

# A selective overview of deep learning

Jianqing Fan, Cong Ma and Yiqiao Zhong

Princeton University and Stanford University

*Abstract.* Deep learning has achieved tremendous success in recent years. In simple words, deep learning uses the composition of many nonlinear functions to model the complex dependency between input features and labels. While neural networks have a long history, recent advances have significantly improved their empirical performance in computer vision, natural language processing, and other predictive tasks. From the statistical and scientific perspective, it is natural to ask: What is deep learning? What are the new characteristics of deep learning, compared with classical statistical methods? What are the theoretical foundations of deep learning?

To answer these questions, we introduce common neural network models (e.g., convolutional neural nets, recurrent neural nets, generative adversarial nets) and training techniques (e.g., stochastic gradient descent, dropout, batch normalization) from a statistical point of view. Along the way, we highlight new characteristics of deep learning (including depth and over-parametrization) and explain their practical and theoretical benefits. We also sample recent results on theories of deep learning, many of which are only suggestive. While a complete understanding of deep learning remains elusive, we hope that our perspectives and discussions serve as a stimulus for new statistical research.

*Key words and phrases:* neural networks, over-parametrization, stochastic gradient descent, approximation theory, generalization error.

*MSC 2010 subject classifications:* Primary 62-01, 62-02; secondary 62H-30.

## 1. INTRODUCTION

Modern machine learning and statistics deal with the problem of *learning from data*: given a training dataset  $\{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$  where  $\mathbf{x}_i \in \mathbb{R}^d$  contains the input features and  $y_i \in \mathbb{R}$  is the output<sup>1</sup>, one seeks a function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  from a certain function class  $\mathcal{F}$  that has good prediction performance on test data. This problem is of fundamental significance and finds applications in numerous

---

*Department of ORFE, Princeton University, Princeton, NJ, 08544 (e-mail: [jqfan@princeton.edu](mailto:jqfan@princeton.edu); [congm@princeton.edu](mailto:congm@princeton.edu); [yiqiaoz@stanford.edu](mailto:yiqiaoz@stanford.edu))*

\*J. Fan is supported in part by the NSF grants DMS-1712591 and DMS-1662139, the NIH grant R01-GM072611 and Cong Ma by the ONR grant N00014-19-1-2120.

<sup>1</sup>When the label  $y$  is given, this problem is often known as *supervised learning*. We mainly focus on this paradigm throughout this paper and remark sparingly on its counterpart, *unsupervised learning*, where  $y$  is not given.

contexts. For instance, in image recognition, the input  $\mathbf{x}$  (reps. the output  $y$ ) corresponds to the raw image (reps. its category) and the goal is to find a mapping  $f(\cdot)$  that can classify future images accurately. Decades of research efforts in statistical machine learning have been devoted to developing methods to find  $f(\cdot)$  efficiently with provable guarantees. Prominent examples include linear classifiers (e.g., linear / logistic regression, linear discriminant analysis), kernel methods (e.g., support vector machines), tree-based methods (e.g., decision trees, random forests), nonparametric regression (e.g., nearest neighbors, local kernel smoothing), etc. Roughly speaking, each aforementioned method corresponds to a different function class  $\mathcal{F}$  from which the final classifier  $f(\cdot)$  is chosen.

Deep learning (LeCun, Bengio and Hinton, 2015), in its simplest form, proposes the following *compositional* function class:

$$(1.1) \quad \{f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}_L \sigma_L(\mathbf{W}_{L-1} \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}))) \mid \boldsymbol{\theta} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}\}.$$

Here, the set of matrices  $\boldsymbol{\theta} = \{\mathbf{W}_1, \dots, \mathbf{W}_L\}$  are the parameters of the model, and for each  $1 \leq l \leq L$ ,  $\sigma_l(\cdot)$  is some fixed nonlinear function. Though simple, deep learning has made significant progress towards addressing the problem of learning from data over the past decade. Specifically, it has performed close to or better than humans in various important tasks in artificial intelligence, including image recognition (He et al., 2016a), game playing (Silver et al., 2017), and machine translation (Wu et al., 2016). Owing to its great promise, the impact of deep learning is also growing rapidly in areas beyond artificial intelligence; examples include statistics (Bauer et al., 2019; Schmidt-Hieber, 2017; Liang, 2017; Romano, Sesia and Candès, 2018; Gao et al., 2018), applied mathematics (Weinan, Han and Jentzen, 2017; Chen et al., 2018), clinical research (De Fauw et al., 2018), etc.

TABLE 1  
Winning models for ILSVRC image classification challenge.

| Model      | Year   | # Layers | # Params | Top-5 error |
|------------|--------|----------|----------|-------------|
| Shallow    | < 2012 | —        | —        | > 25%       |
| AlexNet    | 2012   | 8        | 61M      | 16.4%       |
| VGG19      | 2014   | 19       | 144M     | 7.3%        |
| GoogleNet  | 2014   | 22       | 7M       | 6.7%        |
| ResNet-152 | 2015   | 152      | 60M      | 3.6%        |

To get a better idea of the success of deep learning, let us take the ImageNet Challenge (Russakovsky et al., 2015) (also known as ILSVRC) as an example. In the classification task, one is given a training dataset consisting of 1.2 million color images with 1000 categories, and the goal is to classify images based on the input pixels. The performance of a classifier is then evaluated on a test dataset of  $10^5$  images, and in the end the top-5 error<sup>2</sup> is reported. Table 1 highlights a few popular models and their corresponding performance. As can be seen, deep learning models (the second to the last rows) have a clear edge over shallow models (the first row) that fit linear models / tree-based models on handcrafted features. This significant improvement raises three foundational questions:

*Why is deep learning better than classical methods on tasks like image recognition?*

---

<sup>2</sup>The algorithm makes an error if the true label is not contained in the top-5 predictions made by the algorithm.

*Why doesn't over-parametrization hurt predictions?*

*How are such high-dimensional highly nonconvex functions optimized?*

### 1.1 Intriguing new characteristics of deep learning

It is widely acknowledged that two indispensable factors contribute to the success of deep learning, namely (1) huge datasets that often contain millions of samples and (2) immense computing power resulting from clusters of graphics processing units (GPUs). Admittedly, these resources are only recently available: The latter allows us to train larger and deeper neural networks which reduces biases and the former enables variance reduction. However, these two alone are not sufficient to explain the mystery of deep learning due to some of its intriguing new characteristics.<sup>3</sup>

*1.1.1 Depth and approximation.* Neural networks are closely akin to the classical projection pursuit regression model (Friedman and Stuetzle, 1981), where one postulates

$$(1.2) \quad f(\mathbf{x}) = \sum_{k=1}^K g_k(\beta_k^\top \mathbf{x}),$$

with  $\{\beta_k\}$  being unit vectors and  $\{g_k(\cdot)\}$  being univariate smooth functions. One can essentially view the projection pursuit regression model (1.2) as a neural network with one hidden layer, where  $\{g_k(\cdot)\}$  constitutes the nonlinear activation functions.<sup>4</sup> In contrast, deep learning models the complicated nonlinearity through composing many simple nonlinear functions (cf. (1.1)), namely the depth is much larger than one. While it is known, theoretically, that neural networks with one hidden layer is already a *universal approximator*, i.e., it can approximate any continuous function if the number of hidden units is taken arbitrarily large, empirically, deeper neural networks generally have better performance; see Table 1 for an example. The heuristic explanation for this multilayer structure is that, in many real-world datasets such as images, there are different levels of features and lower-level features are building blocks of higher-level ones. See Yosinski et al. (2015) for a visualization of trained features of convolutional neural nets. This is also supported by empirical results from physiology and neuroscience (Hubel and Wiesel, 1962; Abbasi-Asl et al., 2018). Nevertheless, the understanding of the benefits of depth is far from complete.

*1.1.2 Over-parametrization and generalization.* Modern deep learning models often have far more parameters than the number of samples. Such a large degree of *over-parametrization* allows the model to perfectly fit the training data — a phenomenon often observed in practice (Zhang et al., 2016). This seems to be at odds with traditional statistical wisdom. Indeed, the bias-variance tradeoff tells us that the test / generalization error usually exhibits a U-shaped curve as model complexity increases, and in particular, a model that easily overfits training data is prone to worse performance on the test data. Surprisingly, it is found

<sup>3</sup>We note that these characteristics are interleaved in an intricate way. For instance, depth also plays an important role in the optimization aspect of deep learning. Here, we choose the following decomposition to simplify the narrative.

<sup>4</sup>Note that in the the projection pursuit regression model, the activation functions  $\{g_k(\cdot)\}$  are unspecified and fully nonparametric.

that neural networks can exhibit the “double descent phenomenon”, i.e., the test error further decreases once the model complexity is beyond a certain critical threshold (Belkin et al., 2019; Hastie et al., 2019). This curious phenomenon motivates many statistical questions: What is a good measure of model complexity for over-parametrized networks? How do practical algorithms such as stochastic gradient descent find a low-complexity classifier that generalizes well?

*1.1.3 Nonconvexity and optimization.* Training deep learning models, e.g., estimating the parameters  $\theta$  in (1.1) via maximum likelihood estimation, often results in nonconvex optimization problems that are computationally intractable to solve. In general, we cannot hope for an algorithm that can return to us the global solution and different algorithms might end with different estimators in the parameter space. Therefore, the statistical performance of neural networks (e.g., test accuracy) depends heavily on the particular optimization algorithm used for training (Wilson et al., 2017). This is drastically different from many classical statistical problems, where the related optimization problems are less complicated. For instance, when the associated optimization problem has a relatively simple structure (e.g., convex objective functions, linear constraints), the solution to the optimization problem can often be unambiguously computed and analyzed. However, in deep neural networks, due to over-parametrization, there are usually many local minima with different statistical performance (Li et al., 2018a). Nevertheless, running stochastic gradient descent with random initialization often finds model parameters with good out-of-sample performance.

*1.1.4 Implicit representation learning.* Although deep learning is often used as discriminative models to predict the outcome  $y$  based on the features  $\mathbf{x}$ , its power of obtaining useful representations of the input data  $\mathbf{x}$  has been well documented in the literature. For instance, after training on a large amount of labeled images, the hidden units of deep neural networks can represent features such as edges, corners, wheels, eyes, etc.; see Yosinski et al. (2015). Moreover, people can use a part of the trained models to achieve better performance on different but related datasets than training from scratch (a.k.a. transfer learning) (Yosinski et al., 2014). In the statistical language, this suggests that deep learning not only performs well in finding the conditional density  $p(y|\mathbf{x})$ , but also implicitly informs us of the prior  $p(\mathbf{x})$  even though no explicit unsupervised learning is involved.

## 1.2 Towards the theory of deep learning

Despite the empirical success, the theoretical support for deep learning is still in its infancy. Setting the stage, for any classifier  $f$ , denote by<sup>5</sup>  $R(f)$  the expected risk on fresh sample (a.k.a. test error, prediction error or generalization error), and by<sup>6</sup>  $R_n(f)$  the empirical risk / training error averaged over the training dataset. Arguably, the key theoretical question in deep learning is

*why is  $R(\hat{f}_n)$  small, where  $\hat{f}_n$  is the classifier returned by the training algorithm?*

We follow the conventional approximation-estimation decomposition (sometimes, also bias-variance tradeoff) to decompose the term  $R(\hat{f}_n)$  into two parts. Let  $\mathcal{F}$  be the function class expressible by a family of neural nets. Define  $f^* \triangleq$

<sup>5</sup>In the literature,  $\mathbb{E}(f)$  is also used to denote this quantity.

<sup>6</sup>Parallel to the previous footnote,  $\mathbb{E}_n(f)$  is also used in the literature to denote this quantity.

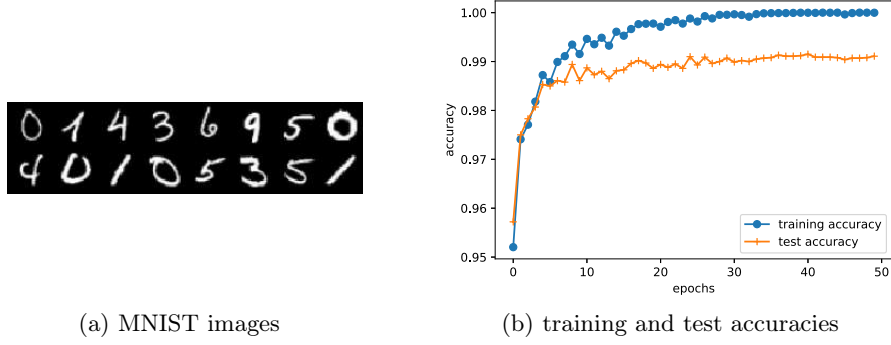


Fig 1: (a) shows the images in the public dataset MNIST; and (b) depicts the training and test accuracies along the training dynamics. Note that the training accuracy is approaching 100% and the test accuracy is still high.

$\operatorname{argmin}_f R(f)$  to be the best possible classifier and  $f_{\mathcal{F}}^* \triangleq \operatorname{argmin}_{f \in \mathcal{F}} R(f)$  to be the best classifier in  $\mathcal{F}$ . Then, we can decompose the excess risk  $\mathcal{E} \triangleq R(\hat{f}_n) - R(f^*)$  into two parts:

$$(1.3) \quad \mathcal{E} = \underbrace{R(f_{\mathcal{F}}^*) - R(f^*)}_{\text{approximation error}} + \underbrace{R(\hat{f}_n) - R(f_{\mathcal{F}}^*)}_{\text{estimation error}}.$$

Both errors can be small for deep learning (cf. Figure 1), which we explain below.

- The *approximation error* is determined by the function class  $\mathcal{F}$ . Intuitively, the larger the class, the smaller the approximation error. Deep learning models use many layers of nonlinear functions (cf. Figure 2) that drive this error small. Indeed, in Section 5, we provide recent theoretical progress of its representation power. For example, deep models allow an efficient representation of the interactions among variable while shallow models cannot.
- The *estimation error* reflects the generalization power, which is influenced by both the complexity of the function class  $\mathcal{F}$  and the properties of the training algorithms. Interestingly, for *over-parametrized* deep neural nets, stochastic gradient descent typically results in a near-zero training error (i.e.,  $R_n(\hat{f}_n) \approx 0$ ; see e.g., the right panel of Figure 1). Moreover, its generalization error remains small or moderate. This “counterintuitive” behavior suggests that for over-parametrized models, gradient-based algorithms enjoy benign statistical properties; we shall see in Section 7 that gradient descent enjoys *implicit regularization* in the over-parametrized regime even without explicit regularization (e.g.,  $\ell_2$  regularization).

The above two points lead to the following heuristic explanation of the success of deep learning models. The large depth of deep neural nets and heavy over-parametrization lead to small or zero training errors, even when running simple algorithms with a moderate number of iterations. In addition, these simple algorithms with a moderate number of iterations do not explore the entire function space and thus have limited complexities, which results in small generalization error with a large sample size.

### 1.3 A roadmap of the paper

We first introduce basic deep learning models in Sections 2–4, and then examine their representation power via the lens of approximation theory in Section 5. Section 6 is devoted to training algorithms and their ability to drive the training error small. In Section 7, we sample recent theoretical progress towards demystifying the generalization power of deep learning. Along the way, we provide our own perspectives, and in the end, we identify a few interesting questions for future research in Section 8. The goal of this paper is to present suggestive methods and results, rather than giving conclusive arguments (which is currently unlikely) or a comprehensive survey. We hope that our discussion serves as a stimulus for new statistics research.

## 2. FEED-FORWARD NEURAL NETWORKS

Before introducing the vanilla feed-forward neural nets, let us set up necessary notations for the rest of this section. We focus primarily on classification problems, as regression problems can be addressed similarly. Given the training dataset  $\{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$  where  $y_i \in [K] \triangleq \{1, 2, \dots, K\}$  and  $\mathbf{x}_i \in \mathbb{R}^d$  are independent across  $i \in [n]$ , supervised learning aims at finding a (possibly random) function  $\hat{f}(\mathbf{x})$  that predicts the outcome  $y$  for a new input  $\mathbf{x}$ , assuming that  $(y, \mathbf{x})$  follows the same distribution as  $\{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$ . In the terminology of machine learning, the input  $\mathbf{x}_i$  is often called the *feature*, the output  $y_i$  called the *label*, and the pair  $(y_i, \mathbf{x}_i)$  is an *example*. They correspond to the statistical terminologies covariate (or predictor), categorical response, and sample. The function  $\hat{f}$  is called the *classifier*, and estimation of  $\hat{f}$  is *training* or *learning*, which is the same as an estimator in statistics. The performance of  $\hat{f}$  is evaluated through the mis-classification error  $\mathbb{P}(y \neq \hat{f}(\mathbf{x}))$ , which can be often estimated from a separate test dataset.

As with classical statistical estimation, for each  $k \in [K]$ , a classifier approximates the conditional probability  $\mathbb{P}(y = k | \mathbf{x})$  using a function  $f_k(\mathbf{x}; \boldsymbol{\theta}_k)$  parametrized by  $\boldsymbol{\theta}_k$ . Then the category with the highest probability is predicted. Thus, learning is essentially estimating the parameters  $\boldsymbol{\theta}_k$ . In statistics, one of the most popular methods is (multinomial) logistic regression, which stipulates a specific form for the functions  $f_k(\mathbf{x}; \boldsymbol{\theta}_k)$ : let  $z_k = \mathbf{x}^\top \boldsymbol{\beta}_k + \alpha_k$  and  $f_k(\mathbf{x}; \boldsymbol{\theta}_k) = Z^{-1} \exp(z_k)$  where  $Z \triangleq \sum_{k=1}^K \exp(z_k)$  is a normalization factor to make  $\{f_k(\mathbf{x}; \boldsymbol{\theta}_k)\}_{1 \leq k \leq K}$  a valid probability distribution. It is clear that logistic regression induces linear decision boundaries in  $\mathbb{R}^d$ , and hence it is restrictive in modeling the nonlinear dependency between  $y$  and  $\mathbf{x}$ . The deep neural networks we introduce below provide a flexible framework for modeling nonlinearity in a fairly general way.

### 2.1 Model setups

From the high level, deep neural networks (DNNs) use the composition of a series of simple nonlinear functions to model nonlinearity

$$\mathbf{h}^{(L)} = \mathbf{g}^{(L)} \circ \mathbf{g}^{(L-1)} \circ \dots \circ \mathbf{g}^{(1)}(\mathbf{x}),$$

where  $\circ$  denotes the composition of two functions and  $L$  is the number of hidden layers, which is usually called the *depth* of a neural network (NN) model. Letting  $\mathbf{h}^{(0)} \triangleq \mathbf{x}$ , one can recursively define  $\mathbf{h}^{(\ell)} \triangleq \mathbf{g}^{(\ell)}(\mathbf{h}^{(\ell-1)})$  for all  $\ell = 1, 2, \dots, L$ . The *feed-forward neural networks*, also called the *multilayer perceptrons* (MLPs), are



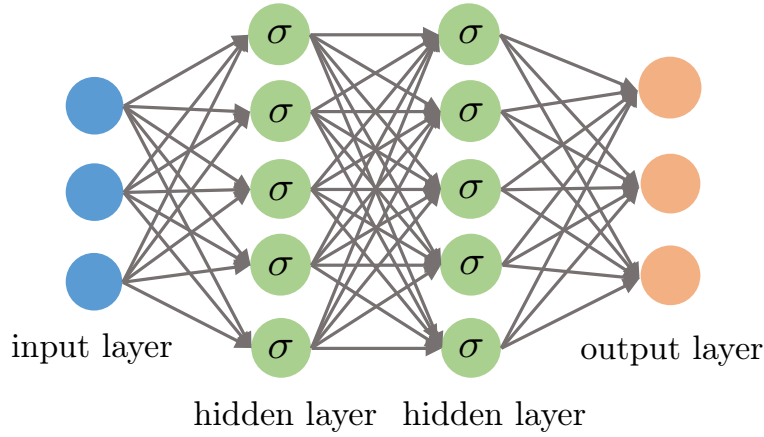


Fig 2: A feed-forward neural network with an input layer, two hidden layers and an output layer. The input layer represents raw features  $\mathbf{x}_i$ . Both hidden layers compute an affine transform (a.k.a. indices) of the input and then apply an element-wise activation function  $\sigma(\cdot)$ . Finally, the output returns a linear transform followed by the softmax activation (resp. simply a linear transform) of the hidden layers for the classification (resp. regression) problem.

neural nets with a specific choice of  $\mathbf{g}^{(\ell)}$ : for  $\ell = 1, \dots, L$ , define

$$(2.1) \quad \mathbf{h}^{(\ell)} = \mathbf{g}^{(\ell)}(\mathbf{h}^{(\ell-1)}) \triangleq \sigma(\mathbf{W}^{(\ell)}\mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)}),$$

where  $\mathbf{W}^{(\ell)}$  and  $\mathbf{b}^{(\ell)}$  are the weight matrix and the bias / intercept, respectively, associated with the  $\ell$ -th layer, and  $\sigma(\cdot)$  is usually a simple fixed nonlinear function called the *activation function*. In statistical terms, we create affine transformations (or multiple indices) using matrix  $\mathbf{W}^{(\ell)}$  with intercept  $\mathbf{b}^{(\ell)}$ , based on the previous layer of input vector  $\mathbf{h}^{(\ell-1)}$ . In words, in each layer  $\ell$ , the input vector  $\mathbf{h}^{(\ell-1)}$  goes through an affine transformation first and then passes through a fixed nonlinear function  $\sigma(\cdot)$ . See Figure 2 for an illustration of a simple feed-forward neural network with two hidden layers. The activation function  $\sigma(\cdot)$  is usually applied element-wise, and a popular choice is the ReLU (Rectified Linear Unit) function:

$$(2.2) \quad [\sigma(\mathbf{z})]_j = \max\{z_j, 0\},$$

which is the positive part of  $z_j$ . Other choices of activation functions include leaky ReLU (Maas, Hannun and Ng, 2013), tanh function and the classical sigmoid function  $(1 + e^{-z})^{-1}$ , which are less used nowadays, due to the computation and the convergence/divergence of gradients, as to be explained later.

