

Mathematics of Deep Learning

Lecture 8: Over-parametrization, Convergence and Lazy Training Regime; PAC Learning and Generalization Bounds

Prof. Predrag R. Jelenković
Time: Tuesday 4:10-6:40pm

Dept. of Electrical Engineering
Columbia University , NY 10027, USA
Office: 812 Schapiro Research Bldg.
Phone: (212) 854-8174
Email: predrag@ee.columbia.edu
URL: <http://www.ee.columbia.edu/~predrag>

New Idea: Over-Parametrization

- ▶ Previous idea: **Composition of functions**
- ▶ New idea: **Over-parametrization**
 - ▶ Mathematically speaking: Passes the width $m \rightarrow \infty$
(this is the opposite of the previous work on expressiveness, where the goal was to minimize the # of neurons)
 - ▶ Mathematical tools: **Laws of large numbers**
 - ▶ "Smooth out" the layer functions
 - ▶ This idea was used in a number of recent papers to:
 - ▶ Connect NNs to Kernels
 - ▶ Show that all local minima are global
 - ▶ Show convergence to global minima
 - ▶ **Show generalization bounds**

Neural Tangent Kernel

- ▶ Today, we'll focus on the following paper:
[Neural Tangent Kernel: Convergence and Generalization in Neural Networks](#), by A. Jacot, F. Gabriel and C. Hongler, NIPS 2018.
- ▶ This paper considered similar questions of convergence as Du et al. (2019).
- ▶ They introduced a terminology of Neural Tangent Kernel (NTK)
- ▶ We'll follow the notation from the paper:
 - ▶ N data points: $(\mathbf{x}_i, y_i), \mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}, 1 \leq i \leq N$.
 - ▶ Find the best fit of a cost $C(\cdot)$ over a function class, say quadratic

$$\min_{f \in \mathcal{F}} C(f) = \min_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N (f(\mathbf{x}_i) - y_i)^2.$$

Linear Model

- ▶ Choose P features $f^{(p)} \in \mathcal{F}, 1 \leq p \leq P$.
- ▶ Let a parametrized function $F : \mathbb{R}^P \rightarrow \mathcal{F}$

$$f(\theta) = \frac{1}{\sqrt{P}} \sum_{p=1}^P \theta_p f^{(p)}$$

- ▶ We can choose $f^{(p)}$ are i.i.d. such that a kernel is given by $\mathbb{E}_P[f^{(P)}(x)f^{(P)}(y)] = K(x, y)$
- ▶ Then, we optimize the composition $C \circ F$, using, say gradient descent, to find optimal linear parameters θ .
- ▶ Convex problem when C is convex.

Neural Networks: Nonlinear Model

- ▶ L layers, $l = 0$ input and $l = L$ output; each hidden layer having $n_l, 1 \leq l \leq L - 1$ neurons
- ▶ Activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, say ReLU, $\sigma(x) = x_+$
- ▶ $W^{(\ell)}$ are $n_\ell \times n_{\ell+1}$ matrices connecting the layers, $\ell = 0, \dots, L - 1$.
- ▶ Total parameters: $\theta = (W^{(0)}, \dots, W^{(L-1)})$.
- ▶ Non-linear parametrization, network function $f_\theta(x)$

$$\alpha^{(0)}(x; \theta) = x$$

$$\tilde{\alpha}^{(\ell+1)}(x; \theta) = \frac{1}{\sqrt{n_\ell}} W^{(\ell)} \alpha^{(\ell)}(x; \theta) + \beta b^{(\ell)} \quad (\beta \in \mathbb{R}_+) \quad (\text{preactivations})$$

$$\alpha^{(\ell)}(x; \theta) = \sigma(\tilde{\alpha}^{(\ell)}(x; \theta)) \quad (\text{activations})$$

$$f_\theta = \tilde{\alpha}^{(L)}(x; \theta)$$

- ▶ Now, the loss $C \circ F^{(L)}$ is **not convex**
- ▶ Linearize the model by studying the infinite-width limit
 $(n_1, \dots, n_{L-1}) \rightarrow \infty$

Gaussian Initialization

- Initialize the parameters $\theta \sim \mathcal{N}(\mathbf{0}, I_P)$, where P is the number of parameters $P = \sum_{\ell=0}^{L-1} (n_\ell + 1)n_{\ell+1}$

Proposition 1 In the infinite width limit $(n_1, \dots, n_{L-1}) \rightarrow \infty$ the preactivations $\tilde{\alpha}_i^{(\ell)}(\cdot; \theta) : \mathbb{R}^{n_\ell} \rightarrow \mathbb{R}$ are i.i.d. Gaussian processes with covariance $\Sigma^{(\ell)}$ given by

$$\Sigma^{(1)}(x, y) = \frac{1}{n_0} x \cdot y + \beta^2$$

$$\Sigma^{(\ell+1)}(x, y) = \mathbb{E}_{\alpha \sim \mathcal{N}(\mathbf{0}, \Sigma^{(\ell)})} [\sigma(\alpha(x))\sigma(\alpha(y))] + \beta^2$$

In particular, f_θ is a Gaussian process with covariance $\Sigma^{(L)}$

The **proof** follows by induction, conditioning and then passing width to infinity, i.e., conditioning on $\tilde{\alpha}^{(L)}$

$$\begin{aligned} \tilde{\Sigma}^{(L+1)}(x, y) &= \frac{1}{n_L} \alpha^{(L)}(x; \theta) \cdot \alpha^{(L)}(y; \theta) + \beta^2 \\ &\rightarrow \mathbb{E}_{\alpha \sim \mathcal{N}(\mathbf{0}, \Sigma^{(L)})} [\sigma(\alpha(x))\sigma(\alpha(y))] + \beta^2 \end{aligned}$$

Training Neural Tangent Kernel (NTK)

- ▶ Gradient descent: path of steepest descent for quadratic loss

$$\partial_t \theta_p = -\partial_\theta (C \circ F^{(L)}) = \frac{2}{N} \sum_{i=1}^N (y_i - f_\theta(x_i)) \partial_{\theta_p} f_\theta(x_i)$$

N - number of data points, for short $\partial_v = \partial / \partial v$.

