

Overfitting and Generalization Performance

October 12, 2020

Introduction

General Aim

Given training sample

$$(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$$

learn a predictor $h_n : \mathbb{R}^d \rightarrow \mathbb{R}$ that predicts y given new x .

Empirical Risk Minimization (ERM)

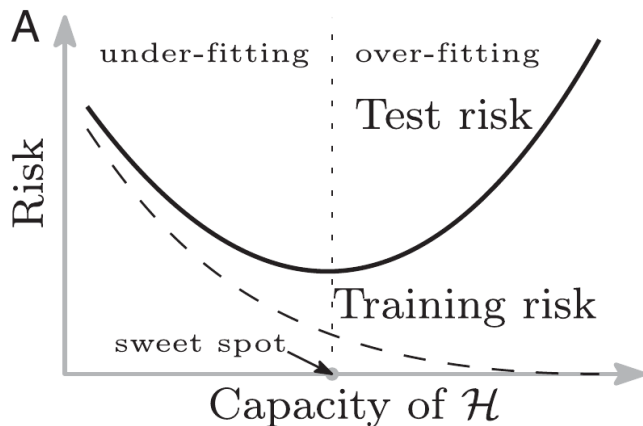
Minimize training risk: $\frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$ given a loss function ℓ .

- Find h_n that performs well on unseen data.
- Minimize true risk: $E[\ell(h(x), y)]$ where (x, y) drawn independently from P .

"Classical" thinking

- Finding a balance between underfitting and overfitting.
- "Bias-Variance Tradeoff"
- 0 training error does not tend to generalize well.
- Control function class \mathcal{H} implicitly or explicitly.

Generalization of performance



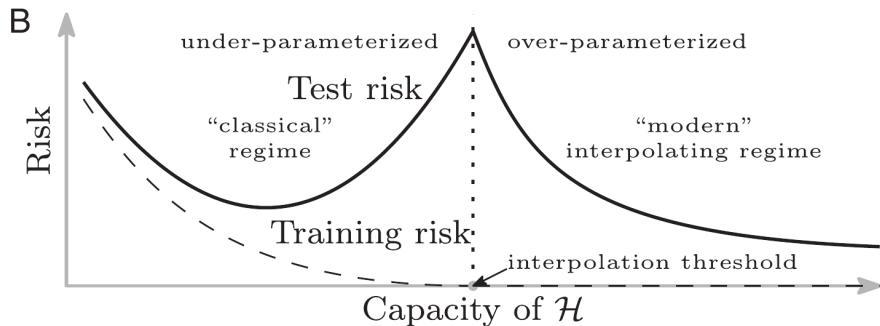
Classical curve from bias variance tradeoff.

- Modern ML methods such as large neural networks and other non-linear predictors have very low to no training risk
- NN architectures chosen such that interpolation can be achieved.
- Works even when training data have high levels of noise.

"Double Descent"

- "Double Descent" curve proposed and empirically observed to some extent.
- Curve the extends beyond the point of interpolation
- Risk decreases beyond this point, typically surpassing performance of classical stopping point.

Double Descent Curve



Explanations on Double Descent

Why should the test risk decrease even when empirical risk stays the same?

- Capacity of function class needs not suit the appropriate inductive bias for the problem.
- By having a larger function class, may find a function that matches the inductive bias better.
- Eg., smoother function, smaller norm, larger margin.

Empirical Evidence - Random Fourier Features

Let the function class \mathcal{H}_∞ be the Reproducing Kernel Hilbert Space (RKHS) corresponding to the Gaussian kernel.

We consider the following non-linear parametric model:

Random Fourier Features (RFF)

Let the function class \mathcal{H}_N consist of functions $h_n : \mathbb{R}^d \rightarrow \mathbb{C}$ of the form:

$$h(\cdot) = \sum_{k=1}^N a_k \phi(\cdot, v_k)$$

where $\phi(\cdot, v_k) := e^{i\langle v_k, \cdot \rangle}$

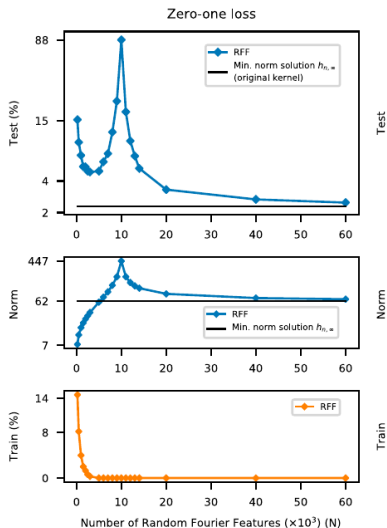
H_N has N parameters in \mathbb{C} , $\{a_1, \dots, a_n\}$.

As $N \rightarrow \infty$, H_N becomes a closer approximation to \mathcal{H}_∞

Empirical Evidence - Learning Procedure

- Given training sample $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$.
- Minimize empirical risk: $\frac{1}{n} \sum_{j=1}^n (h(x_j) - y_j)^2$ for $h \in \mathcal{H}_N$.
- When minimizer not unique ($N > n$), choose the minimizer with the coefficients (a_1, \dots, a_N) that have the smallest ℓ_2 norm.
- Let this predictor be: $h_{n,N} \in \mathcal{H}_N$.

Empirical Evidence - Results



Neural Networks. (Might be hard to explain why SGD is the inductive bias.)

Appendix on Approximation Theorem

On why they choose a function with a smaller norm in RKHS.