# 7

## Representation and GAN

## 7.1 Reviewing

1. When do neural-nets have bad local-minima?

   - Classical results in (Auer *et al.*, 1996) show that for unrealizable case, non-linear activation can easily create bad local minima;

   - Recent results in (Li *et al.*, 2018) show that even for over-parameterized case (thus also realizable), the bad local minima do exist for a class of non-linear activations.

2. People claims that "over-parametrization" smooths the landscape, any rigorous result in this claim?

   There are a few rigorous results for this claim. For instance, (Li *et al.*, 2018) shows that the loss function of over-parametrized networks has no bad basin (or more precisely, it is a "weakly global" function); for many smooth enough activations, over-paramterized networks satisfy a stronger geometrical property $PT$ (i.e. at any point, after a tiny generic perturbation, there is a strictly decreasing path from the perturbed point to a global minimum).

3. How to empirically check nice or bad landscape for any problem (continuous optimization)?

   Check values along paths connecting interesting points

4. When does the neural-nets have enough representation power?

   The most important factor is the activation function. It should be bounded and non-constant; although ReLU is not bounded, it also makes the neural-nets have enough representation power; but linear/quadratic activation does not.

## 7.2   Representation: depth separation

### 7.2.1   A simple proof of threhold activation has enough representation power

Consider the dimension $d = 1$ first. The non-linear activation is $\phi(t) = 1\{t \geq 0\}$. It suffices to show that

$$\overline{\text{span}\{\phi(at + b)\}} = \mathcal{C}(\mathcal{D})$$

for any compact domain $\mathcal{D} \in \mathbb{R}$. Define the pulse function $\psi(t) = 1\{0 \leq t < 1\}$, which can be expressed as $\psi = \phi(t) - \phi(1 - t)$. It suffices to use the pulse function to approximate any continuous function. The general idea is shown in the Figure. 7.1.
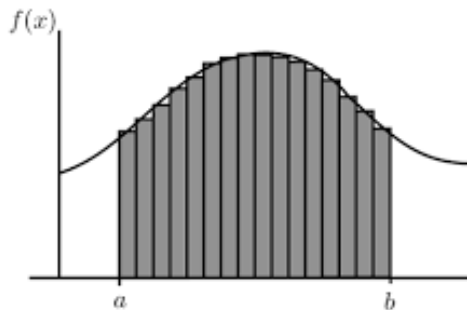


**Figure 7.1:** The pulse function can approximate any continuous function.

We can see that the insights are very similar to those in Riemann integration, and the proof follows the similar idea as well (Wang, 2019a, Theorem (6.4)).

*Proof.* Suppose that $\mathcal{D} = [0, 1]$ and our target function $f^* \in \mathcal{C}[0, 1]$. The continuity of $f^*$ together with the compactness of $\mathcal{D}$ implies that the function $f^*$ is uniformly continuous, i.e., $\forall \varepsilon > 0$, there exists $\delta$ such that for any $|x - x^*| < \delta$,

$$|f - f^*| < \varepsilon.$$

Pick a partition

$$\mathcal{P} = \{b_0 := 0, b_1 = h, b_2 = 2h, \ldots, b_K := Kh := 1\}, \quad \text{with} \quad \frac{1}{K} < \Delta.$$

Therefore, define the approximation function

$$f(x) = \sum_{i=1}^{K} a_i \psi \left( \frac{x - b_i}{b_{i+1} - b_i} \right) \in \text{span}\{\phi(at + b)\},$$

where $a_i \triangleq f^*(b_i), i = 0, \ldots, K - 1$. Then it's easy to verify that $|f(x) - f^*(x)| < \varepsilon$ for any $x \in [0, 1]$. $\qquad \square$

**Remark 7.1.** The first step in the proof explains why we need to define the *compact* domain.

**Bibliography** Then we discuss the representation power for other kinds of activations. For sigmoid function $\phi(t) = \frac{1}{1+e^{-t}}$, it suffices to show that it can approximate the threshold function very well; for other types of functions such as switch function, some techniques from function analysis are needed. The paper (Cybenko, 1989) shows that the sigmoidal-type activation has enough representation power by using arguments from real analysis; the paper (Barron, 1994) further gives an mean integrated squared error between the estimated network and a target function $f$, in terms of number of neurons and the input dimension; the Kolmogorov–Arnold representation theorem actually has solved this problem by using that every multivariate continuous function can be represented as a superposition of continuous functions of one variable,

which is also related to Hilbert's thirteenth problem. The proof in this representation theorem contains the multi-resolution idea, and VCG/Receptor has the similar idea.