- ▶ Evolution of $f_\theta(x)$:

$$\begin{aligned} \partial_t f_\theta(x) &= \sum_{p=1}^P \partial_{\theta_p} f_\theta(x) \partial_t \theta_p \quad (\text{chain rule}) \\ &= \frac{2}{N} \sum_{i=1}^N (y_i - f_\theta(x_i)) \left(\sum_{p=1}^P \partial_{\theta_p} f_\theta(x_i) \partial_{\theta_p} f_\theta(x) \right) \end{aligned}$$

- ▶ Neural Tangent Kernel (NTK):

$$\Theta^{(L)}(x, x') := \sum_{p=1}^P \partial_{\theta_p} f_\theta(x) \partial_{\theta_p} f_\theta(x')$$

Note that this quantity was used implicitly in Du et al., and was denoted as $H(t)$

Infinite Width Limit of NTK

For finite width, $\Theta^{(L)}(t, x, x')$ is random, but it is deterministic in the infinite width limit. (This limit was H^∞ in Du et al.)

Theorem 1 For any $t < T$, as $(n_1, \dots, n_{L-1}) \rightarrow \infty$,

$$\Theta^{(L)}(t) \rightarrow \Theta_\infty^{(L)},$$

where

$$\Theta_\infty^{(L)}(x, x') = \sum_{\ell=1}^L \Sigma^{(\ell)}(x, x') \dot{\Sigma}^{(\ell+1)}(x, x') \cdots \dot{\Sigma}^{(L)}(x, x')$$

and

$$\dot{\Sigma}^{(\ell)}(x, x') = \mathbb{E}_{\alpha \sim \mathcal{N}(\mathbf{0}, \Sigma^{(\ell-1)})} [\dot{\sigma}(\alpha(x)) \dot{\sigma}(\alpha(x'))]$$

Key characteristics: (Proposition 1) For non-polynomial Lipschitz activation σ , and for any input dimension n_0 , the limiting NTK, $\Theta_\infty^{(L)}$, is positive definite on the unit sphere $S^{n_0-1} = \{x \in \mathbb{R}^{n_0} : x \cdot x = 1\}$.

NTK: Main Result

Theorem 2 Assume non-polynomial activation σ , which is twice differentiable with bounded second derivative. Then, for least squares regression (and more), for any fixed T , the dynamics of gradient descent for $t \in [0, T]$ converges to the infinite width limit.

Key ingredient in the proof - Lemma 1 The weights $W^{(\ell)}(t)$ during training do not move much, or more precisely

$$\lim_{n_L \rightarrow \infty} \cdots \lim_{n_1 \rightarrow \infty} \sup_{t \in [0, T]} \left\| \frac{1}{\sqrt{n_\ell}} \left(W^{(\ell)}(t) - W^{(\ell)}(0) \right) \right\| = 0,$$

where $\| \cdot \|$ is the operator norm.

Reflection On Recent Convergence Papers

We have seen several papers that prove convergence, where

- ▶ Key ingredient is the weights do not move during training, i.e.,

$$W(t) \approx W(0)$$

- ▶ In essence, these models behave like linear models corresponding to the infinite with Gaussian kernel - Neural Tangent Kernel, or more precisely random sampling from this infinite width kernel.
- ▶ Are we missing something?
 - ▶ The main characteristics of NNs is that they are nonlinear models. If we are reducing them to linear Kernel models, why not work with kernels in the first place?
- ▶ What is the main property of the considered models that is constraining the weights from moving far?
 - ▶ Is it over-parametrization?

Lazy Training

The following paper pinpoints the key ingredient responsible for weights remaining close to the initial values during training

- [A Note on Lazy Training in Supervised Differentiable Programming](#), by L. Chizat and F. Bach, Dec 2018.

See also an updated version:

[On Lazy Training Differentiable Programming](#), by L. Chizat and F. Bach, NIPS 2019.

This paper puts in perspective many of the papers that we covered recently:

Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In Advances in neural information processing systems, 2018.

Simon S. Du, Xiyu Zhai, Barnabás Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In International Conference on Learning Representations, 2019.

Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. arXiv preprint arXiv:1904.11955, 2019.

Simon S. Du, Lee Jason D., Li Haochuan, Wang Liwei, and Zhai Xiyu. Gradient descent finds global minima of deep neural networks. In International Conference on Machine Learning (ICML), 2019.

Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In Advances in Neural Information Processing Systems, pages 81678176, 2018.

Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Proceedings of the 36th International Conference on Machine Learning, volume 97, pages 242252, 2019. Difan

Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep ReLU networks. Machine Learning Journal, 2019.

The Answer Is In Scaling

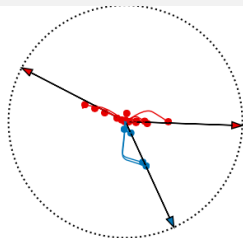
They show that the key part of the model which ensures that the weights do not move far from the initial position ($w_{ij}(0)$) is the **scaling**

$\alpha = 1/\sqrt{m}$, where m is the width of the network.

