



# **Bootstrap and Resampling Methods**

## **Lecture 11: Bootstrap in Statistical Learning**

Elia Lapenta (ENSAE)

# Introduction

The background features a minimalist design with large, overlapping triangles in various shades of teal and light blue. The triangles are arranged in a way that creates a sense of depth and movement, with some pointing upwards and others downwards, meeting at sharp angles.

# Overview of today's class

- ▶ To learn an unknown function we often need to select some **tuning parameters**
  - ▶ **The bootstrap** can be used for this task
- ▶ **Boosting**: method that can reduce the "bias" of a weak learner
  - ▶ We will explore the connection between the **Wild Bootstrap** and the  $L_2$ -boosting for *linear* estimators
- ▶ **Bagging**: method based on the bootstrap to stabilize a high-variance estimator

Throughout this class we will focus on a prototype estimation method: **Smoothing Splines**

# Overview of Smoothing Spline Estimation

The slide features a white background with the title text centered. At the bottom, there are two large, overlapping teal-colored geometric shapes that form a wide 'V' or chevron pattern, pointing towards the center of the slide.

# Smoothing Spline Estimation

We observe  $\{Y_i, X_i\}_{i=1}^n$  with  $X_i \in [0, 1]$ . The model is

$$Y_i = f_0(X_i) + \varepsilon_i \text{ with } \mathbb{E}\{\varepsilon|X\} = 0$$

We want to estimate  $f_0$  **nonparametrically**, i.e. without assuming that  $f$  has a specific functional form (e.g., it is linear in  $X$ ). We just assume that  $f_0$  is a *smooth* function.

The **Smoothing Spline** method estimates  $f$  as

$$\hat{f}_\lambda := \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n [Y_i - f(X_i)]^2 + \lambda \int_0^1 |f^{(2)}(x)|^2 dx$$

where  $\lambda$  is a penalization parameter,  $f^{(2)}$  is the second derivative of  $f$ , and

$\mathcal{F} :=$  Class of functions that are twice continuously differentiable

# Smoothing Spline Estimation: The Role of $\lambda$

$$\hat{f}_\lambda := \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n [Y_i - f(X_i)]^2 + \lambda \int_0^1 |f^{(2)}(x)|^2 dx$$

$\lambda$  is a **penalization parameter**

- ▶ The larger  $\lambda$  the higher the penalty from *variations* of  $f$
- ▶  $\lambda = 0 \Rightarrow$  no penalty, and we will have  $Y_i = \hat{f}_\lambda(X_i)$  (**overfitting**)
- ▶  $\lambda = \infty \Rightarrow \hat{f}_\lambda$  is a linear function

$\lambda$  avoids overfitting **but** introduces a **regularization bias**: How do we select  $\lambda$  ?

## Selection of $\lambda$ by Bootstrap

# An Ideal Objective Function

Given the available data  $\{Y_i, X_i\}_{i=1}^n$  we define the  $L_2$  risk as

$$\widehat{R}(\lambda) := \mathbb{E}_X \left[ f_0(X) - \widehat{f}_\lambda(X) \right]^2 + \sigma_\varepsilon^2 = \mathbb{E}_{Y,X} \left[ Y - \widehat{f}_\lambda(X) \right]^2$$

where  $\mathbb{E}_X$  considers as random **only**  $X$  but **not**  $\widehat{f}_\lambda$ ,  $\mathbb{E}_{Y,X}$  considers as random only  $(Y, X)$  but not  $\widehat{f}_\lambda$ , and  $\sigma_\varepsilon^2$  is the variance of  $\varepsilon$ .