## 7.2.2   Depth Separation (Analysis for ReLU Activation)

It's a common belief that *deep* neural network usually gains better performance. We want to analysis this claim from the perspective of representation power. To show the power of depth, one way is to construct a function represented "*deep*"-net, then show this function is difficult to be represented by shallow networks.

1. Consider a function $\psi$ frequently studied in the dynamical systems literature, which can be represented with the ReLU activation $\phi$:

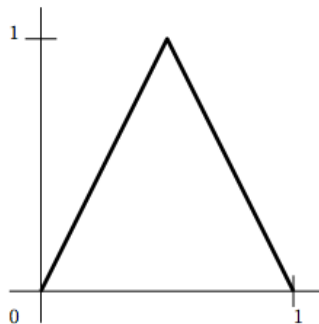$$\psi(x) = \phi(2\phi(x) - 4\phi(x - 0.5)).$$



**Figure 7.2:** The function $\psi$, which has one "peak"

2. Construct a function $f^*(x) = \psi^{(L)}(x)$, a the composition of $L$ $\psi$ functions, which has $2^{L-1}$ peaks. See $\psi^{(2)}$ for instance:
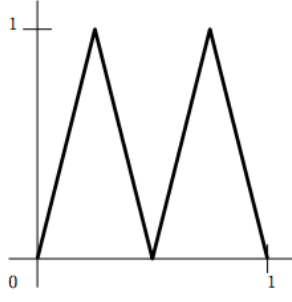
**Figure 7.3:** The function $\psi^{(2)}$, which has two "peaks"

It suffices to show $f^*(x) = \psi^{(L)}(x)$ can be represented by deep neural-nets, but it is difficult to be represented by a shallow network, i.e., we need $\mathcal{O}(2^L)$ neurons of a shallow network for representation. The intuition is that the depth (i.e., function composition) increases oscillation exponentially; while width (i.e., linear combination) increases oscillation linearly.

**Definition 7.1.** We say that $f$ is $K$-sawtooth if $f$ is piecewise affine with $K$ pieces. For example, the ReLU function is 2-sawtooth, and $\psi$ is 4-sawtooth.

We can show that function composition is stronger to produce more "sawtooth" than addition:

**Lemma 7.1.** If $f$ is $a$-sawtooth, $g$ is $b$-sawtooth, then $f + g$ is at most $(a + b)$-sawtooth, $f \circ g$ is at most $ab$-sawtooth.

By using this lemma, we can show the converse error bound on the representation power of shallow network:

**Theorem 7.2.** Given an underlying function $F$ and data points $\{(x_i, y_i \triangleq F(x_i))\}_{i=1}^n, y_i \in \{0, 1\}$, define the classification error for the approximation function $f$:

$$R(f) = \frac{1}{n} \sum_{i=1}^{n} 1\left\{\text{sign}(f(x_i) - 1/2) \neq y_i\right\}.$$

Construct the data points $x_i = \frac{i}{2^{L^*}}$, $y_i = \psi^{(L^*)}(x_i)$, $i = 1, \ldots, 2^{L^*}$. As a result, $y = (0, 1, 0, 1, \ldots)$, and $F(x) = \psi^{(L^*)}(x)$. If a ReLU neural network $f$ has $L$ layers, and width $m < 2^{(L^*-k)/L-1}$, then $R(f) > \frac{1}{2} - \frac{1}{3}\frac{1}{2^{k-1}}$.

**Corollary 7.3.** If $f$ has $L$ layers with width $m < 2^{(L^*-1)/L-1}$, then the error function is lower bounded by a constant: $R(f) > 1/6$.

**Corollary 7.4.** If $f$ has no more than $\sqrt{L^*}$ layers, then we need at least $m > 2^{\mathcal{O}(\sqrt{L^*})}$ neurons to get the error less than $1/6$.

**Remark 7.2.** There is an Implicit assumption on Theorem (7.2), i.e., the neural network is fully connected feedforward. It does not apply to ResNet and RNN. The paper (Lin and Jegelka, 2018) shows the representation power for ResNet.

**Remark 7.3.** The representation power on other kinds of neural-nets is a popular problem, such as graph neural-nets and meta-learning.