Given an output  $\mathbf{h}^{(L)}$  from the final hidden layer and a label  $y$ , we can define a loss function to minimize. A common loss function for classification problems is the multinomial logistic loss. Using the terminology of deep learning, we say that  $\mathbf{h}^{(L)}$  goes through an affine transformation and then the *soft-max* function:

$$f_k(\mathbf{x}; \boldsymbol{\theta}) \triangleq \frac{\exp(z_k)}{\sum_k \exp(z_k)}, \quad \forall k \in [K], \quad \text{where } \mathbf{z} = \mathbf{W}^{(L+1)}\mathbf{h}^{(L)} + \mathbf{b}^{(L+1)} \in \mathbb{R}^K.$$

Then the loss is defined to be the cross-entropy between the label  $y$  (in the form of an indicator vector) and the score vector  $(f_1(\mathbf{x}; \boldsymbol{\theta}), \dots, f_K(\mathbf{x}; \boldsymbol{\theta}))^\top$ , which is

exactly the negative log-likelihood of the multinomial logistic regression model:

$$(2.3) \quad \mathcal{L}(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = - \sum_{k=1}^K \mathbb{1}\{y = k\} \log f_k(\mathbf{x}; \boldsymbol{\theta}),$$

where  $\boldsymbol{\theta} \triangleq \{\mathbf{W}^{(\ell)}, \mathbf{b}^{(\ell)} : 1 \leq \ell \leq L+1\}$ . We emphasize that all the hidden nodes  $\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L)}$  depend on the parameters  $\boldsymbol{\theta}$  through (2.1). As a final remark, the number of parameters scales with both the depth  $L$  and the width (i.e., the dimensionality of  $\mathbf{W}^{(\ell)}$ ), and hence it can be quite large for deep neural nets.

## 2.2 Back-propagation in computational graphs

Training neural networks follows the *empirical risk minimization* paradigm that minimizes the loss (e.g., (2.3)) over all the training data. This minimization is usually done via *stochastic gradient descent* (SGD). In a way similar to gradient descent, SGD starts from a certain initial value  $\boldsymbol{\theta}^0$  and then iteratively updates the parameters  $\boldsymbol{\theta}^t$  by moving it in the direction of the negative gradient. The difference is that, in each update, a small subsample  $\mathcal{B} \subset [n]$  called a *mini-batch*—which is typically of size 32–512—is randomly drawn and the gradient is calculated only on  $\mathcal{B}$  instead of the full batch  $[n]$ . This considerably saves the computational cost in the calculation of gradients. By the law of large numbers, this stochastic gradient should be close to the full sample one, albeit with some random fluctuations. A pass of the whole training set is called an *epoch*. Usually, after several or tens of epochs, the error on a validation set levels off and training is complete. See Section 6 for more details and variants on training algorithms.

The key to the above training procedure, namely SGD, is the calculation of the gradient  $\nabla \ell_{\mathcal{B}}(\boldsymbol{\theta})$ , where

$$(2.4) \quad \ell_{\mathcal{B}}(\boldsymbol{\theta}) \triangleq |\mathcal{B}|^{-1} \sum_{i \in \mathcal{B}} \mathcal{L}(\mathbf{f}(\mathbf{x}_i; \boldsymbol{\theta}), y_i).$$

Gradient computation, however, is in general nontrivial for complex models, and it is susceptible to numerical instability for a model with large depth. Here, we introduce an efficient approach, namely *back-propagation*, for computing gradients in neural networks.

Back-propagation (Rumelhart, Hinton and Williams, 1985) is a direct application of the chain rule in networks. As the name suggests, the calculation is performed in a backward fashion: one first computes  $\partial \ell_{\mathcal{B}} / \partial \mathbf{h}^{(L)}$ , then  $\partial \ell_{\mathcal{B}} / \partial \mathbf{h}^{(L-1)}$ ,  $\dots$ , and finally  $\partial \ell_{\mathcal{B}} / \partial \mathbf{h}^{(1)}$ . For example, in the case of the ReLU activation function<sup>7</sup>, we have the following recursive / backward relation

$$(2.5) \quad \frac{\partial \ell_{\mathcal{B}}}{\partial \mathbf{h}^{(\ell-1)}} = \frac{\partial \mathbf{h}^{(\ell)}}{\partial \mathbf{h}^{(\ell-1)}} \cdot \frac{\partial \ell_{\mathcal{B}}}{\partial \mathbf{h}^{(\ell)}} = (\mathbf{W}^{(\ell)})^{\top} \text{diag} \left( \mathbb{1}\{\mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)} \geq \mathbf{0}\} \right) \frac{\partial \ell_{\mathcal{B}}}{\partial \mathbf{h}^{(\ell)}}$$

where  $\text{diag}(\cdot)$  denotes a diagonal matrix with elements given by the argument. Note that the calculation of  $\partial \ell_{\mathcal{B}} / \partial \mathbf{h}^{(\ell-1)}$  depends on  $\partial \ell_{\mathcal{B}} / \partial \mathbf{h}^{(\ell)}$ , which is the partial derivatives from the next layer. In this way, the derivatives are “back-propagated” from the last layer to the first layer. These derivatives  $\{\partial \ell_{\mathcal{B}} / \partial \mathbf{h}^{(\ell)}\}$  are then used to update the parameters. For instance, the gradient update for

<sup>7</sup>The issue of non-differentiability at the origin is often ignored in implementation.



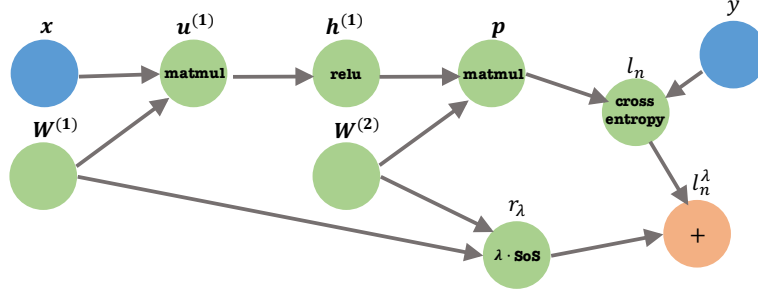


Fig 3: The computational graph illustrates the loss (2.7). For simplicity, we omit the bias terms. Symbols inside nodes represent functions, and symbols outside nodes represent function outputs (vectors/scalars). **matmul** is matrix multiplication, **relu** is the ReLU activation, **cross entropy** is the cross entropy loss, and **SoS** is the sum of squares.

$\mathbf{W}^{(\ell)}$  is given by

$$(2.6) \quad \mathbf{W}^{(\ell)} \leftarrow \mathbf{W}^{(\ell)} - \eta \frac{\partial \ell_{\mathcal{B}}}{\partial \mathbf{W}^{(\ell)}}, \quad \text{where} \quad \frac{\partial \ell_{\mathcal{B}}}{\partial W_{jm}^{(\ell)}} = \frac{\partial \ell_{\mathcal{B}}}{\partial h_j^{(\ell)}} \cdot \sigma' \cdot h_m^{(\ell-1)},$$

where  $\sigma' = 1$  if the  $j$ -th element of  $\mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)}$  is nonnegative, and  $\sigma' = 0$  otherwise. The step size  $\eta > 0$ , also called the *learning rate*, controls how much parameters are changed in a single update.

A more general way to think about neural network models and training is to consider *computational graphs*. Computational graphs are directed acyclic graphs that represent functional relations between variables. They are very convenient and flexible to represent function composition, and moreover, they also allow an efficient way of computing gradients. Consider a feed-forward neural network with a single hidden layer and an  $\ell_2$  regularization:

$$(2.7) \quad \ell_{\mathcal{B}}^{\lambda}(\boldsymbol{\theta}) = \ell_{\mathcal{B}}(\boldsymbol{\theta}) + r_{\lambda}(\boldsymbol{\theta}) = \ell_{\mathcal{B}}(\boldsymbol{\theta}) + \lambda \left( \sum_{j,j'} (W_{j,j'}^{(1)})^2 + \sum_{j,j'} (W_{j,j'}^{(2)})^2 \right),$$

where  $\ell_{\mathcal{B}}(\boldsymbol{\theta})$  is the same as (2.4), and  $\lambda \geq 0$  is a tuning parameter. A similar example is considered in Goodfellow, Bengio and Courville (2016). The corresponding computational graph is shown in Figure 3. Each node represents a function (inside a circle), which is associated with an output of that function (outside a circle). For example, we view the term  $\ell_{\mathcal{B}}(\boldsymbol{\theta})$  as a result of four compositions: first the input data  $\mathbf{x}$  multiplies the weight matrix  $\mathbf{W}^{(1)}$  resulting in  $\mathbf{u}^{(1)}$ , then it goes through the ReLU activation function **relu** resulting in  $\mathbf{h}^{(1)}$ , then it multiplies another weight matrix  $\mathbf{W}^{(2)}$  leading to  $\mathbf{p}$ , and finally it produces the cross-entropy with label  $y$  as in (2.3). The regularization term is incorporated in the graph similarly.

A forward pass is complete when all nodes are evaluated starting from the input  $\mathbf{x}$ . A backward pass then calculates the gradients of  $\ell_{\mathcal{B}}^{\lambda}$  with respect to all other nodes in the reverse direction. Due to the chain rule, the gradient calculation for a variable (say,  $\partial \ell_{\mathcal{B}} / \partial \mathbf{u}^{(1)}$ ) is simple: it only depends on the gradient value of the variables ( $\partial \ell_{\mathcal{B}} / \partial \mathbf{h}$ ) the current node points to, and the function derivative evaluated at the current variable value ( $\sigma'(\mathbf{u}^{(1)})$ ). Thus, in each iteration, a computational graph only needs to (1) calculate and store the function evaluations

at each node in the forward pass, and then (2) calculate all the derivatives in the backward pass.

Back-propagation in computational graphs forms the foundations of popular deep learning programming software, including TensorFlow (Abadi et al., 2015) and PyTorch (Paszke et al., 2017), which allows more efficient building and training of complex neural net models.

### 3. POPULAR MODELS

Moving beyond vanilla feed-forward neural networks, we introduce two other popular deep learning models, namely, the convolutional neural networks (CNNs) and the recurrent neural networks (RNNs). One important characteristic shared by the two models is *weight sharing*, that is, some model parameters are identical across locations in CNNs or across time in RNNs. This is related to the notion of translational invariance in CNNs and stationarity in RNNs. At the end of this section, we introduce a modular thinking for constructing more flexible neural nets.

#### 3.1 Convolutional neural networks

The convolutional neural network (CNN) (LeCun et al., 1998; Fukushima and Miyake, 1982) is a special type of feed-forward neural networks that is tailored for image processing. More generally, it is suitable for analyzing data with salient spatial structures. In this subsection, we focus on image classification using CNNs, where the raw input (i.e., image pixels) and features of each hidden layer are represented by a 3D tensor  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$ . Here, the first two dimensions  $d_1, d_2$  of  $\mathbf{X}$  indicate spatial coordinates of an image while the third  $d_3$  indicates the number of channels. For instance,  $d_3$  is 3 for the raw inputs corresponding to the red, green and blue channels, and  $d_3$  can be much larger (say, 256) for hidden layers. Each channel is also called a *feature map*, because each feature map is specialized to detect the same feature at different locations of the input, which we will soon explain.

Before giving the formal introduction, let us think intuitively what network structures best fit image data. First, objects in images are related to their locations (spatial coordinates), so it would be better if the weights  $\mathbf{W}$  are only receptive to different local patches. This dramatically reduces the number of parameters compared with the feed-forward networks. Second, we expect invariance of objects when translating them, or flipping them in images, so naturally, some weights should share the same values. Third, we hope the network could summarize the image information as we move to higher layers.

These desired properties lead to the following building blocks of CNNs, namely the convolutional layer and the pooling layer.

1. *Convolutional layer (CONV)*. A convolutional layer has the same functionality as described in (2.1), where the input feature  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  goes through an affine transformation first and then an element-wise nonlinear activation. The difference lies in the specific form of the affine transformation. A convolutional layer uses a number of *filters* to extract local features from the previous input. More precisely, each filter is represented by a 3D tensor  $\mathbf{F}_k \in \mathbb{R}^{w \times w \times d_3}$  ( $1 \leq k \leq \tilde{d}_3$ ), where  $w$  is the size of the filter (typically 3 or 5) and  $\tilde{d}_3$  denotes the total number of filters. Note that the third dimension  $d_3$  of  $\mathbf{F}_k$  is equal to that of the

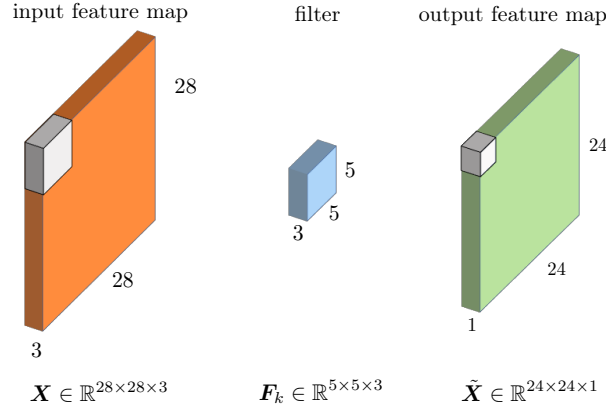


Fig 4:  $\mathbf{X} \in \mathbb{R}^{28 \times 28 \times 3}$  represents the input feature consisting of  $28 \times 28$  spatial coordinates in a total number of 3 channels / feature maps.  $\mathbf{F}_k \in \mathbb{R}^{5 \times 5 \times 3}$  denotes the  $k$ -th filter with size  $5 \times 5$ . The third dimension 3 of the filter automatically matches the number 3 of channels in the previous input. Every 3D patch of  $\mathbf{X}$  gets convolved with the filter  $\mathbf{F}_k$  and this as a whole results in a single output feature map  $\tilde{\mathbf{X}}_{:, :, k}$  with size  $24 \times 24 \times 1$ . Stacking the outputs of all the filters  $\{\mathbf{F}_k\}_{1 \leq k \leq K}$  will lead to the output feature with size  $24 \times 24 \times K$ .

input feature  $\mathbf{X}$ . For this reason, one usually says that the filter has size  $w \times w$ , while suppressing the third dimension  $d_3$ . Each filter  $\mathbf{F}_k$  then convolves with the input feature  $\mathbf{X}$  to obtain one single feature map  $\mathbf{O}^k \in \mathbb{R}^{(d_1-w+1) \times (d_1-w+1)}$ , where<sup>8</sup>

$$(3.1) \quad O_{ij}^k = \langle [\mathbf{X}]_{ij}, \mathbf{F}_k \rangle = \sum_{i'=1}^w \sum_{j'=1}^w \sum_{l=1}^{d_3} [\mathbf{X}]_{i+i'-1, j+j'-1, l} [\mathbf{F}_k]_{i', j', l}.$$

Here  $[\mathbf{X}]_{ij} \in \mathbb{R}^{w \times w \times d_3}$  is a small “patch” of  $\mathbf{X}$  starting at location  $(i, j)$ . See Figure 4 for an illustration of the convolution operation. If we view the 3D tensors  $[\mathbf{X}]_{ij}$  and  $\mathbf{F}_k$  as vectors, then each filter essentially computes their inner product with a part of  $\mathbf{X}$  indexed by  $i, j$  (which can be also viewed as convolution, as its name suggests). One then packs the resulted feature maps  $\{\mathbf{O}^k\}$  into a 3D tensor  $\mathbf{O}$  with size  $(d_1 - w + 1) \times (d_1 - w + 1) \times \tilde{d}_3$ , where

$$(3.2) \quad [\mathbf{O}]_{ijk} = [\mathbf{O}^k]_{ij}.$$

The outputs of convolutional layers are then followed by nonlinear activation functions. One commonly used activation function is the ReLU function, in which case we have

$$(3.3) \quad \tilde{X}_{ijk} = \sigma(O_{ijk}), \quad \forall i \in [d_1 - w + 1], j \in [d_2 - w + 1], k \in [\tilde{d}_3].$$

The convolution operation (3.1) and the ReLU activation (3.3) work together to extract features  $\tilde{\mathbf{X}}$  from the input  $\mathbf{X}$  (a positive value represent presence of a feature). Different from feed-forward neural nets, the filters  $\mathbf{F}_k$  are shared across all locations  $(i, j)$ . A patch  $[\mathbf{X}]_{ij}$  of an input responds strongly (that is, producing a large value) to a filter  $\mathbf{F}_k$  if they are positively correlated. Therefore

<sup>8</sup>To simplify notation, we omit the bias/intercept term associated with each filter.



Fig 5: A  $2 \times 2$  max pooling layer extracts the maximum of 2 by 2 neighboring pixels / features across the spatial dimension.

intuitively, each filter  $\mathbf{F}_k$  serves to extract features similar to  $\mathbf{F}_k$ . As a side note, after the convolution (3.1), the spatial size  $d_1 \times d_2$  of the input  $\mathbf{X}$  shrinks to  $(d_1 - w + 1) \times (d_2 - w + 1)$  of  $\tilde{\mathbf{X}}$  due to the edge effect. However, one may want the spatial size unchanged. This can be achieved via *padding*, where one appends zeros to the margins of the input  $\mathbf{X}$  to enlarge the spatial size to  $d_1 \times d_2$ . In addition, a *stride* in the convolutional layer determines the gap  $i' - i$  and  $j' - j$  between two patches  $\mathbf{X}_{ij}$  and  $\mathbf{X}_{i'j'}$ : in (3.1) the stride is 1, and a larger stride would lead to feature maps with smaller sizes.

2. *Pooling layer (POOL)*. A pooling layer aggregates the information of nearby features into a single one. This downsampling operation reduces the size of the features for subsequent layers and saves computation. One common form of the pooling layer is composed of the  $2 \times 2$  max-pooling filter. It computes  $\max\{X_{i,j,k}, X_{i+1,j,k}, X_{i,j+1,k}, X_{i+1,j+1,k}\}$ , that is, the maximum of the  $2 \times 2$  neighborhood in the spatial coordinates; see Figure 5 for an illustration. Note that the pooling operation is done separately for each feature map  $k$ . As a consequence, a  $2 \times 2$  max-pooling filter acting on  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2 \times d_3}$  will result in an output of size  $d_1/2 \times d_2/2 \times d_3$ . In addition, the pooling layer does not involve any parameters to optimize. Pooling layers serve to reduce redundancy since a small neighborhood around a location  $(i, j)$  in a feature map is likely to contain similar information.

