

# In Defense of Uniform Convergence: Generalization via derandomization with an application to interpolating predictors (Jeffrey Negrea, 2020 ICML)

---

Presented by Insung Kong

2022 1/17

Seoul National University

# Table of Contents

- ① Introduction
- ② Preliminaries
- ③ Glivenko–Cantelli property and surrogate
- ④ Example : overparameterized linear regression

# Introduction

---

- One of the central problems in learning theory is to explain the statistical performance of deep learning algorithms.
- The bulk of recent work on this problem implicitly assumes the classifier learned by SGD belongs to a class for which there is a uniform and tight bound on the generalization error.
- The traditional notions of uniform convergence (Glivenko–Cantelli classes) are may not applicable.

- Paper extend the concept of uniform convergence to the setting of sequences of learning problems of increasing complexity, defined by the structural Glivenko–Cantelli property.
- They introduce a approach to relate sequences of learning problems which are not structural Glivenko–Cantelli to ones that may be.
- Paper applicated it to overparameterized linear regression.

# Preliminaries

---

# Notations

- $Z_1, \dots, Z_n$  : i.i.d. random elements in a space  $\mathcal{S}$  with common distribution  $\mathcal{D}$
- $S = (Z_1, \dots, Z_n)$  : the training sample
- Loss function  $l : \mathcal{H} \times \mathcal{S} \rightarrow \mathbb{R}_+$  for a space  $\mathcal{H}$  of hypotheses.
- $\mathcal{M}_1(\mathcal{H})$  : space of distribution on  $\mathcal{H}$ .
- For  $Q \in \mathcal{M}_1(\mathcal{H})$ , the loss and risk are defined to be

$$\ell(Q, z) = \int \ell(h, z)Q(dh), \quad L_{\mathcal{D}}(Q) = \int \ell(Q, z)\mathcal{D}(dz)$$

- $L_S(Q) = L_{\hat{\mathcal{D}}_n}(Q)$  : empirical risk
- For  $h \in \mathcal{H}$   $L_{\mathcal{D}}(h) = L_{\mathcal{D}}(\delta_h)$  and  $L_S(h) = L_S(\delta_h)$
- $\hat{h} \in \mathcal{H}$  : learned classifier, which is random.

# Glivenko–Cantelli property and surrogate

---



## Glivenko–Cantelli property

### Definition 3.1

Let  $(\mathcal{S}, \mathcal{F}, \mathcal{D})$  be a probability space. Let  $\mathcal{H}$  be a collection of measurable functions on  $(\mathcal{S}, \mathcal{F}, \mathcal{D})$ .

Then  $\mathcal{H}$  has the Glivenko–Cantelli property if

$$\lim_{p \rightarrow \infty} \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \left| \mathcal{D}h - \widehat{\mathcal{D}}_p h \right| \right] = 0$$

where  $Ph = \int h(x)P(dx)$  and  $\widehat{\mathcal{D}}_p$  is the empirical distribution of an IID sample of size  $p$  from  $\mathcal{D}$ .

# Structural Glivenko–Cantelli property

## Definition 3.2

Let  $\{(\mathcal{S}^{(p)}, \mathcal{F}^{(p)}, \mathcal{D}^{(p)})\}_{p \in \mathbb{N}}$  be a sequence of probability space. Let  $\mathcal{H}^{(p)}$  be a collection of measurable functions on  $(\mathcal{S}^{(p)}, \mathcal{F}^{(p)}, \mathcal{D}^{(p)})$ .

Then  $\mathcal{H}^{(\cdot)}$  has the structural  $(\mathcal{D}^{(\cdot)}, n_{(\cdot)})$ -Glivenko–Cantelli property, if

$$\lim_{p \rightarrow \infty} \mathbb{E} \left[ \sup_{h \in \mathcal{H}^{(p)}} \left| \mathcal{D}^{(p)} h - \widehat{\mathcal{D}^{(p)}}_{n_p} h \right| \right] = 0$$

where  $Ph = \int h(x)P(dx)$  and  $\widehat{\mathcal{D}^{(p)}}_p$  is the empirical distribution of an IID sample of size  $p$  from  $\mathcal{D}^{(p)}$ .

