



I. Mysteries of deep learning

Yaoyu Zhang

Institute of Natural Sciences & School of Mathematical Sciences

Shanghai Jiao Tong University

FAU MoD Course 饮水思源•爱国荣校

Deep learning is no longer a black-box



Friedrich-Alexander-Universität **Research Center for** Mathematics of Data | MoD

FAU MoD Course



Towards a mathematical foundation of Deep Learning: From phenomena to theory

Yaoyu Zhang

SHANGHAI JIAO TONG UNIVERSITY



Fri.-Thu. May 2-8, 2025 10:00H (Berlin time)

On-site / Online

Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU) Room H11 / H16 Felix-Klein building Cauerstraße 11, 91058 Erlangen. Bavaria, Germany

Live-streaming: www.fau.tv/fau-mod-livestream-2025

Towards a Mathematical Foundation of Deep Learning: From Phenomena to Theory

Date

Fri. – Thu. May 2 – 8, 2025

Session Titles

- Mysteries of Deep Learning
- 2. Frequency Principle/Spectral Bias
- 3. Condensation Phenomenon
- 4. From Condensation to Loss Landscape Analysis
- From Condensation to Generalization 5

Session Titles:

Analysis

1. Mysteries of Deep Learning 2. Frequency Principle/Spectral Bias

3. Condensation Phenomenon

its advancement and application.

4. From Condensation to Loss Landscape

5. From Condensation to Generalization Theory

Overall, this course serves as a gateway to the vibrant field of deep learning theory, inspiring

participants to contribute fresh perspectives to

Establishing a mathematical foundation for deep learning is a significant and challenging endeavor in mathematics. Recent theoretical advancements are transforming deep learning from a black box into a more transparent and understandable framework. This course offers an in-depth exploration of these developments. emphasizing a promising phenomenological approach. It is designed for those seeking an intuitive understanding of how neural networks learn from data, as well as an appreciation of their theoretical underpinnings. (...)

Unimaginable achievements of Al





Image recognition





^{30%} 25% 20% 15% Rate in 10% Human Performance Zone 5% 0% XRCE AlexNet ZFNet NEC-UIUC GoogLeNet ResNet SENet (2010)(2011)(2012)(2013)(2014)(2015)(2017)

https://www.linkedin.com/pulse/must-read-pathbreaking-papers-image-classification-muktabh-mayank

Krizhevsky, et al, 2012



Go playing























T1037 / 6vr4 90.7 GDT (RNA polymerase domain)

T1049 / 6y4f 93.3 GDT (adhesin tip)

Experimental resultComputational prediction

The Nobel Prize in Chemistry 2024



"for computational protein design"



© Nobel Prize Outreach. Photo: Clément Morin

Demis Hassabis

"for protein structure prediction"



© Nobel Prize Outreach. Photo: Clément Morin

John Jumper

"for protein structure prediction"



© Nobel Prize Outreach. Photo: Clément Morin







Jumper et al., 2021



Image&video generation

/v5_upscale

Woman with a Cactus Hat artwork by Edmund Dulac and Christian Schloe

17 hrs ago

(a) yaros89



alexzz











https://openai.com/index/sora/

https://www.midjourney.com/





https://bootcamp.uxdesign.cc/how-stable-diffusion-works-explained-for-non-technical-people-be6aa674fa1d

Large language model—milestone towards AGI



https://openai.com/

https://cshub.in/what-is-turing-test/





The engine of AI: deep learning

Empirical risk:
$$R_S(\theta) = \frac{1}{n} \sum_{i=1}^n l(f(x_i, \theta), y_i)$$

Model: $f(x, \theta)$
Data: $S = \{(x_i, y_i)\}_{i=1}^n$

Just this?

Common Models:

Linear models: polynomial models, random feature models, … Neural networks: fully-connected, convolutional, ResNet, Transformer, …

Common loss function:

Mean-squared error (I2) loss: $l(y,y') = ||y-y'||_2^2$, Cross entropy, Hinge loss, ...

Common training algorithm:

Gradient decent (GD): $\theta^{t+1} = \theta^t - \eta \nabla R_S(\theta^t)$, Stochastic gradient descent (SGD), Adam, ...



