

# 机器学习理论简明手册

Theoretical Machine Learning:  
A Handbook for Everyone

滕佳烨、王伯元、张臻

[www.tengjiaye.com/mlbook](http://www.tengjiaye.com/mlbook)

2025 年 10 月 6 日      第一版



# 目录

<b>第一章 Week 1 - Representation</b>	<b>1</b>
1.1 The Universal Approximation Theorem (UAT)	1
1.1.1 The Core Idea	1
1.1.2 Key Takeaways from UAT	2
1.2 Representation and the Power of Depth	2
1.2.1 Shallow vs. Deep Representation	2
1.2.2 The Theoretical Advantage of Depth	3
1.2.3 Summary Table	3
<b>第二章 Week 2 - Representation</b>	<b>5</b>
2.1 The Problem: Adversarial Examples	5
2.2 A Universal Law of Robustness via Isoperimetry	6
2.2.1 Connecting Isoperimetry to Machine Learning	6
2.2.2 The Intuition	6
2.3 Representation Requirements for a Robust Neural Network	7
<b>第三章 Week 3 - Representation</b>	<b>9</b>
3.1 Graph Neural Networks (GNNs)	9
3.2 The Expressive Power of GNNs and the 1-WL Test	10
3.3 Rethinking Expressiveness via Graph Biconnectivity	10
3.3.1 Key Definitions	11
3.3.2 The Main Result	11
3.3.3 Intuition: How GNNs Achieve This	11
3.4 Implications for Representation	12

<b>第四章 Week 4 - Optimization</b>	<b>13</b>
4.1 The World of Convex Training . . . . .	13
4.1.1 What is a Convex Function? . . . . .	13
4.1.2 Why Convexity is Desirable ✓ . . . . .	13
4.2 The Reality of Non-Convex Training . . . . .	14
4.2.1 The Treacherous Non-Convex Landscape . . . . .	14
4.2.2 Optimization Algorithms for Deep Learning . . . . .	15
4.3 Conclusion . . . . .	16
<b>第五章 Week 5 - Optimization</b>	<b>17</b>
5.1 The Challenge of Saddle Points in High Dimensions . . . . .	17
5.2 Why Standard Gradient Descent Fails . . . . .	18
5.3 The Power of Noise: Escaping with Stochasticity . . . . .	18
5.3.1 The Escape Mechanism . . . . .	18
5.3.2 Perturbed Gradient Descent . . . . .	19
5.4 Implications for Training Neural Networks . . . . .	19
<b>第六章 Week 6 - Optimization</b>	<b>21</b>
6.1 The Mystery of Overparameterization . . . . .	21
6.2 Neural Network Dynamics as a Linear Model . . . . .	22
6.3 The Neural Tangent Kernel (NTK) . . . . .	22
6.3.1 The Infinite-Width Limit . . . . .	22
6.4 Optimization and Generalization in the NTK Regime . . . . .	23
6.4.1 Convex Optimization in Disguise . . . . .	23
6.4.2 Implicit Bias and Generalization . . . . .	23
6.5 Implications and Limitations . . . . .	23
<b>第七章 Week 7 - Generalization</b>	<b>25</b>
7.1 The Generalization Problem . . . . .	25
7.2 Uniform Convergence . . . . .	25
7.3 Measuring Complexity I: VC Dimension . . . . .	26
7.3.1 Shattering and VC Dimension . . . . .	26
7.3.2 The VC Generalization Bound . . . . .	26
7.4 Measuring Complexity II: Rademacher Complexity . . . . .	27
7.4.1 Fitting to Random Noise . . . . .	27

7.4.2	The Rademacher Generalization Bound . . . . .	27
7.5	Conclusion and The Modern Dilemma . . . . .	27
<b>第八章</b>	<b>Week 8 - Generalization</b>	<b>29</b>
8.1	The Failure of Classical Bounds for Deep Networks . . . . .	29
8.2	A New Approach: Norm-Based Complexity . . . . .	29
8.3	Spectrally-Normalized Margin Bounds . . . . .	30
8.3.1	Controlling the Lipschitz Constant . . . . .	30
8.3.2	The Generalization Bound . . . . .	30
8.3.3	Interpreting the Bound . . . . .	31
8.4	Implications and Conclusion . . . . .	31
<b>第九章</b>	<b>Week 9 - Generalization</b>	<b>33</b>
9.1	Recap: The Story So Far . . . . .	33
9.2	The Empirical Challenge: Rethinking Generalization . . . . .	33
9.2.1	Key Finding: Deep Networks Can Fit Random Noise . . . . .	34
9.3	The Theoretical Challenge: A Provable Separation . . . . .	34
9.3.1	The Main Result . . . . .	34
9.3.2	Why Does This Happen? . . . . .	35
9.4	Conclusion: Beyond Uniform Convergence . . . . .	35
<b>第十章</b>	<b>Week 10 - Generalization</b>	<b>37</b>
10.1	A New Perspective on Generalization . . . . .	37
10.1.1	Formalizing Stability . . . . .	37
10.1.2	The Stability-Generalization Theorem . . . . .	38
10.2	Stability for Convex and Lipschitz Losses . . . . .	38
10.3	Stability in Deep Learning: The Non-Convex Case . . . . .	39
10.3.1	Key Results . . . . .	39
10.4	Conclusion . . . . .	39
<b>第十一章</b>	<b>Week 11 - Generalization</b>	<b>41</b>
11.1	The PAC-Bayes Framework . . . . .	41
11.1.1	Prior and Posterior Distributions . . . . .	41
11.1.2	The PAC-Bayes Bound . . . . .	41
11.1.3	KL Divergence and Occam's Razor . . . . .	42

11.2 Information-Theoretic Bounds . . . . .	42
11.2.1 Mutual Information and Generalization . . . . .	43
11.2.2 The Information-Theoretic Bound . . . . .	43
11.3 Conclusion . . . . .	43
<b>第十二章 Week 12 - Generalization</b>	<b>45</b>
12.1 The Core Idea: What is Implicit Bias? . . . . .	45
12.2 The Simplest Case: Min-Norm Bias in Linear Models . . . . .	46
12.3 Implicit Bias in Classification: The Max-Margin Bias . . . . .	46
12.3.1 The Setup . . . . .	46
12.4 Conclusion: The Story Comes Together . . . . .	47
<b>第十三章 Week 13 - LLM</b>	<b>49</b>
13.1 The Double Descent Phenomenon . . . . .	49
13.1.1 The Modern Overparameterized Regime . . . . .	49
13.1.2 Why Does the Error Decrease Again? . . . . .	50
13.2 Scaling Laws for Large Models . . . . .	50
13.2.1 The Three Key Factors . . . . .	50
13.2.2 The Power-Law Relationship . . . . .	51
13.3 Conclusion: The Era of Scale . . . . .	51
<b>第十四章 Week 14 - LLM</b>	<b>53</b>
14.1 Emergent Abilities and In-Context Learning . . . . .	53
14.2 What is Chain-of-Thought (CoT) Prompting? . . . . .	53
14.2.1 Standard Prompting vs. CoT Prompting . . . . .	53
14.3 A Theoretical Perspective on CoT . . . . .	54
14.3.1 Transforming "Slow" Problems into "Fast" Steps . . . . .	54
14.3.2 The Latent Skill Model . . . . .	55
14.4 Conclusion: From Scaling to Reasoning . . . . .	55
<b>第十五章 Week 15 - LLM</b>	<b>57</b>
15.1 The Surprise of Scale: What is Emergence? . . . . .	57
15.2 A Theoretical Explanation for Emergence . . . . .	57
15.2.1 The Main Argument: Emergence as a Perceptual Artifact . . . . .	58
15.2.2 A Model for Complex Task Performance . . . . .	58

15.2.3 The Illusion of Sharpness . . . . .	58
15.3 Conclusion: The End of the Beginning . . . . .	59





# 第一章 Week 1 - Representation

## 1.1 The Universal Approximation Theorem (UAT)

The **Universal Approximation Theorem** is a cornerstone of neural network theory. In simple terms, it states that a standard, feed-forward neural network with just **one hidden layer** can, in principle, approximate any continuous function to any desired degree of accuracy, provided the hidden layer has enough neurons.

This is a powerful existence proof. It tells us that shallow networks are, theoretically, incredibly capable.

### 1.1.1 The Core Idea

Think of the activation function, like a sigmoid or a ReLU, as a basic building block. A neuron applies a non-linear activation function to a linear combination of its inputs. The theorem's magic lies in combining these simple non-linear building blocks.

By adding more neurons to the hidden layer, you are essentially adding more of these basic functions (e.g., "S" curves for sigmoid, "hinges" for ReLU). With enough of them, you can piece them together to create an approximation of any arbitrarily complex, continuous function.

A single-hidden-layer network's output can be expressed as a weighted sum of these activation functions:

$$f(x) \approx \sum_{i=1}^N c_i \sigma(w_i \cdot x + b_i)$$

Here:

- $f(x)$  is the continuous function we want to approximate.
- $N$  is the number of neurons in the hidden layer (the **width** of the network).

- $\sigma$  is a non-linear activation function (e.g., sigmoid, tanh, ReLU).
- $w_i \in \mathbb{R}^d$ ,  $b_i \in \mathbb{R}$ , and  $c_i \in \mathbb{R}$  are the weights and biases of the network that we learn during training.

### 1.1.2 Key Takeaways from UAT

1. **Capability:** A single hidden layer is sufficient to represent any continuous function.
2. **Activation is Key:** The theorem requires a **non-linear** activation function. A network with only linear activations can only represent linear functions, no matter how many layers it has.
3. **Existence, Not Construction:** The UAT guarantees that a set of weights *exists*, but it doesn't tell us how to find them (that's the job of optimization algorithms like gradient descent) or how many neurons ( $N$ ) we will need. For some functions,  $N$  might need to be astronomically large.

## 1.2 Representation and the Power of Depth

The UAT tells us shallow networks *can* do the job. So why do we use **deep** neural networks? The answer is **efficiency and compositionality**.

While a shallow network *can* learn any function, it may need an *exponentially* large number of neurons to do so. A deep network can often represent the same function with far fewer total parameters.

### 1.2.1 Shallow vs. Deep Representation

Let's think about a complex task like identifying a face in an image.

**A Shallow Network's Approach:** A shallow network with one massive hidden layer would have to learn a direct mapping from raw pixel values to the concept of a "face." Each neuron would have to learn a very complex, global template. To distinguish between thousands of different faces and non-faces, it would need an enormous number of these template-matching neurons. This is incredibly inefficient.

