(X, w_{ij}^L, b_i^L)$ 

Goal: choose  $w_{ij}^L$  and  $b_i^L$  such that  $Y_{pred} \approx Y$ 

### Interpretation of Neural Network 2d Ising Model



- Correlation Probing Net interpolates between a general NN and a minimal optimal NN which has the same performance
- Interpretation by reducing the NN capacity in an ordered manner until one observes a performance drop
- > Inspired by intensive/extensive quantities (averaging layer probes for translational invariance of the quantity Q(S))

### Interpretation of Neural Network 2d Ising Model

### Decision functions $F(S) = \operatorname{sigmoid}(w Q(S) + b)$

$$\succ Q(S) = |1/N\sum_{i} s_i|$$

$$\Rightarrow Q(S) = \frac{1}{N} \sum_{\langle i,j \rangle_{nn}} s_i s_j$$

Deduction visually confirmed:

#### Note:

1x2 Network also has the Magnetization minimum which is easier to find!

Receptive Field Size	Train Loss	Validation Loss
$28 \times 28$	6.1588e - 04	0.0232
<b>1</b> imes <b>2</b>	$1.2559\mathrm{e}\text{-}04$	$1.2105\mathrm{e} extsf{-}07$
<b>1</b> imes <b>1</b>	0.2015	0.1886
baseline	0.6931	0.6931

#### Magnetization

#### Expected Energy per site



# SU(2) Lattice Gauge Theory



Quarks on heavy static lattice sites.

Gluons on the connections between lattice sites are described by Matrices



# SU(2) Lattice Gauge Theory

Data: Monte Carlo samples

$$S_{\text{Wilson}}[U] = \beta_{\text{latt}} \sum_{x} \sum_{\mu < \nu} \text{Re tr} \left( 1 - U_{\mu\nu}^x \right)$$

- Training at well known points in phase diagram
- Labels: Phase

Find phase transition close to lattice calculation



#### Interpretation of Neural Network SU(2) Gauge Theory



Polyakov Loop

### (Variational) Autoencoder 2d Ising Model



**Objective: Minimize Reconstruction error** 

$$MSE = \frac{1}{N} \sum_{k} \left\| x_k - F(x_k) \right\|^2$$

- > Data: Monte Carlo samples
- > Train everywhere in phase diagram
- Labels: None



### (Variational) Autoencoder 2d Ising Model



Ferromagnetic Ising model on the square lattice

Wetzel, PRE 2017

- Latent parameter corresponds to magnetization
- Identification of phases: Latent representations are clustered
- Location of phases: Magnetization, latent parameter and reconstruction loss show a steep change at the phase transition.

# Siamese Neural Networks



- Input : Pair of data points
- Label : same / different
- Network pair contains identical neural networks with shared weights

# Machine Learning Multi Class Classification

"Machine learning is the subfield of computer science that gives computers the ability to learn without being explicitly programmed." - Wikipedia



# Machine Learning Infinite Class Classification

Reformulation of the Problem:

Teach a maching learning algorithm if two pictures show the same class.



### Siamese Neural Networks Particle in Gravitational Potential

Problem:

Given two observations of positions and velocities, do they belong to the same particle trajectory?



SNN Solution:

Prepare Dataset of positive data where the pair is connected by solving the equations of motion

$$((x, y, v_x, v_y), (x', y', v'_x, v'_y))$$

- Prepare Negative Dataset by permuting positive dataset
- > Train SNN to distinguish between positive and negative pairs

### Siamese Neural Networks Particle in Gravitational Potential

**Results:** 

)

Training accuracy : 98% Test accuracy : 97%

Interpretation by polynomial regression on latent representation:

$$f(\mathbf{x}) \approx -403.71xv_y - 4.85x - 0.58xy -0.17xv_x - 0.02v_y^2 - 0.01v_xv_y +0.00v_y^2 + 0.01v_y + 0.02v_x +0.45x^2 + 0.66y^2 + 0.74 +0.99yv_y + 1.24y + 402.44yv_x \approx -403(xv_y - yv_x) = L_z$$



400

200

-10

-400

-200

0

intermediate output

Network has learned the angular momentum to infer its prediction.

### Siamese Neural Networks Lorentz Transformation of Electromagnetic Fields

Problem:

Given two field configurations, can they be transformed into each other by a Lorentz transformation?



SNN Solution:

 Prepare Dataset of positive data where the pair is connected by a Lorentz Transformation

 $((E_x, E_y, E_z, B_x, B_y, B_z), (E'_x, E'_y, E'_z, B'_x, B'_y, B'_z))$ 

- Prepare Negative Dataset by permuting pointive dataset
- > Train SNN to distinguish between positive and negative pairs

#### Siamese Neural Networks Lorentz Transformation of Electromagnetic Fields

**Results:** 

Training accuracy : 95% Test accuracy : 94%

Interpretation by polynomial regression on latent representation:

$$f(\mathbf{x}) \approx -170.53E_2B_2 - 170.22E_1B_1 - 170.20E_3B_3$$
$$-4.13B_3^2 + \dots + 4.92E_2^2 + 53.43$$
$$\approx -170\underbrace{(E_1B_1 + E_2B_2 + E_3B_3)}_{=E \cdot B} + 53$$



Network has learned the determinant of the field strength tensor to infer its prediction.

# Summary

- \* Interpretation of Artificial Neural Networks is hard because information is distributed among many layers and neurons
- \* Interpretation is possible by identifying bottlenecks and performing regression
- \* Interpretation is constructive and can give insight into the underlying physics:

Neural Networks applied to phase recognition learn order parameters or energies

Siamese Networks for similarity detection learn invariants or conserved quantities



# **Twin Neural Network Regression**



Solution of the Original Regression Problem:



## **Bias-Variance Tradeoff**



#### **TNN** implicit ensemble

$$y_i^{pred} = \frac{1}{m} \sum_{j=1}^m F(x_i, x_j^{train}) + y_j^{train} = \frac{1}{m} \sum_{j=1}^m \frac{1}{2} F(x_i, x_j^{train}) - \frac{1}{2} F(x_j^{train}, x_i) + y_j^{train}$$

- Get huge ensemble of twice the training data set size
- Ensemble is relatively uncorrelated, since the predicted differences are different by construction

# **Uncertainty Signal**

#### Do ensemble members agree?

$$y_i^{pred} = \frac{1}{m} \sum_{j=1}^m F(x_i, x_j^{train}) + y_j^{train} = \frac{1}{m} \sum_{j=1}^m \frac{1}{2} F(x_i, x_j^{train}) - \frac{1}{2} F(x_j^{train}, x_i) + y_j^{train}$$

- Uncorrelated predictions make different mistakes
- Measure ensemble standard deviation



(additional uncertainty signal based on loop consistencies)

# **Semi-Supervised Learning**



- Train to enforce loop consistency during training
- Loops can be used as training data even if the data points within them are unlabelled

$$0 = F(x_i, x_j) + F(x_j, x_k) + F(x_k, x_i)$$

It can be viewed as two predictions provide a suggested label for the third.