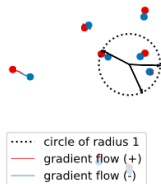
- ▶ Hence, the key to NTK approximation works because of scaling, rather than over-parametrization.
 - ▶ They argue: any model can be scaled so that weights do not move much
- Example: τ - variance of $w_{ij}(0)$; "double trick" $f(x, w(0)) = 0$

Ground truth: x uniform on a unit sphere, y output of NN with 3 neurons

Model: $m = 20$ neurons - not a wide network



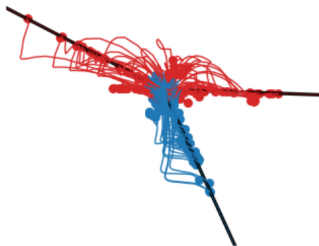
(a) "Active" training ($\tau = 0.1$)



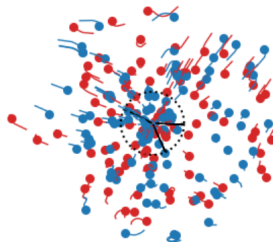
(b) Lazy training ($\tau = 2$)

Many Neurons Example

- ▶ Ground truth: x uniform on a unit sphere, y output of NN with 3 neurons
- ▶ $n = 200$ data samples
- ▶ Model wide network: $m = 200$ neurons
- ▶ τ - variance of $w_{ij}(0)$



(a) “Active” training ($\tau = 0.1$)



(b) Lazy ($\tau = 2$)

Hence, the [key to NTK regime is the scaling](#), rather than over-parametrization.

Recall Tangent Model

Consider any parametric model $f(w, x)$, $w \in \mathbb{R}^p, x \in \mathbb{R}^d$. Then, assuming f is differentiable, we can linearize this model around the initial parameters w_0 as:

$$f_0(w, x) := f(w_0, x) + (w - w_0) \cdot \nabla_w f(w_0, x), \quad (\text{Tangent model}). \quad (1)$$

Consider n data points $(x_i, y_i) \in \mathbb{R}^{d+1}, 1 \leq i \leq n$, and assume a quadratic loss, then training $f_0(w, x)$ is equivalent to solving a dual problem in the dot-product space according to tangent kernel (NTK):

$$K(x, y) = \nabla_w f(w_0, x) \cdot \nabla_w f(w_0, y).$$

Question: When is training $f_0(w, x)$ going to produce approximately the same results as $f(w, x)$?

Obviously, if $f(w, x)$ is linear, then $f_0(w, x) = f(w, x)$.

When Does Lazy Training Occur?

In other words, when can we expect that the trained \mathbf{w}^* close to the initial \mathbf{w}_0 , $\mathbf{w}^* \approx \mathbf{w}_0$.

- ▶ Let $F(\mathbf{w}) := L(f(\mathbf{w}))$, where L is the loss function. e.g., quadratic.
- ▶ Assume $F(\mathbf{w}_0) > 0$, $\nabla F(\mathbf{w}_0) \neq 0$, and consider a GD step

$$\mathbf{w}_1 = \mathbf{w}_0 - \eta \nabla F(\mathbf{w}_0). \quad (2)$$

- ▶ Conditions for lazy training: when the differential of f does not change much while the loss changes significantly, i.e.,

$$\Delta(F) := \frac{|F(\mathbf{w}_1) - F(\mathbf{w}_0)|}{F(\mathbf{w}_0)} \gg \Delta(Df) := \frac{\|Df(\mathbf{w}_1) - Df(\mathbf{w}_0)\|}{\|Df(\mathbf{w}_0)\|} \quad (3)$$

- ▶ Next, using Equation (2), we obtain

$$F(\mathbf{w}_1) \approx F(\mathbf{w}_0) + (\mathbf{w}_1 - \mathbf{w}_0) \cdot \nabla F(\mathbf{w}_0) = F(\mathbf{w}_0) - \eta \|\nabla F(\mathbf{w}_0)\|^2 \quad (4)$$

When Does Lazy Training Occur?

- ▶ Similarly,

$$\begin{aligned} Df(\mathbf{w}_1) &\approx Df(\mathbf{w}_0) + (w_1 - w_0) \cdot D^2 f(\mathbf{w}_1) \\ &= Df(\mathbf{w}_0) - \eta \nabla F(\mathbf{w}_0) \cdot D^2 f(\mathbf{w}_1), \end{aligned} \quad (5)$$

where Df is the Gradient/Jacobian, and $D^2 f$ is a Hessian matrix.

- ▶ Finally, replacing Equation (5) and Equation (4) in Equation (3), yields

$$\begin{aligned} \Delta(F) &\approx \eta \frac{\|\nabla F(\mathbf{w}_0)\|^2}{F(\mathbf{w}_0)} \gg \Delta(Df) \approx \eta \frac{\|\nabla F(\mathbf{w}_0)\| \|D^2 f(\mathbf{w}_0)\|}{\|Df(\mathbf{w}_0)\|} \\ \frac{\|\nabla F(\mathbf{w}_0)\|}{F(\mathbf{w}_0)} &\gg \frac{\|D^2 f(\mathbf{w}_0)\|}{\|Df(\mathbf{w}_0)\|} \end{aligned}$$

Lazy Training: Quadratic Loss

- For quadratic function $L(f) = \|f - y\|^2/2$, the previous condition simplifies to

$$\kappa_f(\mathbf{w}_0) := \|f(\mathbf{w}_0) - y\| \frac{\|D^2 f(\mathbf{w}_0)\|}{\|Df(\mathbf{w}_0)\|^2} \ll 1$$

using the approximation

$$\|\nabla F(\mathbf{w}_0)\| = \|Df(\mathbf{w}_0)^\top (f(\mathbf{w}_0) - y)\| \approx \|Df(\mathbf{w}_0)\| \|f(\mathbf{w}_0) - y\|$$

- $\kappa_f(\mathbf{w}_0)$ could be called **relative scale** of model f at \mathbf{w}_0 .

