

The Mathematical Theory of Neural Network-Based Machine Learning

Prof. Weinan E
Department of Mathematics
Princeton University

Typed by Yuelin Li

2019 APMA colloquium

1 Approximating Functions Using Samples

Given a set of samples from μ , $\{\mathbf{x}_j\}_{j=1}^n$, and $\{y_j = f^*(\mathbf{x}_j)\}_{j=1}^n$

Approximate f^* using $S = \{(\mathbf{x}_j, y_j)\}_{j=1}^n$.

Strategy:

Define some "hypothesis space" (set of functions) \mathcal{H}_m ($m \sim \dim$ of \mathcal{H}_m) :

$$\mathcal{H}_m = \{f(\cdot, \theta)\}$$

Best approximation minimizes the "population risk":

$$\mathcal{R}(\theta) = \mathbb{E} (f(\mathbf{x}, \theta) - f^*(\mathbf{x}))^2 = \int_{\mathbb{R}^d} (f(\mathbf{x}, \theta) - f^*(\mathbf{x}))^2 d\mu$$

In practice: Minimize the "empirical risk":

$$\mathcal{R}_n(\theta) = \frac{1}{n} \sum_j (f(\mathbf{x}_j, \theta) - y_j)^2 = \frac{1}{n} \sum_j (f(\mathbf{x}_j, \theta) - f^*(\mathbf{x}_j))^2$$

To choose the hypothesis space:

$f(\mathbf{x}) = \beta \cdot \mathbf{x} + \beta_0$. Generalized linear models (feature-based models): $f(\mathbf{x}) = \sum_{k=1}^m c_k \phi_k(\mathbf{x})$, where $\{\phi_k\}$ are linearly independent functions.

Simple neural networks: $f(\mathbf{x}) = \sum_k a_k \sigma(\mathbf{b}_k \cdot \mathbf{x} + c_k)$, where σ is some nonlinear function.

- $\sigma(x) = \max(x, 0)$, the ReLU (rectified linear units) function.

- $\sigma(x) = (1 + e^{-x})^{-1}$.

- $\sigma(x) = \cos(x)$

- $\sigma(x) = \tanh(x)$

Deep neural networks (DNN): compositions of functions of the form above. (σ is a scalar function)

2 Dynamical System Viewpoint to Deep Learning

Constructing nonlinear approximations through the flow map of a dynamical system:

$$\frac{dz(\mathbf{x}, t)}{dt} = \mathbf{F}(z(\mathbf{x}, t)), \quad z(0, \mathbf{x}) = \mathbf{V}\mathbf{x}$$

The flow map $x \rightarrow z(x, 1)$ is a nonlinear mapping. Simplest choice of nonlinear \mathbf{F} :

$$\mathbf{F}(z; \mathbf{U}, \mathbf{W}) = \mathbf{U}\sigma \circ (\mathbf{W}z)$$

Choose the optimal $\mathbf{U}, \mathbf{W}(\cdot), \alpha$ to approximate f^* by

$$f^*(\mathbf{x}) \sim \alpha \cdot z(\mathbf{x}, 1)$$

Application: High-dimensional control and PDEs high-dimensional stochastic control, nonlinear parabolic PDE...

Application: Molecular dynamics traditional dilemma cost VS accuracy.

$$E = E(x_1, x_2, \dots, x_i, \dots, x_N) \\ m_i \frac{d^2 x_i}{dt^2} = F_i = -\nabla_{x_i} E$$

Two ways to calculate E and F : 1) Empirical potentials: efficient but unreliable. The Lennard-Jones potential:

$$V_{ij} = 4\epsilon \left[\left(\frac{\sigma}{r_{ij}} \right)^{12} - \left(\frac{\sigma}{r_{ij}} \right)^6 \right], E = \frac{1}{2} \sum_{i \neq j} V_{ij}$$