► Notice that  $\widehat{R}(\lambda)$  is random: it depends on the data because  $\widehat{f}_\lambda$  depends on  $\{Y_i, X_i\}_{i=1}^n$

Our **ideal** selection of  $\lambda$  would minimize the expectation of the  $L_2$  risk:

$$\mathbb{E}_{\text{sample}} \widehat{R}(\lambda)$$

where  $\mathbb{E}_{\text{sample}}$  is the expectation with respect to the sample  $\{Y_i, X_i\}_{i=1}^n$

**The Bootstrap principle will help us creating an estimator of  $\mathbb{E}_{\text{sample}} \widehat{R}$**



# Bootstrap Counterpart of the $L_2$ Risk

$$\mathbb{E}_{sample} \widehat{R}(\lambda) = \mathbb{E}_{sample} \mathbb{E}_{Y,X} \left[ Y - \widehat{f}_\lambda(X) \right]^2$$

- ▶  $(Y, X) \sim P^{Y,X}$ ,  $(Y_i, X_i) \sim \text{iid } P^{Y,X}$ , and  $\widehat{f}_\lambda$  depends on the sample data  $\{Y_i, X_i\}_{i=1}^n$
- ▶ By the bootstrap principle we replace the population with our sample  $\{Y_i, X_i\}_{i=1}^n \Rightarrow P^{Y,X}$  is replaced by

$$\mathbb{P}_n(A) = \frac{1}{n} \sum_{i=1}^n 1\{(Y_i, X_i) \in A\}$$

**Hence:**  $(Y^*, X^*) \sim \mathbb{P}_n$ ,  $(Y_i^*, X_i^*) \sim \text{iid } \mathbb{P}_n$ ,  $\widehat{f}_\lambda^*$  depends on  $\{Y_i^*, X_i^*\}_{i=1}^n$

$$\mathbb{E}_{sample}^* \mathbb{E}_{Y^*, X^*}^* \left[ Y^* - \widehat{f}_\lambda^*(X^*) \right]^2 \text{ where } \mathbb{E}_{Y^*, X^*}^* \left[ Y^* - \widehat{f}_\lambda^*(X^*) \right]^2 = \frac{1}{n} \sum_{i=1}^n \left[ Y_i - \widehat{f}_\lambda^*(X_i) \right]^2$$

## Bootstrap Counterpart of the $L_2$ Risk , ct'ed

$$\mathbb{E}_{sample}^* \mathbb{E}_{Y^*, X^*}^* \left[ Y^* - \hat{f}_\lambda^*(X^*) \right]^2 = \mathbb{E}_{sample}^* \frac{1}{n} \sum_{i=1}^n \left[ Y_i - \hat{f}_\lambda^*(X_i) \right]^2$$

with  $\hat{f}_\lambda^*$  that depends on  $\{Y_i^*, X_i^*\}_{i=1}^n$  and  $\mathbb{E}_{sample}^*$  expectation considering  $\{Y_i^*, X_i^*\}_{i=1}^n$  iid drawn with replacement from  $\{Y_i, X_i\}_{i=1}^n$ .

**We approximate  $\mathbb{E}_{sample}^*$  by the Monte-Carlo Algorithm:**

- ▶ Extract with replacement  $n$  observations from  $\{Y_i, X_i\}_{i=1}^n$  to get  $\{Y_i^*, X_i^*\}_{i=1}^n$
- ▶ Compute  $\hat{f}_\lambda^*$  and  $\hat{R}_b^*(\lambda) = (1/n) \sum_{i=1}^n \left[ Y_i - \hat{f}_\lambda^*(X_i) \right]^2$
- ▶ Repeat the above steps  $B$  times to get  $\{\hat{R}_b^*(\lambda) : b = 1, \dots, B\}$  and approximate

$$\mathbb{E}_{sample}^* \mathbb{E}_{Y^*, X^*}^* \left[ Y^* - \hat{f}_\lambda^*(X^*) \right]^2 \approx \frac{1}{B} \sum_{b=1}^B \hat{R}_b^*(\lambda). \quad \text{So} \quad \lambda^* = \arg \min_{\lambda} \frac{1}{B} \sum_{b=1}^B \hat{R}_b^*(\lambda)$$

# $L_2$ boosting and Wild Bootstrap

# $L_2$ Boosting and Smoothing Splines

$L_2$  **boosting** is a method used to reduce the bias of an estimator. It does so by extracting information from the residuals in an iterative fashion.

$$\hat{f}^Y := \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n [Y_i - f(X_i)]^2 + \lambda \int_0^1 \left| f^{(2)}(x) \right|^2 dx$$

## $L_2$ boosting algorithm

**Initial** Compute  $\hat{f}^Y$  (weak learner) and set  $\hat{f}_0 = \hat{f}^Y$ .

For  $m = 1, \dots, M$

**Step 1** Compute the residuals  $\hat{\varepsilon}_i = Y_i - \hat{f}_{m-1}(X_i)$

**Step 2** Compute  $\hat{f}^{\hat{\varepsilon}}$  by minimizing the above objective function where  $Y_i$  is replaced by  $\hat{\varepsilon}_i$ .