Lazy Training Examples

- **Rescaled model.** Consider a rescaled quadratic loss function

$$\frac{\|\alpha f - y\|^2}{2\alpha^2},$$

which yields

$$\begin{aligned}\kappa_{\alpha f}(\mathbf{w}_0) &= \frac{1}{\alpha} \|\alpha f(\mathbf{w}_0) - y\| \frac{\|D^2 f(\mathbf{w}_0)\|}{\|Df(\mathbf{w}_0)\|^2} \\ &= O\left(\frac{1}{\alpha}\right), \quad (\text{if } \|\alpha f(\mathbf{w}_0) - y\| = O(1)) \\ &\ll 1, \quad (\text{as } \alpha \uparrow \infty)\end{aligned}$$

- $\|\alpha f(\mathbf{w}_0) - y\| = O(1)$ can be ensured by setting $f(\mathbf{w}_0) = 0$, which in NNs can be attained using a "doubling trick", i.e., repeating each neuron in the output layer with a negative linear weight.
- Hence, NN can be trained in the lazy regime even if networks are not wide. **Large α is equivalent to large variance for \mathbf{w}_0 .**

Lazy Training Examples

- **One hidden layer NN.** Consider

$$\begin{aligned}f(x) &= \alpha(m) \sum_{j=1}^m a_j \sigma(w_j \cdot x) \\ &=: \alpha(m) \sum_{j=1}^m \phi(\theta_j, x),\end{aligned}$$

where $\theta_j = (a_j, w_j)$, $\alpha(m) > 0$, $\mathbb{E}\phi(\theta_j, x) = 0$. Recall that the papers we considered used $\alpha(m) = 1/\sqrt{m}$, $\sigma(x) = x^+$ and $a_j \in \{0, 1\}$.

- Since $\mathbb{E}\phi(\theta_j, x) = 0$,

$$\mathbb{E}\|f(\mathbf{w}_0, x)\|^2 = m\alpha(m)^2 \mathbb{E}\|\phi(\theta)\|^2.$$

- For the differential, using the law of large numbers,

$$\frac{1}{m\alpha(m)^2} Df(\mathbf{w}_0) Df(\mathbf{w}_0)^\top = \frac{1}{m} \sum_{j=1}^m D\phi(\theta_j) D\phi(\theta_j)^\top \rightarrow \mathbb{E}[D\phi(\theta) D\phi(\theta)^\top]$$

Lazy Training Examples

- ▶ Hence

$$\mathbb{E}\|Df(\mathbf{w}_0)\|^2 = \mathbb{E}[Df(\mathbf{w}_0)Df(\mathbf{w}_0)^\top] \sim m\alpha(m)^2\mathbb{E}[D\phi(\theta)D\phi(\theta)^\top]$$

- ▶ Also, we can estimate the operator norm of Df as

$$\begin{aligned}\|D^2f(\mathbf{w}_0)\| &= \sup_{u \in \mathbb{R}^{dm}, \|u\| \leq 1} \alpha(m) \sum_{j=1}^m u_j^\top D^2\phi(\theta_j)u_j \\ &\leq \alpha(m) \sup_{\theta_j} \|D^2\phi(\theta_j)\| \leq \alpha(m)\text{Lip}(D\phi),\end{aligned}$$

where $\text{Lip}(D\phi)$ is the Lipschitz constant of $D\phi$.

- ▶ Using the above and

$$\|f(\mathbf{w}_0) - y\| \leq \|f(\mathbf{w}_0)\| + \|y\| = O(\alpha\sqrt{m}) + O(\sqrt{n})$$

$$\mathbb{E}[\kappa_f(\mathbf{w}_0)] \approx O\left((\alpha\sqrt{m}) + \sqrt{n}\right) \frac{\alpha}{\alpha^2 m} = O(m^{-1/2} + (m\alpha(m))^{-1})$$

Lazy Training Examples

Hence, if $\alpha m \rightarrow \infty$, we have a Lazy Training regime.

In particular, in wide networks, we can consider two interesting regimes

$$\alpha(m) = \frac{1}{\sqrt{m}} \Rightarrow \mathbb{E}[\kappa_f(\mathbf{w}_0)] = O(m^{-1/2}) \ll 1, \quad (\text{NTK/lazy regime})$$

$$\alpha(m) = \frac{1}{m} \Rightarrow \mathbb{E}[\kappa_f(\mathbf{w}_0)] = O(1), \quad (\text{Mean field/active regime})$$

In the mean field/active regime

- ▶ The weights \mathbf{w} move considerably during training
- ▶ Much harder problem: requires extending gradient flow/descent to infinite dimensions
 - ▶ We might consider in the future some papers in this regime

General Result On Lazy Training

Theorem 3.2 from Chizat and Bach (2019): We state a [version from Section 8 in Telgarsky's note](#). Additional notation:

- ▶ Consider n data points (x_i, y_i) , $1 \leq i \leq n$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, and let

$$f(w) = (f(x_1; w), \dots, f(x_n; w))^\top \in \mathbb{R}^n.$$

- ▶ Assume quadratic loss

$$L(\alpha f(w)) := \frac{1}{2} \|\alpha f(w) - y\|^2, \quad \alpha > 0.$$

- ▶ Then, the gradient flow evolves as

$$\dot{w}(t) = -\nabla_w L(\alpha f(w(t))) = -\alpha J_t^\top \nabla L(\alpha f(w(t))),$$

where $J_t = Df(w(t))$ is the Jacobian

$$J_t = J_{w(t)} = (\nabla f(x_1; w(t)), \dots, \nabla f(x_n; w(t)))^\top \in \mathbb{R}^{n \times p}$$

- ▶ Linear tangent model flow, with the same initialization $w(0)$

$$f_0(u) = f(w(0)) + J_0(u - w(0))$$

$$\dot{u}(t) = \nabla_w L(\alpha f_0(u(t))) = -\alpha J_0^\top \nabla L(\alpha f_0(u(t)))$$