**Summarization**   There are three criteria for the performance of neural-nets:

$$\begin{cases} \text{Representation Error} \\ \text{Optimization Error} \\ \text{Generalization Error} \end{cases}$$

The important factors for the success of neural-nets are as follows:

1. The depth of neural-nets relates to the representation error;

2. The width of neural-nets relates to the landscape of neural-nets, which further influences the optimization error;

3. The initialization and normalization techniques relate to the convergence performance of optimization;

4. The architecture design influences the representation error; and the optimization error;

5. The SGD algorithm influences the speed for convergence during the optimization process, and people believe that it also tends to give a solution with low generalization error

**Remark 7.4.** Why the over-parametrization of neural-nets usually do not lead to over-fitting? Prof. Ruoyu Sun gives his understanding of this question. Consider true data $\{(x_i, y_i)\}_{i=1}^n$ generated by $f^*(x_i) = y_i$. We want to approximate $f^*$ with $f$, by using these $n$ data points. Let $W^*$ denote the representation-power threshold, and $n^*$ denotes the threshold for the number of data points, under which the approximation is likely to be bad. The number of parameters of $f$ is more than $W^*$ will not cause over-fitting, but when $n < n^*$, it is likely to cause over-fitting.

## 7.3  GAN

Now we turn from supervised learning to unsupervised learning. Prof. Ruoyu Sun will give basic formulation about GAN this lecture, and Prof. Mingyi Hong will provide the introduction to Adversarial Attack & Defense in the next lecture.

### Motivation

> Richard Feymann: What I cannot create, I do not understand

We wish to learn the data distribution $\mathbb{P}_d$, e.g., a style of writing of articles. In order to do so, we build a generative model which generates $\mathbb{P}_g$, e.g., imitates writing articles; and a classifier which judges whether the received sample comes from $\mathbb{P}_d$ or $\mathbb{P}_g$, e.g., give comments to the written samples. Finally, we want $\mathbb{P}_g \approx \mathbb{P}_d$, e.g., the generated article has a similar style of the original one.

We are interested in solving the optimization problem

$$\min_{\mathbb{P}_g} \Phi(\mathbb{P}_d, \mathbb{P}_g)$$

The question is that what distance measure should the $\Phi$ be? Statiscians tend to pick $\Phi(P, Q) = \text{KL}(P, Q)$ or $\Phi(P, Q) = \text{JS}(P, Q)$ empirically. However, it is not clear whether these distance metrics are good metrics. An ideal metric should satisfy the following property: if two images are from the same class, then their distance is small; if they are from different classes, then their distance is large. Is the JS distance between

the images of two cats is smaller than the distance between a cat and a dog? This is not clear. The major issue here is that common distances may not capture the right representation of the images. To explain the solution, next, we use an example of fake paintings.

**Discussion**  Suppose we want to generate an appropriate painting, denoted as $X$. Given an artist paint something, denoted as $\hat{X}$; and hire a critic to judge whether it is good or bad. We use $D(x)$ to represent the probability that the input $x$ is thought by the critic to come from the data $P_d$ rather than $P_g$. The evaluation score can be modeled as

$$L^{\text{GAN}}(\mathbb{P}_d, D) = \mathbb{E}_{x \sim \mathbb{P}_d}[\log D(x)] + \mathbb{E}_{\hat{x} \sim \mathbb{P}_g}[\log(1 - D(\hat{x}))]$$

Pick the *best* critic, i.e., the most strict critic, the distance measure is

$$\Phi(\mathbb{P}_d, \mathbb{P}_g) = \max_D L^{\text{GAN}}(\mathbb{P}_d, D)$$

Therefore, the optimization for GAN is a minimax problem:

$$\min_{\mathbb{P}_g} \max_D \mathbb{E}_{x \sim \mathbb{P}_d}[\log D(x)] + \mathbb{E}_{\hat{x} \sim \mathbb{P}_g}[\log(1 - D(\hat{x}))] \qquad (7.1)$$

When proposing a new mode, the sanity-check is needed, i.e., ensure that the global optimum equals whatever we want, i.e., the optimal solution $\mathbb{P}_g^* = \mathbb{P}_d$.

**Theorem 7.5** ((Goodfellow *et al.*, 2014)).  The global minimum of the problem (7.1) is achieved if and only if $\mathbb{P}_g^* = \mathbb{P}_d$. Moreover, this optimization problem is equivalent to minimizing the Jensen-Shannon divergence

$$\Phi(\mathbb{P}_g, \mathbb{P}_d) = -\log 4 + 2\text{JSD}(\mathbb{P}_d \| \mathbb{P}_g).$$

*Proof.*  Finding the global minima of the problem consists of two steps: first, we need to specify the range of the objective function; second, we need to identify some points that achieve the extreme of the range.