In addition, we also use fully-connected layers as building blocks, which we have already introduced in Section 2. Each fully-connected layer treats input tensor  $\mathbf{X}$  as a vector  $\text{Vec}(\mathbf{X})$ , and computes  $\tilde{\mathbf{X}} = \sigma(\mathbf{W}\text{Vec}(\mathbf{X}))$ . A fully-connected layer

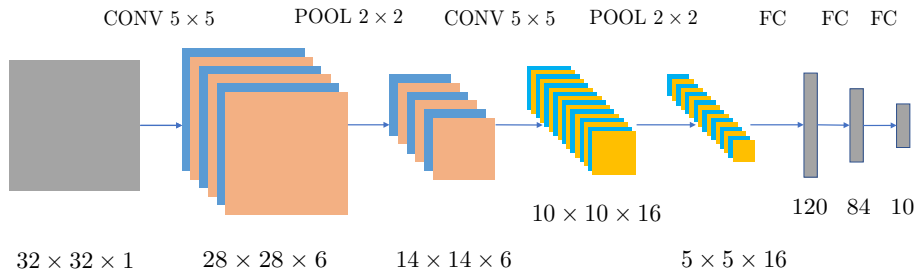


Fig 6: LeNet is composed of an input layer, two convolutional layers, two pooling layers and three fully-connected layers. Both convolutions use  $5 \times 5$  filters. In addition, the two pooling layers use  $2 \times 2$  average pooling.

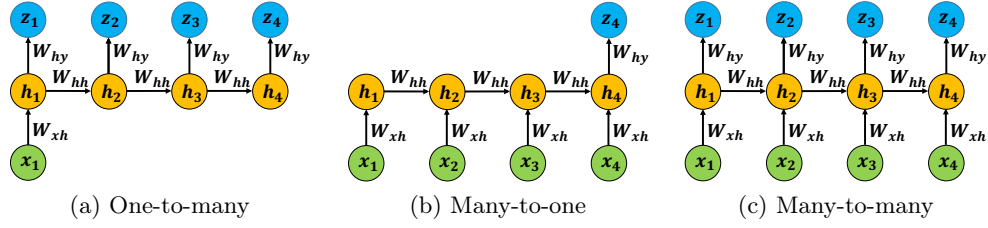


Fig 7: Vanilla RNNs with different inputs/outputs settings. (a) has one input but multiple outputs; (b) has multiple inputs but one output; (c) has multiple inputs and outputs. Note that the parameters are shared across time steps.

does not use weight sharing and is often used in the last few layers of a CNN. As an example, Figure 6 depicts the well-known LeNet 5 (LeCun et al., 1998), which is composed of two sets of CONV-POOL layers and three fully-connected layers.

### 3.2 Recurrent neural networks

Recurrent neural nets (RNNs) are another family of powerful models, which are designed to process time series data and other sequence data. RNNs have successful applications in speech recognition (Sak, Senior and Beaufays, 2014), machine translation (Wu et al., 2016), genome sequencing (Cao et al., 2018), etc. The structure of an RNN naturally forms a computational graph and can be easily combined with other structures such as CNNs to build larger computational graph models for complex tasks. Here we introduce vanilla RNNs and improved variants such as long short-term memory (LSTM).

*3.2.1 Vanilla RNNs.* Suppose we have vector-valued time series inputs  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ . A vanilla RNN models the “hidden state” at time  $t$  by a vector  $\mathbf{h}_t$ , which is subject to the recursive formula

$$(3.4) \quad \mathbf{h}_t = \mathbf{f}_{\boldsymbol{\theta}}(\mathbf{h}_{t-1}, \mathbf{x}_t).$$

Here,  $\mathbf{f}_{\boldsymbol{\theta}}$  is generally a nonlinear function parametrized by  $\boldsymbol{\theta}$ . Concretely, a vanilla RNN with one hidden layer has the following form<sup>9</sup>

$$\begin{aligned} \mathbf{h}_t &= \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t + \mathbf{b}_h), & \text{where } \tanh(a) &= \frac{e^{2a}-1}{e^{2a}+1}, \\ \mathbf{z}_t &= \sigma(\mathbf{W}_{hy}\mathbf{h}_t + \mathbf{b}_z), \end{aligned}$$

where  $\mathbf{W}_{hh}, \mathbf{W}_{xh}, \mathbf{W}_{hy}$  are trainable weight matrices,  $\mathbf{b}_h, \mathbf{b}_z$  are trainable bias vectors, and  $\mathbf{z}_t$  is the output at time  $t$ . Like many classical time series models, those parameters are shared across time. Note that in different applications, we may have different input/output settings (cf. Figure 7). Examples include

- **One-to-many:** a single input with multiple outputs; see Figure 7(a). A typical application is image captioning, where the input is an image and outputs are a series of words.
- **Many-to-one:** multiple inputs with a single output; see Figure 7(b). One application is text sentiment classification, where the input is a series of words in a sentence and the output is a label (e.g., positive vs. negative).

<sup>9</sup>Similar to the activation function  $\sigma(\cdot)$ , the function  $\tanh(\cdot)$  means element-wise operations.

- **Many-to-many:** multiple inputs and outputs; see Figure 7(c). This is adopted in machine translation, where inputs are words of a source language (say Chinese) and outputs are words of a target language (say English).

As the case with feed-forward neural nets, we minimize a loss function using back-propagation, where the loss is typically

$$\ell_{\mathcal{T}}(\boldsymbol{\theta}) = \sum_{t \in \mathcal{T}} \mathcal{L}(y_t, \mathbf{z}_t) = - \sum_{t \in \mathcal{T}} \sum_{k=1}^K \mathbb{1}\{y_t = k\} \log \left( \frac{\exp([\mathbf{z}_t]_k)}{\sum_k \exp([\mathbf{z}_t]_k)} \right),$$

where  $K$  is the number of categories for classification (e.g., size of the vocabulary in machine translation), and  $\mathcal{T} \subset [T]$  is the length of the output sequence. During the training, the gradients  $\partial \ell_{\mathcal{T}} / \partial \mathbf{h}_t$  are computed in the reverse time order (from  $T$  to  $t$ ). For this reason, the training process is often called *back-propagation through time*.

One notable drawback of vanilla RNNs is that they have difficulty in capturing long-range dependencies in sequence data when the length of the sequence is large. This is sometimes due to the phenomenon of *exploding/vanishing gradients*. Take Figure 7(c) as an example. Computing  $\partial \ell_{\mathcal{T}} / \partial \mathbf{h}_1$  involves the product  $\prod_{t=1}^3 (\partial \mathbf{h}_{t+1} / \partial \mathbf{h}_t)$  by the chain rule. However, if the sequence is long, the product will be the multiplication of many Jacobian matrices, which usually results in exponentially large or small singular values. To alleviate this issue, in practice, the forward pass and backward pass are implemented in a shorter sliding window  $\{t_1, t_1 + 1, \dots, t_2\}$ , instead of the full sequence  $\{1, 2, \dots, T\}$ . Though effective in some cases, this technique alone does not fully address the issue of long-term dependency.

**3.2.2 GRUs and LSTM.** There are two improved variants that alleviate the above issue: gated recurrent units (GRUs) (Cho et al., 2014) and long short-term memory (LSTM) (Hochreiter and Schmidhuber, 1997). The motivating intuition behind these specialized structures is the characteristics of natural languages: at a give point of the sequence, a good model should maintain some ‘states’ about past inputs, e.g., the tense of a sentence, the clause structure, semantics, etc. These states can be changed, turned ‘on’ or ‘off’ with newly arriving inputs; this idea leads to the use of ‘gates’.

- A **GRU** refines the recursive formula (3.4) by introducing *gates*, which are vectors of the same length as  $\mathbf{h}_t$ . The gates, which take values in  $[0, 1]$  elementwise, multiply with  $\mathbf{h}_{t-1}$  elementwise and determine how much they keep the old hidden states.
- An **LSTM** similarly uses gates in the recursive formula. In addition to  $\mathbf{h}_t$ , an LSTM maintains a *cell state*, which takes values in  $\mathbb{R}$  elementwise and are analogous to counters.

Here we only discuss LSTM in detail. Denote by  $\odot$  the element-wise multiplication. We have a recursive formula in replace of (3.4):

$$\begin{pmatrix} \mathbf{i}_t \\ \mathbf{f}_t \\ \mathbf{o}_t \\ \mathbf{g}_t \end{pmatrix} = \begin{pmatrix} \boldsymbol{\sigma} \\ \boldsymbol{\sigma} \\ \boldsymbol{\sigma} \\ \tanh \end{pmatrix} \mathbf{W} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \\ 1 \end{pmatrix},$$



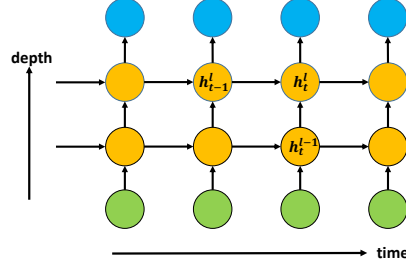


Fig 8: A vanilla RNN with two hidden layers. Higher-level hidden states  $\mathbf{h}_t^\ell$  are determined by the old states  $\mathbf{h}_{t-1}^\ell$  and lower-level hidden states  $\mathbf{h}_t^{\ell-1}$ . Multilayer RNNs generalize both feed-forward neural nets and one-hidden-layer RNNs.

$$\begin{aligned} \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t, \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t), \end{aligned}$$

where  $\mathbf{W}$  is a big weight matrix with appropriate dimensions. The cell state vector  $\mathbf{c}_t$  carries information of the sequence (e.g., singular/plural form in a sentence). The forget gate  $\mathbf{f}_t$  determines how much the values of  $\mathbf{c}_{t-1}$  are kept for time  $t$ , the input gate  $\mathbf{i}_t$  controls the amount of update to the cell state, and the output gate  $\mathbf{o}_t$  gives how much  $\mathbf{c}_t$  reveals to  $\mathbf{h}_t$ . Ideally, the elements of these gates have nearly binary values. For example, an element of  $\mathbf{f}_t$  being close to 1 may suggest the presence of a feature in the sequence data. Similar to the skip connections in residual nets, the cell state  $\mathbf{c}_t$  has an additive recursive formula, which helps back-propagation and thus captures long-range dependencies.

**3.2.3 Multilayer RNNs.** Multilayer RNNs are generalization of the one-hidden-layer RNN discussed above. Figure 8 shows a vanilla RNN with two hidden layers. In place of (3.4), the recursive formula for an RNN with  $L$  hidden layers now reads

$$\mathbf{h}_t^\ell = \tanh \left[ \mathbf{W}^\ell \begin{pmatrix} \mathbf{h}_t^{\ell-1} \\ \mathbf{h}_{t-1}^\ell \\ 1 \end{pmatrix} \right], \quad \text{for all } \ell \in [L], \quad \mathbf{h}_t^0 \triangleq \mathbf{x}_t.$$

Note that a multilayer RNN has two dimensions: the sequence length  $T$  and depth  $L$ . Two special cases are the feed-forward neural nets (where  $T = 1$ ) introduced in Section 2, and RNNs with one hidden layer (where  $L = 1$ ). Multilayer RNNs usually do not have very large depth (e.g., 2–5), since  $T$  is already very large.

Finally, we remark that CNNs, RNNs, and other neural nets can be easily combined to tackle tasks that involve different sources of input data. For example, in image captioning, the images are first processed through a CNN, and then the high-level features are fed into an RNN as inputs. These neural nets combined together form a large computational graph, so they can be trained using back-propagation. This generic training method provides much flexibility in various applications.

### 3.3 Modules

Deep neural nets are essentially the composition of many nonlinear functions. A component function may be designed to have specific properties in a given task, and it can be itself resulted from composing a few simpler functions. In LSTM,

we have seen that the building block consists of several intermediate variables, including cell states and forget gates that can capture long-term dependency and alleviate numerical issues.

This leads to the idea of designing *modules* for building more complex neural net models. Desirable modules usually have low computational costs, alleviate numerical issues in training, and lead to good statistical accuracy. Since modules and the resulting neural net models form computational graphs, training follows the same principle briefly described in Section 2.

Here, we use the examples of *Inception* and *skip connections* to illustrate the ideas behind modules. Figure 9(a) is an example of “Inception” modules used in GoogleNet (Szegedy et al., 2015). As before, all the convolutional layers are followed by the ReLU activation function. The concatenation of information from filters with different sizes gives the model great flexibility to capture spatial information. Note that  $1 \times 1$  filters is a  $1 \times 1 \times d_3$  tensor (where  $d_3$  is the number of feature maps), so its convolutional operation does not interact with other spatial coordinates, only serving to aggregate information from different feature maps at the same coordinate. This reduces the number of parameters and speeds up the computation. Similar ideas appear in other work (Lin, Chen and Yan, 2013; Iandola et al., 2016).

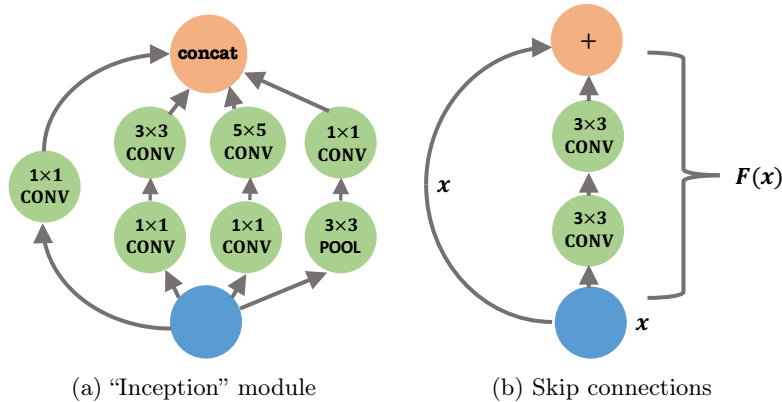


Fig 9: (a) The “Inception” module from GoogleNet. **Concat** means combining all features maps into a tensor. (b) Skip connections are added every two layers in ResNets.

Another module, usually called *skip connections*, is widely used to alleviate numerical issues in very deep neural nets, with additional benefits in optimization efficiency and statistical accuracy. Training very deep neural nets are generally more difficult, but the introduction of skip connections in *residual networks* (He et al., 2016a,b) has greatly eased the task.

The high level idea of skip connections is to add an identity map to an existing nonlinear function. Let  $\mathbf{F}(x)$  be an arbitrary nonlinear function represented by a (fragment of) neural net, then the idea of skip connections is simply replacing  $\mathbf{F}(x)$  with  $x + \mathbf{F}(x)$ . The motivation can be explained from the angle of numerical stability: if the magnitude of  $\mathbf{F}(x)$  is small, then the spectra (singular values) of Jacobian of  $x + \mathbf{F}(x)$  is close to 1, so we expect better numerical stability. Figure 9(b) shows a well-known structure from residual networks (He et al.,

2016a)—for every two layers, an identity map is added:

$$(3.5) \quad \mathbf{x} \mapsto \sigma(\mathbf{x} + \mathbf{F}(\mathbf{x})) = \sigma(\mathbf{x} + \mathbf{W}'\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) + \mathbf{b}'),$$

where  $\mathbf{x}$  can be hidden nodes from any layer and  $\mathbf{W}, \mathbf{W}', \mathbf{b}, \mathbf{b}'$  are corresponding parameters. By repeating (namely composing) this structure throughout all layers, He et al. (2016a,b) are able to train neural nets with hundreds of layers easily, which overcome well-observed training difficulties in deep neural nets. Moreover, deep residual networks also improve statistical accuracy, as the relative classification error on ImageNet challenge was reduced by 46% from 2014 to 2015. As a side note, skip connections can be used flexibly. They are not restricted to the form in (3.5), and can be used between any pair of layers  $\ell, \ell'$  (Huang et al., 2017).

## 4. DEEP UNSUPERVISED LEARNING

In supervised learning, given labeled training set  $\{(y_i, \mathbf{x}_i)\}$ , we focus on discriminative models, which essentially represents  $\mathbb{P}(y|\mathbf{x})$  by a deep neural net  $f(\mathbf{x}; \boldsymbol{\theta})$  with parameters  $\boldsymbol{\theta}$ . Unsupervised learning, in contrast, aims at extracting *information* from *unlabeled* data  $\{\mathbf{x}_i\}$ , where the labels  $\{y_i\}$  are absent. In regard to this information, it can be a low-dimensional embedding of the data  $\{\mathbf{x}_i\}$  or a generative model with latent variables to approximate the distribution  $\mathbb{P}_{\mathbf{X}}(\mathbf{x})$ . To achieve these goals, we introduce two popular unsupervised deep learning models, namely, autoencoders and generative adversarial networks (GANs). The first one can be viewed as a dimension reduction technique and the second as a high-dimensional density estimation. DNNs are the key elements for both of these two models.

### 4.1 Autoencoders

Recall that in dimension reduction, the goal is to reduce the dimensionality of the data and at the same time preserve its salient features. In particular, in principal component analysis (PCA), the goal is to embed the data  $\{\mathbf{x}_i\}_{1 \leq i \leq n}$  into a low-dimensional space via a linear function  $\mathbf{f}$  such that maximum variance can be explained. Equivalently, we want to find linear functions  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^k$  and  $\mathbf{g} : \mathbb{R}^k \rightarrow \mathbb{R}^d$  ( $k \leq d$ ) such that the difference between  $\mathbf{x}_i$  and  $\mathbf{g}(\mathbf{f}(\mathbf{x}_i))$  is minimized. Formally, we let

$$\mathbf{f}(\mathbf{x}) = \mathbf{W}_f \mathbf{x} \triangleq \mathbf{h} \quad \text{and} \quad \mathbf{g}(\mathbf{h}) = \mathbf{W}_g \mathbf{h}, \quad \text{where} \quad \mathbf{W}_f \in \mathbb{R}^{k \times d} \text{ and } \mathbf{W}_g \in \mathbb{R}^{d \times k}.$$

Here, for simplicity, we assume that the intercept/bias terms for  $\mathbf{f}$  and  $\mathbf{g}$  are zero. Then, PCA amounts to minimizing the quadratic loss function

$$(4.1) \quad \text{minimize}_{\mathbf{W}_f, \mathbf{W}_g} \quad \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{W}_f \mathbf{W}_g \mathbf{x}_i\|_2^2.$$

If we write  $\mathbf{W} = \mathbf{W}_f \mathbf{W}_g$ , then the above problem is the same as minimizing  $\|\mathbf{X} - \mathbf{W}\mathbf{X}\|_F^2$  subject to  $\text{rank}(\mathbf{W}) \leq k$ , where  $\mathbf{X} \in \mathbb{R}^{p \times n}$  is the design matrix. The solution is given by the singular value decomposition of  $\mathbf{X}$  (Golub and Van Loan, 2013, Thm. 2.4.8), which is exactly what PCA does. It turns out that PCA is a special case of autoencoders, which is often known as the *undercomplete linear autoencoder*.