Overparametrized NTK/Lazy Regime

- ▶ How close is the nonlinear gradient flow to the linear one?
- ▶ Assumptions

$$\text{rank}(J_0) = n$$

$$\sigma_{\min} = \sigma_{\min}(J_0) = \sqrt{\lambda_{\min}(J_0 J_0^\top)} > 0$$

$$\|J_w - J_v\| \leq \beta \|w - v\|$$

$$\alpha \geq \frac{\beta \sigma_{\max} \sqrt{1152 L_0}}{\sigma_{\min}^3}, \quad L_0 = \frac{1}{2} \|\alpha f(w(0)) - y\|^2$$

Theorem 8.1 (Telgarsky, also Theorem 3.2 in Chizat&Bach, 2019)
Under the preceding assumptions,

$$\max(L(\alpha f(w(t))), L(\alpha f_0(u(t)))) \leq L_0 \exp(-t\alpha^2 \sigma_{\min}^2/2)$$

$$\max(\|w(t) - w(0)\|, \|u(t) - w(0)\|) \leq \frac{3\sigma_{\max} \sqrt{8L_0}}{\alpha \sigma_{\min}^2}$$

Proof next time.

Back to Generalization Error

The main objective of statistical learning is to predict well on future data.

- ▶ How do we measure the error?
 - ▶ Regression: Quadratic error/loss

$$\ell(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$$

- ▶ Classification: $\ell(\hat{y}, y) = 1_{\{\hat{y} \neq y\}}$
- ▶ $\{\mathcal{H}\}$: Hypothesis class of functions, e.g., all functions, $h(x, w)$, that can be generated by a NN of a certain architecture.
- ▶ **Empirical Risk/Loss** - total training error: For a data sample $S = \{(x_i, y_i)\}_{i=1}^n$ and $h \in \mathcal{H}$

$$\hat{L}_n = \hat{L}_n(h) \equiv L_S(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Shai's book [ML2014] uses $L_S(h)$ notation; we'll use these interchangeably.

Empirical Risk Minimization and True Error

During training one typically **minimizes the empirical risk (ERM)**, and obtain $\hat{h}_n \equiv h_S$

$$\hat{h}_n \in \arg \min_{h \in \mathcal{H}} \hat{L}(h)$$

True Risk=Population Risk

- ▶ Let $x \in \mathcal{X}, y \in \mathcal{Y}$, say $\mathcal{X} \subset \mathbb{R}^d, \mathcal{Y} \subset \mathbb{R}, z = (x, y) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$
- ▶ Define a probability measure on \mathcal{Z} , denoted by \mathcal{D} in [ML2014] book
- ▶ Let $\{z_i = (x_i, y_i)\}_{i=1}^n$ be i.i.d. random variables on a product probability space \mathcal{Z}^n .
- ▶ Then, for $h \in \mathcal{H}$, we define a true risk/population risk as

$$L(h) \equiv L_{\mathcal{D}}(h) := \mathbb{E}_{x,y} \ell(h(x), y)$$

Probably Approximately Correct (PAC) Learning

The ultimate goal is to find h that minimizes the true risk, i.e.,

$$h \in \arg \min_{h \in \mathcal{H}} L(h)$$

But this is often impossible, leading to the more relaxed definition

Definition (PAC Learnability) A hypothesis class \mathcal{H} is PAC learnable with respect to a set $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ and a loss function $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow \mathbb{R}_+$, if there exists a learning algorithm, which returns $\hat{h} \equiv \hat{h}(S)$ for a sample S of size $|S| = m$ with the following property: for every $0 < \epsilon, \delta < 1$ there exists $m(\epsilon, \delta)$, such that for all $m \geq m(\epsilon, \delta)$,

$$\mathbb{P} \left[\left\{ S \in \mathcal{Z}^m : L(\hat{h}(S)) \leq \min_{h \in \mathcal{H}} L(h) + \epsilon \right\} \right] \geq 1 - \delta$$

- PAC learning was introduced by Leslie Valiant (1984), who won for it the Turing Award in 2010.

Representative Sample

Definition (ϵ -representative sample)

A training set S is called ϵ -representative if

$$\forall h \in \mathcal{H}, \quad |\hat{L}_S(h) - L(h)| \leq \epsilon .$$

- ▶ Typically, we can only expect this to be true with high probability.

Representative Sample

Lemma

Assume that a training set S is $\frac{\epsilon}{2}$ -representative. Then, any output of $\text{ERM}_{\mathcal{H}}(S)$, namely any $\hat{h}_S \in \operatorname{argmin}_{h \in \mathcal{H}} \hat{L}_S(h)$, satisfies

$$L(\hat{h}_S) \leq \min_{h \in \mathcal{H}} L(h) + \epsilon .$$

Proof: For every $h \in \mathcal{H}$,

$$L(\hat{h}_S) \leq \hat{L}_S(\hat{h}_S) + \frac{\epsilon}{2} \leq \hat{L}_S(h) + \frac{\epsilon}{2} \leq L(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} = L(h) + \epsilon$$

- ▶ Hence, to ensure that ERM algorithm is a PAC (good) learner, we need to ensure that with probability $1 - \delta$ a training set S is ϵ -representative

Uniform Convergence is Sufficient for Learnability

Definition (uniform convergence)

\mathcal{H} has the *uniform convergence property* if there exists $m(\epsilon, \delta)$, such that for every $m \geq m(\epsilon, \delta)$, every distribution \mathcal{D} , and any i.i.d. sample $S \equiv S_m$ with distribution \mathcal{D} , S is ϵ -representative with probability $1 - \delta$, i.e., if $m \geq m(\epsilon, \delta)$

$$\mathbb{P} \left[\left\{ S \in \mathcal{Z}^m : \sup_{h \in \mathcal{H}} |\hat{L}_S(h) - L(h)| \leq \epsilon \right\} \right] \geq 1 - \delta$$

- ▶ A direct consequence of the preceding definition and lemma:

Corollary

- ▶ If \mathcal{H} has the uniform convergence property with a function $m(\epsilon, \delta)$ then \mathcal{H} is PAC learnable with the sample complexity $m(\epsilon, \delta) \leq m(\epsilon/2, \delta)$.
- ▶ Furthermore, in that case, the $\text{ERM}_{\mathcal{H}}$ paradigm is a successful agnostic PAC learner for \mathcal{H} .