- Question 1: What is the range of the objectvie function?

  We find that $D(x) \in (0, 1)$, and therefore $L^{GAN} \in (-\infty, 0)$. It seems that it is meanless to solve an optimization problem with negative infinite value. In fact, the objective is lower bounded since the maximum criteria help.

- Question 2: Check when does the objective achieve the optimum.

  Consider the finite support distribution for simplicity. Denote the pmf from $\mathbb{P}_d$ and $\mathbb{P}_g$ as $\{q_1, \ldots, q_n\}$, $\{p_1, \ldots, p_n\}$, respectively. It suffices to solve

  $$\min_{p \in \mathcal{P}^n} \max_{d_i \in (0,1)} \quad \sum_{i=1}^n q_i \log d_i + \sum_{i=1}^n p_i \log(1 - d_i)$$
  with
  $$\mathcal{P}^n = \{p \mid \textstyle\sum_i p_i = 1, p_i \geq 0\}$$

  Consider the maximum optimization first:

  $$\max_{d_i \in (0,1)} \sum_{i=1}^n q_i \log d_i + \sum_{i=1}^n p_i \log(1 - d_i)$$

  It is decomposable in terms of $i$. For each single problem $\max_{d_i} q_i \log d_i + p_i \log(1 - d_i)$, we find the optimal solution is $d_i = \frac{q_i}{q_i + p_i}$. Substituting this solution into $\Phi(\mathbb{P}_g, \mathbb{P}_d)$, we imply

  $$\Phi(q, p) = \sum_i q_i \log \frac{q_i}{q_i + p_i} + \sum_i p_i \log \frac{p_i}{q_i + p_i}$$
  $$= \mathrm{JSD}(p \| q) - 2 \log 2$$

  We find that it suffices to minimize $\Phi(q, p) \in (-2 \log 2, 0)$, which is a valid probem now. After solving this minimization problem, we obtain the optimal $d$:

  $$d_i^*(p_i) = \frac{q_i}{q_i + p_i} = \begin{cases} 1, & \text{if } p_i = 0, \text{ juage as a bad generator} \\ 0, & \text{if } q_i = 0, p_i > 0, \text{ juage as invalid} \\ 1/2, & \text{juage as a good generator} \end{cases},$$

  i.e., for certain data point $i$, the discrimintator returns a probability $q_i/(q_i + p_i)$. At optimal $p^* = q$, $d_i^*(p^*) = 1/2, \forall i$.

  $\square$

**Remark 7.5.** This result is misleading somehow. For instance, images are continuous distributions, so it is impossible to expect the generated image exactly match the original image, i.e., we can never achieve values for $d_i$ other than $\{0, 1\}$.

**Remark 7.6.** This proof justifies GAN by relating it to Jensen–Shannon divergence, but in the beginning we think that this distance is not good.

**Motivation of W-GAN**    The Jensen–Shannon divergence is not a good metric in some settings. For instance, it is impossible to measure the distance between two distributions with the different supporting set, but Wasserstein distance givens a reasonable measure. The $p$-th Wasserstein distance between two probability measures $\mu, \nu$ is defined as

$$W_p(\mu, \nu) = \min_{p \sim \Gamma(\mu,\nu)} \left( \mathbb{E}_{(x,y) \sim \mathcal{P}} |x - y|^p \right)^{1/p}$$

where $\Gamma(\mu, \nu)$ denotes the set of all couplings of $\mu$ and $\nu$. When $p = 1$, finding the Wasserstein distance reduces to solving an LP problem. Moreover, the $W_1$ distance can be re-expressed using duality of LP:

$$W_1(\mu, \nu) = \sup_{|f|_L \leq 1} \mathbb{E}_{x \sim \mu}[f(x)] - \mathbb{E}_{y \sim \nu}[f(y)]$$

where the supremum is taken over all the 1-Lipschitz functions $f$. The W-GAN solves the following problem:

$$\min_{\mathbb{P}_g} \max_{|f|_L \leq 1} \mathbb{E}_{x \sim \mu}[f(x)] - \mathbb{E}_{\hat{x} \sim \nu}[f(\hat{x})]$$