More broadly, autoencoders are neural network models for (nonlinear) dimension reduction, which generalize PCA. An autoencoder has two key components,

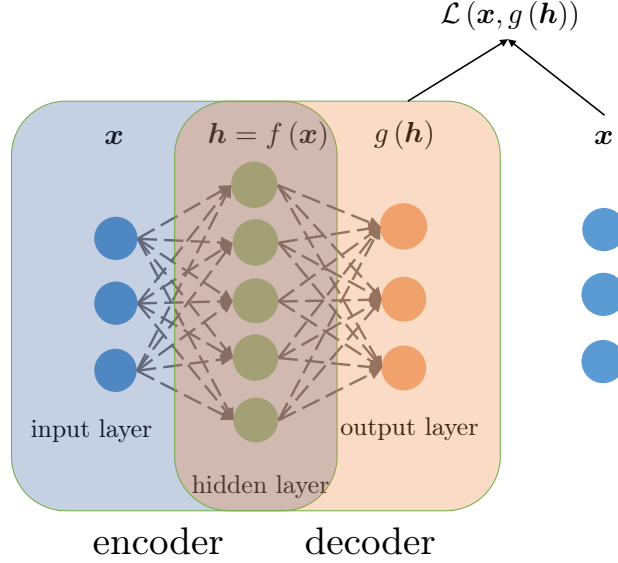


Fig 10: A diagram of a (sparse) autoencoder. First an input  $x$  goes through the encoder  $f(\cdot)$ , and we obtain its hidden representation  $h = f(x)$ . Then, we use the decoder  $g(\cdot)$  to get  $g(h)$  as a reconstruction of  $x$ . Finally, the loss is determined from the difference between the original input  $x$  and its reconstruction  $g(f(x))$ . A regularizer is used if the dimension of  $h$  is larger than that of  $x$ .

namely, the encoder function  $f(\cdot)$ , which maps the input  $x \in \mathbb{R}^d$  to a hidden code/representation  $h \triangleq f(x) \in \mathbb{R}^k$ , and the decoder function  $g(\cdot)$ , which maps the hidden representation  $h$  to a point  $g(h) \in \mathbb{R}^d$ . Both functions can be multi-layer neural networks as (2.1). See Figure 10 for an illustration of autoencoders. Let  $\mathcal{L}(x_1, x_2)$  be a loss function that measures the difference between  $x_1$  and  $x_2$  in  $\mathbb{R}^d$ . Similar to PCA, an autoencoder is used to find the encoder  $f$  and decoder  $g$  such that  $\mathcal{L}(x, g(f(x)))$  is as small as possible. Mathematically, this amounts to solving the following minimization problem

$$(4.2) \quad \text{minimize}_{f,g} \quad \frac{1}{n} \sum_{i=1}^n \mathcal{L}(x_i, g(h_i)) \quad \text{with} \quad h_i = f(x_i), \quad \text{for all } i \in [n].$$

One needs to make structural assumptions on the functions  $f$  and  $g$  in order to find useful representations of the data, which leads to different types of autoencoders. One way is to model  $f$  and  $g$  by neural network functions. If no assumption is made, choosing  $f$  and  $g$  to be identity functions clearly minimizes the above optimization problem. To avoid this trivial solution, one natural way is to require that the encoder  $f$  maps the data onto a space with a smaller dimension, i.e.,  $k < d$ . This is the *undercomplete autoencoder* that includes PCA as a special case. There are other structured autoencoders which add desired properties to the model such as sparsity or robustness, mainly through regularization terms. Below we present two other common types of autoencoders.

- *Sparse autoencoders.* One may believe that the dimension  $k$  of the hidden code  $h_i$  is larger than the input dimension  $d$ , and that  $h_i$  admits a sparse representation. As with LASSO (Tibshirani, 1996) or SCAD (Fan and Li, 2001),

one may add a regularization term to the reconstruction loss  $\mathcal{L}$  in (4.2) to encourage sparsity (Poultney et al., 2007). A sparse autoencoder solves

$$\min_{\mathbf{f}, \mathbf{g}} \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, \mathbf{g}(\mathbf{h}_i))}_{\text{loss}} + \underbrace{\lambda \|\mathbf{h}_i\|_1}_{\text{regularizer}} \quad \text{with} \quad \mathbf{h}_i = \mathbf{f}(\mathbf{x}_i), \text{ for all } i \in [n].$$

This is similar to *dictionary learning*, where one aims at finding a sparse representation of input data on an overcomplete basis. Due to the imposed sparsity, the model can potentially learn useful features of the data.

- *Denoising autoencoders.* One may hope that the model is robust to noise in the data: even if the input data  $\mathbf{x}_i$  are corrupted by small noise  $\boldsymbol{\xi}_i$  or miss some components (the noise level or the missing probability is typically small), an ideal autoencoder should faithfully recover the original data. A denoising autoencoder (Vincent et al., 2008) achieves this robustness by explicitly building a noisy data  $\tilde{\mathbf{x}}_i = \mathbf{x}_i + \boldsymbol{\xi}_i$  as the new input, and then solves an optimization problem similar to (4.2) where  $\mathcal{L}(\mathbf{x}_i, \mathbf{g}(\mathbf{h}_i))$  is replaced by  $\mathcal{L}(\mathbf{x}_i, \mathbf{g}(\mathbf{f}(\tilde{\mathbf{x}}_i)))$ . A denoising autoencoder encourages the encoder/decoder to be stable in the neighborhood of an input, which is generally a good statistical property. An alternative way could be constraining  $\mathbf{f}$  and  $\mathbf{g}$  in the optimization problem, but that would be very difficult to optimize. Instead, sampling by adding small perturbations in the input provides a simple implementation. We shall see similar ideas in Section 6.3.3.

## 4.2 Generative adversarial networks

Given unlabeled data  $\{\mathbf{x}_i\}_{1 \leq i \leq n}$ , density estimation aims to estimate the underlying probability density function  $\mathbb{P}_{\mathbf{X}}$  from which the data is generated. Both parametric and nonparametric estimators (Silverman, 1998) have been proposed and studied under various assumptions on the underlying distribution. Different from these classical density estimators, where the density function is explicitly defined in relatively low dimension, generative adversarial networks (GANs) (Goodfellow et al., 2014) can be categorized as an *implicit* density estimator in much higher dimension. The reasons are twofold: (1) GANs put more emphasis on sampling from the distribution  $\mathbb{P}_{\mathbf{X}}$  than estimation; (2) GANs define the density estimation implicitly through a source distribution  $\mathbb{P}_{\mathbf{Z}}$  and a generator function  $g(\cdot)$ , which is usually a deep neural network. We introduce GANs from the perspective of sampling from  $\mathbb{P}_{\mathbf{X}}$  and later we will generalize the vanilla GANs using its relation to density estimators.

*4.2.1 Sampling view of GANs.* Suppose the data  $\{\mathbf{x}_i\}_{1 \leq i \leq n}$  at hand are all real images, and we want to generate *new* natural images. With this goal in mind, GAN models a *zero-sum* game between two players, namely, the generator  $\mathcal{G}$  and the discriminator  $\mathcal{D}$ . The generator  $\mathcal{G}$  tries to generate fake images akin to the true images  $\{\mathbf{x}_i\}_{1 \leq i \leq n}$  while the discriminator  $\mathcal{D}$  aims at differentiating the fake ones from the true ones. Intuitively, one hopes to train a generator  $\mathcal{G}$  to generate images where the *best* discriminator  $\mathcal{D}$  cannot distinguish. Therefore the payoff is higher for the generator  $\mathcal{G}$  if the probability of the discriminator  $\mathcal{D}$  getting wrong is higher, and correspondingly the payoff for the discriminator correlates positively with its ability to tell the truth.

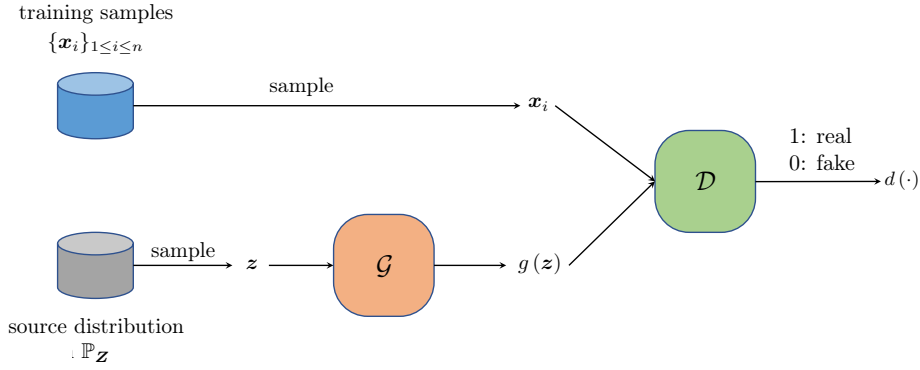


Fig 11: GANs consist of two components, a generator  $\mathcal{G}$  which generates fake samples and a discriminator  $\mathcal{D}$  which differentiate the true ones from the fake ones.

Mathematically, the generator  $\mathcal{G}$  consists of two components, an source distribution  $\mathbb{P}_Z$  (usually a standard multivariate Gaussian distribution with hundreds of dimensions) and a function  $\mathbf{g}(\cdot)$  which maps a sample  $\mathbf{z}$  from  $\mathbb{P}_Z$  to a point  $\mathbf{g}(\mathbf{z})$  living in the same space as  $\mathbf{x}$ . Note that any continuous random vector  $\mathbf{X}$  in  $\mathbb{R}^d$  can be generated by  $\mathbf{X} = \mathbf{g}(\mathbf{Z})$  for the standard Gaussian random vector  $\mathbf{Z}$  in  $\mathbb{R}^d$  and some function  $\mathbf{g}(\cdot)$ . To achieve low-dimensional embedding for a class of images, we take  $\mathbf{Z}$  in much smaller space, often in the order of hundreds. For generating images,  $\mathbf{g}(\mathbf{z})$  would be a 3D tensor. Here  $\mathbf{g}(\mathbf{z})$  is the fake sample generated from  $\mathcal{G}$ . Similarly the discriminator  $\mathcal{D}$  is composed of one function which takes an image  $\mathbf{x}$  (real or fake) and return a number  $d(\mathbf{x}) \in [0, 1]$ , the probability of  $\mathbf{x}$  being a real sample from  $\mathbb{P}_X$  or not. Oftentimes, both the generating function  $\mathbf{g}(\cdot)$  and the discriminating function  $d(\cdot)$  are constructed by deep neural networks, e.g., CNNs introduced in Section 3.1. See Figure 11 for an illustration for GANs. Denote  $\theta_{\mathcal{G}}$  and  $\theta_{\mathcal{D}}$  the parameters in  $\mathbf{g}(\cdot)$  and  $d(\cdot)$ , respectively. Then GAN tries to solve the following *min-max* problem:

$$(4.3) \quad \min_{\theta_{\mathcal{G}}} \max_{\theta_{\mathcal{D}}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_X} [\log(d(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim \mathbb{P}_Z} [\log(1 - d(\mathbf{g}(\mathbf{z})))] .$$

Recall that  $d(\mathbf{x})$  models the belief / probability that the discriminator thinks that  $\mathbf{x}$  is a true sample. Fix the parameters  $\theta_{\mathcal{G}}$  and hence the generator  $\mathcal{G}$  and consider the inner maximization problem. We can see that the goal of the discriminator is to maximize its ability of differentiation. Similarly, if we fix  $\theta_{\mathcal{D}}$  (and hence the discriminator), the generator tries to generate more realistic images  $\mathbf{g}(\mathbf{z})$  to fool the discriminator.

**4.2.2 Density estimation view of GANs.** Let us now take a density-estimation view of GANs. Fixing the source distribution  $\mathbb{P}_Z$ , any generator  $\mathcal{G}$  induces a distribution  $\mathbb{P}_G$  over the space of images. Removing the restrictions on  $d(\cdot)$  and minimizing more generally  $\mathbb{P}_G$  over a class of distributions, one can write (4.3) as

$$(4.4) \quad \min_{\mathbb{P}_G} \max_{d(\cdot)} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_X} [\log(d(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_G} [\log(1 - d(\mathbf{x}))] .$$



Observe that the inner maximization problem is solved by the likelihood ratio, i.e.

$$d^*(\mathbf{x}) = \frac{\mathbb{P}_{\mathbf{X}}(\mathbf{x})}{\mathbb{P}_{\mathbf{X}}(\mathbf{x}) + \mathbb{P}_{\mathcal{G}}(\mathbf{x})}.$$

As a result, (4.4) can be simplified as

$$(4.5) \quad \min_{\mathbb{P}_{\mathcal{G}}} \quad \text{JS}(\mathbb{P}_{\mathbf{X}} \parallel \mathbb{P}_{\mathcal{G}}),$$

where  $\text{JS}(\cdot \parallel \cdot)$  denotes the Jensen–Shannon divergence between two distributions

$$\text{JS}(\mathbb{P}_{\mathbf{X}} \parallel \mathbb{P}_{\mathcal{G}}) = \frac{1}{2} \text{KL}(\mathbb{P}_{\mathbf{X}} \parallel \frac{\mathbb{P}_{\mathbf{X}} + \mathbb{P}_{\mathcal{G}}}{2}) + \frac{1}{2} \text{KL}(\mathbb{P}_{\mathcal{G}} \parallel \frac{\mathbb{P}_{\mathbf{X}} + \mathbb{P}_{\mathcal{G}}}{2}).$$

In words, the vanilla GAN (4.3) seeks a density  $\mathbb{P}_{\mathcal{G}}$  that is closest to  $\mathbb{P}_{\mathbf{X}}$  in terms of the Jensen–Shannon divergence. This view allows to generalize GANs to other variants, by changing the distance metric. Examples include f-GAN (Nowozin, Cseke and Tomioka, 2016), Wasserstein GAN (W-GAN) (Arjovsky, Chintala and Bottou, 2017), MMD GAN (Li, Swersky and Zemel, 2015), etc. We single out the Wasserstein GAN (W-GAN) (Arjovsky, Chintala and Bottou, 2017) to introduce due to its popularity. As the name suggests, it minimizes the Wasserstein distance between  $\mathbb{P}_{\mathbf{X}}$  and  $\mathbb{P}_{\mathcal{G}}$ :

$$(4.6) \quad \min_{\theta_{\mathcal{G}}} \quad \text{WS}(\mathbb{P}_{\mathbf{X}} \parallel \mathbb{P}_{\mathcal{G}}) = \min_{\theta_{\mathcal{G}}} \sup_{f: f \text{ 1-Lipschitz}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\mathbf{X}}} [f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\mathcal{G}}} [f(\mathbf{x})],$$

where  $f(\cdot)$  is taken over all Lipschitz functions with coefficient 1 and  $\mathbb{P}_{\mathcal{G}}$  is generated by a neural network model  $\mathbf{g}_{\theta_{\mathcal{G}}}$ . Comparing W-GAN (4.6) with the original formulation of GAN (4.3), one finds that the Lipschitz function  $f$  in (4.6) corresponds to the discriminator  $\mathcal{D}$  in (4.3) in the sense that they share similar objectives to differentiate the true distribution  $\mathbb{P}_{\mathbf{X}}$  from the fake one  $\mathbb{P}_{\mathcal{G}}$ . In the end, we would like to mention that GANs are more difficult to train than supervised deep learning models such as CNNs (Salimans et al., 2016). Apart from the training difficulty, how to evaluate GANs objectively and effectively is an ongoing research.

## 5. REPRESENTATION POWER: APPROXIMATION THEORY

Having seen the building blocks of deep learning models in the previous sections, it is natural to ask: what are the benefits of composing multiple layers of nonlinear functions? In this section, we address this question from an approximation theoretical point of view. Mathematically, letting  $\mathcal{H}$  be the space of functions representable by NNs, how well can a function  $f$  (with certain properties) be approximated by functions in  $\mathcal{H}$ ? We first revisit universal approximation theories, which are mostly developed for shallow neural nets (neural nets with a single hidden layer), and then provide recent results that demonstrate the benefits of depth in neural nets. Other notable works include Kolmogorov–Arnold superposition theorem (Arnold, 2009; Sprecher, 1965), and circuit complexity for neural nets (Parberry, 1994).

### 5.1 Universal approximation theory for shallow NNs

The universal approximation theories study the approximation of  $f$  in a space  $\mathcal{F}$  by a function represented by a one-hidden-layer neural net

$$(5.1) \quad g(\mathbf{x}) = \sum_{j=1}^N c_j \sigma_*(\mathbf{w}_j^\top \mathbf{x} - b_j),$$

where  $\sigma_* : \mathbb{R} \rightarrow \mathbb{R}$  is certain activation function and  $N$  is the number of hidden units in the neural net. For different space  $\mathcal{F}$  and activation function  $\sigma_*$ , there are upper bounds and lower bounds on the approximation error  $\|f - g\|$ . See [Pinkus \(1999\)](#) for a comprehensive overview. Here we present representative results.

First, as  $N \rightarrow \infty$ , any continuous function  $f$  can be approximated by some  $g$  under mild conditions. Loosely speaking, this is because each component  $\sigma_*(\mathbf{w}_j^\top \mathbf{x} - b_j)$  behaves like a basis function and functions in a suitable space  $\mathcal{F}$  admits a basis expansion. Given the above heuristics, the next natural question is: what is the rate of approximation for a finite  $N$ ?

Let us restrict the domain of  $\mathbf{x}$  to a unit ball  $B^d$  in  $\mathbb{R}^d$ . For  $p \in [1, \infty)$  and integer  $m \geq 1$ , consider the  $L^p$  space and the Sobolev space with standard norms

$$\|f\|_p = \left[ \int_{B^n} |g(\mathbf{x})|^p d\mathbf{x} \right]^{1/p}, \quad \|f\|_{m,p} = \left[ \sum_{0 \leq |\mathbf{k}| \leq m} \|D^{\mathbf{k}} f\|_p^p \right]^{1/p},$$

where  $D^{\mathbf{k}} f$  denotes partial derivatives indexed by  $\mathbf{k} \in \mathbb{Z}_+^d$ . Let  $\mathcal{F} \triangleq \mathcal{F}_p^m$  be the space of functions  $f$  in the Sobolev space with  $\|f\|_{m,p} \leq 1$ . Note that functions in  $\mathcal{F}$  have bounded derivatives up to  $m$ -th order, and that smoothness of functions is controlled by  $m$  (larger  $m$  means smoother). Denote by  $\mathcal{H}_N$  the space of functions with the form (5.1). The following general upper bound is due to [Mhaskar \(1996\)](#).

**THEOREM 5.1** (Theorem 2.1 in [Mhaskar \(1996\)](#)). *Assume  $\sigma_* : \mathbb{R} \rightarrow \mathbb{R}$  is such that  $\sigma_*$  has arbitrary order derivatives in an open interval  $I$ , and that  $\sigma_*$  is not a polynomial on  $I$ . Then, for any  $p \in [1, \infty)$ ,  $d \geq 2$ , and integer  $m \geq 1$ ,*

$$\sup_{f \in \mathcal{F}_p^m} \inf_{g \in \mathcal{H}_N} \|f - g\|_p \leq C_{d,m,p} N^{-m/d},$$

where  $C_{d,m,p}$  is independent of  $N$ , the number of hidden units.

In the above theorem, the condition on  $\sigma_*(\cdot)$  is mainly technical. Note, however, that the ReLU function does not satisfy this condition. This upper bound is useful when the dimension  $d$  is not large. It clearly implies that the one-hidden-layer neural net is able to approximate any smooth function with enough hidden units. However, it is unclear how to find a good approximator  $g$ ; nor do we have control over the magnitude of the parameters (huge weights are impractical). While increasing the number of hidden units  $N$  leads to better approximation, the exponent  $-m/d$  suggests the presence of the *curse of dimensionality*. The following (nearly) matching lower bound is stated in [Maiorov and Meir \(2000\)](#).

**THEOREM 5.2** (Theorem 5 in [Maiorov and Meir \(2000\)](#)). *Let  $p \geq 1$ ,  $m \geq 1$  and  $N \geq 2$ . If the activation function is the standard sigmoid function  $\sigma(t) =$*

$(1 + e^{-t})^{-1}$ , then

$$(5.2) \quad \sup_{f \in \mathcal{F}_p^m} \inf_{g \in \mathcal{H}_N} \|f - g\|_p \geq C'_{d,m,p} (N \log N)^{-m/d},$$

where  $C'_{d,m,p}$  is independent of  $N$ .

Results for other activation functions are also obtained by [Maierov and Meir \(2000\)](#). Moreover, the term  $\log N$  can be removed if we assume an additional continuity condition ([Mhaskar, 1996](#)).

