

# OPTIMIZATION FOR MACHINE LEARNING

Recall that ERM rule aims:

$$\hat{f}_{\text{ERM}} \in \underset{f \in \mathcal{H}}{\operatorname{argmin}} \hat{\mathcal{R}}_S(f) := \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

**Parameterization:** This is the traditional methodology in statistics.

For a compact set  $\Theta \subset \mathbb{R}^d$ , a class of parameterized predictors is considered:

$$\mathcal{H}_\Theta = \{x \mapsto f_\theta(x) : \theta \in \Theta, x \in \mathcal{X}\}.$$

Then,

$$\hat{\theta}_{\text{ERM}} \in \underset{\theta \in \Theta}{\operatorname{argmin}} \hat{\mathcal{R}}_S(f_\theta),$$

leads to the empirical risk minimizer  $f_{\hat{\theta}_{\text{ERM}}}$ .

EXAMPLE (ReLU neural networks)

$$\Theta \subset \mathbb{R}^{m \times d} \times \mathbb{R}^m, \quad \mathcal{X} = \mathbb{R}^d, \quad \mathcal{Y} = \mathbb{R}$$

$$\mathcal{H}_\Theta = \left\{ x \mapsto \sum_{i=1}^m c_i \sigma(w_i^T x) : (w, c) \in \Theta, x \in \mathbb{R}^d \right\}$$

where

$\sigma(z) = \max\{0, z\}$  is the so-called ReLU activation function.

$$w = \begin{bmatrix} - & w_1^T & - \\ & \vdots & \\ - & w_m^T & - \end{bmatrix} \in \mathbb{R}^{m \times d}, \quad \rightarrow \text{hidden-layer weights}$$

$$c \in \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} \in \mathbb{R}^m \quad \rightarrow \text{output layer weights.}$$

$m \in \mathbb{Z}_+$  is the network-width.

The trainable (or learnable) parameter is  $\theta = (w, c) \in \Theta$ .

Thus, given a training set  $S = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, 2, \dots, n\}$ ,

$$\hat{\theta}_{\text{ERM}} \in \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \left( y_i - f_\theta(x_i) \right)^2,$$

$$f_\theta(x) := \sum_{i=1}^m c_i \sigma(w_i^T x).$$

**Vital question:** Can we find  $\hat{\theta}_{\text{ERM}}$  in a computationally efficient way?

Thus, we cast ERM as an optimization problem:

$$g(\theta) \triangleq \hat{R}_S(f_\theta), \theta \in \mathcal{H}$$

$$\hat{\theta}_{\text{ERM}} \in \underset{\theta \in \mathcal{H}}{\operatorname{argmin}} g(\theta)$$

## Uniform Grid Method and Curse of Dimensionality:

Consider  $\mathcal{H} = \{\theta \in \mathbb{R}^P : 0 \leq \theta_i \leq 1, i=1,2,\dots,P\}$  [hypercube].

### ALG1 (Uniform Grid)

input:  $q \in \mathbb{Z}_+$

Step 0: Form  $q^P$  points:

$$\theta_\alpha = \left( \frac{2i_1-1}{2q}, \frac{2i_2-1}{2q}, \dots, \frac{2i_P-1}{2q} \right), \alpha = (i_1, i_2, \dots, i_P) \in \{1, 2, \dots, q\}^P.$$

Step 1: Find the optimal point in the grid:

$$\alpha^* \in \underset{\alpha \in \{1, \dots, q\}^P}{\operatorname{argmin}} g(\theta_\alpha)$$

Step 2: Return  $(\theta_{\alpha^*}, g(\theta_{\alpha^*}))$ .

### Assumption 1 (Lipschitz continuity)

$g: \mathbb{R}^P \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous on  $\mathcal{H}$ :

$$|g(\theta) - g(\theta')| \leq L \cdot \|\theta - \theta'\|_\infty, \forall \theta, \theta' \in \mathcal{H}.$$

### THEOREM 1 (Oracle complexity of Uniform Grid)

Let  $g^*$  be a global optimum value of  $g$ . Then,

$$g(\theta_{\alpha^*}) - g^* \leq \frac{L}{2q}, \text{ for any } q \in \mathbb{Z}_+.$$

Then,  $\left( \left\lfloor \frac{L}{2\varepsilon} \right\rfloor + 1 \right)^P$  function evaluations should be performed

to achieve  $g(\theta_{\alpha^*}) - g^* \leq \varepsilon$  for  $\varepsilon > 0$ .

Proof: For a multi-index  $\alpha = (i_1, i_2, \dots, i_P)$ , define

$$X_\alpha = \{\theta \in \mathbb{R}^P : \|\theta - \theta_\alpha\|_\infty \leq \frac{1}{2q}\}.$$

Then,  $\bigcup_{\alpha \in \{1, \dots, q\}^P} X_\alpha = \mathcal{H}$ . Let  $\theta^* \in \mathcal{H}$  be the globally

optimum solution. Then,  $\exists \bar{\alpha} \in \{1, \dots, q\}^P$  s.t.  $\theta^* \in X_{\bar{\alpha}}$ .

$$\Rightarrow \|\theta^* - \theta_{\bar{\alpha}}\|_\infty \leq \frac{1}{2q}.$$

$$\Rightarrow g(\theta_{\alpha^*}) - g^* \leq g(\theta_{\bar{\alpha}}) - g^* \leq \|\theta_{\bar{\alpha}} - \theta^*\|_\infty \leq \frac{1}{2q}.$$

### Remark (Dimension)

For  $\epsilon$ -optimality, a naive optimization method needs

$$O\left(\frac{1}{\epsilon^p}\right) \text{ computations if there are } p$$

parameters. In our ERM search, this would mean

$$O\left(\frac{n}{\epsilon^p}\right) \text{ computations.}$$

In modern deep learning,  $p$  is huge. Especially in the overparameterized regime, where  $p \gg n$ ,  $O\left(\frac{n}{\epsilon^p}\right)$  would imply completely impractical optimization.

For  $\epsilon = 0.001$  error,  $O(n \cdot 10^{3p})$  computations.

Solutions: For tractability, inspired by convex optimization theory, local iterative methods (e.g., gradient descent) are used. As we will see, they will lead to dimension-free (i.e., independent of  $p$ ) iteration complexities.

Furthermore, since  $n$  is very large in deep learning, stochastic gradient methods are developed.

As we will study, these will bring elimination of  $n$  in complexity bounds.