**Remark 7.7.** The origin GAN using Jensen–Shannon divergence also works, but using Wasserstein distance is better. However, Wasserstein distance does has some disadvantages. For example, it is not generalizable, i.e., to approximate $W_1(\mu, \nu)$, we require $\exp(d)$ samples from $(\mu, \nu)$, where $d$ is the supporting dimension of $\mu$ and $\nu$, even when they are Gassuain distributions. One solution is to realize that the objective functions in the original GAN and W-GAN are actually using different distance metrics called "neural network distance". Again, we emphasize that neural-network distance is not a new distance, but a distance with good generalization property, and is used by everyone although the distance metric is not explicitly defined before.

# References

Auer, P., M. Herbster, and M. K. Warmuth (1996). "Exponentially many local minima for single neurons". In: *Advances in Neural Information Processing Systems 8*. Ed. by D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo. MIT Press. 316–322. URL: http://papers.nips.cc/paper/ 1028-exponentially-many-local-minima-for-single-neurons.pdf.

Baldi, P. and K. Hornik (1989). "Neural networks and principal component analysis: Learning from examples without local minima". *Neural Networks*. 2(1): 53–58. ISSN: 0893-6080. DOI: https://doi.org/ 10.1016/0893-6080(89)90014-2. URL: http://www.sciencedirect.com/ science/article/pii/0893608089900142.

Balduzzi, D., M. Frean, L. Leary, J. P. Lewis, K. W.-D. Ma, and B. McWilliams (2017). "The Shattered Gradients Problem: If Resnets Are the Answer, then What is the Question?" In: *Proceedings of the 34th International Conference on Machine Learning - Volume 70. ICML'17*. Sydney, NSW, Australia: JMLR.org. 342–350. URL: http://dl.acm.org/citation.cfm?id=3305381.3305417.

Barron, A. R. (1994). "Approximation and estimation bounds for artificial neural networks". *Machine Learning*. 14(1): 115–133.

Billingsley, P. (1986). *Probability and Measure*. Second. John Wiley and Sons.

Carlini, N. and D. Wagner (2017). "Towards Evaluating the Robustness of Neural Networks". In: *2017 IEEE Symposium on Security and Privacy (SP)*. 39–57. DOI: 10.1109/SP.2017.49.

Chen, P.-Y., H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh (2017). "ZOO: Zeroth Order Optimization Based Black-box Attacks to Deep Neural Networks Without Training Substitute Models". In: *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security. AISec'17*. Dallas, Texas, USA: ACM. 15–26. DOI: 10.1145/3128572.3140448.

Cybenko, G. (1989). "Approximation by superpositions of a sigmoidal function". *Mathematics of Control, Signals, and Systems (MCSS)*. 2(4): 303–314. ISSN: 0932-4194. DOI: 10.1007/BF02551274. URL: http://dx.doi.org/10.1007/BF02551274.

Frankle, J. and M. Carbin (2019). "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks". In: *International Conference on Learning Representations*. URL: https://openreview.net/forum?id=rJl-b3RcF7.

Garipov, T., P. Izmailov, D. Podoprikhin, D. P. Vetrov, and A. G. Wilson (2018). "Loss Surfaces, Mode Connectivity, and Fast Ensembling of DNNs". In: *Advances in Neural Information Processing Systems 31*. Ed. by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett. Curran Associates, Inc. 8789–8798. URL: http://papers.nips.cc/paper/8095-loss-surfaces-mode-connectivity-and-fast-ensembling-of-dnns.pdf.

Gilboa, D., B. Chang, M. Chen, G. Yang, S. S. Schoenholz, E. H. Chi, and J. Pennington (2019). "Dynamical Isometry and a Mean Field Theory of LSTMs and GRUs". *CoRR*. abs/1901.08987. arXiv: 1901.08987. URL: http://arxiv.org/abs/1901.08987.

Glorot, X. and Y. Bengio (2010). "Understanding the difficulty of training deep feedforward neural networks". In: *In Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS?10). Society for Artificial Intelligence and Statistics*.

Glorot, X., A. Bordes, and Y. Bengio (2010). "Deep Sparse Rectifier Neural Networks". In: vol. 15.

Goodfellow, I., J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio (2014). "Generative Adversarial Nets". In: *Advances in Neural Information Processing Systems 27*. Ed. by Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger. Curran Associates, Inc. 2672–2680. URL: http://papers.nips.cc/paper/5423-generative-adversarial-nets.pdf.