**A Deep Network’s Approach (Hierarchical Representation):** A deep network learns a **hierarchy of features**. This is also known as **compositionality**.

- **Layer 1 (Low-Level Features):** The first layer might learn to recognize simple patterns like edges, corners, and color gradients from the raw pixels.
- **Layer 2 (Mid-Level Features):** The second layer takes the features from Layer 1 as input and learns to combine them into more complex patterns like eyes, noses, and mouths.
- **Layer 3 (High-Level Features):** The third layer combines these mid-level features to detect face-like structures.
- **Final Layer:** The output layer takes these high-level face representations and makes a final classification.

This hierarchical approach is exponentially more efficient. Instead of every neuron trying to understand a whole face from pixels, the network reuses and combines simpler patterns to build up more abstract concepts.

### 1.2.2 The Theoretical Advantage of Depth

Many functions that occur in the real world have a compositional structure. For these functions, theory shows that deep networks have a fundamental advantage.

**Key Idea:** There exists a family of functions that can be represented efficiently by a moderately-sized deep network, but would require a shallow network of *exponential width* to represent with the same level of accuracy.

For example, consider a function that involves many nested operations, such as  $f(x) = g_L(\dots g_2(g_1(x)) \dots)$ . A deep network naturally mirrors this structure, with each layer representing one of the  $g_i$  functions. A shallow network would have to “flatten” this entire compositional structure, requiring a huge number of neurons to capture all the possible interactions.

### 1.2.3 Summary Table

In conclusion, while the Universal Approximation Theorem motivates why even simple neural networks are powerful, the **efficiency of representation** provided by **depth** is what makes modern deep learning so effective in practice.

Feature	Shallow Network (1 Hidden Layer)	Deep Network (>1 Hidden Layer)
<b>Universal Approximation</b>	Yes ✓	Yes ✓
<b>Representation Strategy</b>	Sum of simple functions (e.g., sigmoids).	Hierarchical composition of features.
<b>Efficiency</b>	Can be very inefficient; may require an exponential number of neurons (width).	Far more parameter-efficient for many real-world, compositional problems.
<b>Feature Learning</b>	Learns global, non-compositional features.	Learns a hierarchy of features, from simple to abstract. More powerful.
<b>Practicality</b>	Primarily a theoretical concept.	The standard for state-of-the-art results in nearly all domains.

表 1.1: Comparison of shallow and deep network representations.

## 第二章 Week 2 - Representation

Last week, we established that neural networks are universal approximators. However, this theoretical power comes with a catch: standard networks are often incredibly **brittle**. This week, we'll explore what it means for a network to be **robust** and what representation properties are required to achieve it.

### 2.1 The Problem: Adversarial Examples

The brittleness of neural networks is most famously demonstrated by **adversarial examples**. These are inputs that are intentionally perturbed with a small, often human-imperceptible amount of noise to cause a model to misclassify them.

- **Original Image** ( $x$ ): Classified correctly as “panda” with high confidence.
- **Perturbation** ( $\delta$ ): A carefully crafted noise pattern, invisible to the human eye.
- **Adversarial Example** ( $x + \delta$ ): Classified as “gibbon” with high confidence.

This reveals a critical flaw: the function  $f(x)$  learned by the network can be highly unstable. A **robust** neural network is one that is resistant to such perturbations. Formally, for a classifier  $f$  and an input  $x$ , we want:

$$f(x) = f(x + \delta) \quad \text{for all } \|\delta\|_p \leq \epsilon$$

Here,  $\epsilon$  defines the “robustness radius” around the input  $x$  (measured in some  $\ell_p$ -norm) within which the classification must remain constant.

## 2.2 A Universal Law of Robustness via Isoperimetry

Why does this brittleness exist? Is it a flaw in our training methods, or something more fundamental? The paper “A Universal Law of Robustness via Isoperimetry” by Bubeck, Sellke, et al. provides a profound answer rooted in geometry.

The core idea comes from the **isoperimetric inequality**, a classical mathematical principle. In simple terms, it states that among all shapes with a given perimeter, the circle encloses the maximum area. This principle connects the “size” of a set (its volume) to the size of its “boundary.”

### 2.2.1 Connecting Isoperimetry to Machine Learning

Let’s translate this to a classification problem:

- **Sets:** The regions in the high-dimensional input space corresponding to each class. For example, the set  $\mathcal{A}$  contains all images that are correctly classified as “cat.”
- **Volume of a set:** The probability mass of that class,  $\mathbb{P}(\mathcal{A})$ .
- **Boundary of a set:** The decision boundary separating one class from another.

The paper shows that a fundamental trade-off exists between a classifier’s accuracy on the original data and its robustness to adversarial examples. The “Universal Law” can be expressed conceptually as:

$$\text{Standard Accuracy} + \text{Adversarial Robustness} \leq 1$$

(Note: This is a simplification of the formal result, which relates standard error and robust error.)

### 2.2.2 The Intuition

Imagine two classes whose data points are highly intertwined in the input space. To achieve high **standard accuracy**, the decision boundary must be very complex and “wiggly” to correctly separate them. This complex boundary, by its very nature, will pass very close to many data points from both classes. These points are inherently **non-robust**. A tiny push is all it takes to shove them over to the other side of the boundary.

To achieve high **robustness**, the classifier would need to draw a simple, smooth boundary that stays far away from the data. But if the classes are naturally tangled, such a boundary

would inevitably misclassify many points, leading to low **standard accuracy**. This illustrates the intrinsic tension between the two objectives.

## 2.3 Representation Requirements for a Robust Neural Network

The Universal Law tells us that robustness is only possible if the data representations are “nice.” A robust neural network must learn an internal representation that transforms the data to make it more separable.

1. **Separated Data Manifolds:** The primary task of the network’s feature extractor (the hidden layers) is to learn a mapping  $g(x)$  that takes the raw, entangled input data and projects it into a new feature space where the classes are **well-separated**. The goal is to transform the tangled data clouds into distinct, compact clusters that are far from each other. Robustness is impossible if the learned representations of different classes overlap significantly.
2. **Feature Smoothing (Low Lipschitz Constant):** The learned function must be smooth. A function is “smooth” if small changes in the input lead to small changes in the output. This is mathematically captured by the **Lipschitz constant**,  $L$ , of the function, which bounds the function’s steepness. For a function  $g$ :

$$\|g(x_1) - g(x_2)\| \leq L\|x_1 - x_2\|$$

A robust network must have a small Lipschitz constant. This prevents a small perturbation  $\delta$  in the input from creating a massive change in the network’s internal representation or final output, thus preserving the classification.

3. **Simplification over Memorization:** A robust model must prioritize learning the essential, abstract features that define a class, rather than memorizing the noisy details of the training set. This is because the points that are most vulnerable to adversarial attacks are often those where the model has overfitted to noisy patterns near the decision boundary. Achieving robustness forces the model to have a “simpler” view of the world, which might mean sacrificing a few points of accuracy on the “clean” training data. This is the accuracy-robustness trade-off in action.

In summary, for a neural network to be robust, it can't just be a powerful approximator. It must learn a representation that is **geometrically well-structured**, mapping complex, high-dimensional data into a feature space where the classes are simple and far apart.



## 第三章 Week 3 - Representation

In the previous weeks, we discussed representation in standard feed-forward networks. This week, we turn our attention to data with inherent relational structure: **graphs**. We will explore the expressive power of Graph Neural Networks (GNNs), their fundamental limitations, and a more nuanced perspective on what they can represent, guided by the paper “Rethinking the Expressive Power of GNNs via Graph Biconnectivity.”

### 3.1 Graph Neural Networks (GNNs)

Graph Neural Networks are a class of models designed to learn representations of nodes, edges, or entire graphs. The most prevalent framework is the **Message Passing Neural Network (MPNN)**.

In this framework, each node  $v$  in a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  maintains a hidden state or feature vector  $h_v$ . The network updates these states over  $K$  layers (or time steps) by iteratively aggregating information from neighboring nodes. The update rule for node  $v$  at layer  $k$  is:

$$h_v^{(k)} = \text{UPDATE}^{(k)} \left( h_v^{(k-1)}, \text{AGGREGATE}^{(k)} \left( \{h_u^{(k-1)} : u \in \mathcal{N}(v)\} \right) \right)$$

where:

- $h_v^{(k)}$  is the feature vector of node  $v$  at layer  $k$ .
- $\mathcal{N}(v)$  is the set of neighbors of node  $v$ .
- AGGREGATE is a permutation-invariant function (e.g., SUM, MEAN, MAX) that combines messages from neighbors.
- UPDATE is a function (often a small neural network) that combines the aggregated message with the node’s own previous state to compute its new state.

After  $K$  iterations, the final node representations  $h_v^{(K)}$  can be used for node-level tasks (like node classification) or aggregated to form a graph-level representation for graph classification.

### 3.2 The Expressive Power of GNNs and the 1-WL Test

A central question is: what is the **expressive power** of GNNs? In this context, expressiveness refers to the ability to distinguish between two non-isomorphic graphs. If a GNN assigns the same graph-level representation to two different graphs, it cannot distinguish them.

It turns out that the power of standard MPNNs is fundamentally limited by a classical graph algorithm: the **1-Weisfeiler-Lehman (1-WL) test** for graph isomorphism.

**The 1-WL Test:** This test is an iterative node coloring algorithm.

1. Initially, all nodes are assigned the same color (or a color based on their initial features).
2. In each round, every node  $v$  collects the multiset of colors from its neighbors  $\{\{c_u : u \in \mathcal{N}(v)\}\}$ .
3. It then updates its own color based on its current color and the multiset of neighbor colors using an injective hash function.
4. The process repeats until the set of colors in the graph stabilizes. Two graphs are considered indistinguishable by the 1-WL test if their final color histograms are identical.

**Theorem:** The expressive power of any MPNN is at most that of the 1-WL test.

This is because the message-passing update is a neural version of the 1-WL color update. Consequently, any two graphs that the 1-WL test cannot distinguish, a standard GNN cannot distinguish either. For example, a GNN cannot distinguish a 6-node cycle graph from two disjoint 3-node cycle graphs, as both are 2-regular graphs and appear identical from a local neighborhood perspective.

### 3.3 Rethinking Expressiveness via Graph Biconnectivity

The 1-WL limit suggests that GNNs are weak at capturing global graph structure. However, the paper “Rethinking the Expressive Power of GNNs via Graph Biconnectivity” argues that this view is incomplete. It shows that GNNs are surprisingly powerful at identifying crucial topological structures related to connectivity.