For the natural space  $\mathcal{F}_p^m$  of smooth functions, the exponential dependence on  $d$  in the upper and lower bounds may look unappealing. However, [Barron \(1993\)](#) showed that for a different function space, there is a good dimension-free approximation by the neural nets. Suppose that a function  $f : \mathbb{R}^d \mapsto \mathbb{R}$  has a Fourier representation

$$(5.3) \quad f(\mathbf{x}) = \int_{\mathbb{R}^d} e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} \tilde{f}(\boldsymbol{\omega}) d\boldsymbol{\omega},$$

where  $\tilde{f}(\boldsymbol{\omega}) \in \mathbb{C}$ . Assume that  $f(\mathbf{0}) = 0$  and that the following quantity is finite

$$(5.4) \quad C_f = \int_{\mathbb{R}^d} \|\boldsymbol{\omega}\|_2 |\tilde{f}(\boldsymbol{\omega})| d\boldsymbol{\omega}.$$

[Barron \(1993\)](#) uncovers the following dimension-free approximation guarantee.

**THEOREM 5.3** (Proposition 1 in [Barron \(1993\)](#)). *Fix a  $C > 0$  and an arbitrary probability measure  $\mu$  on the unit ball  $B^d$  in  $\mathbb{R}^d$ . For every function  $f$  with  $C_f \leq C$  and every  $N \geq 1$ , there exists some  $g \in \mathcal{H}_N$  such that*

$$\left[ \int_{B^d} (f(\mathbf{x}) - g(\mathbf{x}))^2 \mu(d\mathbf{x}) \right]^{1/2} \leq \frac{2C}{\sqrt{N}}.$$

Moreover, the coefficients of  $g$  may be restricted to satisfy  $\sum_{j=1}^N |c_j| \leq 2C$ .

The upper bound is now independent of the dimension  $d$ . However,  $C_f$  may implicitly depend on  $d$ , as the formula in (5.4) involves an integration over  $\mathbb{R}^d$  (so for some functions  $C_f$  may depend exponentially on  $d$ ). Nevertheless, this theorem does characterize an interesting function space with an improved upper bound. Details of the function space are discussed by [Barron \(1993\)](#). This theorem can be generalized; see [Makovoz \(1996\)](#) for an example. Earlier papers on universal approximation include [Gybenko \(1989\)](#) and [Hornik \(1991\)](#).

To help understand why a dimensionality-free approximation holds, let us appeal to a heuristic argument given by Monte Carlo simulations. It is well-known that Monte Carlo approximation errors are independent of dimensionality in evaluation of high-dimensional integrals. Generally, let us generate  $\{\boldsymbol{\omega}_j\}_{1 \leq j \leq N}$  randomly from a given density  $p(\cdot)$  in  $\mathbb{R}^d$ , e.g., one can take  $p(\cdot) \propto \tilde{f}$ . Consider the approximation to (5.3) by

$$(5.5) \quad g_N(\mathbf{x}) = \frac{1}{N} \sum_{j=1}^N c_j e^{i\langle \boldsymbol{\omega}_j, \mathbf{x} \rangle}, \quad c_j = \frac{\tilde{f}(\boldsymbol{\omega}_j)}{p(\boldsymbol{\omega}_j)}.$$

Then,  $g_N(\mathbf{x})$  is a one-hidden-layer neural network with  $N$  units and the sinusoid

activation function. Note that  $\mathbb{E}g_N(\mathbf{x}) = f(\mathbf{x})$ , where the expectation is taken with respect to randomness  $\{\omega_j\}$ . Now, by independence, we have

$$\mathbb{E}(g_N(\mathbf{x}) - f(\mathbf{x}))^2 = \frac{1}{N} \text{Var}(c_j e^{i\langle \omega_j, \mathbf{x} \rangle}) \leq \frac{1}{N} \mathbb{E}c_j^2,$$

if  $\mathbb{E}c_j^2 < \infty$ . Therefore, the rate is independent of the dimensionality  $d$ , though the constant can be.

## 5.2 Approximation theory for multi-layer NNs

The approximation theory for multilayer neural nets is less understood compared with neural nets with one hidden layer. Driven by the success of deep learning, there are many recent papers focusing on the expressivity of deep neural nets. As studied by [Telgarsky \(2016\)](#); [Eldan and Shamir \(2016\)](#); [Mhaskar, Liao and Poggio \(2016\)](#); [Poggio et al. \(2017\)](#); [Bauer et al. \(2019\)](#); [Schmidt-Hieber \(2017\)](#); [Lin, Tegmark and Rolnick \(2017\)](#); [Rolnick and Tegmark \(2017\)](#), deep neural nets excel at representing *composition* of functions. This is perhaps not surprising since deep neural nets are themselves defined by composing layers of functions. Nevertheless, it points to a new territory rarely studied in statistics before. While deep neural nets are not panacea for all problems, it is meaningful to study specific examples where they succeed. Below we present two results, one based on [Lin, Tegmark and Rolnick \(2017\)](#); [Rolnick and Tegmark \(2017\)](#) and the other based on [Bauer et al. \(2019\)](#).

Suppose that the inputs  $\mathbf{x}$  have a bounded domain  $[-1, 1]^d$  for simplicity. As before, let  $\sigma_* : \mathbb{R} \rightarrow \mathbb{R}$  be a generic function, and  $\sigma_* = (\sigma_*, \dots, \sigma_*)^\top$  be element-wise application of  $\sigma_*$ . Consider a neural net which is similar to (2.1) but with scalar output:  $g(\mathbf{x}) = \mathbf{W}_\ell \sigma_*(\dots \sigma_*(\mathbf{W}_2 \sigma_*(\mathbf{W}_1 \mathbf{x})) \dots)$ . A unit or neuron refers to an element of vectors  $\sigma_*(\mathbf{W}_k \dots \sigma_*(\mathbf{W}_2 \sigma_*(\mathbf{W}_1 \mathbf{x})) \dots)$  for any  $k = 1, \dots, \ell - 1$ . For a multivariate polynomial  $p$ , define  $m_k(p)$  to be the smallest integer such that, for any  $\epsilon > 0$ , there exists a neural net  $g(\mathbf{x})$  satisfying  $\sup_{\mathbf{x}} |p(\mathbf{x}) - g(\mathbf{x})| < \epsilon$ , with  $k$  hidden layers (i.e.,  $\ell = k + 1$ ) and no more than  $m_k(p)$  neurons in total. Essentially,  $m_k(p)$  is the minimum number of neurons required to approximate  $p$  arbitrarily well.

**THEOREM 5.4** (Theorem 4.1 in [Rolnick and Tegmark \(2017\)](#)). *Let  $p(\mathbf{x})$  be a monomial  $x_1^{r_1} x_2^{r_2} \dots x_d^{r_d}$  with  $q = \sum_{j=1}^d r_j$ . Suppose that  $\sigma_*$  has derivatives of order  $2q$  at the origin, and that they are nonzero. Then,*

- (i)  $m_1(p) = \prod_{j=1}^d (r_j + 1)$ ;
- (ii)  $\min_k m_k(p) \leq \sum_{j=1}^d (7 \lceil \log_2(r_j) \rceil + 4)$ .

This theorem reveals a sharp distinction between shallow networks (one hidden layer) and deep networks. To represent a monomial function, a shallow network requires exponentially many neurons in terms of the dimension  $d$ , whereas linearly many neurons suffice for a deep network (with bounded  $r_j$ ). The exponential dependence on  $d$ , as shown in Theorem 5.4(i), is resonant with the curse of dimensionality widely seen in many fields; see [Donoho \(2000\)](#). One may ask: how does depth help? Depth circumvents this issue, at least for certain functions, by allowing us to represent function composition efficiently. Indeed, Theorem 5.4(ii) offers a nice result with clear intuitions: it is known that the product of two scalar

inputs can be represented using 4 neurons (Lin, Tegmark and Rolnick, 2017), so by composing multiple products, we can express monomials with  $O(d)$  neurons.

Recent advances in nonparametric regressions also support the idea that deep neural nets excel at representing composition of functions (Bauer et al., 2019; Schmidt-Hieber, 2017). In classical nonparametric regression setting, we have i.i.d. data  $\mathcal{D}_n = \{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$ . Well-known results (Stone, 1982) show that the optimal minimax rate of convergence for the regression function  $f(\mathbf{x}) = \mathbb{E}(y_i | \mathbf{x}_i = \mathbf{x})$  under the  $L^2$  error is  $O(n^{-\frac{2p}{2p+d}})$  where  $p$  is a measure of smoothness of  $f$ .

Without structural assumptions, the quality of estimates will suffer from a large dimension  $d$ . Subsequent works focused on restricted families of regression functions whose intrinsic dimensionality is small, which include additive models (Stone et al., 1985), single-index models (Hardle et al., 1993), functions with low-order interactions (Stone et al., 1994), etc. Recently, Bauer et al. (2019) considered a general form of hierarchical structure where the true regression function can be represented by a tree where each node has at most  $d^*$  children. In a very simplified fashion, the functions under consideration are represented by compositions  $\mathbf{G}_L(\mathbf{G}_{L-1}(\dots \mathbf{G}_0(\mathbf{a}_1^\top \mathbf{x}, \dots, \mathbf{a}_{d^*}^\top \mathbf{x})))$ , where each  $\mathbf{G}_\ell : \mathbb{R}^{d^*} \rightarrow \mathbb{R}^{d^*}$  ( $\ell = 0, \dots, L-1$ ) is a map in lower-dimensional spaces, and it satisfies certain smoothness condition. Then, it is shown that there is a function estimator  $\hat{f}_n$  representable by neural networks such that the  $L^2$  error satisfies

$$(5.6) \quad \mathbb{E}_{\mathcal{D}_n} \mathbb{E}_{\mathbf{x}} \left| \hat{f}_n(\mathbf{x}) - f(\mathbf{x}) \right|^2 = O((\log n)^3 n^{-\frac{2p}{2p+d^*}}).$$

As a relief, this result implies that the intrinsic dimension  $d^*$ , rather than the ambient dimension  $d$ , determines the convergence rate.

This result provides another justification for deep neural nets: if data are truly hierarchical, then the quality of approximators by deep neural nets depends on the intrinsic dimensionality, which avoids the curse of dimensionality.

We point out that the approximation theory for deep learning is far from complete. For example, in Theorem 5.4, the condition on  $\sigma_*$  excludes the widely used ReLU activation function, and there are no constraints on the magnitude of the weights (so they can be unreasonably large). In (5.6), there is no guarantee that training neural nets with SGD would find such  $\hat{f}_n$ .

## 6. TRAINING DEEP NEURAL NETS

The *existence* of a good function approximator in the NN function class does not explain why in practice we can easily *find* them. In this section, we introduce standard methods, namely *stochastic gradient descent* (SGD) and its variants, to train deep neural networks (or to find such a good approximator). As with many statistical machine learning tasks, training DNNs follows the *empirical risk minimization* (ERM) paradigm which solves the following optimization problem

$$(6.1) \quad \text{minimize}_{\boldsymbol{\theta} \in \mathbb{R}^p} \quad \ell_n(\boldsymbol{\theta}) \triangleq \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i).$$

Here  $\mathcal{L}(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$  measures the discrepancy between the prediction  $f(\mathbf{x}_i; \boldsymbol{\theta})$  of the neural network and the true label  $y_i$ . Correspondingly, denote by  $\ell(\boldsymbol{\theta}) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\mathcal{L}(f(\mathbf{x}; \boldsymbol{\theta}), y)]$  the out-of-sample error, where  $\mathcal{D}$  is the joint distribution over  $(y, \mathbf{x})$ . Solving ERM (6.1) for deep neural nets faces various challenges that

roughly fall into the following three categories.

- *Scalability and nonconvexity.* Both the sample size  $n$  and the number of parameters  $p$  can be huge for modern deep learning applications, as we have seen in Table 1. Many optimization algorithms are not practical due to the computational costs and memory constraints. What is worse, the empirical loss function  $\ell_n(\boldsymbol{\theta})$  in deep learning is often nonconvex. It is *a priori* not clear whether an optimization algorithm can drive the empirical loss (6.1) small.
- *Numerical stability.* With a large number of layers in DNNs, the magnitudes of the hidden nodes can be drastically different, which may result in the “exploding gradients” or “vanishing gradients” issue during the training process. This is because the recursive relations across layers often lead to exponentially increasing / decreasing values in both forward passes and backward passes.
- *Generalization performance.* Our ultimate goal is to find a parameter  $\hat{\boldsymbol{\theta}}$  such that the out-of-sample error  $\ell(\hat{\boldsymbol{\theta}})$  is small. However, in the over-parametrized regime where  $p$  is much larger than  $n$ , the underlying neural network has the potential to fit the training data perfectly while performing poorly on the test data. To avoid this overfitting issue, proper regularization, whether explicit or implicit, is needed in the training process for the neural nets to generalize.

In the following three subsections, we discuss practical solutions / proposals to address these challenges.

### 6.1 Stochastic gradient descent

Stochastic gradient descent (SGD) (Robbins and Monro, 1951) is by far the most popular optimization algorithm to solve ERM (6.1) for large-scale problems. It has the following simple update rule:

$$(6.2) \quad \boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_t G(\boldsymbol{\theta}^t) \quad \text{with} \quad G(\boldsymbol{\theta}^t) = \nabla \mathcal{L}(f(\mathbf{x}_{i_t}; \boldsymbol{\theta}^t), y_{i_t})$$

for  $t = 0, 1, 2, \dots$ , where  $\eta_t > 0$  is the step size (or learning rate),  $\boldsymbol{\theta}^0 \in \mathbb{R}^p$  is an initial point and  $i_t$  is chosen randomly from  $\{1, 2, \dots, n\}$ . It is easy to verify that  $G(\boldsymbol{\theta}^t)$  is an unbiased estimate of  $\nabla \ell_n(\boldsymbol{\theta}^t)$ . The advantage of SGD is clear: compared with the full gradient descent, which goes over the entire dataset to compute the average gradient in every update, SGD uses one sample to compute the gradient in each update and hence is considerably more efficient in terms of both computation and memory (especially in the first few iterations).

Apart from practical benefits of SGD, how well does SGD perform theoretically in terms of minimizing  $\ell_n(\boldsymbol{\theta})$ ? We begin with the convex case, i.e., the case where the loss function is convex w.r.t.  $\boldsymbol{\theta}$ . It is well understood in literature that with proper choices of the step sizes  $\{\eta_t\}$ , SGD is guaranteed to achieve both *consistency* and *asymptotic normality*.

- *Consistency.* If  $\ell(\boldsymbol{\theta})$  is a strongly convex function<sup>10</sup>, then under some mild

<sup>10</sup>For results on consistency and asymptotic normality, we consider the case where in each step of SGD, the stochastic gradient is computed using a fresh sample  $(y, \mathbf{x})$  from  $\mathcal{D}$ . This allows to view SGD as an optimization algorithm to minimize the population loss  $\ell(\boldsymbol{\theta})$ .



conditions<sup>11</sup>, learning rates that satisfy

$$(6.3) \quad \sum_{t=0}^{\infty} \eta_t = +\infty \quad \text{and} \quad \sum_{t=0}^{\infty} \eta_t^2 < +\infty$$

guarantee almost sure convergence to the unique minimizer  $\boldsymbol{\theta}^* \triangleq \operatorname{argmin}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})$ , i.e.,  $\boldsymbol{\theta}^t \xrightarrow{\text{a.s.}} \boldsymbol{\theta}^*$  as  $t \rightarrow \infty$  (Robbins and Monro, 1951; Kiefer et al., 1952; Bottou, 1998; Kushner and Yin, 2003). The requirements in (6.3) can be viewed from the perspective of bias-variance tradeoff: the first condition ensures that the iterates can reach the minimizer (controlled bias), and the second ensures that stochasticity does not prevent convergence (controlled variance).

- *Asymptotic normality.* It is proved by Polyak and Tsykin (1979) that for robust linear regression with fixed dimension  $p$ , under the choice  $\eta_t = t^{-1}$ ,  $\sqrt{t}(\boldsymbol{\theta}^t - \boldsymbol{\theta}^*)$  is asymptotically normal under some regularity conditions (but  $\boldsymbol{\theta}^t$  is not asymptotically efficient in general). Moreover, by averaging the iterates of SGD, Polyak and Juditsky (1992) proved that even with a *larger* step size  $\eta_t \propto t^{-\alpha}$ ,  $\alpha \in (1/2, 1)$ , the averaged iterate  $\bar{\boldsymbol{\theta}}^t = t^{-1} \sum_{s=1}^t \boldsymbol{\theta}^s$  is asymptotic efficient for robust linear regression. These strong results show that SGD with averaging performs as well as the MLE asymptotically, in addition to its computational efficiency.

These classical results, however, fail to explain the effectiveness of SGD when dealing with nonconvex loss functions in deep learning. Admittedly, finding global minima of nonconvex functions is computationally infeasible in the worst case. Nevertheless, recent work (Allen-Zhu, Li and Song, 2018; Du et al., 2018) bypasses the worst-case scenario by focusing on losses incurred by over-parametrized deep learning models. In particular, they show that (stochastic) gradient descent converges linearly towards the *global* minimizer of  $\ell_n(\boldsymbol{\theta})$  as long as the neural network is sufficiently *over-parametrized*. This phenomenon is formalized below.

**THEOREM 6.1** (Theorem 2 in Allen-Zhu, Li and Song, 2018). *Let  $\{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$  be a training set satisfying  $\min_{i,j:i \neq j} \|\mathbf{x}_i - \mathbf{x}_j\|_2 \geq \delta > 0$ . Consider fitting the data using a feed-forward neural network (1.1) with ReLU activations. Denote by  $L$  (resp.  $W$ ) the depth (resp. width) of the network. Suppose that the neural network is sufficiently over-parametrized, i.e.,*

$$(6.4) \quad W \gg \text{poly} \left( n, L, \frac{1}{\delta} \right),$$

*where **poly** means a polynomial function. Then with high probability, running SGD (6.2) with certain random initialization and properly chosen step sizes yields  $\ell_n(\boldsymbol{\theta}^t) \leq \varepsilon$  in  $t \asymp \log \frac{1}{\varepsilon}$  iterations.*

Two notable features are worth mentioning: (1) first, the network under consideration is sufficiently over-parametrized (cf. (6.4)) in which the number of parameters is *much* larger than the number of samples, and (2) one needs to initialize the weight matrices to be in near-isometry such that the magnitudes of the hidden nodes do not blow up or vanish. In a nutshell, *over-parametrization* and *random initialization* together ensure that the loss function (6.1) has a benign landscape<sup>12</sup>

<sup>11</sup>One example of such condition can be constraining the second moment of the gradients:  $\mathbb{E} [\|\nabla \mathcal{L}(\mathbf{x}_i, y_i; \boldsymbol{\theta}^t)\|_2^2] \leq C_1 + C_2 \|\boldsymbol{\theta}^t - \boldsymbol{\theta}^*\|_2^2$  for some  $C_1, C_2 > 0$ . See Bottou (1998) for details.

<sup>12</sup>In Allen-Zhu, Li and Song (2018), the loss function  $\ell_n(\boldsymbol{\theta})$  satisfies the PL condition.

around the initial point, which in turn implies fast convergence of SGD iterates.

There are certainly other challenges for vanilla SGD to train deep neural nets: (1) training algorithms are often implemented in GPUs, and therefore it is important to tailor the algorithm to the infrastructure, (2) the vanilla SGD might converge very slowly for deep neural networks, albeit good theoretical guarantees for well-behaved problems, and (3) the learning rates  $\{\eta_t\}$  can be difficult to tune in practice. To address the aforementioned challenges, three important variants of SGD, namely *mini-batch SGD*, *momentum-based SGD*, and *SGD with adaptive learning rates* are introduced.