Goodfellow, I., J. Shlens, and C. Szegedy (2015a). "Explaining and Harnessing Adversarial Examples". In: *International Conference on Learning Representations*. URL: http://arxiv.org/abs/1412.6572.

Goodfellow, I., O. Vinyals, and A. Saxe (2015b). "Qualitatively Characterizing Neural Network Optimization Problems". In: *International Conference on Learning Representations*. URL: http://arxiv.org/abs/1412.6544.

Gotmare, A., N. Shirish Keskar, C. Xiong, and R. Socher (2018). *Using Mode Connectivity for Loss Landscape Analysis*.

Han, S., J. Pool, J. Tran, and W. Dally (2015). "Learning both Weights and Connections for Efficient Neural Network". In: *Advances in Neural Information Processing Systems 28*. Ed. by C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett. Curran Associates, Inc. 1135–1143. URL: http://papers.nips.cc/paper/5784-learning - both - weights - and - connections - for - efficient - neural - network.pdf.

Hanin, B. and D. Rolnick (2018). "How to Start Training: The Effect of Initialization and Architecture". In: *Advances in Neural Information Processing Systems 31*. Ed. by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett. Curran Associates, Inc. 571–581. URL: http://papers.nips.cc/paper/7338-how-to-start-training-the-effect-of-initialization-and-architecture.pdf.

He, K., X. Zhang, S. Ren, and J. Sun (2015). "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification". In: *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV). ICCV '15*. Washington, DC, USA: IEEE Computer Society. 1026–1034. ISBN: 978-1-4673-8391-2. DOI: 10.1109/ICCV.2015.123. URL: http://dx.doi.org/10.1109/ICCV.2015.123.

He, K., X. Zhang, S. Ren, and J. Sun (2016). "Deep Residual Learning for Image Recognition". In: 770–778. DOI: 10.1109/CVPR.2016.90.

Hornik, K. (1991). "Approximation Capabilities of Multilayer Feedforward Networks". *Neural Netw.* 4(2): 251–257. ISSN: 0893-6080. DOI: 10.1016/0893-6080(91)90009-T. URL: http://dx.doi.org/10.1016/0893-6080(91)90009-T.

"How to comment the paper "The Lottery Ticket Hypothesis"" (n.d.). https://www.zhihu.com/question/323214798. Accessed: 2019-08-14.

Ilyas, A., L. Engstrom, A. Athalye, and J. Lin (2018). "Black-box Adversarial Attacks with Limited Queries and Information". In: *Proceedings of the 35th International Conference on Machine Learning.* Vol. 80. *Proceedings of Machine Learning Research.* PMLR. 2137–2146.

Kawaguchi, K. (2016). "Deep Learning without Poor Local Minima". In: *Advances in Neural Information Processing Systems 29.* Ed. by D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett. Curran Associates, Inc. 586–594. URL: http://papers.nips.cc/paper/6112-deep-learning-without-poor-local-minima.pdf.

Kingma, D. P. and J. Ba (2015). "Adam: A method for stochastic optimization". In: *International Conference on Learning Representations (ICLR).*

Kurach, K., M. Lucic, X. Zhai, M. Michalski, and S. Gelly (2018). "The GAN Landscape: Losses, Architectures, Regularization, and Normalization". *CoRR.* abs/1807.04720. arXiv: 1807.04720. URL: http://arxiv.org/abs/1807.04720.

Lee, J. D., M. Simchowitz, M. I. Jordan, and B. Recht (2016). "Gradient Descent Only Converges to Minimizers". In: *29th Annual Conference on Learning Theory.* Ed. by V. Feldman, A. Rakhlin, and O. Shamir. Vol. 49. *Proceedings of Machine Learning Research.* Columbia University, New York, New York, USA: PMLR. 1246–1257. URL: http://proceedings.mlr.press/v49/lee16.html.

Li, D., T. Ding, and R. Sun (2018). *Over-Parameterized Deep Neural Networks Have No Strict Local Minima For Any Continuous Activations.*

Li, P. and P.-M. Nguyen (2019). "On Random Deep Weight-Tied Autoencoders: Exact Asymptotic Analysis, Phase Transitions, and Implications to Training". In: *International Conference on Learning Representations*.