### 3.3.1 Key Definitions

- **Bridge (or Cut-Edge):** An edge whose removal increases the number of connected components in the graph. Bridges represent critical vulnerabilities or bottlenecks.
- **Articulation Point (or Cut-Vertex):** A vertex whose removal increases the number of connected components.
- **Biconnected Component (BCC):** A maximal subgraph that remains connected even after removing any single vertex. Inside a BCC, there are at least two vertex-disjoint paths between any two nodes.

### 3.3.2 The Main Result

The paper’s central finding is a positive result on GNN expressiveness:

**Theorem:** An MPNN with a sufficiently powerful UPDATE function and an injective AGGREGATE function (like SUM) can learn to count the number of biconnected components that any given node or edge belongs to.

As a direct consequence, these GNNs can **identify all bridges, articulation points, and biconnected components** in a graph.

### 3.3.3 Intuition: How GNNs Achieve This

The key lies in the structure of the computation graphs (or trees) formed during message passing.

- For an edge  $(u, v)$  that is a **bridge**, the messages flowing from  $u$ ’s side to  $v$  are completely distinct from those flowing from  $v$ ’s side to  $u$ . The GNN sees two disjoint computational subtrees, allowing it to learn a unique representation for this “bottleneck” structure.
- For an edge inside a **biconnected component** (e.g., a cycle), there are at least two paths for information to travel between its endpoints. The messages each node receives are therefore redundant and mixed from different paths. This creates a fundamentally different message pattern than that of a bridge.

The GNN’s node update functions can be trained to act as pattern recognizers, distinguishing these distinct message-passing patterns and thereby identifying the underlying topological structures.

### 3.4 Implications for Representation

This result provides a more optimistic view of what GNNs learn.

1. **Beyond Local Neighborhoods:** While limited by 1-WL, GNNs capture topological information that is non-local and crucial for understanding a graph's global structure and robustness.
2. **Hierarchical Understanding:** The ability to identify BCCs means GNNs can implicitly learn a graph's *block-cut tree* structure—a hierarchical representation of its connectivity.
3. **Explaining Practical Success:** This helps explain why GNNs work so well in domains like chemistry and social networks. Molecules often have bridge-like bonds, and social networks have communities (dense components) connected by weak links. The GNN's ability to represent these structures is key to its empirical success.

In conclusion, while standard GNNs are not universal graph classifiers, their representations are rich enough to capture fundamental properties of graph connectivity, making them powerful tools for many real-world problems.

## 第四章 Week 4 - Optimization

Having discussed what neural networks can *represent*, we now turn to the question of *how* we find the parameters (weights and biases) that achieve a good representation. This is the domain of **optimization**. The goal of training a model is to find parameters  $\theta$  that minimize a loss function  $L(\theta)$  over a dataset.

The general form of this problem is called **Empirical Risk Minimization (ERM)**:

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x_i; \theta), y_i)$$

where  $f(x_i; \theta)$  is our model's prediction and  $\mathcal{L}$  is a loss function (e.g., Mean Squared Error, Cross-Entropy). The landscape of this loss function determines the difficulty of the optimization problem.

### 4.1 The World of Convex Training

In optimization, convexity is the dividing line between problems we can solve reliably and those we cannot.

#### 4.1.1 What is a Convex Function?

A function  $L(\theta)$  is **convex** if, for any two points  $\theta_1, \theta_2$  in its domain, the line segment connecting  $(\theta_1, L(\theta_1))$  and  $(\theta_2, L(\theta_2))$  lies on or above the graph of the function.

Mathematically, for any  $\lambda \in [0, 1]$ :

$$L(\lambda\theta_1 + (1 - \lambda)\theta_2) \leq \lambda L(\theta_1) + (1 - \lambda)L(\theta_2)$$

#### 4.1.2 Why Convexity is Desirable ✓

Convex optimization problems are “easy” because they possess a critical property:

**The Golden Rule of Convexity:** Any local minimum is also a global minimum.

This means we don't have to worry about getting stuck in a suboptimal solution. If an algorithm finds a point where the gradient is zero, it has found the best possible solution.

For convex problems, simple algorithms like **Gradient Descent (GD)** are guaranteed to converge to the global minimum (given an appropriate step size  $\eta$ ). The update rule for GD is:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t)$$

This involves repeatedly taking steps in the direction of the steepest descent of the loss function.

**Examples in Machine Learning:** Several classic models have convex loss functions, including:

- Linear Regression (with Mean Squared Error loss)
- Logistic Regression
- Support Vector Machines (SVMs)

## 4.2 The Reality of Non-Convex Training

The loss landscapes of deep neural networks are highly **non-convex**. This makes training them a far greater challenge.

### 4.2.1 The Treacherous Non-Convex Landscape

The loss function of a neural network contains many pathological features that do not exist in convex problems:

- **Multiple Local Minima:** There are many points that are locally optimal but globally suboptimal. An optimizer can easily get stuck in a “valley” that is not the deepest one.
- **Saddle Points:** These are points where the gradient is zero, but they are not minima. They are a minimum along some directions and a maximum along others. In the high-dimensional spaces of neural networks, saddle points are exponentially more common than local minima and are a major obstacle for optimization.

- **Plateaus:** Large, flat regions where the gradient is very close to zero. An optimizer can slow down to a crawl in these regions, making training progress extremely slow.
- **Sharp vs. Flat Minima:** Not all minima are created equal. Empirical evidence suggests that solutions found in wide, **flat minima** tend to generalize better to unseen data than solutions in sharp, narrow minima. A flat minimum is more robust to small shifts between the training and test data distributions.

### 4.2.2 Optimization Algorithms for Deep Learning

Standard Gradient Descent is often ineffective in this landscape. Instead, we rely on more sophisticated, stochastic algorithms.

**Stochastic Gradient Descent (SGD):** Instead of computing the gradient over the entire dataset (which is computationally expensive), SGD computes it on a small, random mini-batch of data. The update rule is:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} L(\theta_t; x_{\text{batch}}, y_{\text{batch}})$$

The noise introduced by the mini-batch sampling helps the optimizer escape from sharp local minima and navigate saddle points more effectively.

**Momentum:** To overcome plateaus and oscillations, we can add a **momentum** term. This acts like a “heavy ball” that accumulates velocity in directions of persistent gradient descent.

$$v_{t+1} = \beta v_t + \eta \nabla_{\theta} L(\theta_t)$$

$$\theta_{t+1} = \theta_t - v_{t+1}$$

The momentum term  $\beta$  (e.g., 0.9) helps the optimizer power through flat regions and dampens updates in directions where the gradient changes frequently.

**Adaptive Methods (e.g., Adam):** Algorithms like **Adam** (Adaptive Moment Estimation) take this a step further by maintaining a per-parameter learning rate that is adapted based on the history of gradients (both first and second moments). This allows for faster convergence and makes them less sensitive to the initial choice of learning rate.

### 4.3 Conclusion

While the theory of **convex optimization** provides strong guarantees and a solid foundation, the practice of training deep neural networks is fundamentally a **non-convex** problem. There are no guarantees of finding the global minimum. However, through a combination of stochastic methods, momentum, and adaptive learning rates, we have developed powerful heuristic algorithms that can navigate these complex landscapes to find solutions that, while not provably optimal, are remarkably effective and generalize well in practice.



## 第五章 Week 5 - Optimization

Last week, we established that training neural networks is a non-convex optimization problem, characterized by a complex loss landscape filled with local minima, plateaus, and saddle points. This week, we will delve deeper into one of the most significant challenges in high-dimensional optimization: escaping from **saddle points**.

### 5.1 The Challenge of Saddle Points in High Dimensions

A **critical point** of a loss function  $L(\theta)$  is any point  $\theta^*$  where the gradient is zero, i.e.,  $\nabla L(\theta^*) = \mathbf{0}$ . These points can be local minima, local maxima, or saddle points.

To distinguish them, we must look at the curvature of the loss function, which is described by the **Hessian matrix**,  $\mathbf{H} = \nabla^2 L(\theta)$ . The eigenvalues of the Hessian tell us about the shape of the landscape around the critical point.

- **Local Minimum:** All eigenvalues of  $\mathbf{H}$  are strictly positive. The surface curves upwards in all directions.
- **Local Maximum:** All eigenvalues of  $\mathbf{H}$  are strictly negative. The surface curves downwards in all directions.
- **Saddle Point:** The Hessian  $\mathbf{H}$  has at least one positive eigenvalue and at least one negative eigenvalue. The surface curves up in some directions and down in others.

In the high-dimensional parameter spaces of neural networks, saddle points are exponentially more common than local minima. Therefore, an effective optimizer *must* be able to escape them.

## 5.2 Why Standard Gradient Descent Fails

Standard (full-batch) Gradient Descent (GD) is defined by the update rule  $\theta_{t+1} = \theta_t - \eta \nabla L(\theta_t)$ .

Consider the behavior of GD near a saddle point  $\theta^*$ . The direction corresponding to a negative eigenvalue of the Hessian is an “escape route”—moving in this direction would decrease the loss. However, if the optimization trajectory lands on a point where the gradient component along this escape direction is exactly zero, GD will not move in that direction. It gets stuck on a “ridge” and converges to the saddle point, failing to make further progress. This is a critical failure for a non-convex optimizer.

## 5.3 The Power of Noise: Escaping with Stochasticity

The key insight, explored in papers like “Escaping From Saddle Points – Online Stochastic Gradient for Tensor Decomposition” and others, is that **noise is a feature, not a bug**. Stochastic algorithms, like SGD, use noise to their advantage.

The gradient in Stochastic Gradient Descent (SGD) is an unbiased but noisy estimate of the true gradient:

$$\nabla L(\theta_t; \text{batch}) = \nabla L(\theta_t) + \epsilon_t$$

where the noise term  $\epsilon_t$  has zero mean but non-zero variance.

### 5.3.1 The Escape Mechanism

At a saddle point, the true gradient  $\nabla L(\theta_t)$  might have a zero component along the escape direction (the eigenvector  $\mathbf{v}_{\min}$  corresponding to the minimum eigenvalue  $\lambda_{\min} < 0$ ).

However, the stochastic gradient’s noise component  $\epsilon_t$  is random and isotropic (it points in all directions with equal probability). Therefore, it is almost guaranteed that the noise will have a non-zero component along the escape direction  $\mathbf{v}_{\min}$ .

$$\text{Component of stochastic gradient along } \mathbf{v}_{\min} = (\nabla L(\theta_t) + \epsilon_t) \cdot \mathbf{v}_{\min} \approx \epsilon_t \cdot \mathbf{v}_{\min} \neq 0$$

This small, random “kick” from the stochastic update is enough to push the parameter iterate  $\theta_t$  off the ridge and into the region of negative curvature. Once it is in this region, the true gradient  $\nabla L(\theta)$  will point further downhill, and the optimizer will successfully escape the saddle point’s vicinity.