Finite Classes are PAC Learnable

We will prove the following:

Theorem

Assume \mathcal{H} is finite and the range of the loss function is $[0, 1]$.

Then, \mathcal{H} is agnostically PAC learnable using the $\text{ERM}_{\mathcal{H}}$ algorithm with sample complexity

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2 \log(2|\mathcal{H}|/\delta)}{\epsilon^2} \right\rceil .$$

Proof: It suffices to show that \mathcal{H} has the uniform convergence property with

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2} \right\rceil .$$

Proof (cont.)

- To show uniform convergence, we need to show

$$\mathbb{P}(\{S : \sup_{h \in \mathcal{H}} |\hat{L}_S(h) - L(h)| > \epsilon\}) < \delta .$$

- Using the union bound:

$$\begin{aligned} \mathbb{P}(\max_{h \in \mathcal{H}} |\hat{L}_S(h) - L(h)| > \epsilon) &= \mathbb{P}(\cup_{h \in \mathcal{H}} \{|\hat{L}_S(h) - L(h)| > \epsilon\}) \\ &\leq \sum_{h \in \mathcal{H}} \mathbb{P}(\{|\hat{L}_S(h) - L(h)| > \epsilon\}) . \end{aligned}$$

Proof (cont.)

- ▶ Recall: $z = (x, y), z_i = (x_i, y_i)$
- ▶ Recall: $L(h) = \mathbb{E}_z[\ell(h, z)]$ and $\hat{L}_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$.
- ▶ Denote $\theta_i = \ell(h, z_i)$.
- ▶ Then, for all i , $\mathbb{E}[\theta_i] = L(h)$

Then, by Hoeffding's inequality, since $0 \leq \theta_i \leq 1$, this implies:

$$\mathbb{P}[|\hat{L}_S(h) - L(h)| > \epsilon] \leq 2 \exp(-2m\epsilon^2) .$$

Which, in conjunction with the union bound, results in

$$\mathbb{P}[\sup_{h \in \mathcal{H}} |\hat{L}_S(h) - L(h)| > \epsilon] \leq 2|\mathcal{H}| \exp(-2m\epsilon^2)$$

So, if $m \geq \frac{\log(2|\mathcal{H}|/\delta)}{2\epsilon^2}$ then the right-hand side is at most δ as required. □

The Discretization Trick

- ▶ Suppose \mathcal{H} is parameterized by d numbers
- ▶ Suppose we are happy with a representation of each number using b bits (say, $b = 32$)
- ▶ Then $|\mathcal{H}| \leq 2^{db}$, and so

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{2db + 2 \log(2/\delta)}{\epsilon^2} \right\rceil .$$

- ▶ While not very elegant, it's a great tool for upper bounding sample complexity

Learning Infinite Hypothesis Classes

For most hypothesis classes $|\mathcal{H}| = \infty$: What can be done here?

Common measures of complexity of hypothesis classes

- ▶ Vapnik-Chervonenkis (VC) dimension
Chapter 6 in Shai's [ML2014] book
- ▶ Rademacher Complexity
Chapter 26 in [ML2014] book; this was recently used in
 - ▶ [A Priori Estimates For Two-layer Neural Networks](#), by Weinan et al., Jan 2019.
 - ▶ [Fine-Grained Analysis of Optimization and Generalization for Overparameterized Two-Layer Neural Networks](#), by Arora et al., Jan 2019.
- ▶ PAC - Bayes: Chapter 31 in [ML20014]
used recently in several NN papers, e.g.: see Neyshabur et al.
- ▶ Compression Bounds: Chapter 30 in [ML2014]

Rademacher Complexity: Motivation

- ▶ To simplify the notation, let $f(z) \equiv f(h, z) = \ell(h, z)$ and let \mathcal{F} be the set of these functions
- ▶ Then, for $f \in \mathcal{F}$, the population/true risk and empirical risk are equal to

$$L(f) = \mathbb{E}_z[f(z)], \quad \hat{L}_S(f) = \frac{1}{m} \sum_{i=1}^m f(z_i)$$

- ▶ Recall ϵ -representative sample: A training set S is called ϵ -representative if

$$\forall h \in \mathcal{H}, \quad |\hat{L}_S(h) - L(h)| \leq \epsilon .$$

- ▶ Which motivates the following definition

Definition *Representativeness of S* : is the largest gap between the true error and empirical error

$$\text{Rep}(\mathcal{F}, S) = \sup_{f \in \mathcal{F}} |\hat{L}_S(f) - L(f)|$$

Rademacher Complexity: Motivation

- ▶ But we don't know the true error $L(f)$
- ▶ Replace $L(f)$ by another empirical error
- ▶ Partition the training data S into two disjoint sets: S_1, S_2 and estimate the representativeness of S by

$$\sup_{f \in \mathcal{F}} (\hat{L}_{S_1}(f) - \hat{L}_{S_2}(f))$$

- ▶ Let $\sigma_i = 1$ if $z_i \in S_1$ and $\sigma_i = -1$, otherwise. Then

$$\sup_{f \in \mathcal{F}} (\hat{L}_{S_1}(f) - \hat{L}_{S_2}(f)) = \frac{1}{m} \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(z_i)$$

which motivates the following definition

Definition Let σ_i be i.i.d. with $\mathbb{P}[\sigma_i = 1] = \mathbb{P}[\sigma_i = -1] = 1/2$. Then, the *Rademacher Complexity* of \mathcal{F} with respect to sample S is defined as

$$R(\mathcal{F}, S) := \frac{1}{m} \mathbb{E}_{\sigma} \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(z_i)$$