2) Multi-scale modeling: computing the inter-atomic forces on the fly using QM, e.g. the Car-Parrinello MD. Accurate but expensive:

$$E = \langle \Psi_0 | H_e^{KS} | \Psi_0 \rangle, \mu \ddot{\phi}_i = H_e^{KS} \phi_i + \sum_j \Lambda_{ij} \phi_j$$

Issues: Basic questions we want to understand: high dimensionality; models are highly over-parametrized, classical machine learning theory would suggest overfitting; models are non-convex, yet simple gradient algorithms seem to work (compare with structural optimization in material science and chemistry)

More advanced questions: why deep networks seem to perform better than shallow ones? Why stochastic gradient descent seems to perform better than gradient descent? etc.

3 Classical Numerical Analysis

We define a well-posed math model (the hypothesis space, the loss function...) splines: hypothesis space = C^1 piecewise cubic polynomials the data

$$I_n(f) = \frac{1}{n} \sum_{j=1}^n (f(x_j) - y_j)^2 + \lambda \int |D^2 f(x)|^2 dx$$

finite elements: hypothesis space = C^0 piecewise polynomials.

We identify the right function spaces: direct and inverse approximation theorem (Bernstein and Jackson type theorems): f can be approximated by trig polynomials in L^2 to order s iff $f \in H^s$, $\|f\|_{H^s}^2 = \sum_{k=0}^s \|\nabla^k f\|_{L^2}^2$. Functions of interest are in the right spaces (PDE theory, real analysis, etc).

Optimal error estimates

- A priori estimates (for piecewise linear finite elements, $\alpha = 1/d, s = 2$)

$$\|f_m - f^*\|_{H^1} \leq C m^{-\alpha} \|f^*\|_{H^s}$$

- A posteriori estimates (say in finite elements):

$$\|f_m - f^*\|_{H^1} \leq C m^{-\alpha} \|f_m\|_h$$

Convergence rate is $O(n^{-1/2})$ independent of dimension.

Analyzing the error The empirical risk

$$\mathcal{R}_n(\theta) = \frac{1}{n} \sum_i (f(\mathbf{x}_i, \theta) - f^*(\mathbf{x}_i))^2$$

To estimate the population risk:

$$\mathcal{R}(\hat{\theta}) = \mathbb{E} \left(f(\mathbf{x}, \hat{\theta}) - f^*(\mathbf{x}) \right)^2, \quad \hat{\theta} = \operatorname{argmin} \mathcal{R}_n(\theta)$$

(Yaim Cooper: Such $\hat{\theta}$'s form a $m - n$ dimensional (smooth) manifold).

Difficulty: $\hat{\theta}$ is highly correlated with the data set S Study the worse case situation in the hypothesis space. Expect the optimal error rate to be $O(1/m) + O(1/\sqrt{n})$, both are Monte Carlo rates.

Generalization Gap

$$\mathcal{R}(\hat{\theta}) - \mathcal{R}_n(\hat{\theta}) = I(g) - I_n(g), \quad g(\mathbf{x}) = (f(\mathbf{x}, \theta) - f^*(\mathbf{x}))^2$$

$$I(g) = \int_{\mathbb{R}^d} g(\mathbf{x}) d\mu, \quad I_n(g) = \frac{1}{n} \sum_j g(\mathbf{x}_j)$$

- For fixed $g = h$, we have

$$|I(h) - I_n(h)| \sim \frac{1}{\sqrt{n}}$$

- For Lipschitz functions (Wasserstein distance)

$$\sup_{\|h\|_{Lip} \leq 1} |I(h) - I_n(h)| \sim \frac{1}{n^{1/d}}$$

- For functions in Barron space, to be defined later

$$\sup_{\|h\|_{\mathcal{B}} \leq 1} |I(h) - I_n(h)| \sim \frac{1}{\sqrt{n}}$$

4 Rademacher Complexity

Let \mathcal{H} be a set of functions, and $S = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ be a set of data points. Then, the Rademacher complexity of \mathcal{H} with respect to S is defined as

$$\hat{R}_S(\mathcal{H}) = \frac{1}{n} \mathbb{E}_\xi \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^n \xi_i h(\mathbf{x}_i) \right]$$

where $\{\xi_i\}_{i=1}^n$ i.i.d. random variables taking values ± 1 with equal probability.