**Step 3** Update  $\hat{f}_m = \hat{f}_{m-1} + \hat{f}^{\hat{\varepsilon}}$

# $L_2$ boosting and the Wild Bootstrap

$L_2$  boosting can be seen as a **bias correction** based on the **Wild Bootstrap**

- First define  $\hat{f}^n := (\hat{f}(X_1), \dots, \hat{f}(X_n))^T$ . From the theory of smoothing splines (*Reproducing Kernel Hilbert Spaces/Support Vector Machines*) we know that

$$\hat{f}^n = SY^n$$

where  $Y^n := (Y_1, \dots, Y_n)^T$  and  $S$  is a matrix that depends only on the data  $\{X_i\}_{i=1}^n$

- So, the 1st iteration in the  $L_2$  boosting is

$$\hat{f}_1^n = SY^n + S(Y^n - SY^n) = SY^n + S(I - S)Y^n$$

## $L_2$ boosting and the Wild Bootstrap , ct'ed

- Now, generate an artificial sample by the **Wild Bootstrap**

$$Y_i^* = \hat{f}(X_i) + \xi_i \hat{\varepsilon}_i$$

where  $\hat{\varepsilon}_i = Y_i - \hat{f}(X_i)$  and  $\{\xi_i\}_{i=1}^n$  are iid bootstrap weights independent from the sample with  $\mathbb{E}\xi = 0$  and  $\mathbb{E}\xi^2 = 1$ . In vector form

$$Y^{*n} = \hat{f}^n + (\xi \hat{\varepsilon})^n$$

with  $Y^{*n} = (Y_1^*, \dots, Y_n^*)^T$  and  $(\xi \hat{\varepsilon})^n = (\xi_1 \hat{\varepsilon}_1, \dots, \xi_n \hat{\varepsilon}_n)^T$ .

- Since  $\hat{f}^n = SY^n$ , the smoothing spline estimator in the bootstrap world is

$$\hat{f}^{*n} = SY^{*n} = S\hat{f}^n + S(\xi \hat{\varepsilon})^n = SSY^n + S(\xi \hat{\varepsilon})^n$$

- So, the **Bias in the bootstrap world** is

$$E_\xi^* \hat{f}^{*n} - \hat{f}^n = SSY^n - SY^n = -S(I - S)Y^n$$

because  $E_\xi^* S(\xi \hat{\varepsilon})^n = 0$  (as  $E_\xi^* \xi = 0$ ).

## $L_2$ boosting and the Wild Bootstrap , ct'ed

- We can use the **bias in the bootstrap world** to estimate  $Bias(\hat{f}^n)$ . So, the *bias corrected* estimator is

$$\hat{f}^n - \widehat{Bias}(\hat{f}^n) = \mathcal{S}Y^n + \mathcal{S}(I - \mathcal{S})Y^n$$

which corresponds to the **boosted estimator after one boosting step**.

- So, the **boosted estimator** after one boosting iterations corresponds to a **bootstrap bias-corrected estimator** where the bias is estimated by the Wild bootstrap
- We can also generalize this to **any number of boosting steps**!

**Note** that the above procedure is valid for any type of estimator that is linear in  $Y^n$ . Such a correspondence does not hold for estimators that are not linear in  $Y^n$  (e.g., neural nets)

# The Role of $M$ (=number of iterations)

## Idea behind the $L_2$ Boosting:

1. Extract information from the residuals (i.e. estimate the model that considers the residuals as a response variable)
2. Add this piece of information to the "initial" estimate

How does the estimator behave when we increase the booting iterations?

- If we iterate too many times we run in **overfitting!**

**Let us analyze the Mean Squared Error of  $\hat{f}$  as a function of  $M$ .** Let  $\{\lambda_k : j = 1, \dots, K\}$  be the eigenvalues of the matrix  $\mathcal{S}$ . We have  $\lambda_k \in [0, 1]$ .