Rademacher Complexity

Lemma 26.2 The expected value of the representativeness of S is bounded by

$$\mathbb{E}_S[\text{Rep}(\mathcal{F}, S)] = \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} |\hat{L}_S(f) - L(f)| \right] \leq \frac{2}{m} \mathbb{E} \sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(z_i) = 2 \mathbb{E}_S[R(\mathcal{F}, S)]$$

Proof: Let $S' = \{z'_i\}_{i=1}^m$ be another i.i.d. sample independent of S .
Then

$$L(f) - \hat{L}_S(f) = \mathbb{E}_{S'}[\hat{L}_{S'}(f) - \hat{L}_S(f)]$$

which implies (note $\hat{L}_{S'}(f) - \hat{L}_S(f) \leq \sup_{f \in \mathcal{F}} (\hat{L}_{S'}(f) - \hat{L}_S(f))$)

$$\sup_{f \in \mathcal{F}} (L(f) - \hat{L}_S(f)) = \sup_{f \in \mathcal{F}} \mathbb{E}_{S'}[\hat{L}_{S'}(f) - \hat{L}_S(f)] \leq \mathbb{E}_{S'} \left[\sup_{f \in \mathcal{F}} (\hat{L}_{S'}(f) - \hat{L}_S(f)) \right]$$

Rademacher Complexity

Proof: (continued) Taking the expectation over S

$$\begin{aligned} \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} (L(f) - \hat{L}_S(f)) \right] &\leq \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} (\hat{L}_{S'}(f) - \hat{L}_S(f)) \right] \\ &= \frac{1}{m} \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m (f(z'_i) - f(z_i)) \right] \\ &= \frac{1}{m} \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left(f(z'_j) - f(z_j) + \sum_{i \neq j} (f(z'_i) - f(z_i)) \right) \right] \\ &= \frac{1}{m} \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left(f(z_j) - f(z'_j) + \sum_{i \neq j} (f(z'_i) - f(z_i)) \right) \right] \\ &= \frac{1}{m} \mathbb{E}_{S, S', \sigma_j} \left[\sup_{f \in \mathcal{F}} \left(\sigma_j(f(z'_j) - f(z_j)) + \sum_{i \neq j} (f(z'_i) - f(z_i)) \right) \right] \\ &= \frac{1}{m} \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in \mathcal{F}} \left(\sum_i \sigma_i(f(z'_i) - f(z_i)) \right) \right] \end{aligned}$$

Rademacher Complexity

Proof: (continued) Taking the expectation over S

$$\begin{aligned}\mathbb{E}_S \left[\sup_{f \in \mathcal{F}} (L(f) - \hat{L}_S(f)) \right] &\leq \frac{1}{m} \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in \mathcal{F}} \left(\sum_i \sigma_i (f(z'_i) - f(z_i)) \right) \right] \\ &\leq \frac{1}{m} \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in \mathcal{F}} \left(\sum_i \sigma_i f(z'_i) \right) + \sup_{f \in \mathcal{F}} \left(\sum_i -\sigma_i f(z_i) \right) \right] \\ &= 2 \mathbb{E}_S [R(\mathcal{F}, S)]\end{aligned}$$

since σ_i and $-\sigma_i$ have the same distribution. This completes the proof.

To proceed further we need a more general inequality.

Need a More General Inequality

Generalizations of Hoeffding's Inequality:

- ▶ Azuma's Inequality
 - ▶ Generalization to **martingales with bounded differences**
- ▶ McDiarmid's Inequality: the one we need
 - ▶ Basically a corollary of Azuma's
 - ▶ Can be proved directly using the same proof as for Azuma's inequality
 - ▶ The proof is a combination of martingale arguments and the ones we used for Hoeffding's inequality

McDiarmid's Inequality

Lemma

McDiarmid's Inequality Consider independent random variables $X_1, X_2, \dots, X_n \in \mathcal{X}$ and a function $f : \mathcal{X}^n \rightarrow \mathbb{R}$. If for all $i \in \{1, \dots, n\}$ and all $x_1, \dots, x_n, x'_i \in \mathcal{X}$, the function f satisfies

$$|f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c,$$

then

$$\mathbb{P}[|f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)]| > t] \leq 2 \exp\left(-\frac{2t^2}{nc^2}\right)$$

or, equivalently, with probability at least $1 - \delta$,

$$|f(X_1, \dots, X_n) - \mathbb{E}[f(X_1, \dots, X_n)]| \leq c \sqrt{\frac{2}{n} \log\left(\frac{2}{\delta}\right)}.$$

McDiarmid's Inequality

Proof: To simplify the notation, let's denote $f \equiv f(X_1, \dots, X_n)$. Note that $M_i = \mathbb{E}[f|X_1, \dots, X_i]$, $M_0 = \mathbb{E}[f]$, is a martingale and consider martingale differences

$$V_i = \mathbb{E}[f|X_1, \dots, X_i] - \mathbb{E}[f|X_1, \dots, X_{i-1}].$$

Next, observe that $\mathbb{E}[V_i] = 0$ and

$$\sum_{i=1}^n V_i = \mathbb{E}[f|X_1, \dots, X_n] - \mathbb{E}[f] = f(X_1, \dots, X_n) - \mathbb{E}[f].$$