### 5.3.2 Perturbed Gradient Descent

This idea has been formalized in the analysis of **Perturbed Gradient Descent (PGD)**. Theoretical results show that if you add a small amount of artificial noise to the parameters whenever the gradient becomes small, you can provably escape saddle points and converge to a **second-order stationary point** (a point that is, with high probability, a local minimum).

The algorithm is conceptually simple:

1. Run standard gradient descent.
2. If  $\|\nabla L(\theta_t)\|$  becomes small, it implies we are near a critical point.
3. Add a small, spherical random noise vector to the parameters:  $\theta_t \leftarrow \theta_t + \xi$ .
4. Continue gradient descent from this perturbed point.

This analysis confirms that noise is a sufficient condition for escaping saddle points.

## 5.4 Implications for Training Neural Networks

- **SGD is a Natural Fit:** The inherent noise in SGD from mini-batch sampling serves the same role as the explicit noise injection in PGD. This is a primary theoretical reason for SGD’s remarkable success in deep learning. It’s not just a faster, approximate version of GD; its stochastic nature is fundamental to its ability to navigate non-convex landscapes.
- **Focus Shifts from “Getting Stuck”:** Because saddle points are easily escaped by modern optimizers, the main challenge in deep learning optimization is no longer about avoiding getting stuck. Instead, the focus has shifted to understanding the properties of the local minima that are found.
- **The Generalization Puzzle:** The current frontier of research is to understand why the local minima found by SGD tend to be “flat” and why these flat minima generalize well to unseen data. The dynamics of escaping saddles and settling into wide valleys is a key part of this ongoing investigation.

In summary, the problem of saddle points, once thought to be a major barrier to training deep networks, is effectively solved in practice by the stochasticity of our standard optimization algorithms.



## 第六章 Week 6 - Optimization

We've seen that the optimization landscape of neural networks is highly non-convex, yet simple algorithms like SGD consistently find good solutions. This week, we explore one of the most significant theoretical breakthroughs in understanding this phenomenon: the **Neural Tangent Kernel (NTK)**. The NTK provides a powerful lens for analyzing the training dynamics of very wide neural networks.

### 6.1 The Mystery of Overparameterization

A central puzzle in modern deep learning is the success of **overparameterization**. We routinely train models with far more parameters than training samples (e.g., millions of parameters for thousands of images).

- **Classical View:** From a classical statistics perspective, such models should overfit catastrophically. They have enough capacity to simply memorize the entire training set, including its noise, and should fail to generalize to new, unseen data.
- **Deep Learning Reality:** Empirically, we observe the opposite. Massively overparameterized networks, trained with gradient-based methods, not only achieve zero training error but also often generalize remarkably well.

The NTK provides a compelling explanation for how this is possible, connecting the complex, non-convex optimization of neural networks to the simpler, well-understood world of kernel methods.

## 6.2 Neural Network Dynamics as a Linear Model

Let's consider a neural network's output function  $f(\boldsymbol{\theta}, x)$  for an input  $x$  with parameters  $\boldsymbol{\theta}$ . During training, the parameters evolve from an initialization  $\boldsymbol{\theta}_0$ . If the parameter updates are small, we can analyze the function's behavior using a first-order Taylor expansion around  $\boldsymbol{\theta}_0$ :

$$f(\boldsymbol{\theta}, x) \approx f(\boldsymbol{\theta}_0, x) + \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0, x)^T (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

This equation is key. It shows that for small changes in parameters  $(\boldsymbol{\theta} - \boldsymbol{\theta}_0)$ , the neural network function  $f(\boldsymbol{\theta}, x)$  behaves like a **linear model** with respect to its parameters. The “features” of this linear model are the gradients of the network output with respect to the initial parameters,  $\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0, x)$ .

This is called the “lazy training” or “NTK” regime, where the model is so overparameterized that it can fit the data without moving its parameters far from their random initialization.

## 6.3 The Neural Tangent Kernel (NTK)

The dynamics of training with gradient descent can be studied in function space instead of parameter space. The change in the function's output over time is governed by a special kernel.

The **Neural Tangent Kernel (NTK)** is defined as the inner product of these gradient-features:

$$\Theta(\boldsymbol{\theta}, x, x') = \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, x)^T \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}, x')$$

This kernel measures the similarity between two inputs  $x$  and  $x'$  in terms of how a change in the model's parameters affects their respective outputs.

### 6.3.1 The Infinite-Width Limit

The groundbreaking result, formalized by Jacot et al. (2018) and explored in detail by papers like the reference “Fine-Grained Analysis...”, is what happens in the limit as the width of the network layers goes to infinity.

**Theorem (Informal):** As the width of a neural network tends to infinity, its Neural Tangent Kernel  $\Theta(\boldsymbol{\theta}, \cdot, \cdot)$  converges to a deterministic, fixed kernel  $\Theta_{\infty}(\cdot, \cdot)$  at initialization, and crucially, **it stays constant throughout training**.

This means that for an infinitely wide network, the complex dynamics simplify dramatically. The network behaves as if its features,  $\nabla_{\theta} f(\theta, x)$ , are frozen at their initial values.

## 6.4 Optimization and Generalization in the NTK Regime

### 6.4.1 Convex Optimization in Disguise

Because the kernel is fixed during training, the optimization problem becomes much simpler. Training a very wide neural network with gradient descent on a Mean Squared Error (MSE) loss is equivalent to performing **kernel regression** with the fixed kernel  $\Theta_{\infty}$ .

Kernel regression is a **convex optimization problem**. This provides a stunning result: the highly non-convex optimization of an infinitely wide neural network is effectively convex. This explains why gradient descent can find a global minimum and achieve zero training error.

### 6.4.2 Implicit Bias and Generalization

The NTK framework also helps explain generalization. When gradient descent minimizes the training loss, it doesn't just find any solution; it finds a specific one. This is known as **implicit bias**.

In the NTK regime, gradient descent finds the solution that has the minimum norm in the Reproducing Kernel Hilbert Space (RKHS) associated with the kernel  $\Theta_{\infty}$ . This minimum-norm solution is inherently “simple” or “smooth,” which helps it generalize well to new data. The properties of the kernel  $\Theta_{\infty}$  thus determine the generalization performance of the trained network.

## 6.5 Implications and Limitations

- **A Powerful Explanatory Framework:** The NTK provides the first rigorous mathematical theory that explains both the successful optimization and the good generalization of massively overparameterized neural networks.
- **Limitations – The “Lazy Training” View:** The NTK describes a regime where the network is “lazy”—it fits the data by making very small adjustments to its initial function. This perspective has some crucial limitations:

1. It doesn't capture **feature learning**. A key belief in deep learning is that networks succeed by learning hierarchical, meaningful features in their hidden layers. The NTK, by assuming the kernel is fixed, does not account for this representation learning.
2. The theory holds in the infinite-width limit, and while it provides insights for very wide finite networks, it may not accurately describe the behavior of narrower, more practical architectures.

In conclusion, the Neural Tangent Kernel provides a profound theoretical lens, revealing that in the overparameterized limit, neural networks can behave like simple, linear kernel machines, allowing us to offer guarantees for their training and generalization.



## 第七章 Week 7 - Generalization

So far, we have focused on representation and optimization: what functions can our models represent, and how do we find good parameters? This week, we address the most fundamental question in machine learning: why should a model trained on a finite dataset work well on new, unseen data? This is the problem of **generalization**.

### 7.1 The Generalization Problem

When we train a model, we minimize the loss on the training data. This is the **empirical risk**:

$$R_{\text{emp}}(f) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x_i), y_i)$$

However, our true goal is to minimize the **true risk** (or generalization error), which is the expected loss over the entire data distribution  $\mathcal{D}$ :

$$R_{\text{true}}(f) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\mathcal{L}(f(x), y)]$$

The **generalization gap** is the difference  $|R_{\text{true}}(f) - R_{\text{emp}}(f)|$ . We can only trust our trained model if this gap is small.

### 7.2 Uniform Convergence

How can we guarantee a small generalization gap? A powerful approach is through the principle of **uniform convergence**. Instead of just bounding the gap for the single function  $f_{\text{train}}$  that we found, we seek a much stronger guarantee. We want to bound the worst-case gap over the *entire class of functions*  $\mathcal{F}$  that our model could possibly learn.

We aim to bound the supremum of the generalization gap:

$$\mathbb{P} \left( \sup_{f \in \mathcal{F}} |R_{\text{true}}(f) - R_{\text{emp}}(f)| > \epsilon \right) \leq \delta$$

If we can show that this supremum is small with high probability, it means that the empirical risk landscape is a good proxy for the true risk landscape across our entire function class. Consequently, the function that minimizes the empirical risk will also have a true risk that is close to the minimum possible true risk within our class  $\mathcal{F}$ .

The size of this gap is determined by the **complexity** or **capacity** of the function class  $\mathcal{F}$ . A more complex class can generate more varied functions, making it easier to fit the specific noise of our training sample, which leads to overfitting and a larger generalization gap.

### 7.3 Measuring Complexity I: VC Dimension

For binary classification problems, the Vapnik-Chervonenkis (VC) dimension is a classic, combinatorial measure of complexity.

#### 7.3.1 Shattering and VC Dimension

A hypothesis class  $\mathcal{F}$  can **shatter** a set of  $n$  points  $\{x_1, \dots, x_n\}$  if, for every one of the  $2^n$  possible binary labelings of these points, there exists a function  $f \in \mathcal{F}$  that produces that exact labeling.

**Definition:** The **VC dimension** of a class  $\mathcal{F}$ , denoted  $\text{VCdim}(\mathcal{F})$ , is the size of the largest set of points that  $\mathcal{F}$  can shatter.

For example, the set of 2D linear classifiers can shatter any 3 non-collinear points, but cannot shatter any set of 4 points. Thus, its VC dimension is 3.

#### 7.3.2 The VC Generalization Bound

The VC dimension allows us to bound the generalization gap. For any  $f \in \mathcal{F}$ , with probability at least  $1 - \delta$ :

$$R_{\text{true}}(f) \leq R_{\text{emp}}(f) + \sqrt{\frac{C_1 \cdot \text{VCdim}(\mathcal{F}) (\log(n/\text{VCdim}(\mathcal{F}))) + C_2 \log(1/\delta)}{n}}$$

This bound tells us that the generalization gap shrinks as  $n \rightarrow \infty$  and grows with the model's complexity, as measured by its VC dimension.