**6.1.1 Mini-batch SGD.** Modern computational infrastructures (e.g., GPUs) can evaluate the gradient on a number (say 64) of examples (a sub-sample) as efficiently as evaluating that on a single example. To utilize this advantage, mini-batch SGD with batch size  $K \geq 1$  forms the stochastic gradient through  $K$  random samples:

$$(6.5) \quad \boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_t G(\boldsymbol{\theta}^t) \quad \text{with} \quad G(\boldsymbol{\theta}^t) = \frac{1}{K} \sum_{k=1}^K \nabla \mathcal{L}(f(\mathbf{x}_{i_t^k}; \boldsymbol{\theta}^t), y_{i_t^k}),$$

where for each  $1 \leq k \leq K$ ,  $i_t^k$  is sampled uniformly from  $\{1, 2, \dots, n\}$ . Mini-batch SGD, which is an “interpolation” between gradient descent and stochastic gradient descent, achieves the best of both worlds: (1) using  $1 \ll K \ll n$  samples to estimate the gradient, one effectively reduces the variance and hence accelerates the convergence, and (2) by taking the batch size  $K$  appropriately (say 64 or 128), the stochastic gradient  $G(\boldsymbol{\theta}^t)$  can be efficiently computed using the matrix computation toolboxes on GPUs.

**6.1.2 Momentum-based SGD.** While mini-batch SGD forms the foundation of training neural networks, it can sometimes be slow to converge due to its oscillation behavior (Sutskever et al., 2013). Optimization community has long investigated how to accelerate the convergence of gradient descent, which results in a beautiful technique called *momentum methods* (Polyak, 1964; Nesterov, 1983). Similar to gradient descent with moment, *momentum-based SGD*, instead of moving the iterate  $\boldsymbol{\theta}^t$  in the direction of the current stochastic gradient  $G(\boldsymbol{\theta}^t)$ , smooth the past (stochastic) gradients  $\{G(\boldsymbol{\theta}^t)\}$  to stabilize the update directions. Mathematically, let  $\mathbf{v}^t \in \mathbb{R}^p$  be the direction of update in the  $t$ th iteration, i.e.,

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_t \mathbf{v}^t.$$

Here  $\mathbf{v}^0 = G(\boldsymbol{\theta}^0)$  and for  $t = 1, 2, \dots$

$$(6.6) \quad \mathbf{v}^t = \rho \mathbf{v}^{t-1} + G(\boldsymbol{\theta}^t)$$

with  $0 < \rho < 1$ . A typical choice of  $\rho$  is 0.9. Notice that  $\rho = 0$  recovers the mini-batch SGD (6.5), where no past information of gradients is used. A simple unrolling of  $\mathbf{v}^t$  reveals that  $\mathbf{v}^t$  is actually an exponential smoothing (un-normalized) of the past gradients, i.e.,  $\mathbf{v}^t = \sum_{j=0}^t \rho^{t-j} G(\boldsymbol{\theta}^j)$ . Compared with vanilla mini-batch SGD, the inclusion of the momentum “smooths” the oscillation direction and accumulates the persistent descent direction. We want to emphasize that theoretical justifications of momentum in the *stochastic* setting is not fully understood (Kidambi et al., 2018; Jain et al., 2017).

*6.1.3 SGD with adaptive learning rates.* In optimization, *preconditioning* is often used to accelerate first-order optimization algorithms. In principle, one can apply this to SGD, which yields the following update rule:

$$(6.7) \quad \boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta_t \mathbf{P}_t^{-1} G(\boldsymbol{\theta}^t)$$

with  $\mathbf{P}_t \in \mathbb{R}^{p \times p}$  being a preconditioner at the  $t$ -th step. Newton's method can be viewed as one type of preconditioning where  $\mathbf{P}_t = \nabla^2 \ell_n(\boldsymbol{\theta}^t)$ . The advantages of preconditioning are two-fold: first, a good preconditioner reduces the condition number by changing the local geometry to be more homogeneous, which is amenable to fast convergence; second, a good preconditioner frees practitioners from laboring tuning of the step sizes, as is the case with Newton's method. AdaGrad, an adaptive gradient method proposed by [Duchi, Hazan and Singer \(2011\)](#), builds a preconditioner  $\mathbf{P}_t$  based on information of the past gradients:

$$(6.8) \quad \mathbf{P}_t = \left\{ \text{diag} \left( \sum_{j=0}^t G(\boldsymbol{\theta}^j) G(\boldsymbol{\theta}^j)^\top \right) \right\}^{1/2}.$$

Since we only require the diagonal part, this preconditioner (and its inverse) can be efficiently computed in practice. In addition, investigating (6.7) and (6.8), one can see that AdaGrad adapts to the importance of each coordinate of the parameters, which receive small learning rates if updated frequently and vice versa. In practice, one adds a small quantity  $\delta > 0$  (say  $10^{-8}$ ) to the diagonal entries to avoid singularity (numerical underflow). A notable drawback of AdaGrad is that the effective learning rate vanishes quickly along the learning process. This is because the historical sum of the gradients can only increase with time. RMSProp ([Hinton, Srivastava and Swersky, 2012](#)) is a popular remedy for this problem which incorporates the idea of exponential smoothing:

$$(6.9) \quad \mathbf{P}_t = \left\{ \text{diag} \left( \rho \mathbf{P}_{t-1} + (1 - \rho) G(\boldsymbol{\theta}^t) G(\boldsymbol{\theta}^t)^\top \right) \right\}^{1/2}.$$

Again, the decaying parameter  $\rho$  is usually set to be 0.9. Later, Adam ([Kingma and Ba, 2014](#); [Reddi, Kale and Kumar, 2018](#)) combines the momentum method and adaptive learning rate and becomes the default training algorithms in many deep learning applications.

## 6.2 Easing numerical instability

For very deep neural networks or RNNs with long dependencies, training difficulties often arise when the values of nodes have different magnitudes or when the gradients “vanish” or “explode” during back-propagation. Here we discuss three partial solutions to alleviate this problem. Other proposals for promoting numerical stability include AdaNet ([Cortes et al., 2017](#)) for example.

*6.2.1 ReLU activation function.* One useful characteristic of the ReLU function is that its derivative is either 0 or 1, and the derivative remains 1 even for a large input. This is in sharp contrast with the standard sigmoid function  $(1 + e^{-t})^{-1}$  which results in a very small derivative when inputs have large magnitude. The consequence of small derivatives across many layers is that gradients tend to be “killed”, which means that gradients become approximately zero in deep nets.

The popularity of the ReLU activation function and its variants (e.g., leaky ReLU) is largely attributable to the above reason. It has been well observed that

the ReLU activation function has superior training performance over the sigmoid function (Krizhevsky, Sutskever and Hinton, 2012; Maas, Hannun and Ng, 2013).

**6.2.2 Skip connections.** We have introduced skip connections in Section 3.3. Why are skip connections helpful for reducing numerical instability? This structure does not introduce a larger function space, since the identity map can be also represented with ReLU activations:  $\mathbf{x} = \sigma(\mathbf{x}) - \sigma(-\mathbf{x})$ .

One explanation is that skip connections bring ease to the training / optimization process. Suppose that we have a general nonlinear function  $\mathbf{F}(\mathbf{x}_\ell; \boldsymbol{\theta}_\ell)$ . With a skip connection, we represent the map as  $\mathbf{x}_{\ell+1} = \mathbf{x}_\ell + \mathbf{F}(\mathbf{x}_\ell; \boldsymbol{\theta}_\ell)$  instead. Now the gradient  $\partial \mathbf{x}_{\ell+1} / \partial \mathbf{x}_\ell$  becomes

$$(6.10) \quad \frac{\partial \mathbf{x}_{\ell+1}}{\partial \mathbf{x}_\ell} = \mathbf{I} + \frac{\partial \mathbf{F}(\mathbf{x}_\ell; \boldsymbol{\theta}_\ell)}{\partial \mathbf{x}_\ell} \quad \text{instead of} \quad \frac{\partial \mathbf{F}(\mathbf{x}_\ell; \boldsymbol{\theta}_\ell)}{\partial \mathbf{x}_\ell},$$

where  $\mathbf{I}$  is an identity matrix. By the chain rule, gradient update requires computing products of many components, e.g.,  $\frac{\partial \mathbf{x}_L}{\partial \mathbf{x}_1} = \prod_{\ell=1}^{L-1} \frac{\partial \mathbf{x}_{\ell+1}}{\partial \mathbf{x}_\ell}$ , so it is desirable to keep the spectra (singular values) of each component  $\frac{\partial \mathbf{x}_{\ell+1}}{\partial \mathbf{x}_\ell}$  close to 1. In neural nets, with skip connections, this is easily achieved if the parameters have small values; otherwise, this may not be achievable even with careful initialization and tuning. Notably, training neural nets with hundreds of layers is possible with the help of skip connections.

**6.2.3 Batch normalization.** Recall that in regression analysis, one often standardizes the design matrix so that the features have zero mean and unit variance. Batch normalization extends this standardization procedure from the input layer to all the hidden layers. Mathematically, fix a mini-batch of input data  $\{(\mathbf{x}_i, y_i)\}_{i \in \mathcal{B}}$ , where  $\mathcal{B} \subset [n]$ . Let  $\mathbf{h}_i^{(\ell)}$  be the feature of the  $i$ -th example in the  $\ell$ -th layer ( $\ell = 0$  corresponds to the input  $\mathbf{x}_i$ ). The batch normalization layer computes the normalized version of  $\mathbf{h}_i^{(\ell)}$  via the following steps:

$$\boldsymbol{\mu} \triangleq \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \mathbf{h}_i^{(\ell)}, \quad \boldsymbol{\sigma}^2 \triangleq \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} (\mathbf{h}_i^{(\ell)} - \boldsymbol{\mu})^2 \quad \text{and} \quad \mathbf{h}_{i,\text{norm}}^{(l)} \triangleq \frac{\mathbf{h}_i^{(\ell)} - \boldsymbol{\mu}}{\boldsymbol{\sigma}}.$$

Here all the operations are element-wise. In words, batch normalization computes the z-score for each feature over the mini-batch  $\mathcal{B}$  and uses that as inputs to subsequent layers. To make it more versatile, a typical batch normalization layer has two additional learnable parameters  $\boldsymbol{\gamma}^{(\ell)}$  and  $\boldsymbol{\beta}^{(\ell)}$  such that

$$\mathbf{h}_{i,\text{new}}^{(l)} = \boldsymbol{\gamma}^{(l)} \odot \mathbf{h}_{i,\text{norm}}^{(l)} + \boldsymbol{\beta}^{(l)}.$$

Again  $\odot$  denotes the element-wise multiplication. As can be seen,  $\boldsymbol{\gamma}^{(\ell)}$  and  $\boldsymbol{\beta}^{(\ell)}$  set the new feature  $\mathbf{h}_{i,\text{new}}^{(l)}$  to have mean  $\boldsymbol{\beta}^{(\ell)}$  and standard deviation  $\boldsymbol{\gamma}^{(\ell)}$ . The introduction of batch normalization makes the training of neural networks much easier and smoother. More importantly, it allows the neural nets to perform well over a large family of hyper-parameters including the number of layers, the number of hidden units, etc. At test time, the batch normalization layer needs more care. For brevity we omit the details and refer to Ioffe and Szegedy (2015).

### 6.3 Regularization techniques

So far we have focused on training techniques to drive the empirical loss (6.1) small efficiently. Here we proceed to discuss common practice to improve the

generalization power of trained neural nets.

**6.3.1 Weight decay.** One natural regularization idea is to add an  $\ell_2$  penalty to the loss function. This regularization technique is known as the weight decay in deep learning. We have seen one example in (2.7). For general deep neural nets, the loss to optimize is  $\ell_n^\lambda(\boldsymbol{\theta}) = \ell_n(\boldsymbol{\theta}) + r_\lambda(\boldsymbol{\theta})$  where

$$r_\lambda(\boldsymbol{\theta}) = \lambda \sum_{\ell=1}^L \sum_{j,j'} [W_{j,j'}^{(\ell)}]^2.$$

Note that the bias (intercept) terms are not penalized. If  $\ell_n(\boldsymbol{\theta})$  is a least square loss, then regularization with weight decay gives precisely ridge regression. The penalty  $r_\lambda(\boldsymbol{\theta})$  is a smooth function and thus it can be also implemented efficiently with back-propagation.

**6.3.2 Dropout.** Dropout, introduced by Hinton et al. (2012), prevents overfitting by randomly dropping out subsets of features during training. Take the  $l$ -th layer of the feed-forward neural network as an example. Instead of propagating all the features in  $\mathbf{h}^{(\ell)}$  for later computations, dropout randomly omits some of its entries by

$$\mathbf{h}_{\text{drop}}^{(\ell)} = \mathbf{h}^{(\ell)} \odot \text{mask}^\ell,$$

where  $\odot$  denotes element-wise multiplication as before, and  $\text{mask}^\ell$  is a vector of Bernoulli variables with success probability  $p$ . It is sometimes useful to rescale the features  $\mathbf{h}_{\text{inv drop}}^{(\ell)} = \mathbf{h}_{\text{drop}}^{(\ell)} / p$ , which is called *inverted dropout*. During training,  $\text{mask}^\ell$  are i.i.d. vectors across mini-batches and layers. However, when testing on fresh samples, dropout is disabled and the original features  $\mathbf{h}^{(\ell)}$  are used to compute the output label  $y$ . It has been nicely shown by Wager, Wang and Liang (2013) that for generalized linear models, dropout serves as adaptive regularization: the loss function of the perturbed inputs on average equals the loss of the unperturbed inputs plus certain regularization term of the parameters. In particular, Wager, Wang and Liang (2013) shows that in the simplest case of linear regression, it is equivalent to  $\ell_2$  regularization. Another possible way to understand the regularization effect of dropout is through the lens of bagging (Goodfellow, Bengio and Courville, 2016). Since different mini-batches have different masks, dropout can be viewed as training a large ensemble of classifiers at the same time, with a further constraint that the parameters are shared. Theoretical justification remains elusive.

**6.3.3 Data augmentation.** Data augmentation is a technique of enlarging the dataset when we have knowledge about the invariance structure of data. It implicitly increases the sample size and usually regularizes the model effectively. For example, in image classification, we have strong prior knowledge about what invariance properties a good classifier should possess. The label of an image should not be affected by translation, rotation, flipping, and even crops of the image. Hence one can augment the dataset by randomly translating, rotating and cropping the images in the original dataset.

Formally, during training we want to minimize the loss  $\ell_n(\boldsymbol{\theta}) = \sum_i \mathcal{L}(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$  w.r.t. parameters  $\boldsymbol{\theta}$ , and we know a priori that certain transformation  $T \in \mathcal{T}$  where  $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$  (e.g., affine transformation) should not change the category / label

of a training sample. In principle, if computation costs were not a consideration, we could convert this knowledge to a constraint  $f_{\theta}(T\mathbf{x}_i) = f_{\theta}(\mathbf{x}_i), \forall T \in \mathcal{T}$  in the minimization formulation. Instead of solving a constrained optimization problem, data augmentation enlarges the training dataset by sampling  $T \in \mathcal{T}$  and generating new data  $\{(T\mathbf{x}_i, y_i)\}$ . In this sense, data augmentation induces invariance properties through sampling, which results in a much bigger dataset than the original one (Chen, Dobriban and Lee, 2019).

## 7. GENERALIZATION POWER

Section 6 has focused on the in-sample / training error obtained via SGD, but this alone does not guarantee good performance with respect to the out-of-sample / test error. The gap between the in-sample error and the out-of-sample error, namely the *generalization gap*, has been the focus of statistical learning theory since its birth; see Shalev-Shwartz and Ben-David (2014) for an excellent introduction to this topic.

While understanding the generalization power of deep neural nets is difficult (Zhang et al., 2016), we sample recent endeavors in this section. From a high-level point of view, these approaches can be divided into two categories, namely *algorithm-independent controls* and *algorithm-dependent controls*. More specifically, algorithm-independent controls focus solely on bounding the *complexity* of the function class represented by certain deep neural networks. In contrast, algorithm-dependent controls take into account the algorithm (e.g., SGD) used to train the neural network.

### 7.1 Algorithm-independent controls: uniform convergence

The key to algorithm-independent controls is the notion of *complexity* of the function class parametrized by certain neural networks. Informally, as long as the complexity is not too large, the generalization gap of *any* function in the function class is well-controlled. However, the standard complexity measure (e.g., VC dimension (Vapnik and Chervonenkis, 1971)) is at least proportional to the number of weights in a neural network (Anthony and Bartlett, 2009; Shalev-Shwartz and Ben-David, 2014), which fails to explain the practical success of deep learning. The caveat here is that the function class under consideration is *all* the functions realized by certain neural networks, with *no* restrictions on the size of the weights at all. On the other hand, for the class of linear functions with bounded norm, i.e.,  $\{\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} \mid \|\mathbf{w}\|_2 \leq M\}$ , it is well understood that the complexity of this function class (measured in terms of the empirical Rademacher complexity) with respect to a random sample  $\{\mathbf{x}_i\}_{1 \leq i \leq n}$  is upper bounded by  $\max_i \|\mathbf{x}_i\|_2 M / \sqrt{n}$ , which is independent of the number of parameters in  $\mathbf{w}$ . This motivates researchers to investigate the complexity of *norm-controlled* deep neural networks<sup>13</sup> (Neyshabur, Tomioka and Srebro, 2015; Bartlett, Foster and Telgarsky, 2017; Golowich, Rakhlin and Shamir, 2017; Li et al., 2018b). Setting the stage, we introduce a few necessary notations and facts. The key object under study is the function class parametrized by the following fully-connected neural network with depth  $L$ :

$$(7.1) \quad \mathcal{F}_L \triangleq \left\{ \mathbf{x} \mapsto \mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\cdots \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))) \mid (\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathcal{W} \right\}.$$

<sup>13</sup>Such attempts have been made in the seminal work Bartlett (1998).



Here  $(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_L) \in \mathcal{W}$  represents a certain constraint on the parameters. For instance, one can restrict the Frobenius norm of each parameter  $\mathbf{W}_l$  through the constraint  $\|\mathbf{W}_l\|_F \leq M_F(l)$ , where  $M_F(l)$  is some positive quantity. With regard to the complexity measure, it is standard to use *Rademacher complexity* to control the capacity of the function class of interest.

DEFINITION 7.1 (Empirical Rademacher complexity). *The empirical Rademacher complexity of a function class  $\mathcal{F}$  w.r.t. a dataset  $S \triangleq \{\mathbf{x}_i\}_{1 \leq i \leq n}$  is defined as*

$$(7.2) \quad \mathcal{R}_S(\mathcal{F}) = \mathbb{E}_{\boldsymbol{\varepsilon}} \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\mathbf{x}_i) \right],$$

where  $\boldsymbol{\varepsilon} \triangleq (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)$  is composed of i.i.d. Rademacher random variables, i.e.,  $\mathbb{P}(\varepsilon_i = 1) = \mathbb{P}(\varepsilon_i = -1) = 1/2$ .

In words, Rademacher complexity measures the ability of the function class to fit the random noise represented by  $\boldsymbol{\varepsilon}$ . Intuitively, a function class with a larger Rademacher complexity is more prone to overfitting. We now formalize the connection between the empirical Rademacher complexity and the out-of-sample error; see Chapter 24 in [Shalev-Shwartz and Ben-David \(2014\)](#).

THEOREM 7.1. *Assume that for all  $f \in \mathcal{F}$  and all  $(y, \mathbf{x})$  we have  $|\mathcal{L}(f(\mathbf{x}), y)| \leq 1$ . In addition, assume that for any fixed  $y$ , the univariate function  $\mathcal{L}(\cdot, y)$  is Lipschitz with constant 1. Then with probability at least  $1 - \delta$  over the sample  $S \triangleq \{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$ , one has for all  $f \in \mathcal{F}$*

$$\underbrace{\mathbb{E}_{(y, \mathbf{x}) \sim \mathcal{D}} [\mathcal{L}(f(\mathbf{x}), y)]}_{\text{out-of-sample error}} \leq \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(\mathbf{x}_i), y_i)}_{\text{in-sample error}} + 2\mathcal{R}_S(\mathcal{F}) + 4\sqrt{\frac{\log(4/\delta)}{n}}.$$

In words, the generalization gap of any function  $f$  that lies in  $\mathcal{F}$  is well-controlled as long as the Rademacher complexity of  $\mathcal{F}$  is not too large. With this connection in place, we single out the following complexity bound.