Then, (see Buhlmann and Yu , 2003)

$$MSE = Bias^2(\hat{f}_M) + Var(\hat{f}_M)$$

$$\text{with } Bias(\hat{f}_M) \sim diag((1 - \lambda_k)^{2M}) \text{ and } Var(\hat{f}_M) \sim [1 - (1 - \lambda_k)^{M+1}]^2$$



# Bagging



# Bagging as a Variance Reduction Technique

**Bagging**= bootstrap aggregating.

**Main Idea:** Extract bootstrap samples from the original sample, compute your estimator on each bootstrap sample, and then average out the results.

- ▶ Draw with replacement  $n$  observations from your sample  $\{Y_i, X_i\}_{i=1}^n$  to get  $\{Y_i^*, X_i^*\}_{i=1}^n$
- ▶ compute the estimator  $\hat{f}^*$  on the bootstrap sample
- ▶ Repeat the above steps  $B$  times to obtain  $\{\hat{f}_b^* : b = 1 \dots, B\}$
- ▶ Average the results:

$$\hat{f}_{Bagged} = \frac{1}{B} \sum_{b=1}^B \hat{f}_b^*$$

# Why Bagging Works?

If we could, it would be better to use  $\mathbb{E}_{sam}\hat{f}(x)$  instead of  $\hat{f}(x)$ .

- In fact, as already noticed

$$L_2 \text{ risk of } \hat{f} = \mathbb{E}_{Y,X} \left[ f(X) - \hat{f}(X) \right]^2 + \sigma_\varepsilon^2 = \mathbb{E}_{Y,X} \left[ Y - \hat{f}(X) \right]^2$$

- Next, notice that

$$\begin{aligned} \mathbb{E}_{sam} \mathbb{E}_{Y,X} \left[ Y - \hat{f}(X) \right]^2 &= \mathbb{E}_{sam} \mathbb{E}_{Y,X} \left[ Y - \mathbb{E}_{sam} \hat{f}(X) + \mathbb{E}_{sam} \hat{f}(X) - \hat{f}(X) \right]^2 \\ &= \mathbb{E}_{sam} \mathbb{E}_{Y,X} \left[ Y - \mathbb{E}_{sam} \hat{f}(X) \right]^2 + \mathbb{E}_{sam} \mathbb{E}_{Y,X} \left[ \mathbb{E}_{sam} \hat{f}(X) - \hat{f}(X) \right]^2 \\ &\geq \mathbb{E}_{Y,X} \left[ Y - \mathbb{E}_{sam} \hat{f}(X) \right]^2 = L_2 \text{ risk of } \mathbb{E}_{sample} \hat{f} \end{aligned}$$

where the cross-product in the second equality is

$$\mathbb{E}_{sam} \mathbb{E}_{Y,X} \left\{ \left[ Y - \mathbb{E}_{sam} \hat{f}(X) \right] \left[ \mathbb{E}_{sam} \hat{f}(X) - \hat{f}(X) \right] \right\} = 0$$

# Why Bagging Works? Ct'ed

- ▶ So, if we could use  $\mathbb{E}_{sam} \hat{f}(X)$  we would get a smaller risk ( $L_2$  error) than if we used  $\hat{f}(x)$
- ▶ In practice we cannot compute  $\mathbb{E}_{sam}$ , so we approximate it by the bootstrap:

$$\mathbb{E}_{sam}^* \hat{f}^*(x)$$

where  $\mathbb{E}_{sam}^*$  is the expectation that considers each observation in  $\{Y_i^*, X_i^*\}_{i=1}$  drawn with replacement from  $\mathbb{P}_n$  and  $\hat{f}^*$  is the estimator based on the bootstrapped sample  $\{Y_i^*, X_i^*\}$ .

- ▶ Similarly as in Slide 6, we approximate  $\mathbb{E}_{sam}^* \hat{f}^*(x)$  by the Monte-Carlo algorithm:

$$\mathbb{E}_{sam}^* \hat{f}^*(x) \approx \frac{1}{B} \sum_{b=1}^B \hat{f}_b^*(x)$$