

# Introductions of Approximation Optimization and Generalization Errors

Shijun Zhang\*

January 4, 2021

In this note, we introduce approximation, optimization, and generalization errors in order to measure the discrepancy between the target function and the final network attained by a numerical training or optimization method.

Let  $\phi(\mathbf{x}; \boldsymbol{\theta})$  denote a function computed by a (fully-connected) network with  $\boldsymbol{\theta}$  as the set of parameters. See Figure 1 for an example of a  $\sigma$ -activated network with width 5 and depth 2.

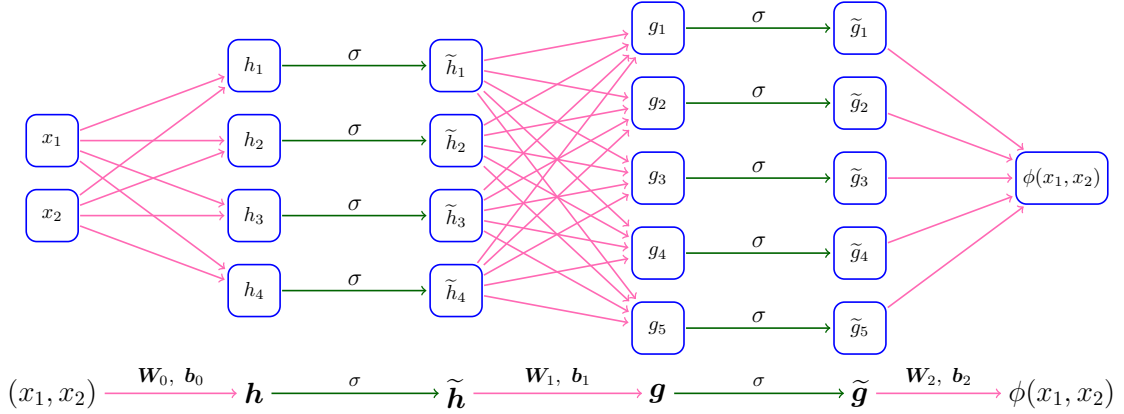


Figure 1: An example of a  $\sigma$ -activated network with width 5 and depth 2. In this example,  $\boldsymbol{\theta} = (\mathbf{W}_0, \mathbf{b}_0, \mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2)$ .

Given a target function  $f$ , consider the expected error/risk of  $\phi(\mathbf{x}; \boldsymbol{\theta})$

$$R_{\mathcal{D}}(\boldsymbol{\theta}) := \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}), f(\mathbf{x}))]$$

with a loss function typically taken as  $\ell(y, y') = \frac{1}{2}|y - y'|^2$ , where  $U(\mathcal{X})$  is an unknown data **distribution** over  $\mathcal{X}$ . For example, when  $\ell(y, y') = \frac{1}{2}|y - y'|^2$  and  $U$  is a uniform distribution over  $\mathcal{X} = [0, 1]^d$ ,

$$R_{\mathcal{D}}(\boldsymbol{\theta}) = \int_{[0,1]^d} \frac{1}{2} |\phi(\mathbf{x}; \boldsymbol{\theta}) - f(\mathbf{x})|^2 d\mathbf{x}.$$

---

\*Department of Mathematics, National University of Singapore (zhangshijun@u.nus.edu).

17 The expected risk minimizer  $\theta_{\mathcal{D}}$  is defined as

$$18 \quad \theta_{\mathcal{D}} := \arg \min_{\theta} R_{\mathcal{D}}(\theta).$$

19 It is unachievable in practice since  $f$  and  $U(\mathcal{X})$  are not available.

20 In practice, for given **samples**  $\{(\mathbf{x}_i, f(\mathbf{x}_i))\}_{i=1}^n$ , we use the empirical risk

$$21 \quad R_{\mathcal{S}}(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(\phi(\mathbf{x}_i; \theta), f(\mathbf{x}_i)).$$

22 to approximate/model the expected risk  $R_{\mathcal{D}}(\theta)$ . The goal of supervised learning is to  
23 identify the empirical risk minimizer

$$24 \quad \theta_{\mathcal{S}} := \arg \min_{\theta} R_{\mathcal{S}}(\theta), \quad (1)$$

25 obtain  $\phi(\mathbf{x}; \theta_{\mathcal{S}}) \approx f(\mathbf{x})$ . When a numerical optimization method is applied to solve (1),  
26 it may result in a numerical solution (denoted as  $\theta_{\mathcal{N}}$ ) that is not a global minimizer.  
27 Hence, the actually learned function generated by a neural network is  $\phi(\mathbf{x}; \theta_{\mathcal{N}})$ . And the  
28 discrepancy between the target function  $f$  and the actually learned function  $\phi(\mathbf{x}; \theta_{\mathcal{N}})$  is  
29 measured by an inference error

$$30 \quad R_{\mathcal{D}}(\theta_{\mathcal{N}}) = \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \theta_{\mathcal{N}}), f(\mathbf{x}))] \stackrel{e.g.}{=} \int_{[0,1]^d} \frac{1}{2} |\phi(\mathbf{x}; \theta_{\mathcal{N}}) - f(\mathbf{x})|^2 d\mathbf{x},$$