Theorem 1. (Rademacher complexity and the generalization gap) Given a function class \mathcal{H} , for any $\delta \in (0, 1)$, with probability at least $1 - \delta$ over the random samples $S = (x_1, \dots, x_n)$

$$\begin{aligned} \sup_{h \in \mathcal{H}} \left| \mathbb{E}_x[h(\mathbf{x})] - \frac{1}{n} \sum_{i=1}^n h(\mathbf{x}_i) \right| &\leq 2\hat{R}_S(\mathcal{H}) + \sup_{h \in \mathcal{H}} \|h\|_\infty \sqrt{\frac{\log(2/\delta)}{2n}} \\ \sup_{h \in \mathcal{H}} \left| \mathbb{E}_x[h(\mathbf{x})] - \frac{1}{n} \sum_{i=1}^n h(\mathbf{x}_i) \right| &\geq \frac{1}{2}\hat{R}_S(\mathcal{H}) - \sup_{h \in \mathcal{H}} \|h\|_\infty \sqrt{\frac{\log(2/\delta)}{2n}} \end{aligned}$$

We want large space but low complexity. (Donsker space).

If \mathcal{H} contains a single function, then $\hat{R}_S(\mathcal{H}) \sim O(1/\sqrt{n})$

If \mathcal{H} contains functions that can fit any random values on S , then $\hat{R}_S(\mathcal{H}) \sim O(1)$

If \mathcal{H} = unit ball in Barron space: $\hat{R}_S(\mathcal{H}) \sim O(1/\sqrt{n})$

If \mathcal{H} = unit ball in Lipschitz space: $\hat{R}_S(\mathcal{H}) \sim O(1/n^{1/d})$

If \mathcal{H} = unit ball in C^0 : $R_S(\mathcal{H}) \sim O(1)$.

In summary, in low dimension, different approximation methods differ by their order of convergence. In high dimension, different machine learning models differ by the size of the function spaces they are associated with.

Within the right space, we have Monte Carlo convergence rates. Outside the space, we encounter curse of dimensionality.

Identify the right spaces (direct/inverse approx. theorems, complexity control)

- kernel method \rightarrow RKHS space
- shallow networks \rightarrow Barron space
- deep residual networks \rightarrow compositional function space

Well-posed models:

- explicit regularization gives rise to optimal error rates
- current results on implicit regularization severely limits the size of the function spaces

A priori estimates for regularized kernel method: fix any $\lambda > 0$,

$$L(f) = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - f^*(\mathbf{x}_i))^2 + \frac{\lambda}{\sqrt{n}} \|f\|_{\mathcal{H}_k}$$

Theorem 2. Assume $f^* \in \mathcal{H}_k$. Let

$$\hat{f}_n \stackrel{\text{def}}{=} \operatorname{argmin}_{f \in \mathcal{H}_k} L(f)$$

We have that

$$\mathbb{E}_x \left[\left| \hat{f}_n(\mathbf{x}) - f^*(\mathbf{x}) \right|^2 \right] \leq C(\lambda) \frac{\|f^*\|_{\mathcal{H}_k}}{\sqrt{n}}$$

5 Barron spaces

Two-layer neural networks:

$$\frac{1}{m} \sum_{j=1}^m a_j \sigma(\mathbf{b}_j^T \mathbf{x} + c_j)$$

Consider the function $f : D_0 = [-1.1]^d \mapsto \mathbb{R}$ of the following form

$$f(\mathbf{x}) = \int_{\Omega} a \sigma(\mathbf{b}^T \mathbf{x} + c) \rho(da, d\mathbf{b}, dc), \quad \mathbf{x} \in D_0$$