Lin, H. and S. Jegelka (2018). "ResNet with one-neuron hidden layers is a Universal Approximator". In: *Advances in Neural Information Processing Systems 31*. Ed. by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett. Curran Associates, Inc. 6169–6178. URL: http://papers.nips.cc/paper/7855-resnet-with-one-neuron-hidden-layers-is-a-universal-approximator.pdf.

Nesterov, Y. (2011). "Random gradient-free minimization of convex functions". Jan.

Pennington, J., S. S. Schoenholz, and S. Ganguli (2017). "Resurrecting the sigmoid in deep learning through dynamical isometry: theory and practice". In: *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*. 4785–4795.

Pennington, J., S. S. Schoenholz, and S. Ganguli (2018). "The Emergence of Spectral Universality in Deep Networks". In: *AISTATS*.

Poole, B., S. Lahiri, M. Raghu, J. Sohl-Dickstein, and S. Ganguli (2016). "Exponential expressivity in deep neural networks through transient chaos". In: *Advances in Neural Information Processing Systems 29*. Ed. by D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett. Curran Associates, Inc. 3360–3368. URL: http://papers.nips.cc/paper/6322-exponential-expressivity-in-deep-neural-networks-through-transient-chaos.pdf.

Saxe, A. M., J. L. Mcclelland, and S. Ganguli (2014). "Exact solutions to the nonlinear dynamics of learning in deep linear neural network". In: *In International Conference on Learning Representations*.

Srivastava, R. K., K. Greff, and J. Schmidhuber (2015). "Highway Networks". cite arxiv:1505.00387Comment: 6 pages, 2 figures. Presented at ICML 2015 Deep Learning workshop. Full paper is at arXiv:1507.06228. URL: http://arxiv.org/abs/1505.00387.

Szegedy, C., S. Ioffe, and V. Vanhoucke (2016). "Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning". In: *AAAI*.

Szegedy, C., W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Good-fellow, and R. Fergus (2014). "Intriguing properties of neural networks". In: *International Conference on Learning Representations*. URL: http://arxiv.org/abs/1312.6199.

"Understanding nonconvex optimization" (n.d.). http://praneethnetrapalli.org/UnderstandingNonconvexOptimization-V5.pdf. Accessed: 2019-08-18.

Wang, J. (2019a). *MAT2006: Elementary Real Analysis*. Available at the link https://walterbabyrudin.github.io/information/Notes/MAT2006.pdf.

Wang, J. (2019b). *MAT3006: Real Analysis; Lecture 8*. Available at the link https://walterbabyrudin.github.io/information/Updates/MAT3006/Week4_Wednesday.pdf.

Wong, E., F. R. Schmidt, J. H. Metzen, and J. Z. Kolter (2018). "Scaling Provable Adversarial Defenses". In: *Proceedings of the 32Nd International Conference on Neural Information Processing Systems. NIPS'18*. Montr&#233;al, Canada: Curran Associates Inc. 8410–8419. URL: http://dl.acm.org/citation.cfm?id=3327757.3327932.

Wu, Y. and K. He (2018). "Group Normalization". In: *The European Conference on Computer Vision (ECCV)*.

Xiao, L., Y. Bahri, J. Sohl-Dickstein, S. Schoenholz, and J. Pennington (2018). "Dynamical Isometry and a Mean Field Theory of CNNs: How to Train 10,000-Layer Vanilla Convolutional Neural Networks". In: *Proceedings of the 35th International Conference on Machine Learning*. Ed. by J. Dy and A. Krause. Vol. 80. *Proceedings of Machine Learning Research*. Stockholmsmassan, Stockholm Sweden: PMLR. 5393–5402.

Xiao-Hu Yu and Guo-An Chen (1995). "On the local minima free condition of backpropagation learning". *IEEE Transactions on Neural Networks*. 6(5): 1300–1303. ISSN: 1045-9227. DOI: 10.1109/72.410380.

Zhang, H., Y. N. Dauphin, and T. Ma (2019). "Residual Learning Without Normalization via Better Initialization". In: *International Conference on Learning Representations*. URL: https://openreview.net/forum?id=H1gsz30cKX.

Zhang, Y., R. Tapia, and L. Velazquez (2000). "On Convergence of
    Minimization Methods: Attraction, Repulsion, and Selection". *Jour-
    nal of Optimization Theory and Applications*. 107(3): 529–546. ISSN:
    1573-2878. DOI: 10.1023/A:1026443131121. URL: https://doi.org/10.
    1023/A:1026443131121.