Deep learning remains a "black" technology





AlphaGo, AlphaFold, ChatGPT, SORA,

• • •

______SJTUM

https://www.ibm.com/cloud/learn/neural-networks

Synthetic diamond

- Atomic bomb
- The Apollo Program
- ChatGPT
- Quantum computer
- I light-speed spaceship



Bitter lesson for deep learning theory









The Bitter Lesson

Rich Sutton

March 13, 2019

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin. The ultimate reason for this is Moore's law, or rather its generalization of continued exponentially falling cost per unit of

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin."

Leverage computation (learning) instead of human knowledge



http://www.incompleteideas.net/IncIdeas/BitterLesson.htm

Timeline of neural network development



https://beamandrew.github.io/deeplearning/2017/02/23/deep_learning_101_part1.html



- Igentiation 1969 book of Perceptrons (lead to the first winter)
- 1986 Backpropagation (emergence of modern deep learning)
- 1989 Universal approximation theorem
- Igeneralization puzzle proposed (not well solved till now)
 - The Vapnik-Jackel Bet (witnessed by Yann Lecun)
- @2018 Neural Tangent Kernel (lead to a surge in DL theory research)
 - **Frequency principle/Spectral bias**

Despite 40 years of effort, framework for its math foundation yet to emerge



The bet on deep learning theory

The Vapnik-Jackel Bet in 1995

1. Jackel bets (one fancy dinner) that by March 14, 2000 people will understand

quantitatively why big neural nets working on large databases are not so bad. (Understanding means that there will be clear conditions and bounds)

Vapnik bets (one fancy dinner) that Jackel is wrong.

But .. If Vapnik figures out the bounds and conditions, Vapnik still wins the bet.

2. Vapnik bets (one fancy dinner) that by March 14, 2005, no one in his right mind will use neural nets that are essentially like those used in 1995.

Jackel bets (one fancy dinner) that Vapnik is wrong

3/14/95

3/14//95

3/14/95

V. Vapnik

1 m

L. Jackel

Witnessed by Y. LeCun





From Lecun's talk





Intelligent Machines

The Dark Secret at the Heart of Al

No one really knows how the most advanced algorithms do what they do. That could be a problem.

by Will Knight April 11, 2017

L

ast year, a strange self-driving car was released onto the quiet roads of Monmouth County, New Jersey. The experimental vehicle, developed by researchers at the chip maker Nvidia,

didn't look different from other autonomous cars, but it was unlike anything demonstrated by Google, Tesla, or General Motors, and it showed the rising power of artificial intelligence. The car didn't follow a single instruction provided by an engineer or programmer. Instead, it relied entirely on an algorithm that had taught itself to drive by watching a human do it.







Theory of deep learning?



Donoho's PPT, Stats 385 Stanford







Figure : Every theorist who looks at it see what they wish

Donoho's PPT, Stats 385 Stanford

A (personal) bitter lesson:

All previously existing frameworks, irrespective of their origin or demonstrated success, are ineffective for understanding deep learning.

Existing frameworks:

statistical learning theory, numerical analysis, statistical physics, statistics, optimization, neuroscience, psychology, ...



Existing frameworks often mislead



In face of deep learning, all of us are blind men.



A phenomenological methodology for blind men

- Suspension: Suspend the prior and belief one may hold and focus on the facts about the object.
- 2. Cumulation: Discover and cumulate all possible facts about the object. Prioritize the more informative ones.
- 3. Emergence: A new framework shall emerge once enough pieces are uncovered.



https://www.sloww.co/blind-men-elephant/









Suspension

Cumulation

Emergence





Phenomenon as a key family of facts to uncover

- Frequency principle/spectral bias
- Condensation
- Double descent
- Edge of stability
- Lottery ticket
- Neural collapse
- Grokking







Basics of deep learning theory





Single artificial neuron:

$$f_{\theta}(\boldsymbol{x}) = \sigma(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + \boldsymbol{b})$$

Parameters (weights): $\theta = (w, b)$, activation function: $\sigma(\cdot) : \mathbb{R} \to \mathbb{R}$

Illustration:









Deep neural networks:



$$m{ heta} := \left(m{W}^{[1]}, m{b}^{[1]}, \dots, m{W}^{[L]}, m{b}^{[L]}
ight)$$
 $m{f}_{m{ heta}}^{[l]}(m{x}) := \sigma(m{W}^{[l]}m{f}_{m{ heta}}^{[l-1]}(m{x}) + m{b}^{[l]})$



https://www.ibm.com/cloud/learn/neural-networks

Universal Approximation Theorem



Neural networks with a single hidden layer can be used to approximate any continuous function to any desired precision.