THEOREM 7.2 (Theorem 1 in [Golowich, Rakhlin and Shamir, 2017](#)). *Consider the function class  $\mathcal{F}_L$  in (7.1), where each parameter  $\mathbf{W}_l$  has Frobenius norm at most  $M_F(l)$ . Further suppose that the element-wise activation function  $\sigma(\cdot)$  is 1-Lipschitz and positive-homogeneous (i.e.,  $\sigma(c \cdot x) = c\sigma(x)$  for all  $c \geq 0$ ). Then the empirical Rademacher complexity (7.2) w.r.t.  $S \triangleq \{\mathbf{x}_i\}_{1 \leq i \leq n}$  satisfies*

$$(7.3) \quad \mathcal{R}_S(\mathcal{F}_L) \leq \max_i \|\mathbf{x}_i\|_2 \cdot \frac{4\sqrt{L} \prod_{l=1}^L M_F(l)}{\sqrt{n}}.$$

The upper bound of the empirical Rademacher complexity (7.3) is in a similar vein to that of linear functions with bounded norm, i.e.,  $\max_i \|\mathbf{x}_i\|_2 M / \sqrt{n}$ , where  $\sqrt{L} \prod_{l=1}^L M_F(l)$  plays the role of  $M$  in the latter case. Moreover, ignoring the term  $\sqrt{L}$ , the upper bound (7.3) does not depend on the size of the network in an explicit way if  $M_F(l)$  sharply concentrates around 1. This reveals that the capacity of the neural network is well-controlled, regardless of the number of parameters, as long as the Frobenius norm of the parameters is bounded. Extensions to other



norm constraints, e.g., spectral norm constraints, path norm constraints have been considered by Neyshabur, Tomioka and Srebro (2015); Bartlett, Foster and Telgarsky (2017); Li et al. (2018b); Klusowski and Barron (2016); E, Ma and Wang (2019). This line of work improves upon traditional capacity analysis of neural networks in the over-parametrized setting, because the upper bounds derived are often size-independent. Having said this, two important remarks are in order: (1) the upper bounds (e.g.,  $\prod_{l=1}^L M_F(l)$ ) involve implicit dependence on the size of the weight matrix and the depth of the neural network, which is hard to characterize; (2) the upper bound on the Rademacher complexity offers a uniform bound over all functions in the function class, which is a pure statistical result. However, it stays silent about how and why standard training algorithms like SGD can obtain a function whose parameters have small norms.

## 7.2 Algorithm-dependent controls

In this subsection, we bring computational thinking into statistics and investigate the role of algorithms in the generalization power of deep learning. The consideration of algorithms is quite natural and well motivated: (1) local/global minima reached by different algorithms can exhibit totally different generalization behaviors due to extreme nonconvexity, which marks a huge difference from traditional models, (2) the *effective* capacity of neural nets is possibly not large, since a particular algorithm does not explore the entire parameter space.

These demonstrate the fact that on top of the complexity of the function class, the inherent property of the algorithm we use plays an important role in the generalization ability of deep learning. In what follows, we survey three different ways to obtain upper bounds on the generalization errors by exploiting properties of the algorithms.

*7.2.1 Mean field view of neural nets.* As we have emphasized, modern deep learning models are highly over-parametrized. A line of work focuses on the over-parametrized one-hidden-layer neural networks (Mei, Montanari and Nguyen, 2018; Sirignano and Spiliopoulos, 2018; Rotskoff and Vanden-Eijnden, 2018; Chizat and Bach, 2018; Mei, Misiakiewicz and Montanari, 2019; Javanmard, Mondelli and Montanari, 2019). More specifically, let  $f(\mathbf{x}; \boldsymbol{\theta}) = N^{-1} \sum_{i=1}^N \sigma(\boldsymbol{\theta}_i^\top \mathbf{x})$  be a function given by a one-hidden-layer neural net with  $N$  hidden units and parameters  $\boldsymbol{\theta} \triangleq [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N]^\top \in \mathbb{R}^{N \times d}$ . Here,  $\sigma(\cdot)$  is the ReLU activation function. Suppose that one runs SGD on the samples (training examples)  $(\mathbf{x}_k, y_k)$  (where  $k = 1, 2, \dots$ ) to minimize the population risk  $R(\boldsymbol{\theta}) = \mathbb{E}[(y - \hat{f}(\mathbf{x}; \boldsymbol{\theta}))^2]$  w.r.t. the parameters  $\boldsymbol{\theta}$ . The central question here is how good is  $R(\boldsymbol{\theta}^k)$ , where  $\boldsymbol{\theta}^k$  results from running  $k$  steps of SGD?

To answer this question, a key observation is that this population risk depends on the parameters  $\boldsymbol{\theta}$  only through its empirical distribution, i.e.,  $\hat{\rho}_N(\boldsymbol{\theta}) = N^{-1} \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i}$  where  $\delta_{\boldsymbol{\theta}_i}$  is a point mass at  $\boldsymbol{\theta}_i$ . This motivates us to view  $R(\boldsymbol{\theta})$  equivalently as  $R(\hat{\rho}_N(\boldsymbol{\theta}))$  — a mapping from distributions to real numbers. With this observation in place, one can prove that running SGD for  $R(\cdot)$ —in a suitable scaling limit—results in a gradient flow on the space of distributions endowed with the Wasserstein metric that minimizes  $R(\cdot)$ . In addition, it turns out that the empirical distribution  $\hat{\rho}_N(\boldsymbol{\theta}^k)$  is well approximated by the gradient flow, as long as the neural net is over-parametrized (e.g.,  $N \gg d$ ) and the number of steps is not too large. In particular, Mei, Montanari and Nguyen (2018) have

shown that under certain regularity conditions,

$$\sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} \left| R(\hat{\rho}_N(\boldsymbol{\theta}^k)) - R(\rho_{k\varepsilon}) \right| \lesssim e^T \sqrt{\frac{1}{N}} \vee \varepsilon \cdot \sqrt{d + \log \frac{N}{\varepsilon}},$$

where  $\varepsilon > 0$  is an proxy for the step size of SGD and  $\rho_{k\varepsilon}$  is the distribution of the gradient flow at time  $k\varepsilon$ . In words, the out-of-sample error of  $\boldsymbol{\theta}^k$  generated by SGD is well-approximated by that of  $\rho_{k\varepsilon}$ , generated by the gradient flow. [Mei, Montanari and Nguyen \(2018\)](#) further demonstrated that in some simple settings, the out-of-sample error  $R(\rho_{k\varepsilon})$  of the distributional limit can be fully characterized. Passing the finite-dimensional SGD to its limit version — a mean field perspective, greatly simplifies the problem conceptually. Nevertheless, how well does  $R(\rho_{k\varepsilon})$  perform and how fast it converges remain largely open for general problems. Moreover, the current approach seems confined to the one-hidden-layer case and the extension to multi-layer scenario seems highly non-trivial.

**7.2.2 Stability.** A second way to understand the generalization ability of deep learning is through the *stability* of SGD. An algorithm is considered stable if a slight change of the input does not alter the output much. It has long been observed that a stable algorithm has a small generalization gap; examples include  $k$  nearest neighbors ([Rogers and Wagner, 1978](#); [Devroye and Wagner, 1979](#)), bagging ([Breiman, 1996](#); [Breiman et al., 1996](#)), etc. The precise connection between stability and generalization gap is stated by [Bousquet and Elisseeff \(2002\)](#); [Shalev-Shwartz et al. \(2010\)](#). In what follows, we formalize the idea of *stability* and its connection with the generalization gap. Let  $\mathcal{A}$  denote an algorithm (possibly randomized) which takes the training set  $S \triangleq \{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$  of size  $n$  and returns an estimated parameter  $\hat{\boldsymbol{\theta}} \triangleq \mathcal{A}(S)$ . Suppose coordinates of  $\mathbf{x}_i$  and  $y_i$  are bounded. Following [Hardt, Recht and Singer \(2015\)](#), we have the following definition for *stability*.

**DEFINITION 7.2.** *An algorithm (possibly randomized)  $\mathcal{A}$  is  $\varepsilon$ -uniformly stable with respect to the loss function  $\mathcal{L}(\cdot, \cdot)$  if for all datasets  $S, S'$  of size  $n$  which differ in at most one example, one has*

$$\sup_{\mathbf{x}, y} \mathbb{E}_{\mathcal{A}} [\mathcal{L}(f(\mathbf{x}; \mathcal{A}(S)), y) - \mathcal{L}(f(\mathbf{x}; \mathcal{A}(S')), y)] \leq \varepsilon.$$

Here the expectation is taken w.r.t. the randomness in the algorithm  $\mathcal{A}$  and  $\varepsilon$  might depend on  $n$ . The loss function  $\mathcal{L}(\cdot, \cdot)$  takes an example (say  $(\mathbf{x}, y)$ ) and the estimated parameter (say  $\mathcal{A}(S)$ ) as inputs and outputs a real value.

Surprisingly, an  $\varepsilon$ -uniformly stable algorithm incurs small generalization gap in expectation, which is stated in the following lemma.

**LEMMA 7.1** (Theorem 2.2 in [Hardt, Recht and Singer, 2015](#)). *Let  $\mathcal{A}$  be  $\varepsilon$ -uniformly stable. Then the expected generalization gap is no larger than  $\varepsilon$ , i.e.,*

$$(7.4) \quad \left| \mathbb{E}_{\mathcal{A}, S} \left[ \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(\mathbf{x}_i; \mathcal{A}(S)), y_i) - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathcal{L}(f(\mathbf{x}; \mathcal{A}(S)), y)] \right] \right| \leq \varepsilon.$$

With Lemma 7.1 in hand, it suffices to prove stability bound on specific algorithms. It turns out that SGD introduced in Section 6 is uniformly stable when

solving smooth nonconvex functions.

**THEOREM 7.3** (Theorem 3.12 in [Hardt, Recht and Singer \(2015\)](#)). *Assume that for any fixed  $(y, \mathbf{x})$ , the loss function  $\mathcal{L}(f(x; \boldsymbol{\theta}), y)$ , viewed as a function of  $\boldsymbol{\theta}$ , is  $L$ -Lipschitz and  $\beta$ -smooth. Consider running SGD on the empirical loss function with decaying step size  $\alpha_t \leq c/t$ , where  $c$  is some small absolute constant. Then SGD is uniformly stable with*

$$\varepsilon \lesssim \frac{T^{1-\frac{1}{\beta c+1}}}{n},$$

where we have ignored the dependency on  $\beta, c$  and  $L$ .

Theorem 7.3 reveals that SGD operating on nonconvex loss functions is indeed uniformly stable as long as the number of steps  $T$  is not large compared with  $n$ . This together with Lemma 7.1 demonstrates the generalization ability of SGD in expectation. Nevertheless, two important limitations are worth mentioning. First, Lemma 7.1 provides an upper bound on the out-of-sample error *in expectation*, but ideally, instead of an on-average guarantee under  $\mathbb{E}_{\mathcal{A}, S}$ , we would like to have a high probability guarantee as in the convex case ([Feldman and Vondrak, 2019](#)). Second, controlling the generalization gap alone is not enough to achieve a small out-of-sample error, since it is unclear whether SGD can achieve a small training error within  $T$  steps.

**7.2.3 Implicit regularization.** In the presence of over-parametrization (number of parameters larger than the sample size), conventional wisdom informs us that we should apply some regularization techniques (e.g.,  $\ell_1 / \ell_2$  regularization) so that the model will not overfit the data. However, in practice, neural networks without explicit regularization generalize well. This phenomenon motivates researchers to look at the regularization effects introduced by training algorithms (e.g., SGD) in this over-parametrized regime. While there might exist multiple, if not infinite global minima of the empirical loss (6.1), it is possible that practical algorithms tend to converge to solutions with better generalization powers.

Take the underdetermined linear system  $\mathbf{X}\boldsymbol{\theta} = \mathbf{y}$  as a starting point. Here  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and  $\boldsymbol{\theta} \in \mathbb{R}^p$  with  $p$  much larger than  $n$ . Running gradient descent on the loss  $\frac{1}{2}\|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  from the origin (i.e.,  $\boldsymbol{\theta}^0 = \mathbf{0}$ ) results in the solution with the minimum Euclidean norm, that is GD converges to

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \|\boldsymbol{\theta}\|_2 \quad \text{subject to} \quad \mathbf{X}\boldsymbol{\theta} = \mathbf{y}.$$

In words, without any  $\ell_2$  regularization in the loss function, gradient descent automatically finds the solution with the least  $\ell_2$  norm. This phenomenon, often called as *implicit regularization*, not only has been empirically observed in training neural networks, but also has been theoretically understood in some simplified cases, e.g., logistic regression with separable data. In logistic regression, given a training set  $\{(y_i, \mathbf{x}_i)\}_{1 \leq i \leq n}$  with  $\mathbf{x}_i \in \mathbb{R}^p$  and  $y_i \in \{1, -1\}$ , one aims to fit a logistic regression model by solving the following program:

$$(7.5) \quad \min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \ell(y_i \mathbf{x}_i^\top \boldsymbol{\theta}).$$

Here,  $\ell(u) \triangleq \log(1 + e^{-u})$  denotes the logistic loss. Further assume that the data

is separable, i.e., there exists  $\theta^* \in \mathbb{R}^p$  such that  $y_i \theta^{*\top} x_i > 0$  for all  $i$ . Under this condition, the loss function (7.5) can be arbitrarily close to zero for certain  $\theta$  with  $\|\theta\|_2 \rightarrow \infty$ . What happens when we minimize (7.5) using gradient descent? Soudry et al. (2018) uncovers a striking phenomenon.

**THEOREM 7.4** (Theorem 3 in Soudry et al., 2018). *Consider the logistic regression (7.5) with separable data. If we run GD*

$$\theta^{t+1} = \theta^t - \eta \frac{1}{n} \sum_{i=1}^n y_i x_i \ell'(y_i x_i^\top \theta^t)$$

*from any initialization  $\theta^0$  with an appropriate step size  $\eta > 0$ , then normalized  $\theta^t$  converges to a solution with the maximum  $\ell_2$  margin. That is,*

$$(7.6) \quad \lim_{t \rightarrow \infty} \frac{\theta^t}{\|\theta^t\|_2} = \hat{\theta},$$

*where  $\hat{\theta}$  is the solution to the hard margin support vector machine:*

$$(7.7) \quad \hat{\theta} \triangleq \arg \min_{\theta \in \mathbb{R}^p} \|\theta\|_2, \quad \text{subject to } y_i x_i^\top \theta \geq 1 \quad \text{for all } 1 \leq i \leq n.$$

The above theorem reveals that gradient descent, when solving logistic regression with separable data, implicitly regularizes the iterates towards the  $\ell_2$  max margin vector (cf. (7.6)), without any explicit regularization as in (7.7). Similar results have been obtained by Ji and Telgarsky (2018). In addition, Gunasekar et al. (2018a) studied algorithms other than gradient descent and showed that coordinate descent produces a solution with the maximum  $\ell_1$  margin.

Moving beyond logistic regression, which can be viewed as a neural net with no hidden layers, the theoretical understanding of implicit regularization in deeper neural networks is still limited; see Gunasekar et al. (2018b) for an illustration in deep linear convolutional neural networks.

## 8. DISCUSSION

Due to space limitations, we have omitted several important deep learning models; notable examples include deep reinforcement learning (Mnih et al., 2015), deep probabilistic graphical models (Salakhutdinov and Hinton, 2009), variational autoencoders (Kingma and Welling, 2013), transfer learning (Yosinski et al., 2014), etc. Apart from the modeling aspect, interesting theories on generative adversarial networks (Arora et al., 2017; Bai, Ma and Risteski, 2018), recurrent neural networks (Allen-Zhu and Li, 2019), connections with kernel methods (Jacot, Gabriel and Hongler, 2018; Arora et al., 2019) are also emerging. We have also omitted the inverse-problem view of deep learning where the data are assumed to be generated from a certain neural net and the goal is to recover the weights in the NN with as few examples as possible. Various algorithms (e.g., GD with spectral initialization) have been shown to recover the weights successfully in some simplified settings (Zhong et al., 2017; Soltanolkotabi, 2017; Goel, Klivans and Meka, 2018; Mondelli and Montanari, 2018; Chen et al., 2019a; Fu, Chi and Liang, 2018).

In the end, we identify a few important directions for future research.

- *New characterization of data distributions.* The success of deep learning relies on its power of efficiently representing complex functions relevant to real

data. Comparatively, classical methods often have optimal guarantee if a problem has a certain known structure, such as smoothness, sparsity, and low-rankness (Stone, 1982; Donoho and Johnstone, 1994; Candès and Tao, 2009; Chen et al., 2019b), but they are insufficient for complex data such as images. How to characterize the high-dimensional real data that can free us from known barriers, such as the curse of dimensionality, is an interesting open question.

- *Understanding various computational algorithms for deep learning.* As we have emphasized throughout this survey, computational algorithms (e.g., variants of SGD) play a vital role in the success of deep learning. They allow fast training of deep neural nets and probably contribute towards the good generalization behavior of deep learning in practice. Understanding these computational algorithms and devising better ones are crucial components in understanding deep learning.
- *Robustness.* It has been well documented that DNNs are sensitive to small adversarial perturbations that are indistinguishable to humans (Szegedy et al., 2013). This raises serious safety issues once deep learning models are deployed in applications such as self-driving cars, healthcare, etc. It is therefore crucial to refine current training practice to enhance robustness in a principled way (Singh, Murdoch and Yu, 2018).
- *Low signal-to-noise ratio (SNR).* Arguably, for image data and audio data where the signal-to-noise ratio is high, deep learning has achieved great success. In many other statistical problems, the SNR may be very low. For example, in financial applications, the firm characteristic and covariates may only explain a small part of the financial returns; in healthcare systems, the uncertainty of an illness may not be predicted well from a patient’s medical history. How to adapt deep learning models to excel at such tasks is an interesting direction to pursue.

## ACKNOWLEDGEMENTS

We thank Ruying Bao, Yuxin Chen, Chenxi Liu, Qingcan Wang and Pengkun Yang for helpful comments and discussions. We also thank the editor and the reviewers for pointing out references on universal approximation theorems and AdaNet.

## REFERENCES

- ABADI, M. and ET. AL. (2015). TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. Software available from tensorflow.org.
- ABBASI-ASL, R., CHEN, Y., BLONIAZ, A., OLIVER, M., WILLMORE, B. D., GALLANT, J. L. and YU, B. (2018). The DeepTune framework for modeling and characterizing neurons in visual cortex area V4. *bioRxiv* 465534.
- ALLEN-ZHU, Z., LI, Y. and SONG, Z. (2018). A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*.
- ALLEN-ZHU, Z. and LI, Y. (2019). Can SGD Learn Recurrent Neural Networks with Provable Generalization? *ArXiv e-prints* **abs/1902.01028**.
- ANTHONY, M. and BARTLETT, P. L. (2009). *Neural network learning: Theoretical foundations*. cambridge university press.
- ARJOVSKY, M., CHINTALA, S. and BOTTOU, L. (2017). Wasserstein Generative Adversarial Networks. **70** 214–223.
- ARNOLD, V. I. (2009). On functions of three variables. *Collected Works: Representations of Functions, Celestial Mechanics and KAM Theory, 1957–1965* 5–8.