## 7.4 Measuring Complexity II: Rademacher Complexity

VC dimension can be loose and is difficult to apply to models beyond binary classifiers. **Rademacher complexity** is a more flexible and often tighter, data-dependent measure.

### 7.4.1 Fitting to Random Noise

The core idea is to measure how well our function class  $\mathcal{F}$  can correlate with random noise. A complex class has the capacity to fit random patterns, while a simple class does not.

We introduce a vector of independent random variables  $\sigma = (\sigma_1, \dots, \sigma_n)$ , where each  $\sigma_i$  is a Rademacher variable ( $\pm 1$  with probability 0.5).

**Definition:** The **empirical Rademacher complexity** of  $\mathcal{F}$  on a dataset  $S = \{x_1, \dots, x_n\}$  is:

$$\hat{\mathcal{R}}_S(\mathcal{F}) = \mathbb{E}_{\sigma} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right]$$

This measures the average correlation between the functions in  $\mathcal{F}$  and the random noise vector  $\sigma$  on the given data points.

### 7.4.2 The Rademacher Generalization Bound

Rademacher complexity gives us a powerful generalization bound. For any  $f \in \mathcal{F}$ , with probability at least  $1 - \delta$ :

$$R_{\text{true}}(f) \leq R_{\text{emp}}(f) + 2 \cdot \mathcal{R}_n(\mathcal{L} \circ \mathcal{F}) + C \sqrt{\frac{\log(1/\delta)}{n}}$$

where  $\mathcal{R}_n(\mathcal{L} \circ \mathcal{F})$  is the (non-empirical) Rademacher complexity of the class of loss functions. In practice, this is often bounded using the empirical version  $\hat{\mathcal{R}}_S$ .

## 7.5 Conclusion and The Modern Dilemma

Uniform convergence theory, using tools like VC dimension and Rademacher complexity, provides the foundational framework for statistical learning theory. It establishes a clear link between model complexity and generalization.

However, this classical view presents a puzzle for modern deep learning. The VC dimension and Rademacher complexity of large neural networks are often enormous or infinite. According

to these bounds, the models should overfit terribly and have a vacuous generalization gap (a bound greater than 1). Yet, in practice, they generalize well.

This discrepancy shows that while uniform convergence is a cornerstone of learning theory, it does not tell the whole story. It has motivated a new wave of research into other principles of generalization, such as the implicit bias of optimizers, properties of the loss landscape (flat vs. sharp minima), and norm-based controls, which we may explore in the coming weeks.

## 第八章 Week 8 - Generalization

Last week, we explored the classical theory of generalization through uniform convergence, using complexity measures like VC dimension and Rademacher complexity. We concluded with a puzzle: these classical bounds fail to explain the success of modern deep learning, as they predict that massively overparameterized networks should not generalize. This week, we investigate a more modern approach that provides a much clearer picture: **norm-based complexity control**.

### 8.1 The Failure of Classical Bounds for Deep Networks

Classical generalization bounds from uniform convergence theory typically depend polynomially on the number of parameters in the model. For a deep neural network with millions or billions of parameters, this leads to a **vacuous bound**—a guaranteed error rate that is greater than 1, which is meaningless.

This tells us that the raw parameter count is the wrong way to measure the “effective complexity” of a neural network. The key insight of modern theory is that generalization is controlled not by the size of the architecture, but by the properties of the specific solution found by the optimization algorithm.

### 8.2 A New Approach: Norm-Based Complexity

The central idea is that the generalization ability of a neural network  $f$  is governed by the **norms of its weight matrices**. Even if a network is enormous, if the weights are small or have a specific structure, the function it represents can be simple enough to generalize well.

To formalize this, we re-introduce two key concepts:

- **The Margin ( $\gamma$ ):** For a binary classifier, the margin is the measure of confidence in

its predictions on the training data. For a data point  $(x, y)$  where  $y \in \{-1, +1\}$ , the margin is  $y \cdot f(x)$ . We are interested in the minimum margin over the training set,  $\gamma = \min_i y_i f(x_i)$ . A larger margin implies a more stable decision boundary that is far from the training points.

- **The Lipschitz Constant ( $L_f$ ):** This measures the smoothness of the function. A small Lipschitz constant means that small perturbations to the input do not drastically change the output. Intuitively, smoother functions are simpler and should generalize better.

## 8.3 Spectrally-Normalized Margin Bounds

The breakthrough, developed in papers like “Spectrally-normalized margin bounds for neural networks” (Bartlett et al.) and “Improved Sample Complexities...” (Wei et al.), was to connect these ideas by controlling the network’s Lipschitz constant via the norms of its weight matrices.

### 8.3.1 Controlling the Lipschitz Constant

Consider a deep neural network composed of  $L$  layers:  $f(x) = W_L \sigma(\dots \sigma(W_1 x) \dots)$ . The overall Lipschitz constant of the network is bounded by the product of the Lipschitz constants of its individual layers.

- For a linear layer  $h(x) = Wx$ , its Lipschitz constant is the **spectral norm** of the weight matrix,  $\|W\|_\sigma$ . The spectral norm is the largest singular value of the matrix and measures its maximum stretching factor.
- For common activation functions like ReLU, the Lipschitz constant is 1.

This gives us a bound on the network’s overall Lipschitz constant:  $L_f \leq \prod_{i=1}^L \|W_i\|_\sigma$ .

### 8.3.2 The Generalization Bound

Combining Rademacher complexity analysis with this norm-based control of the function’s smoothness leads to powerful generalization bounds. A simplified version of the bound for a network  $f$  is:

With high probability, the true classification error of  $f$  is bounded by:

$$\hat{R}_\gamma(f) + \tilde{O} \left( \frac{\left( \prod_{i=1}^L \|W_i\|_\sigma \right) \left( \sum_{i=1}^L \frac{\|W_i\|_F^2}{\|W_i\|_\sigma^2} \right)^{1/2}}{\gamma \sqrt{n}} \right)$$

where  $\hat{R}_\gamma(f)$  is the empirical margin loss (zero if all training points have margin  $\geq \gamma$ ),  $\|W_i\|_F$  is the Frobenius norm, and  $\tilde{O}$  hides logarithmic factors.

### 8.3.3 Interpreting the Bound

This is a significant departure from classical bounds:

1. **No Parameter Count:** The bound does **not** depend explicitly on the number of parameters (the width or depth) of the network. It depends on the norms of the weights of the *final learned solution*.
2. **Complexity Measure:** The complexity is captured by the product of the spectral norms,  $\prod \|W_i\|_\sigma$ . This term controls the function’s smoothness.
3. **Importance of Margin:** The generalization gap is inversely proportional to the margin  $\gamma$ . This formalizes the long-held intuition that large-margin classifiers generalize well.

The second reference paper, “Improved Sample Complexities...”, refines this further by introducing an **all-layer margin**, suggesting that for good generalization, the representations at every layer of the network should be well-separated.

## 8.4 Implications and Conclusion

- **The Benefit of Overparameterization:** These bounds help explain why overparameterization can be beneficial. Having more parameters may create a smoother optimization landscape, making it easier for SGD to find solutions where the weight norms are small and the margin is large. The network uses its vast number of parameters to find a simple function, rather than a complex one.
- **Implicit Bias of SGD:** This theory provides strong support for the idea of **implicit bias**. SGD does not find just any set of parameters that achieve zero training error. It appears to be implicitly biased towards solutions with small norms, which are precisely the solutions that these bounds certify as having good generalization performance.

In summary, norm-based uniform convergence bounds provide a far more realistic framework for understanding generalization in deep learning. They successfully shift the focus from the architectural size to the properties of the learned function, offering a powerful explanation for why deep learning works in the overparameterized regime.



## 第九章 Week 9 - Generalization

Over the past two weeks, we built up the theory of generalization based on **uniform convergence**. We started with classical VC theory and refined it with modern, norm-based bounds that seemed to offer a promising explanation for why overparameterized networks generalize. This week, we confront a series of powerful arguments and results suggesting that even this refined view of uniform convergence may be fundamentally insufficient to explain generalization in deep learning.

### 9.1 Recap: The Story So Far

- **Week 7:** We learned that classical uniform convergence bounds (VC dimension, Rademacher complexity) fail for deep networks because they depend on the parameter count, which is massive.
- **Week 8:** We improved the theory with norm-based bounds, decoupling generalization from the parameter count and linking it to properties of the learned solution (small norms, large margin). This seemed to resolve the paradox.

The question we address today is: Is uniform convergence, even in its modern form, the right framework for understanding deep learning?

### 9.2 The Empirical Challenge: Rethinking Generalization

The influential paper “Understanding deep learning requires rethinking generalization” (Zhang et al., 2017) presented a series of simple yet profound experiments that challenge the uniform convergence perspective.

### 9.2.1 Key Finding: Deep Networks Can Fit Random Noise

The authors conducted a now-famous experiment:

1. They took a standard image dataset (e.g., CIFAR-10) and a powerful convolutional neural network architecture (e.g., Inception).
2. They then completely randomized the labels, destroying any relationship between the images and their classes.
3. They trained the network on this dataset with random labels.

**The shocking result:** The network was still able to achieve nearly **100% training accuracy**. It had sufficient capacity to perfectly memorize every single image and its associated random label.

**Implication:** The effective capacity of these models is enormous. From a classical perspective (VC dimension, Rademacher complexity), a function class that is powerful enough to fit pure random noise should have the worst possible generalization properties. The fact that the *exact same architecture* generalizes well when trained on true labels but can also perfectly memorize noise suggests that the properties of the function class **alone** cannot explain generalization.

The explanation must lie in the interplay between the model, the optimization algorithm, and the properties of the true data distribution.

## 9.3 The Theoretical Challenge: A Provable Separation

While the Zhang et al. paper provided compelling empirical evidence, the paper “Uniform convergence may be unable to explain generalization in deep learning” (Nagarajan & Kolter, 2019) provided a formal, theoretical proof of this separation.

### 9.3.1 The Main Result

The authors constructed a simple, learnable binary classification problem and an over-parameterized neural network architecture. They proved that for this setup, two things are simultaneously true:

1. The network trained with SGD,  $f_{\text{SGD}}$ , achieves a **low true risk**. It generalizes well.

2. For the function class  $\mathcal{F}$  defined by the network’s architecture, the uniform convergence bound is **vacuous**. The worst-case generalization gap over the class is maximal.