## Lemma 3.3 (Surrogate decomposition)

*For every random element  $Q$  in  $\mathcal{M}_1(\mathcal{H})$*

$$\begin{aligned}\mathbb{E} \left[ L_{\mathcal{D}}(\hat{h}) - L_S(\hat{h}) \right] &= \mathbb{E} \left[ L_{\mathcal{D}}(\hat{h}) - L_{\mathcal{D}}(Q) \right] \\ &\quad + \mathbb{E} \left[ L_{\mathcal{D}}(Q) - L_S(Q) \right] \\ &\quad + \mathbb{E} \left[ L_S(Q) - L_S(\hat{h}) \right],\end{aligned}$$

*provided the three expectations on the r.h.s. are finite*

## Example : overparameterized linear regression

---

- $X_i \stackrel{\text{iid}}{\sim} N_{1 \times d}(0, \Sigma_n)$  : random row vectors
- $X^\top = (X_1^\top, \dots, X_n^\top)^\top$
- $(Y_i | X) \stackrel{\text{ind}}{\sim} N(X_i \beta_n, \sigma^2)$ .
- $\ell(\beta, (x, y)) = (x\beta - y)^2$ .
- $\hat{\beta}(X, Y) = (X^\top X)^+ X^\top Y$  : minimum norm interpolating linear predictor

## Failure of uniform convergence

**Lemma 4.1 (Failure of uniform convergence for overparameterized linear regression)**

*There is no sequence of measurable set  $\{A_n\}_{n \in \mathcal{N}}$  such that  $\mathbb{P}((X, Y) \in A_n) > 2/3$  for all  $n \in \mathcal{N}$  and for which*

$$\limsup_{n \rightarrow \infty} \mathbb{E} \sup_{(\tilde{X}, \tilde{Y}) \in A_n} \left| L_D(\hat{\beta}(\tilde{X}, \tilde{Y})) - L_S(\hat{\beta}(\tilde{X}, \tilde{Y})) \right| \leq \frac{3}{2} L_D(\beta)$$

# The sequence of surrogate hypothesis classes

We will consider the surrogate given by the minimum norm interpolating predictor for the training data with label noise removed. i.e.

$$\hat{\beta}_0 = (X'X)^+ X'X\beta$$

**Lemma 4.2** ((The sequence of surrogate hypothesis classes is SGC)

$\left\{ \hat{\beta}_0(S) : S \in \mathcal{S}^{(n)} \right\}_{n \in \mathbb{N}}$  is  $(\mathcal{D}^{(n)}, n)$  - SGC. Quantitatively, for a universal constant  $C > 0$ ,

$$\mathbb{E} \sup_{(X_0, Y_0) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n} \left| L_D \left( \hat{\beta}_0(X_0, Y_0) \right) - L_S \left( \hat{\beta}_0(X_0, Y_0) \right) \right| \leq C \frac{\sigma^2 + \|\beta_n\|^2 \|\Sigma_n\| \max \left( \sqrt{r_0(\Sigma_n)}, r_0(\Sigma_n) / \sqrt{n} \right)}{\sqrt{n}}$$

## Surrogate decomposition of $\hat{\beta}$

**Lemma 4.3 (Surrogate decomposition of  $\hat{\beta}$ )**

$$\begin{aligned} & L_D(\hat{\beta}) - L_S(\hat{\beta}) \\ &= \left( L_S(\hat{\beta}_0) - L_S(\hat{\beta}) \right) + \left( L_D(\hat{\beta}) - L_D(\hat{\beta}_0) \right) \\ & \quad + \left( L_D(\hat{\beta}_0) - L_S(\hat{\beta}_0) \right), \end{aligned}$$

*with*

$$\begin{aligned} L_S(\hat{\beta}_0) - L_S(\hat{\beta}) &= \frac{1}{n} \|Z\|^2 \\ L_D(\hat{\beta}) - L_D(\hat{\beta}_0) &= \text{Tr} \left( X (X'X)^+ \Sigma (X'X)^+ X' Z Z' \right), \\ L_D(\hat{\beta}_0) - L_S(\hat{\beta}_0) &= \sigma^2 - \frac{\|Z\|^2}{n} + \beta' P(X)^\perp \Sigma P(X)^\perp \beta. \end{aligned}$$



## Theorem 4.4 (Expected risk bound for overparameterized linear regression)

For a universal constant  $C > 0$ ,

$$\mathbb{E}L_D(\hat{\beta}) \leq \sigma^2 + C \frac{\sigma^2 + \|\beta_n\|^2 \|\Sigma_n\| \max\left(\sqrt{r_0(\Sigma_n)}, r_0(\Sigma_n)/\sqrt{n}\right)}{\sqrt{n}} + c\sigma^2 \left(\frac{k_n^*}{n} + \frac{n}{R_{k_n^*}(\Sigma_n)}\right)$$

In particular, if  $\{\Sigma_n\}_{n \in \mathbb{N}}$  and  $\left\{\|\beta_n\|^2 \|\Sigma_n\|\right\}_{n \in \mathbb{N}}$  is bounded then

$$\mathbb{E}L_D(\hat{\beta}) \rightarrow \sigma^2.$$