Now, we show that the martingale differences, V_i , are uniformly bounded. (Recall, having bounded random variables was the key component of Hoeffding's proof.) To this end, define, **for x in i th position**

$$L_i = \inf_x \mathbb{E}[f(X_1, \dots, x, \dots, X_n)|X_1, \dots, X_i] - \mathbb{E}[f|X_1, \dots, X_{i-1}]$$

$$U_i = \sup_x \mathbb{E}[f(X_1, \dots, x, \dots, X_n)|X_1, \dots, X_i] - \mathbb{E}[f|X_1, \dots, X_{i-1}]$$

and note that

$$L_i \leq V_i \leq U_i$$

McDiarmid's Inequality

Proof: (continued) Next, we show that $U_i - L_i$ is uniformly bounded

$$\begin{aligned} U_i - L_i &= \sup_{x, x'} (\mathbb{E}[f(X_1, \dots, x, \dots, X_n) - f(X_1, \dots, x', \dots, X_n) | X_1, \dots, X_{i-1}]) \\ &\leq c \end{aligned}$$

where the last equality is by the assumption on f and independence.

Hence, $V_i \in [L_i, L_i + c]$, where L_i is a function of X_1, \dots, X_{i-1} .

Therefore, L_i can be treated as a constant with respect to the conditional expectation $\mathbb{E}[\cdot | X_1, \dots, X_{i-1}]$

Meaning, we can use exactly the same arguments, the convexity of $e^{\lambda x}$, to prove

$$\mathbb{E} \left[e^{\lambda V_i} | X_1, \dots, X_{i-1} \right] \leq e^{\frac{\lambda^2 c^2}{8}}$$

McDiarmid's Inequality

Proof: (continued) Then, using the preceding inequality, we obtain

$$\begin{aligned}\mathbb{P}[f - \mathbb{E}[f] > t] &= \mathbb{P}\left[\sum_{i=1}^n V_i > t\right] \\ &\leq e^{-\lambda t} \mathbb{E}\left[\prod_{i=1}^n e^{\lambda V_i}\right] \\ &\leq e^{-\lambda t} \mathbb{E}\left[\prod_{i=1}^{n-1} e^{\lambda V_i} \mathbb{E}[e^{\lambda V_n} | X_1, \dots, X_{n-1}]\right] \\ &\leq e^{-\lambda t} e^{\frac{\lambda^2 c^2}{8}} \mathbb{E}\left[\prod_{i=1}^{n-1} e^{\lambda V_i}\right] \\ &\leq \dots \leq e^{-\lambda t} e^{\frac{n\lambda^2 c^2}{8}}\end{aligned}$$

Finally, the proof is completed by optimizing over λ , i.e., setting

$$\lambda = \frac{4t}{nc^2}.$$

Back to Generalization Bounds

Theorem 26.6 (Shais' book) Assume that for all z and $h \in \mathcal{H}$, we have $|\ell(h, z)| \leq c$, and recall $S = \{z_1, \dots, z_m\}$ is the sample. Then, with probability at least $1 - \delta$, for any $h \in \mathcal{H}$ (and in particular for $h = \text{ERM}_{\mathcal{H}}(S)$),

$$L(h) - \hat{L}_S(h) \leq 2 \mathbb{E}_S[R(\mathcal{F}, S)] + c \sqrt{\frac{2}{m} \log \left(\frac{2}{\delta} \right)}$$

Proof: First, note that random variable

$$\text{Rep}(\mathcal{F}, S) = \sup_{h \in \mathcal{H}} (L(h) - \hat{L}_S(h))$$

satisfies the bounded difference condition of McDiarmid's lemma with a constant $2c/m$.

Combining the bound from Lemma 26.2 and McDiarmid's lemma, we obtain that with probability at least $1 - \delta$,

$$\text{Rep}(\mathcal{F}, S) \leq \mathbb{E}[\text{Rep}(\mathcal{F}, S)] + c \sqrt{\frac{2}{m} \log \left(\frac{2}{\delta} \right)} \leq 2 \mathbb{E}_S[R(\mathcal{F}, S)] + c \sqrt{\frac{2}{m} \log \left(\frac{2}{\delta} \right)}$$

Reading

NTK/Lazy training regime:

- ▶ [A Note on Lazy Training in Supervised Differentiable Programming](#), by L. Chizat and F. Bach, Dec 2018.

An updated version:

[On Lazy Training Differentiable Programming](#), by L. Chizat and F. Bach, NIPS 2019.

- ▶ Chapter 8 in the recent monograph [Deep Learning Theory Lecture Notes](#), by Telgarsky, Oct 2021.

This paper puts in perspective many of the papers that we covered recently:

Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In Advances in neural information processing systems, 2018.

Simon S. Du, Xiyu Zhai, Barnabs Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In International Conference on Learning Representations, 2019.

Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. arXiv preprint arXiv:1904.11955, 2019.

Simon S. Du, Lee Jason D., Li Haochuan, Wang Liwei, and Zhai Xiyu. Gradient descent finds global minima of deep neural networks. In International Conference on Machine Learning (ICML), 2019.

Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In Advances in Neural Information Processing Systems, pages 81678176, 2018.

Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Proceedings of the 36th International Conference on Machine Learning, volume 97, pages 242252, 2019.

Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep ReLU networks. Machine Learning Journal, 2019.

Reading

Generalization bounds:

- ▶ PAC Learning and Generalization Theory - [ML2014] book:
PAC learning: Chapters 2-4; Rademacher Complexity: Chapter 26
(In particular, Theorem 26.5 and Lemmas 26.9 & 26.11)
In general, for PAC learning theory see: Chapters: 2-6, 26-31
- ▶ Generalization bounds for NNs
 - ▶ [A Priori Estimates For Two-layer Neural Networks](#), by Weinan et al., Jan 2019.
 - ▶ [Fine-Grained Analysis of Optimization and Generalization for Overparameterized Two-Layer Neural Networks](#), by Arora et al., Jan 2019.
 - ▶ See the references in the preceding papers

Have Fun!