Formally, they show it’s possible to have:

$$R_{\text{true}}(f_{\text{SGD}}) \approx 0 \quad \text{while at the same time} \quad \sup_{f \in \mathcal{F}} |R_{\text{true}}(f) - R_{\text{emp}}(f)| \approx 1$$

### 9.3.2 Why Does This Happen?

The intuition behind their proof is that the function class  $\mathcal{F}$  is extremely rich. It contains both “good” functions (like the one found by SGD, which generalizes) and many “bad” interpolating functions that also achieve zero training error but have terrible true error.

- The **uniform convergence bound** is a worst-case measure over the entire class  $\mathcal{F}$ . Its value is therefore dominated by the presence of these “bad” functions, making the bound loose and uninformative.
- The **SGD algorithm**, when run on structured data, has an **implicit bias** that leads it to select a “good” function from  $\mathcal{F}$ , even though the class itself is pathological from a uniform convergence standpoint.

This result is a formal separation: it proves that there are realistic scenarios where neural networks generalize, yet the uniform convergence framework is provably too weak to explain why.

## 9.4 Conclusion: Beyond Uniform Convergence

The evidence from these papers strongly suggests that our search for an explanation for deep learning generalization cannot end with uniform convergence. The key takeaways are:

1. **Generalization is not just a property of the model class.** The explanation must depend on the algorithm (e.g., SGD’s implicit bias) and the data (e.g., the structure of real-world datasets).
2. **Worst-case bounds are too pessimistic.** Uniform convergence analyzes the worst-case function in a class, but SGD seems to find a function that is far from the worst case.

This motivates the exploration of alternative theoretical frameworks that are more algorithm-dependent, such as:

- **Algorithmic Stability:** Analyzes how much the learned model changes if we perturb the training set by a single example. Stable algorithms tend to generalize well.
- **PAC-Bayes:** A framework that bounds the generalization error of a stochastic classifier (a distribution over hypotheses), which can be specialized to analyze the solution found by SGD.

The failure of uniform convergence to fully capture the phenomenon of deep learning generalization remains one of the most active and exciting areas of theoretical machine learning research.

# 第十章 Week 10 - Generalization

Last week, we saw compelling evidence that the uniform convergence framework, which analyzes properties of the entire function class, is likely insufficient to explain generalization in deep learning. This motivates a shift in perspective. Instead of focusing on the function class, what if we analyze the properties of the **learning algorithm** itself? This leads us to the theory of **algorithmic stability**.

## 10.1 A New Perspective on Generalization

Algorithmic stability offers an intuitive and powerful alternative to uniform convergence. The core idea is simple:

An algorithm that generalizes well should be **stable**. Its output should not change drastically if we make a small change to the training dataset.

If an algorithm's learned hypothesis depends heavily on a single training point, it is likely overfitting to the noise in that specific example. A stable algorithm, by contrast, learns a model that reflects the overall trends in the data distribution, making it robust to the specifics of the sample.

### 10.1.1 Formalizing Stability

Let  $A$  be a learning algorithm that maps a training set  $S$  to a hypothesis  $f_S = A(S)$ . Let  $S = \{(z_1, \dots, z_n)\}$  be a training set of size  $n$ , where  $z_i = (x_i, y_i)$ . Let  $S^{(i)}$  be a dataset identical to  $S$ , but with the  $i$ -th data point  $z_i$  replaced by a new point  $z'_i$ .

**Definition (Uniform Stability):** An algorithm  $A$  is  $\beta$ -uniformly stable if for any

datasets  $S$  and  $S^{(i)}$ , the following holds for any test point  $z$ :

$$\sup_{S,i,z} |\mathcal{L}(A(S), z) - \mathcal{L}(A(S^{(i)}), z)| \leq \beta$$

The stability coefficient  $\beta$  measures the worst-case change in the model's loss on a new point when a single training point is changed.

### 10.1.2 The Stability-Generalization Theorem

The power of this framework comes from a fundamental theorem by Bousquet and Elisseeff (2002), which directly links stability to generalization.

**Theorem:** If an algorithm  $A$  is  $\beta$ -uniformly stable, then its expected generalization gap is bounded:

$$|\mathbb{E}_S[R_{\text{true}}(A(S))] - \mathbb{E}_S[R_{\text{emp}}(A(S))]| \leq \beta$$

This is a remarkable result. To understand generalization, we no longer need to analyze the complexity of the entire function class. We just need to analyze the stability of our algorithm. If we can show  $\beta$  is small, we can guarantee generalization.

## 10.2 Stability for Convex and Lipschitz Losses

The classical setting for stability analysis is for algorithms that minimize a loss function that is both **convex** and **Lipschitz-continuous**. Many classic machine learning models, like SVMs with hinge loss, fit this description.

For Stochastic Gradient Descent (SGD), the stability coefficient  $\beta$  depends on parameters like the learning rate  $\eta$  and the number of training steps  $T$ . As shown in “Train faster, generalize better” (Hardt et al., 2016), for convex problems, the stability behaves roughly as  $\beta \propto \frac{T\eta}{n}$ .

This leads to a crucial insight:

- **Early Stopping as Regularization:** Running SGD for fewer epochs (a smaller  $T$ ) results in a smaller  $\beta$ , which means a smaller generalization gap. This provides a rigorous theoretical justification for the common practice of early stopping. “Training faster” (in epochs) leads to a more stable algorithm that generalizes better.

## 10.3 Stability in Deep Learning: The Non-Convex Case

The setting for deep learning is far more challenging. The loss landscapes are non-convex, and common loss functions like cross-entropy are not globally Lipschitz. This means the gradient can, in theory, become arbitrarily large, potentially allowing a single data point to have an unbounded influence and destroying stability.

The paper “Fine-Grained Analysis of Stability and Generalization for Stochastic Gradient Descent” provides a groundbreaking analysis that extends stability to this difficult setting.

### 10.3.1 Key Results

The authors show that even for non-convex and non-Lipschitz losses, the stability of SGD can be bounded. The key is that while the loss might not be globally well-behaved, the gradients encountered during training on typical data are often bounded.

They prove that if SGD is run for  $T$  steps (total number of iterations) on a function whose gradients are bounded, the algorithm’s stability is on the order of:

$$\beta \approx \mathcal{O}\left(\frac{T}{n}\right)$$

This is a powerful result with significant implications:

1. **The Number of Epochs is Key:** The stability, and thus the generalization gap, is directly controlled by the number of passes over the data ( $E = T/n$ ). This confirms that the number of training epochs is a critical complexity parameter for deep learning.
2. **Generalization in the Interpolation Regime:** Modern networks are often trained to achieve zero training error (they *interpolate* the data). This analysis shows that even in this regime, the model can generalize well, provided it is trained for a small number of epochs. The algorithm can be stable even if it perfectly fits the training data.

## 10.4 Conclusion

Algorithmic stability provides a powerful, algorithm-centric framework for understanding generalization that overcomes many of the limitations of uniform convergence. It successfully formalizes the intuition that algorithms that are less sensitive to individual data points should generalize better.

The key takeaways are:

- Stability directly bounds the generalization gap.
- It provides a rigorous explanation for why common regularization techniques like **early stopping** are effective.
- Modern analysis has extended this framework to the challenging non-convex and non-Lipschitz setting of deep learning, showing that the number of training epochs is a fundamental controller of generalization.

This places the focus squarely on the behavior of the optimization algorithm, providing one of the most compelling theoretical explanations for the success of deep learning to date.



# 第十一章 Week 11 - Generalization

In our quest to understand generalization, we have moved from class-based analyses (Uniform Convergence) to algorithm-based analyses (Algorithmic Stability). This week, we explore a third, powerful perspective that analyzes a *distribution* over hypotheses rather than a single deterministic one. This is the world of **PAC-Bayes** and **Information-Theoretic** bounds.

## 11.1 The PAC-Bayes Framework

The PAC-Bayes framework provides a bound on the generalization error of a *stochastic* learning algorithm. Instead of returning a single hypothesis, the algorithm returns a full distribution over the hypothesis class.

### 11.1.1 Prior and Posterior Distributions

The framework is built on two key concepts from Bayesian statistics:

- **The Prior Distribution ( $P$ ):** This is a distribution over the hypothesis class  $\mathcal{H}$  that is fixed *before* we see any data. It represents our prior beliefs or assumptions about which hypotheses are likely to be good. A common choice is a simple Gaussian distribution over the model's weights, e.g.,  $P = \mathcal{N}(0, \sigma^2 I)$ .
- **The Posterior Distribution ( $Q$ ):** This is a distribution over  $\mathcal{H}$  that the learning algorithm chooses *after* seeing the training set  $S$ . It represents our updated beliefs, incorporating evidence from the data.

### 11.1.2 The PAC-Bayes Bound

The classic PAC-Bayes theorem (McAllester, 1999) provides a bound on the expected true risk of a hypothesis drawn from the posterior  $Q$ . A simplified version of the theorem is as

follows:

**Theorem (PAC-Bayes):** For any prior  $P$  (independent of the data), with probability at least  $1 - \delta$  over the draw of the training set  $S$  of size  $n$ , the following holds for **all** posterior distributions  $Q$ :

$$\mathbb{E}_{f \sim Q}[R_{\text{true}}(f)] \leq \mathbb{E}_{f \sim Q}[R_{\text{emp}}(f)] + \sqrt{\frac{\text{KL}(Q \parallel P) + \log(n/\delta)}{2n}}$$

Let's break down this powerful statement:

- $\mathbb{E}_{f \sim Q}[R_{\text{true}}(f)]$  is the quantity we want to be small: the expected true risk under our posterior.
- $\mathbb{E}_{f \sim Q}[R_{\text{emp}}(f)]$  is the expected empirical risk. Our algorithm can make this small by choosing a posterior  $Q$  that is concentrated on hypotheses that fit the training data well.
- The complexity term is driven by the **Kullback-Leibler (KL) divergence**,  $\text{KL}(Q \parallel P)$ .

### 11.1.3 KL Divergence and Occam's Razor

The KL divergence measures the "information cost" or "distance" of updating our beliefs from the prior  $P$  to the posterior  $Q$ .

$$\text{KL}(Q \parallel P) = \int_{\mathcal{H}} Q(f) \log \frac{Q(f)}{P(f)} df$$

If the posterior  $Q$  is very different from the prior  $P$ , the KL divergence will be large. If they are similar, it will be small.

The PAC-Bayes bound is a beautiful formalization of **Occam's Razor**: to generalize well, we must find a posterior  $Q$  that achieves a trade-off. It must both (1) explain the data well (small empirical risk) and (2) be simple, in the sense that it doesn't deviate too much from our simple prior (small KL divergence).