$\Omega = \mathbb{R}^1 \times \mathbb{R}^d \times \mathbb{R}^1$, ρ is a probability distribution on Ω . Fourier analog: $\rho(da, d\omega) = \delta(a - f(\omega))dad\omega$ (not normalizable).

$$\begin{aligned} f(\mathbf{x}) &= \int_{\mathbb{R}^d} f(\omega) \cos(\omega^T \mathbf{x}) d\omega = \int_{\mathbb{R}^1 \times \mathbb{R}^d} a \cos(\omega^T \mathbf{x}) \rho(da, d\omega) \\ \|f\|_{\mathcal{B}_p} &= \inf_{\rho} (\mathbb{E}_{\rho} [|a|^p (\|\mathbf{b}\|_1 + |c|)^p])^{1/p} \\ \mathcal{B}_p &= \text{completion of } \{f \in \mathcal{S} : \|f\|_{\mathcal{B}_p} < \infty\}, \quad \mathcal{B}_{\infty} \subset \cdots \mathcal{B}_2 \subset \mathcal{B}_1 \end{aligned}$$

It is a good idea because approximation by two layer networks becomes Monte Carlo integration:

$$f(\mathbf{x}) \sim f_m(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\mathbf{b}_j^T \mathbf{x} + c_j)$$

Theorem 3. (Direct Approximation Theorem) There exists an absolute constant C_0 such that

$$\|f - f_m\|_{L^2(D_0)} \leq \frac{C_0 \|f\|_{\mathcal{B}_2}}{\sqrt{m}}$$

Theorem 4. (Inverse Approximation Theorem) For $p > 1$, let

$$\mathcal{N}_{p,C} \stackrel{\text{def}}{=} \left\{ \frac{1}{m} \sum_{k=1}^m a_k \sigma(b_k^T \mathbf{x} + c_k) : \left(\frac{1}{m} \sum_{k=1}^m |a_k|^p (\|b_k\|_1 + c_k)^p \right)^{1/p} \leq C, m \in \mathbb{N}^+ \right\}$$

Let f^* be a continuous function. Assume there exists a constant C and a sequence of functions $f_m \in \mathcal{N}_{p,C}$ such that

$$f_m(\mathbf{x}) \rightarrow f^*(\mathbf{x})$$

for all $\mathbf{x} \in D_0$. Then there exists a probability distribution ρ on Ω , such that

$$f^*(\mathbf{x}) = \int a \sigma(\mathbf{b}^T \mathbf{x} + c) \rho(da, d\mathbf{b}, dc)$$

for all $\mathbf{x} \in D_0$.

Theorem 5. Let $\mathcal{F}_Q = \{f \in \mathcal{B}_1, \|f\|_{\mathcal{B}_1} \leq Q\}$. Then we have

$$\hat{\mathcal{R}}_n(\mathcal{F}_Q) \leq 2Q \sqrt{\frac{2 \ln(2d)}{n}}$$

...

6 Compositional Function Space

$$\begin{aligned} z_{l+1,L}(\mathbf{x}) &= z_{l,L}(\mathbf{x}) + \frac{1}{L} \mathbf{U}_l \sigma \circ (\mathbf{W}_{l^2 l, L}(\mathbf{x})) \\ z_{0,L}(\mathbf{x}) &= \mathbf{V} \mathbf{x}, \quad f(\mathbf{x}, \theta) = \alpha \cdot z_{L,L}(\mathbf{x}) \end{aligned}$$

A discrete flow with RHS in Barron space. Let $\{\rho_t\}$ be a family of prob distributions (for (\mathbf{U}, \mathbf{W})) such that $\mathbb{E}_{\rho_t}, g(\mathbf{U}, \mathbf{W})$ is integrable as a function of t for any continuous function g . Define:

$$\begin{aligned} z(\mathbf{x}, 0) &= \mathbf{V} \mathbf{x} \\ \frac{d}{dt} z(\mathbf{x}, t) &= \mathbb{E}_{(\mathbf{U}, \mathbf{W}) \sim \rho_t} \mathbf{U} \sigma \circ (\mathbf{W} z(\mathbf{x}, t)) \end{aligned}$$