31 where the second equality holds when  $\ell(y, y') = \frac{1}{2}|y - y'|^2$  and  $U$  is a uniform distribution  
32 over  $\mathcal{X} = [0, 1]^d$ ,

33 Since  $R_{\mathcal{D}}(\theta_{\mathcal{N}})$  is the expected inference error over all possible data samples, it can  
34 quantify how good the learned function  $\phi(\mathbf{x}; \theta_{\mathcal{N}})$  is. Note that

$$\begin{aligned} 35 \quad R_{\mathcal{D}}(\theta_{\mathcal{N}}) &= \underbrace{[R_{\mathcal{D}}(\theta_{\mathcal{N}}) - R_{\mathcal{S}}(\theta_{\mathcal{N}})]}_{\text{GE}} + \underbrace{[R_{\mathcal{S}}(\theta_{\mathcal{N}}) - R_{\mathcal{S}}(\theta_{\mathcal{S}})]}_{\text{OE}} + \underbrace{[R_{\mathcal{S}}(\theta_{\mathcal{S}}) - R_{\mathcal{S}}(\theta_{\mathcal{D}})]}_{\leq 0 \text{ by Eq. (1)}} + \underbrace{[R_{\mathcal{S}}(\theta_{\mathcal{D}}) - R_{\mathcal{D}}(\theta_{\mathcal{D}})]}_{\text{GE}} + \underbrace{R_{\mathcal{D}}(\theta_{\mathcal{D}})}_{\text{AE}} \\ 36 \quad &\leq \underbrace{R_{\mathcal{D}}(\theta_{\mathcal{D}})}_{\text{Approximation error (AE)}} + \underbrace{[R_{\mathcal{S}}(\theta_{\mathcal{N}}) - R_{\mathcal{S}}(\theta_{\mathcal{S}})]}_{\text{Optimization error (OE)}} + \underbrace{[R_{\mathcal{D}}(\theta_{\mathcal{N}}) - R_{\mathcal{S}}(\theta_{\mathcal{N}})] + [R_{\mathcal{S}}(\theta_{\mathcal{D}}) - R_{\mathcal{D}}(\theta_{\mathcal{D}})]}_{\text{Generalization error (GE)}}. \quad (2) \\ 37 \quad &\quad \quad \quad \text{Approximation error (AE)} \quad \quad \quad \text{Optimization error (OE)} \quad \quad \quad \text{Generalization error (GE)} \end{aligned}$$

38 where the inequality comes from the fact that  $[R_{\mathcal{S}}(\theta_{\mathcal{S}}) - R_{\mathcal{S}}(\theta_{\mathcal{D}})] \leq 0$  since  $\theta_{\mathcal{S}}$  is a global  
39 minimizer of  $R_{\mathcal{S}}(\theta)$ .

40 Constructive approximation provides an upper bound of  $R_{\mathcal{D}}(\theta_{\mathcal{D}})$  in terms of the  
41 network size, e.g., in terms of the network width and depth, or in terms of the number  
42 of parameters. The second term of Equation (2) is bounded by the optimization error  
43 of the numerical algorithm applied to solve the empirical loss minimization problem in  
44 Equation (1). Note that one only needs to make  $R_{\mathcal{S}}(\theta_{\mathcal{N}}) - R_{\mathcal{S}}(\theta_{\mathcal{S}})$  small, but not  $\theta_{\mathcal{N}} - \theta_{\mathcal{S}}$ .  
45 The study of the bounds for the third and fourth terms is referred to as the generalization  
46 error analysis of neural networks. See Figure 2 for the intuitions of these three errors.

47 One of the key targets in the area of deep learning is to develop algorithms to  
48 reduce  $R_{\mathcal{D}}(\theta_{\mathcal{N}})$ . In [1, 2, 3, 4, 5, 6], we provide upper bounds of the approximation error  
49  $R_{\mathcal{D}}(\theta_{\mathcal{D}})$  for several function spaces, which is crucial to estimate an upper bound of  
50  $R_{\mathcal{D}}(\theta_{\mathcal{N}})$ . Instead of deriving an approximator to attain the approximation error bound,

51 deep learning algorithms aim to identify a solution  $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$  reducing the generalization  
 52 and optimization errors in Equation (2). Solutions minimizing both generalization and  
 53 optimization errors will lead to a good solution only if we also have a good upper bound  
 54 estimate of  $R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{D}})$  as shown in Equation (2). Independent of whether our analysis here  
 55 leads to a good approximator, which is an interesting topic to pursue, the theory here  
 56 does provide a key ingredient in the error analysis of deep learning algorithms.