## 11.2 Information-Theoretic Bounds

A closely related perspective on generalization comes from information theory. This framework connects the generalization gap directly to the amount of information the learned hypothesis contains about the training set.

### 11.2.1 Mutual Information and Generalization

The key quantity here is the **mutual information**  $I(A(S); S)$  between the output of the learning algorithm  $A(S)$  (which is a hypothesis, or in this context, a random variable representing the hypothesis) and the training set  $S$ .

The mutual information measures, in bits, how much we learn about the specific training set  $S$  by observing the algorithm's output.

- **High  $I(A(S); S)$ :** The learned hypothesis is highly dependent on the specifics of the training data. This is characteristic of overfitting.
- **Low  $I(A(S); S)$ :** The learned hypothesis is largely independent of the training data. This is characteristic of a stable algorithm that captures the underlying distribution.

### 11.2.2 The Information-Theoretic Bound

The work of Russo & Zou and Xu & Raginsky provides the following elegant bound on the expected generalization gap:

**Theorem (Information-Theoretic):** The expected generalization gap is bounded by the mutual information:

$$\mathbb{E}[R_{\text{true}}(A(S)) - R_{\text{emp}}(A(S))] \leq \sqrt{\frac{2 \cdot \text{Var}(\mathcal{L}) \cdot I(A(S); S)}{n}}$$

(Here,  $\text{Var}(\mathcal{L})$  is the variance of the loss, often assumed to be bounded.)

This provides a direct link: to ensure good generalization, we must use a learning algorithm that does not "extract" too much information from or "memorize" the training data. This provides a deep connection back to the concept of algorithmic stability. In fact, one can show that low mutual information implies the algorithm is stable on average.

## 11.3 Conclusion

The PAC-Bayes and Information-Theoretic frameworks offer a profound and unified perspective on generalization. They shift the focus from a single hypothesis to a distribution over hypotheses, which is a natural fit for analyzing stochastic algorithms like SGD.

The core principle is that **generalization is inversely related to the amount of information the model learns from the specific training sample**. This "information

cost” can be measured by the KL divergence between a posterior and a prior (PAC-Bayes) or by the mutual information between the hypothesis and the data.

These frameworks provide deep insights into the implicit bias of optimization algorithms and remain a highly active and promising area of research for finally demystifying the success of deep learning.

## 第十二章 Week 12 - Generalization

We have reached the final topic in our journey to understand the foundations of deep learning. We’ve seen that overparameterized models have the capacity to fit random noise, yet they generalize well on real data. We’ve also seen that classical generalization theories are not sufficient to explain this. The missing piece of the puzzle is not in the model architecture or the loss function alone, but in the behavior of the optimization algorithm itself. This is the theory of **implicit bias**.

### 12.1 The Core Idea: What is Implicit Bias?

When a model is overparameterized, there is typically not just one, but an entire space of solutions (parameter settings) that achieve zero training error. For example, in a simple linear model with more features than data points, there are infinitely many solutions.

**Definition (Implicit Bias):** The implicit bias of a learning algorithm is its tendency to converge to a particular solution among all possible solutions that minimize the training loss. This bias is not caused by explicit regularization terms (like  $\ell_2$  penalty) but is an emergent property of the algorithm’s dynamics.

The central claim is that the optimization algorithms we use, like Gradient Descent (GD) and its stochastic variants, are implicitly biased towards “simple” or “good” solutions—the very kind that generalization theories tell us should perform well on unseen data.

## 12.2 The Simplest Case: Min-Norm Bias in Linear Models

Let's start with an underdetermined linear regression problem. We want to find a weight vector  $\mathbf{w} \in \mathbb{R}^d$  that solves  $\mathbf{X}\mathbf{w} = \mathbf{y}$ , where  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $n < d$ . There are infinitely many solutions.

Which one does Gradient Descent find?

**Theorem:** For the underdetermined least squares problem, if Gradient Descent is initialized at  $\mathbf{w}_0 = \mathbf{0}$ , it will converge to the solution of the optimization problem:

$$\lim_{t \rightarrow \infty} \mathbf{w}(t) = \arg \min_{\mathbf{w}} \|\mathbf{w}\|_2^2 \quad \text{subject to} \quad \mathbf{X}\mathbf{w} = \mathbf{y}$$

GD finds the solution with the **minimum Euclidean ( $\ell_2$ ) norm**.

**Intuition:** The gradient updates for the least squares loss are always in the form  $\mathbf{X}^T(\dots)$ . This means that the parameter vector  $\mathbf{w}(t)$  always remains in the row span of the data matrix  $\mathbf{X}$ . The minimum norm solution to  $\mathbf{X}\mathbf{w} = \mathbf{y}$  is precisely the one that lies entirely within this span. Any component of  $\mathbf{w}$  orthogonal to this span would increase the norm without contributing to the solution. GD, by its nature, never explores these useless directions and thus finds the simplest, min-norm solution.

## 12.3 Implicit Bias in Classification: The Max-Margin Bias

A more profound result comes from classification, as detailed in “The Implicit Bias of Gradient Descent on Separable Data” (Soudry et al., 2018).

### 12.3.1 The Setup

Consider a linear classifier trained on **linearly separable data** using a loss function with an exponential tail, such as the logistic loss or cross-entropy loss. Because the data is separable, there are infinitely many separating hyperplanes that can achieve zero classification error.

To drive the logistic loss to its minimum (zero), the logits  $y_i(\mathbf{w}^T \mathbf{x}_i)$  must go to  $+\infty$  for all  $i$ . This requires the norm of the weight vector,  $\|\mathbf{w}\|$ , to grow infinitely large. So, the parameters themselves do not converge.

However, the *direction* of the weight vector,  $\frac{\mathbf{w}}{\|\mathbf{w}\|}$ , does converge. The question is, to what direction?

**Theorem (Informal):** For separable data, Gradient Descent on an exponential-tailed loss converges in direction to the **max-margin** separating hyperplane. It finds the same solution as the hard-margin Support Vector Machine (SVM).

This is a remarkable result. We did not build any notion of margin into our loss function. Yet, the dynamics of gradient descent on a standard loss function implicitly perform margin maximization. Since we know from previous weeks (e.g., norm-based bounds) that large-margin solutions generalize well, this provides a direct causal link: **GD generalizes because it implicitly maximizes the margin.**

## 12.4 Conclusion: The Story Comes Together

Implicit bias provides a powerful, unifying narrative that connects the disparate themes of our course.

- **Representation (Weeks 1-3):** We learned that neural networks are powerful enough to represent almost any function. This overparameterization is what creates the problem of multiple solutions.
- **Optimization (Weeks 4-6):** We saw how algorithms like SGD can navigate the complex loss landscape to find a solution that perfectly fits the training data (a global minimum of the empirical risk).
- **Generalization (Weeks 7-11):** We established that for a model to generalize, its solution should be “simple” in some sense (e.g., small norm, large margin, stable, low information content).
- **Implicit Bias (Week 12):** This is the final link. It tells us that the optimization algorithm doesn’t just find *any* solution; it is implicitly biased towards finding a solution that possesses exactly the simplicity properties required for good generalization.

The implicit bias of the algorithm acts as a form of hidden regularization, guiding the overparameterized model to a solution that is not only accurate on the training data but also robust and predictive on new data. This concept is central to the modern understanding of deep learning and remains a vibrant area of ongoing research.



## 第十三章 Week 13 - LLM

Our course has built a theoretical foundation for understanding machine learning, culminating in the concept of implicit bias. This week, we examine phenomena that are particularly relevant in the modern era of massive, overparameterized models. These observations challenge classical statistical intuition and provide a new roadmap for building state-of-the-art systems.

### 13.1 The Double Descent Phenomenon

For decades, the cornerstone of statistical learning was the **bias-variance trade-off**, which describes a U-shaped curve for test error as model complexity increases.

- **Underparameterized Regime (Classical):** As we add parameters to a simple model, bias decreases while variance increases. The test error is minimized at a “sweet spot” of complexity, beyond which the model starts to overfit.

However, empirical work in deep learning revealed a surprising extension to this curve.

#### 13.1.1 The Modern Overparameterized Regime

The **double descent** phenomenon, detailed in papers like “Reconciling modern machine learning practice and the classical bias–variance trade-off” (Belkin et al.), describes what happens as we continue to increase model capacity far beyond the classical “sweet spot.”

- **The Interpolation Threshold:** There is a critical point where the model has just enough capacity to perfectly fit (or *interpolate*) the training data. At this point, the number of parameters  $p$  is close to the number of samples  $n$ . The test error often **peaks** here. The model is forced into a single, brittle solution that fits the training data perfectly but is highly sensitive to noise, causing a spike in variance.

- **The Overparameterized Regime ( $p \gg n$ ):** As we increase the model’s capacity even further beyond the interpolation threshold, the test error surprisingly begins to **decrease again**. This is the second “descent.”

### 13.1.2 Why Does the Error Decrease Again?

When the model is massively overparameterized ( $p \gg n$ ), there are once again infinitely many solutions that can perfectly fit the training data. As we saw in our lecture on implicit bias, the optimization algorithm (like GD) does not choose just any of these solutions. It is biased towards a “simple” one.

As the model capacity increases, the set of possible interpolating solutions grows. This makes it easier for the algorithm to find a solution that not only fits the training data but is also smooth and simple (e.g., has a small norm). This “simpler” solution generalizes better, causing the test error to fall.

**Implication:** The most dangerous place to be is often at the interpolation threshold. It can be better to use a much larger model that operates deep in the overparameterized regime than one that just barely has enough capacity to fit the data. “Bigger is better” has a theoretical justification.

## 13.2 Scaling Laws for Large Models

While double descent provides a theoretical picture, **scaling laws** provide an empirical, predictive framework for the behavior of large models. Research from labs like OpenAI and DeepMind, summarized in papers like “Scaling Laws for Neural Language Models” (Kaplan et al.), found that model performance improves smoothly and predictably with scale.

### 13.2.1 The Three Key Factors

The performance of a large model (specifically its test loss,  $L$ ) is primarily determined by three factors:

1. **Model Size ( $N$ ):** The number of non-embedding parameters in the model.
2. **Dataset Size ( $D$ ):** The number of tokens in the training dataset.
3. **Compute ( $C$ ):** The total amount of computation used for training, measured in floating-point operations (FLOPs).

### 13.2.2 The Power-Law Relationship

The key finding is that the test loss scales as a **power-law** with respect to these factors, as long as the other two are not bottlenecks.

$$L(N) \approx \left(\frac{N_c}{N}\right)^{\alpha_N} \quad L(D) \approx \left(\frac{D_c}{D}\right)^{\alpha_D}$$

where  $N_c, D_c, \alpha_N, \alpha_D$  are constants that can be estimated by fitting curves to the results of smaller-scale experiments.

