

A Mathematical Perspective on Machine Learning

Weinan E

**Center for Machine Learning Research and
School of Mathematical Sciences**

Peking University

(z)activation function



data

data

data

data

goal

data

data

data

average error for all the data

data

data

data

data

model

m: number of free parameters

training data

test data

training error

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

square error

training error

test error

training error

training error

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

training error

test error

test error

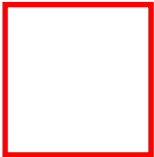
training error

when n is large, use a random



is a property distribution, so can be written as an expectation

approximate, finite sum



doesn't work in high dimention



convergence rate is independent of dimention



“New” approach: Let π be a probability distribution

$$f^*(\mathbf{x}) = \int_{\mathbb{R}^d} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})} \pi(d\boldsymbol{\omega}) = \mathbb{E}_{\boldsymbol{\omega} \sim \pi} a(\boldsymbol{\omega}) e^{i(\boldsymbol{\omega}, \mathbf{x})}$$

repalce $d\boldsymbol{\omega}$ written as an expectation

Let $\{\boldsymbol{\omega}_j\}$ be an **i.i.d. sample of π** , $f_m(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a(\boldsymbol{\omega}_j) e^{i(\boldsymbol{\omega}_j, \mathbf{x})}$

$$\mathbb{E}|f^*(\mathbf{x}) - f_m(\mathbf{x})|^2 = \frac{\text{var}(f)}{m}$$

does not suffer from CoD

Note: $f_m(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m a_j \sigma(\boldsymbol{\omega}_j^T \mathbf{x}) = \text{two-layer neural network}$ with $\sigma(z) = e^{iz}$.

Conclusion:

Functions of the **this type** (i.e. can be expressed as this kind of expectation) can be approximated by two-layer neural networks with a dimension-independent error rate.

Approximation theory for the *random feature model*

ω : feature vector

Ω : sapce of feature vector

- Let $\phi(\cdot; \mathbf{w})$ be a feature function parametrized by $\mathbf{w} \in \Omega$,
e.g. $\phi(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$.
We will **assume** that ϕ is **continuous** and Ω is **compact**.
- Let π_0 be a **fixed distribution** for the **random variable** \mathbf{w} .
- Let $\{\mathbf{w}_j^0\}_{j=1}^m$ be a **set of i.i.d samples** drawn from π_0 .

The random feature model (RFM) associated with the features $\{\phi(\cdot; \mathbf{w}_j^0)\}$ is given by

$$f_m(\mathbf{x}; \mathbf{a}) = \frac{1}{m} \sum_{j=1}^m a_j \phi(\mathbf{x}; \mathbf{w}_j^0).$$

linear combination with a_j

What spaces of functions are “well approximated” (say with the same convergence rate as in Monte Carlo) **by the random feature model?**

- In classical approximation theory, these are the Sobolev or Besov spaces: They are characterized by the convergence behavior for some specific approximation schemes.
- Direct and inverse approximation theorems.

Define the kernel function associated with the random feature model:

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w} \sim \pi_0} [\phi(\mathbf{x}; \mathbf{w}) \phi(\mathbf{x}'; \mathbf{w})]$$

Let \mathcal{H}_k be the reproducing kernel Hilbert space (RKHS) induced by the kernel k .

Probabilistic characterization: $a(\cdot)$: coefficient function

$f \in \mathcal{H}_k$ if and only if there exists $a(\cdot) \in L^2(\pi_0)$ such that

$$f(\mathbf{x}) = \int a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w}) d\pi_0(\mathbf{w}) = \mathbb{E}_{\mathbf{w} \sim \pi_0} a(\mathbf{w}) \phi(\mathbf{x}; \mathbf{w})$$

and

$$\|f\|_{\mathcal{H}_k}^2 = \int a^2(\mathbf{w}) d\pi_0(\mathbf{w}) = \mathbb{E}_{\mathbf{w} \sim \pi_0} a^2(\mathbf{w})$$

sqaure of l2 norm

the rest in : 2022_ICM_A Mathematical Perspective of Machine Learning_note2