Let $f_{\alpha, \{\rho_t\}, V}(\mathbf{x}) = \alpha^T z(\mathbf{x}, 1)$ and define

$$\begin{aligned} \frac{d}{dt} \mathbf{N}_p(t) &= (\mathbb{E}_{\rho_t} |\mathbf{U}|^p |\mathbf{W}|^p)^{1/p} \mathbf{N}_p(t), \quad \mathbf{N}_p(0) = \mathbf{I} \\ \|f\|_{\mathcal{D}_p} &= \inf_{f=f_{\alpha, \{\rho_t\}, V}} \|\alpha\|_p \|\mathbf{N}_p(1)\|_{p,p} \|\mathbf{V}\|_{p,p} \end{aligned}$$

The matrix operations are done entry-wise.

$$\begin{aligned} f(\mathbf{x}) &= \int_{\Omega} a \sigma(\mathbf{b}^T \mathbf{x} + c) \rho(da, d\mathbf{b}, dc) \\ \frac{d}{dt} z(\mathbf{x}, t) &= \mathbb{E}_{\rho(a, b, c)} \begin{bmatrix} a \\ 0 \\ 0 \end{bmatrix} \sigma \circ ([0, b^T, c] z(\mathbf{x}, t)), \\ z(\mathbf{x}, 0) &= \begin{bmatrix} 0 \\ \mathbf{x} \\ 1 \end{bmatrix} \end{aligned}$$

Then, $f(\mathbf{x}) = e_1^T z(\mathbf{x}, 1)$. $\mathcal{B}_2 \subset \mathcal{D}_2$. There exists constant $C > 0$, such that

$$\|f\|_{\mathcal{D}_2} \leq \sqrt{d+1} \|f\|_{\mathcal{B}_2}$$

holds for any $f \in \mathcal{B}_2$...

7 Regularized Model and a Priori Estimates

Regularized loss function:

$$\begin{aligned} J(\theta) &= L(\theta) + \lambda (\|\theta\|_{\mathcal{D}_1} + 1) \sqrt{\frac{2 \log(2d)}{n}} \\ \|\theta\|_P &= \|\alpha\|^T \left(I + \frac{1}{L} |\mathbf{U}_{L-1}| |\mathbf{W}_{L-1}| \right) \cdots \left(I + \frac{1}{L} |\mathbf{U}_1| |\mathbf{W}_1| \right) |\mathbf{V}| \end{aligned}$$

Theorem 6. (A-priori estimate) Assume that $f^* : [-1, 1]^d \rightarrow [-1, 1]$ such that $f^* \in \mathcal{D}_2$. Let

$$\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta)$$

Then there exist fixed constants C_0 and L_0 such that if $\lambda > C_0$ and $L > L_0$, then for any $\delta > 0$, with probability at least $1 - \delta$,

$$L(\hat{\theta}) \lesssim \frac{\|f^*\|_{\mathcal{D}_2}^2}{L} + \lambda (\|f^*\|_{\mathcal{D}_1}^3 + 1) \sqrt{\frac{\log(2d)}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}$$

Gradient descent for unregularized model Questions: can gradient descent algorithm achieve global minima for the empirical risk? can the solution obtained generalize (small population risk)? "implicit regularization". Over-parametrized regime $m \geq O(n^\alpha)$: CLT scaling

$$f_0(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{k=1}^m a_k^0 \sigma(\mathbf{x}^T \mathbf{b}_k^0)$$

Optimization: Empirical risk goes to 0?

With random (normal) initialization, GD can fit any target (random labels).

Generalization: Population risk (real error) small?

Error estimates involve ε or κ