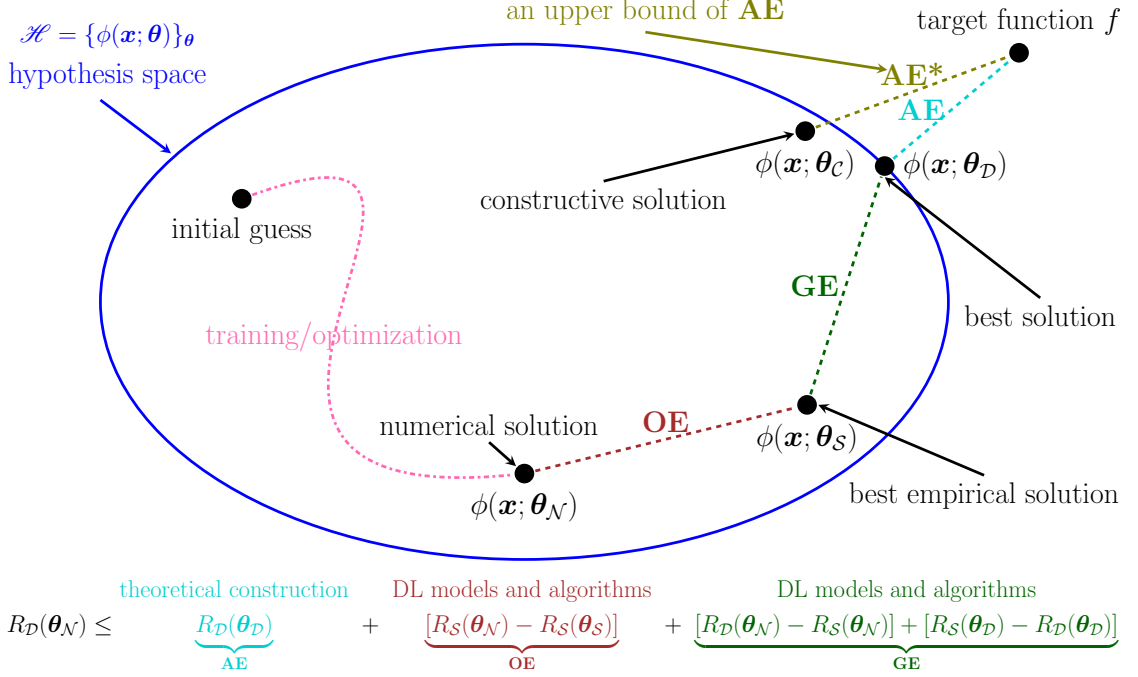


Figure 2: The intuitions of the approximation error (AE), the optimization error (OE), and the generalization error (GE). DL is short of deep learning. One needs to control AE, OE, and GE in order to bound the discrepancy between the target function  $f$  and the numerical solution  $\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}})$  (what we can get in practice), measured by

$$R_{\mathcal{D}}(\boldsymbol{\theta}_{\mathcal{N}}) = \mathbb{E}_{\mathbf{x} \sim U(\mathcal{X})} [\ell(\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}}), f(\mathbf{x}))] \stackrel{e.g.}{=} \int_{[0,1]^d} \frac{1}{2} |\phi(\mathbf{x}; \boldsymbol{\theta}_{\mathcal{N}}) - f(\mathbf{x})|^2 d\mathbf{x}.$$

## 57 References

- 58 [1] Jianfeng Lu, Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approx-  
 59 imation for smooth functions. *arXiv e-prints*, page arXiv:2001.03040, Jan 2020.
- 60 [2] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Nonlinear approximation via com-  
 61 positions. *Neural Networks*, 119:74–84, 2019.
- 62 [3] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network approximation  
 63 characterized by number of neurons. *Communications in Computational Physics*,  
 64 28(5):1768–1811, 2020.
- 65 [4] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Deep network with approximation  
 66 error being reciprocal of width to power of square root of depth. *arXiv e-prints*, page  
 67 arXiv:2006.12231, Jun 2020.

- 68 [5] Zuowei Shen, Haizhao Yang, and Shijun Zhang. Neural network approximation:  
69 Three hidden layers are enough. *arXiv e-prints*, page arXiv:2010.14075, Oct 2020.
- 70 [6] Shijun Zhang. *Deep Neural Network Approximation via Function Compositions*. PhD  
71 thesis, Signapre, 2020.