Cybenko 89, Hornik 89, Hornik 91, Barron 93

Requirement for transfer function:

 $\sigma(z)$ is well-defined as $z \to -\infty$ and $z \to \infty$

$$\left| f(x) - \sum_{j} k_{j} \sigma \left(w_{ij} x_{i} + b_{j} \right) \right| < \epsilon$$

Sketch of a constructive proof:

 Construct Heaviside function from the given transfer function
 Construct "bump" function (1-d) or "tower" function (2-d)
 Approximate the target continuous function with "bump" or "tower" functions

Illustration of constructive proof (three layer)









No Free Lunch Theorem (Wolpert and Macready)

Theorem—Given a finite set V and a finite set S of real numbers, assume that $f: V \to S$ is chosen at random according to uniform distribution on the set S^V of all possible functions from V to S. For the problem of optimizing f over the set V, then no algorithm performs better than blind search.



https://en.wikipedia.org/wiki/No_free_lunch_theorem

Generalization

- Optimization
- Approximation
- Robustness
- Interpretability









Despite strongly nonconvex loss landscape, gradient-based training of large DNNs often find global minima.



What is the geometry of loss landscape?





Some architectures are more parameter efficient than others regarding particular class of tasks.

Ex:CNN vs. FNN for image, Transformer vs. LSTM for language

Problem

How to quantify the difference in parameter efficiency between architectures?





Problem

Output of well-trained DNNs are often susceptible to tiny

adversarial perturbation.



x

"panda" 57.7% confidence



"nematode" 8.2% confidence



 $x + \epsilon sign(\nabla_x J(\theta, x, y))$ "gibbon" 99.3 % confidence

Goodfellow et al.

Why is that? How to improve robustness?





One can hardly obtain an explanation with prediction power.



1: Great_Pyrenees

1:Great_Pyrenees/kuvasz



Problem

When is it possible to obtain explanations with prediction





Generalization puzzle of deep learning





Learning systems with increasingly large size



Suzana Herculano-Houzel, 2009

Parameters of transformer-based language models



¢2023 TECHTARGET, ALL RIGHTS RESERVED TECHTARGET





"With four parameters you can fit an elephant to a curve; with five you can make him wiggle his trunk." -- John von Neumann



Complex models easily overfit.







Large complexity → Large generalization gap Generalization Gap



Occam Razor: Entities should not be multiplied unnecessarily





Leo Breiman

1995

Statistics Department, University of California, Berkeley, CA 94305; e-mail: leo@stat.berkeley.edu

Reflections After Refereeing Papers for NIPS

Our fields would be better off with far fewer theorems, less emphasis on faddish stuff, and much more scientific inquiry and engineering. But the latter requires real thinking.

For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?



Modern verification of generalization mystery

Benjamin Recht[†]

brecht@berkeley.edu





UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

Chiyuan Zhang* Massachusetts Institute of Technology chiyuan@mit.edu

Samy Bengio **Google Brain** bengio@google.com

Moritz Hardt Google Brain mrtz@google.com

University of California, Berkeley

Oriol Vinyals Google DeepMind vinyals@google.com

Cifar10: 60,000 training data

model	# params	random crop	weight decay	train accuracy	test accuracy
Inception	1,649,402	yes yes	yes no ves	100.0 100.0 100.0	89.05 89.31 86.03
		no	no	100.0	85.75
(fitting random labels)		no	no	100.0	9.78

Zhang et al., 2017

Generalization mystery in 1-d interpolation

Find an interpolation of $\mathscr{D}: \{(x_i, y_i)\}_{i=1}^n$ in $\mathscr{H}: \{h(\cdot; \Theta) | \Theta \in \mathbb{R}^m\}$

Example:

 $h(x;\Theta) = \theta_1 + \theta_2 x + \dots + \theta_M x^{m-1}$ with m = n





Thanks!