**Implication:** This is a profoundly useful result for engineering. It turns the development of massive models from a high-risk research endeavor into a more predictable engineering discipline. One can invest a small amount of compute to train a range of smaller models, fit the scaling law, and then use the resulting equation to accurately predict the performance of a model that is 100x or 1000x larger, before committing the massive resources to train it.

These laws also provide guidance on how to best allocate a fixed computational budget. For optimal performance, model size and dataset size must be scaled in tandem. For a given increase in compute, there is an optimal way to distribute that increase between making the model larger and training on more data.

## 13.3 Conclusion: The Era of Scale

The phenomena of double descent and scaling laws have reshaped our understanding of how to build the best machine learning models.

- **Double Descent** revises the classical bias-variance trade-off and gives us confidence that pushing models into the deeply overparameterized regime is a sound strategy.
- **Scaling Laws** provide a practical, quantitative guide for navigating this regime, making the behavior of massive models predictable and allowing us to efficiently allocate resources to achieve state-of-the-art performance.

Together, these concepts form the intellectual bedrock of the current era of large-scale AI, where building ever-larger and more capable models is a central goal.



## 第十四章 Week 14 - LLM

In our final lecture, we explore one of the most surprising and impactful discoveries in the era of large models: their capacity for complex reasoning. Last week, we saw that scaling models leads to predictable improvements and new phenomena like double descent. Today, we focus on an *emergent ability* that is unlocked at scale: the ability to reason, elicited through a technique known as **Chain-of-Thought (CoT) prompting**.

### 14.1 Emergent Abilities and In-Context Learning

As language models grow, they don't just get better at existing tasks; they acquire entirely new capabilities that are absent in smaller models. A prime example is **in-context learning**, where a model can perform a new task by simply being shown a few examples in its prompt, without any updates to its weights.

Chain-of-Thought is an advanced form of in-context learning that is specifically designed to unlock a model's latent reasoning abilities.

### 14.2 What is Chain-of-Thought (CoT) Prompting?

CoT is a prompting technique that guides a Large Language Model (LLM) to solve a multi-step problem by instructing it to generate a series of intermediate, step-by-step reasoning steps before arriving at the final answer.

#### 14.2.1 Standard Prompting vs. CoT Prompting

Consider a simple arithmetic word problem.

**Standard Prompting (few-shot):**

**Q:** Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now? **A:** 11.

**Q:** The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have? **A:** 9.

With a standard prompt, the model is expected to map the question directly to the final answer. For multi-step problems, this is a difficult, often brittle mapping.

**Chain-of-Thought Prompting (few-shot):**

**Q:** Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now? **A:** Roger started with 5 balls. 2 cans of 3 tennis balls each is  $2 * 3 = 6$  balls. So he has  $5 + 6 = 11$  balls. The answer is 11.

**Q:** The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have? **A:**

When presented with the CoT-style examples, the model learns the *pattern of reasoning*. For the new question, it will generate a similar step-by-step thought process:

*The cafeteria started with 23 apples. They used 20, so they had  $23 - 20 = 3$ . Then they bought 6 more, so they have  $3 + 6 = 9$ . The answer is 9.*

By decomposing the problem into a series of simpler steps, CoT dramatically improves the model's performance on tasks requiring arithmetic, commonsense, and symbolic reasoning.

## 14.3 A Theoretical Perspective on CoT

Why does this work? The paper "Towards Revealing the Mystery behind Chain of Thought: A Theoretical Perspective" (Feng et al., 2023) provides a formal framework for understanding this phenomenon.

### 14.3.1 Transforming "Slow" Problems into "Fast" Steps

The paper draws an analogy to Daniel Kahneman's two systems of thought:

- **System 1 (Fast):** Intuitive, automatic, parallel. A standard prompt asks the model to solve the problem in a single forward pass, which is akin to a System 1 task.
- **System 2 (Slow):** Deliberate, sequential, logical. CoT allows the model to engage in a process similar to System 2 thinking.

The core argument is that CoT transforms a single, computationally difficult “slow” reasoning problem into a **sequence of easy, “fast” sub-problems**. Each step in the chain of thought is a simpler task (e.g., “what is  $2 * 3$ ?”) that the model can solve reliably in a single step.

### 14.3.2 The Latent Skill Model

The paper formalizes this by modeling complex tasks as requiring a sequence of **latent skills**. For the tennis ball problem, these skills might be: ‘identify initial quantity’, ‘calculate new quantity’, ‘perform addition’.

- A standard prompt requires the model to learn a direct (and complex) mapping from the question to the answer:  $P(\text{answer}|\text{question})$ .
- A CoT prompt guides the model to learn a sequence of simpler mappings:  $P(\text{skill}_1|\text{question})$ , then  $P(\text{step}_1|\text{skill}_1, \text{question})$ , and so on.

By providing the reasoning chain, we give the model intermediate supervision, guiding it through the latent space of skills needed to solve the problem. This significantly simplifies the learning task.

Furthermore, CoT forces the model to externalize its reasoning, which discourages it from relying on flimsy “semantic shortcuts” (e.g., just multiplying any two numbers in a word problem) and encourages it to follow a more robust, logical process.

## 14.4 Conclusion: From Scaling to Reasoning

Chain-of-Thought is more than just a clever prompting trick; it represents a fundamental shift in how we interact with and understand large models.

- It demonstrates that at sufficient scale, models acquire latent reasoning abilities that can be “activated” through appropriate prompting.

- It highlights that the **interface** to the model is as crucial as its architecture. Prompt engineering is a key component of harnessing the model’s full potential.
- It moves us closer to a more transparent and interpretable AI, where the model can “show its work,” allowing us to verify its reasoning process.

The journey of this course has taken us from the basic building blocks of representation to the frontiers of AI. We began by asking what a network can represent and ended by seeing how, at scale, these representations can support complex, human-like reasoning. The theoretical work to fully explain these emergent phenomena is just beginning, marking one of the most exciting frontiers in the science of artificial intelligence.



# 第十五章 Week 15 - LLM

Welcome to the final lecture of our course. We have traced the theoretical underpinnings of machine learning from the single neuron to the planet-scale models of today. Our journey concludes with one of the most fascinating and debated topics in modern AI: the phenomenon of **emergence**.

## 15.1 The Surprise of Scale: What is Emergence?

We've seen that scaling models leads to predictable, continuous improvements in performance (Scaling Laws). However, scale also brings surprises. Certain complex abilities seem to appear suddenly and unpredictably in large models, while being completely absent in their smaller counterparts.

**Definition (Emergent Ability):** An ability is emergent if it is not present in smaller models but is present in larger models. Performance on a task involving this ability is near-random until a certain critical scale is reached, after which performance rises sharply to well above random.

This is analogous to a **phase transition** in physics, like water suddenly freezing into ice at 0°C. Examples of emergent abilities include multi-step arithmetic, answering questions in unexpected languages, and the very ability to perform Chain-of-Thought reasoning, which we studied last week.

## 15.2 A Theoretical Explanation for Emergence

Is this phenomenon a deep mystery, a true "phase transition" in the model's intelligence? Or is there a simpler explanation? The paper "A Theory for Emergence of Complex Skills

in Language Models” (Schaeffer, Miranda, & Koyejo, 2024) offers a compelling theoretical argument.

### 15.2.1 The Main Argument: Emergence as a Perceptual Artifact

The paper proposes that emergence is not a fundamental property of the model’s learning process, but rather an artifact of the **metrics we use to measure performance**.

The core idea is that a model’s underlying capabilities on basic sub-tasks may be improving smoothly and predictably with scale. However, when we evaluate its performance on a complex task that requires the successful execution of *many* such sub-tasks, this smooth improvement appears to us as a sudden, sharp jump.

### 15.2.2 A Model for Complex Task Performance

Let’s consider a complex task that requires the successful completion of  $k$  independent sub-tasks or skills.

- Let  $p_{\text{sub}}$  be the probability that the model can successfully execute a single sub-task. We can assume that  $p_{\text{sub}}$  improves smoothly with model scale, in line with the scaling laws we studied.
- To succeed at the overall complex task, the model must get all  $k$  sub-tasks correct. The probability of success on the complex task is therefore:

$$P_{\text{complex}} = (p_{\text{sub}})^k$$

### 15.2.3 The Illusion of Sharpness

This simple non-linear relationship,  $P_{\text{complex}} = (p_{\text{sub}})^k$ , is sufficient to create the visual appearance of emergence.

- When the model is small, its capability on each sub-task  $p_{\text{sub}}$  is low (e.g., 0.1). The overall performance,  $(0.1)^k$ , is astronomically small and effectively zero for any moderately large  $k$ .
- As the model scales,  $p_{\text{sub}}$  improves smoothly and linearly. For a while,  $P_{\text{complex}}$  remains near zero.

- However, as  $p_{\text{sub}}$  crosses a certain threshold (e.g., 0.5, 0.6, 0.7), the exponential function  $(p_{\text{sub}})^k$  causes the overall performance to rise extremely rapidly.

This smooth, continuous improvement in the model’s core capabilities ( $p_{\text{sub}}$ ) is perceived by us as a sudden, discontinuous ”emergence” of the complex skill, simply because of our choice of an ”all-or-nothing” success metric.

**Implication:** The mystery of emergence may be less about an unknown phase transition inside the model and more about how we define and measure success. The underlying reality may be one of predictable, continuous improvement governed by scaling laws.

## 15.3 Conclusion: The End of the Beginning

This brings our course to a close. We have journeyed from the foundational question of what a single neuron can represent to the emergent, seemingly intelligent behaviors of planet-scale models.

- We learned that **Representation** is unlocked by depth and non-linearity.
- We saw that **Optimization** is driven by the remarkable ability of SGD to navigate non-convex landscapes, guided by an **Implicit Bias** towards simple solutions.
- We understood that **Generalization** is a complex interplay of stability, implicit regularization, and the structure of data itself.
- Finally, we’ve seen that massive **Scale** is a transformative ingredient, leading to predictable improvements (Scaling Laws), new theoretical regimes (Double Descent), new ways of interacting with models (CoT), and surprising emergent abilities.

The theory of emergence reminds us that as we build more powerful AI systems, we must also build more sophisticated tools for measuring and understanding them. The mysteries that remain are not just about the models, but about our own perception and interpretation of their behavior. The quest to build a complete, rigorous, and predictive theory of intelligence remains one of the greatest scientific challenges of our time.