- ARORA, S., GE, R., LIANG, Y., MA, T. and ZHANG, Y. (2017). Generalization and equilibrium in generative adversarial nets (GANs). In *Proceedings of the 34th International Conference on Machine Learning-Volume 70* 224–232. JMLR. org.
- ARORA, S., DU, S. S., HU, W., LI, Z. and WANG, R. (2019). Fine-Grained Analysis of Optimization and Generalization for Overparameterized Two-Layer Neural Networks. *arXiv preprint arXiv:1901.08584*.
- BAI, Y., MA, T. and RISTESKI, A. (2018). Approximability of discriminators implies diversity in GANs. *arXiv preprint arXiv:1806.10586*.
- BARRON, A. R. (1993). Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information theory* **39** 930–945.
- BARTLETT, P. L. (1998). The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network. *IEEE transactions on Information Theory* **44** 525–536.
- BARTLETT, P. L., FOSTER, D. J. and TELGARSKY, M. J. (2017). Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems 30* (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett, eds.) 6240–6249. Curran Associates, Inc.
- BAUER, B., KOHLER, M. et al. (2019). On deep learning as a remedy for the curse of dimensionality in nonparametric regression. *The Annals of Statistics* **47** 2261–2285.
- BELKIN, M., HSU, D., MA, S. and MANDAL, S. (2019). Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences* **116** 15849–15854.
- BOTTOU, L. (1998). Online learning and stochastic approximations. *On-line learning in neural networks* **17** 142.
- BOUSQUET, O. and ELISSEEFF, A. (2002). Stability and generalization. *Journal of machine learning research* **2** 499–526.
- BREIMAN, L. (1996). Bagging predictors. *Machine learning* **24** 123–140.
- BREIMAN, L. et al. (1996). Heuristics of instability and stabilization in model selection. *The annals of statistics* **24** 2350–2383.
- CANDÈS, E. J. and TAO, T. (2009). The power of convex relaxation: Near-optimal matrix completion. *arXiv preprint arXiv:0903.1476*.
- CAO, C., LIU, F., TAN, H., SONG, D., SHU, W., LI, W., ZHOU, Y., BO, X. and XIE, Z. (2018). Deep learning and its applications in biomedicine. *Genomics, proteomics & bioinformatics* **16** 17–32.
- CHEN, S., DOBRIBAN, E. and LEE, J. H. (2019). Invariance reduces variance: Understanding data augmentation in deep learning and beyond. *arXiv preprint arXiv:1907.10905*.
- CHEN, T., RUBANOVA, Y., BETTENCOURT, J. and DUVENAUD, D. (2018). Neural Ordinary Differential Equations. *arXiv preprint arXiv:1806.07366*.
- CHEN, Y., CHI, Y., FAN, J. and MA, C. (2019a). Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Mathematical Programming* 1–33.
- CHEN, Y., CHI, Y., FAN, J., MA, C. and YAN, Y. (2019b). Noisy Matrix Completion: Understanding Statistical Guarantees for Convex Relaxation via Nonconvex Optimization. *arXiv preprint arXiv:1902.07698*.
- CHIZAT, L. and BACH, F. (2018). On the global convergence of gradient descent for overparameterized models using optimal transport. In *Advances in neural information processing systems* 3040–3050.
- CHO, K., VAN MERRIËNBOER, B., GULCEHRE, C., BAHDANAU, D., BOUGARES, F., SCHWENK, H. and BENGIO, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.
- CORTES, C., GONZALVO, X., KUZNETSOV, V., MOHRI, M. and YANG, S. (2017). Adanet: Adaptive structural learning of artificial neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70* 874–883. JMLR. org.
- DE FAUW, J., LEDSAM, J. R., ROMERA-PAREDES, B., NIKOLOV, S., TOMASEV, N., BLACKWELL, S., ASKHAM, H., GLOROT, X., O’DONOGHUE, B., VISENTIN, D. et al. (2018). Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature medicine* **24** 1342.
- DEVROYE, L. and WAGNER, T. (1979). Distribution-free performance bounds for potential function rules. *IEEE Transactions on Information Theory* **25** 601–604.
- DONOHU, D. L. (2000). High-dimensional data analysis: The curses and blessings of dimension-

- ality. *AMS math challenges lecture* **1** 32.
- DONOHO, D. L. and JOHNSTONE, J. M. (1994). Ideal spatial adaptation by wavelet shrinkage. *biometrika* **81** 425–455.
- DU, S. S., LEE, J. D., LI, H., WANG, L. and ZHAI, X. (2018). Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*.
- DUCHI, J., HAZAN, E. and SINGER, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* **12** 2121–2159.
- E, W., MA, C. and WANG, Q. (2019). A Priori Estimates of the Population Risk for Residual Networks. *arXiv preprint arXiv:1903.02154*.
- ELDAN, R. and SHAMIR, O. (2016). The power of depth for feedforward neural networks. In *Conference on Learning Theory* 907–940.
- FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association* **96** 1348–1360.
- FELDMAN, V. and VONDRAK, J. (2019). High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. *arXiv preprint arXiv:1902.10710*.
- FRIEDMAN, J. H. and STUETZLE, W. (1981). Projection pursuit regression. *Journal of the American statistical Association* **76** 817–823.
- FU, H., CHI, Y. and LIANG, Y. (2018). Local geometry of one-hidden-layer neural networks for logistic regression. *arXiv preprint arXiv:1802.06463*.
- FUKUSHIMA, K. and MIYAKE, S. (1982). Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition. In *Competition and cooperation in neural nets* 267–285. Springer.
- GAO, C., LIU, J., YAO, Y. and ZHU, W. (2018). Robust Estimation and Generative Adversarial Nets. *arXiv preprint arXiv:1810.02030*.
- GOEL, S., KLIVANS, A. and MEKA, R. (2018). Learning one convolutional layer with overlapping patches. *arXiv preprint arXiv:1802.02547*.
- GOLOWICH, N., RAKHLIN, A. and SHAMIR, O. (2017). Size-independent sample complexity of neural networks. *arXiv preprint arXiv:1712.06541*.
- GOLUB, G. H. and VAN LOAN, C. F. (2013). *Matrix computations*, 4 ed. JHU Press.
- GOODFELLOW, I., BENGIO, Y. and COURVILLE, A. (2016). *Deep Learning*. MIT Press.
- GOODFELLOW, I., POUGET-ABADIE, J., MIRZA, M., XU, B., WARDE-FARLEY, D., OZAI, S., COURVILLE, A. and BENGIO, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems* 2672–2680.
- GUNASEKAR, S., LEE, J., SOUDRY, D. and SREBRO, N. (2018a). Characterizing implicit bias in terms of optimization geometry. *arXiv preprint arXiv:1802.08246*.
- GUNASEKAR, S., LEE, J. D., SOUDRY, D. and SREBRO, N. (2018b). Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems* 9482–9491.
- GYBENKO, G. (1989). Approximation by superposition of sigmoidal functions. *Mathematics of Control, Signals and Systems* **2** 303–314.
- HARDLE, W., HALL, P., ICHIMURA, H. et al. (1993). Optimal smoothing in single-index models. *The annals of Statistics* **21** 157–178.
- HARDT, M., RECHT, B. and SINGER, Y. (2015). Train faster, generalize better: Stability of stochastic gradient descent. *arXiv preprint arXiv:1509.01240*.
- HASTIE, T., MONTANARI, A., ROSSET, S. and TIBSHIRANI, R. J. (2019). Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*.
- HE, K., ZHANG, X., REN, S. and SUN, J. (2016a). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* 770–778.
- HE, K., ZHANG, X., REN, S. and SUN, J. (2016b). Identity mappings in deep residual networks. In *European conference on computer vision* 630–645. Springer.
- HINTON, G., SRIVASTAVA, N. and SWERSKY, K. (2012). Neural networks for machine learning lecture 6a overview of mini-batch gradient descent.
- HINTON, G. E., SRIVASTAVA, N., KRIZHEVSKY, A., SUTSKEVER, I. and SALAKHUTDINOV, R. R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*.
- HOCHREITER, S. and SCHMIDHUBER, J. (1997). Long short-term memory. *Neural computation* **9** 1735–1780.
- HORNIK, K. (1991). Approximation capabilities of multilayer feedforward networks. *Neural networks* **4** 251–257.



- HUANG, G., LIU, Z., VAN DER MAATEN, L. and WEINBERGER, K. Q. (2017). Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition* 4700–4708.
- HUBEL, D. H. and WIESEL, T. N. (1962). Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex. *The Journal of physiology* **160** 106–154.
- IANDOLA, F. N., HAN, S., MOSKEWICZ, M. W., ASHRAF, K., DALLY, W. J. and KEUTZER, K. (2016). SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size. *arXiv preprint arXiv:1602.07360*.
- IOFFE, S. and SZEGEDY, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- JACOT, A., GABRIEL, F. and HONGLER, C. (2018). Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems* 8580–8589.
- JAIN, P., KAKADE, S. M., KIDAMBI, R., NETRAPALLI, P. and SIDFORD, A. (2017). Accelerating stochastic gradient descent. *arXiv preprint arXiv:1704.08227*.
- JAVANMARD, A., MONDELLI, M. and MONTANARI, A. (2019). Analysis of a Two-Layer Neural Network via Displacement Convexity. *arXiv preprint arXiv:1901.01375*.
- Ji, Z. and TELGARSKY, M. (2018). Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*.
- KIDAMBI, R., NETRAPALLI, P., JAIN, P. and KAKADE, S. (2018). On the insufficiency of existing momentum schemes for stochastic optimization. In *2018 Information Theory and Applications Workshop (ITA)* 1–9. IEEE.
- KIEFER, J., WOLFOWITZ, J. et al. (1952). Stochastic estimation of the maximum of a regression function. *The Annals of Mathematical Statistics* **23** 462–466.
- KINGMA, D. P. and BA, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- KINGMA, D. P. and WELLING, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- KLUSOWSKI, J. M. and BARRON, A. R. (2016). Risk bounds for high-dimensional ridge function combinations including neural networks. *arXiv preprint arXiv:1607.01434*.
- KRIZHEVSKY, A., SUTSKEVER, I. and HINTON, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* 1097–1105.
- KUSHNER, H. and YIN, G. G. (2003). *Stochastic approximation and recursive algorithms and applications* **35**. Springer Science & Business Media.
- LECUN, Y., BENGIO, Y. and HINTON, G. (2015). Deep learning. *nature* **521** 436.
- LECUN, Y., BOTTOU, L., BENGIO, Y. and HAFFNER, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86** 2278–2324.
- LI, Y., SWERSKY, K. and ZEMEL, R. (2015). Generative moment matching networks. In *International Conference on Machine Learning* 1718–1727.
- LI, H., XU, Z., TAYLOR, G., STUDER, C. and GOLDSTEIN, T. (2018a). Visualizing the loss landscape of neural nets. In *Advances in Neural Information Processing Systems* 6391–6401.
- LI, X., LU, J., WANG, Z., HAUPT, J. and ZHAO, T. (2018b). On tighter generalization bound for deep neural networks: CNNs, ResNets, and beyond. *arXiv preprint arXiv:1806.05159*.
- LIANG, T. (2017). How Well Can Generative Adversarial Networks (GAN) Learn Densities: A Nonparametric View. *arXiv preprint arXiv:1712.08244*.
- LIN, M., CHEN, Q. and YAN, S. (2013). Network in network. *arXiv preprint arXiv:1312.4400*.
- LIN, H. W., TEGMARK, M. and ROLNICK, D. (2017). Why does deep and cheap learning work so well? *Journal of Statistical Physics* **168** 1223–1247.
- MAAS, A. L., HANNUN, A. Y. and NG, A. Y. (2013). Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml* **30** 3.
- MAIOROV, V. and MEIR, R. (2000). On the near optimality of the stochastic approximation of smooth functions by neural networks. *Advances in Computational Mathematics* **13** 79–103.
- MAKOVZ, Y. (1996). Random approximants and neural networks. *Journal of Approximation Theory* **85** 98–109.
- MEI, S., MISIAKIEWICZ, T. and MONTANARI, A. (2019). Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. *arXiv preprint arXiv:1902.06015*.
- MEI, S., MONTANARI, A. and NGUYEN, P.-M. (2018). A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences* **115** E7665–

- E7671.
- MHASKAR, H. N. (1996). Neural networks for optimal approximation of smooth and analytic functions. *Neural computation* **8** 164–177.
- MHASKAR, H., LIAO, Q. and POGGIO, T. (2016). Learning functions: when is deep better than shallow. *arXiv preprint arXiv:1603.00988*.
- MNIH, V., KAVUKCUOGLU, K., SILVER, D., RUSU, A. A., VENESS, J., BELLEMARE, M. G., GRAVES, A., RIEDMILLER, M., FIDJELAND, A. K., OSTROVSKI, G. et al. (2015). Human-level control through deep reinforcement learning. *Nature* **518** 529.
- MONDELLI, M. and MONTANARI, A. (2018). On the connection between learning two-layers neural networks and tensor decomposition. *arXiv preprint arXiv:1802.07301*.
- NESTEROV, Y. E. (1983). A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ . In *Dokl. Akad. Nauk SSSR* **269** 543–547.
- NEYSHABUR, B., TOMIOKA, R. and SREBRO, N. (2015). Norm-based capacity control in neural networks. In *Conference on Learning Theory* 1376–1401.
- NOWOZIN, S., CSEKE, B. and TOMIOKA, R. (2016). f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems* 271–279.
- PARBERRY, I. (1994). *Circuit complexity and neural networks*. MIT press.
- PASZKE, A., GROSS, S., CHINTALA, S., CHANAN, G., YANG, E., DEVITO, Z., LIN, Z., DESMAISON, A., ANTIGA, L. and LERER, A. (2017). Automatic differentiation in PyTorch.
- PINKUS, A. (1999). Approximation theory of the MLP model in neural networks. *Acta numerica* **8** 143–195.
- POGGIO, T., MHASKAR, H., ROSASCO, L., MIRANDA, B. and LIAO, Q. (2017). Why and when can deep-but not shallow-networks avoid the curse of dimensionality: a review. *International Journal of Automation and Computing* **14** 503–519.
- POLYAK, B. T. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics* **4** 1–17.
- POLYAK, B. T. and JUDITSKY, A. B. (1992). Acceleration of stochastic approximation by averaging. *SIAM Journal on Control and Optimization* **30** 838–855.
- POLYAK, B. T. and TSYPKIN, Y. Z. (1979). Adaptive estimation algorithms: convergence, optimality, stability. *Avtomatika i Telemekhanika* **3** 71–84.
- POULTNEY, C., CHOPRA, S., LECUN, Y. et al. (2007). Efficient learning of sparse representations with an energy-based model. In *Advances in neural information processing systems* 1137–1144.
- REDDI, S. J., KALE, S. and KUMAR, S. (2018). On the convergence of adam and beyond.
- ROBBINS, H. and MONRO, S. (1951). A Stochastic Approximation Method. *The Annals of Mathematical Statistics* **22** 400–407.
- ROGERS, W. H. and WAGNER, T. J. (1978). A finite sample distribution-free performance bound for local discrimination rules. *The Annals of Statistics* 506–514.
- ROLNICK, D. and TEGMARK, M. (2017). The power of deeper networks for expressing natural functions. *arXiv preprint arXiv:1705.05502*.
- ROMANO, Y., SESIA, M. and CANDÈS, E. J. (2018). Deep Knockoffs. *arXiv preprint arXiv:1811.06687*.
- ROTSKOFF, G. M. and VANDEN-EIJNDEN, E. (2018). Neural networks as interacting particle systems: Asymptotic convexity of the loss landscape and universal scaling of the approximation error. *arXiv preprint arXiv:1805.00915*.
- RUMELHART, D. E., HINTON, G. E. and WILLIAMS, R. J. (1985). Learning internal representations by error propagation Technical Report, California Univ San Diego La Jolla Inst for Cognitive Science.
- RUSSAKOVSKY, O., DENG, J., SU, H., KRAUSE, J., SATHEESH, S., MA, S., HUANG, Z., KARPATHY, A., KHOSLA, A., BERNSTEIN, M., BERG, A. C. and FEI-FEI, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* **115** 211–252.
- SAK, H., SENIOR, A. and BEAUFAYS, F. (2014). Long short-term memory recurrent neural network architectures for large scale acoustic modeling. In *Fifteenth annual conference of the international speech communication association*.
- SALAKHUTDINOV, R. and HINTON, G. (2009). Deep boltzmann machines. In *Artificial intelligence and statistics* 448–455.
- SALIMANS, T., GOODFELLOW, I., ZAREMBA, W., CHEUNG, V., RADFORD, A. and CHEN, X.

- (2016). Improved techniques for training GANs. In *Advances in Neural Information Processing Systems* 2234–2242.
- SCHMIDT-HIEBER, J. (2017). Nonparametric regression using deep neural networks with ReLU activation function. *arXiv preprint arXiv:1708.06633*.
- SHALEV-SHWARTZ, S. and BEN-DAVID, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge university press.
- SHALEV-SHWARTZ, S., SHAMIR, O., SREBRO, N. and SRIDHARAN, K. (2010). Learnability, stability and uniform convergence. *Journal of Machine Learning Research* **11** 2635–2670.
- SILVER, D., SCHRITTWIESER, J., SIMONYAN, K., ANTONOGLOU, I., HUANG, A., GUEZ, A., HUBERT, T., BAKER, L., LAI, M., BOLTON, A. et al. (2017). Mastering the game of go without human knowledge. *Nature* **550** 354.
- SILVERMAN, B. W. (1998). *Density estimation for statistics and data analysis*. Chapman & Hall, CRC.
- SINGH, C., MURDOCH, W. J. and YU, B. (2018). Hierarchical interpretations for neural network predictions. *arXiv preprint arXiv:1806.05337*.
- SIRIGNANO, J. and SPILIOPOULOS, K. (2018). Mean field analysis of neural networks. *arXiv preprint arXiv:1805.01053*.
- SOLTANOLKOTABI, M. (2017). Learning relus via gradient descent. In *Advances in Neural Information Processing Systems* 2007–2017.
- SOUDRY, D., HOFFER, E., NACSON, M. S., GUNASEKAR, S. and SREBRO, N. (2018). The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research* **19** 2822–2878.
- SPRECHER, D. A. (1965). On the structure of continuous functions of several variables. *Transactions of the American Mathematical Society* **115** 340–355.
- STONE, C. J. (1982). Optimal global rates of convergence for nonparametric regression. *The annals of statistics* 1040–1053.
- STONE, C. J. et al. (1985). Additive regression and other nonparametric models. *The annals of Statistics* **13** 689–705.
- STONE, C. J. et al. (1994). The use of polynomial splines and their tensor products in multivariate function estimation. *The Annals of Statistics* **22** 118–171.
- SUTSKEVER, I., MARTENS, J., DAHL, G. and HINTON, G. (2013). On the importance of initialization and momentum in deep learning. In *International conference on machine learning* 1139–1147.
- SZEGEDY, C., ZAREMBA, W., SUTSKEVER, I., BRUNA, J., ERHAN, D., GOODFELLOW, I. and FERGUS, R. (2013). Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- SZEGEDY, C., LIU, W., JIA, Y., SERMANET, P., REED, S., ANGUELOV, D., ERHAN, D., VANHOUCHE, V. and RABINOVICH, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* 1–9.
- TELGARSKY, M. (2016). Benefits of depth in neural networks. *arXiv preprint arXiv:1602.04485*.
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)* **58** 267–288.
- VAPNIK, V. and CHERVONENKIS, A. Y. (1971). On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Theory of Probability & Its Applications* **16** 264–280.
- VINCENT, P., LAROCHELLE, H., BENGIO, Y. and MANZAGOL, P.-A. (2008). Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning* 1096–1103. ACM.
- WAGER, S., WANG, S. and LIANG, P. S. (2013). Dropout training as adaptive regularization. In *Advances in neural information processing systems* 351–359.
- WEINAN, E., HAN, J. and JENTZEN, A. (2017). Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations. *Communications in Mathematics and Statistics* **5** 349–380.
- WILSON, A. C., ROELOFS, R., STERN, M., SREBRO, N. and RECHT, B. (2017). The Marginal Value of Adaptive Gradient Methods in Machine Learning. In *Advances in Neural Information Processing Systems 30* (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett, eds.) 4148–4158. Curran Associates, Inc.
- WU, Y., SCHUSTER, M., CHEN, Z., LE, Q. V., NOROUZI, M., MACHEREY, W., KRIKUN, M., CAO, Y., GAO, Q., MACHEREY, K. et al. (2016). Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.

- YOSINSKI, J., CLUNE, J., BENGIO, Y. and LIPSON, H. (2014). How transferable are features in deep neural networks? In *Advances in neural information processing systems* 3320–3328.
- YOSINSKI, J., CLUNE, J., NGUYEN, A., FUCHS, T. and LIPSON, H. (2015). Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*.
- ZHANG, C., BENGIO, S., HARDT, M., RECHT, B. and VINYALS, O. (2016). Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*.
- ZHONG, K., SONG, Z., JAIN, P., BARTLETT, P. L. and DHILLON, I. S. (2017). Recovery guarantees for one-hidden-layer neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70* 4140–4149. JMLR. org.