# Supervisory Committee Meeting

Xiaomeng Ju

UBC

18/08/2021

# Overview

- Thesis scope
- Contents of the thesis
- Projects and progress
- Challenges
- Timeline

# Thesis Scope

- Title: Robust boosting for complex data
- Goal: estimate $F$ in a regression model, $Y \in \mathbb{R}$

$$Y = F(X) + \epsilon \tag{1}$$

- Complex data:
    - (a) $X \in \mathbb{R}^p$, $Y$ contains outliers $\rightarrow$ Chapter 2
    - (b) $X \in$ Hilbert space, $Y$ follows (1) $\rightarrow$ Chapter 3
    - (c) $X \in$ Hilbert space, $Y$ contains outliers $\rightarrow$ Chapter 4
    - (d) $X \in L^2(\mathcal{I})$ and evaluated on a sparse grid, possibly contaminated, $Y$ follows (1) $\rightarrow$ Chapter 5

    For (b), (c), and (d), $X$ may also contain real-valued predictors $\in \mathbb{R}^d$

# Robust gradient boosting: MM-estimators

▶ Motivated by MM-estimator for regression
  ▶ S-estimator: highly robust

$$\hat{F}_S = \underset{F}{\operatorname{argmin}}\, \hat{\sigma}_S(F)\,,$$

  where $\hat{\sigma}_S(F)$ is a robust scale of residuals
  ▶ M-estimator: highly efficient, initialized at $\hat{F}_S$

$$\hat{F}_M = \underset{F}{\operatorname{argmin}} \sum_{i \in \mathcal{I}_{\text{train}}} \rho\left(\frac{y_i - F(\mathbf{x}_i)}{\hat{\sigma}_S(\hat{F}_S)}\right)$$

  ▶ MM-estimator: highly robust and highly efficient

# Robust gradient boosting: methodology

## RRBoost

- **Stage 1** : compute an $S$-type boosting estimator $\hat{F}_S$ with high robustness but possibly low efficiency (S-scale as $L$)
- **Stage 2**: compute an $M$-type boosting estimator initialized at the function estimator ($\hat{F}_S$) and scale estimator $\hat{\sigma}_S$ obtained in Stage 1. (bonded loss as $L$)

# Boosting for functional regression: progress

- Ju, Xiaomeng, and Matías Salibián-Barrera. "Robust boosting for regression problems." Computational Statistics & Data Analysis 153 (2021): 107065.
- RRBoost package on CRAN
- Repo: `https://github.com/xmengju/RRBoost`

# Boosting for functional regression: model

- Goal: estimate $F : X \to Y$ in order to make predictions for future observations
- Multi-index model:

$$F(X) = r(\langle X, \alpha_1 \rangle, ..., \langle X, \alpha_p \rangle)$$

Fit complex functions; capture iterations between indices

- Approximation:

$$F(X) \approx r_1(\langle X, \beta_{1,1} \rangle, ..., \langle X, \beta_{1,k} \rangle) + ... + r_T(\langle X, \beta_{T,1} \rangle, ..., \langle X, \beta_{T,K} \rangle),$$

where each $r_j(\langle X, \beta_{j,1} \rangle, ..., \langle X, \beta_{j,K} \rangle)$ is fitted by a functional multi-index tree.

# Boosting for functional regression: methodology

- ▶ Propose a boosting algorithm: TFBoost
- ▶ Input data: $(\mathbf{x}_i, y_i)$, $i \in \{\mathcal{I}_{\text{train}} \cup \mathcal{I}_{\text{val}}\}$
- ▶ Loss function: $L(y_i, F(x_i))$
- ▶ Every boosting iteration:
  calculate negative gradient $\rightarrow$ fit base learner (functional multi-index tree) $\rightarrow$ find step size $(\alpha_t)$ $\rightarrow$ update function

$$\hat{F}(x_i) = \sum_{t=1}^{T} \alpha_t \hat{r}_t(\langle x_i, \hat{\beta}_{t,1} \rangle, ..., \langle x_i, \hat{\beta}_{t,k} \rangle)$$

- ▶ Multi-index tree
  - ▶ Type A tree: optimal indices for the whole tree
  - ▶ Type B tree: optimal index for each split (fast calculation)

# Boosting for functional regression: progress

- Paper draft soon to be submitted (in August)
- TFBoost package completed
- Repo: `https://github.com/xmengju/TFBoost`

# Robust TFBoost: problem description

- Extend TFBoost to data with outliers
- Most proposals are for linear models
- Types of outliers (include a figure):
    - Shape outliers
    - Magnitude outliers (curve, point, or interval)
    - Vertical outliers

# Robust TFBoost: methodology

- TFBoost(LAD): TFBoost with L1 loss
- TFBoost(LAD-M): TFBoost(LAD) $\rightarrow$ residual scale $\rightarrow$ M-type TFBoost
- TFBoost(RR): S-type TFBoost $\rightarrow$ M-type TFBoost

# Robust TFBoost: progress

- Simulation results comparing 3 proposals with competing robust functional regression methods in the literature
- Technical report that describes the methodology and the simulation study
- Deadline: mid September

# Sparse TFBoost: methodology

- ▶ Problem: difficulty to calculate the inner product with sparsely observed functions
- ▶ Idea: borrow strength across functions

$$\tilde{X}(t) = \sum_{j=1}^{K} \xi_j \phi_j(t),$$

where $\phi_j(t)$ are FPC and $\xi_j = \langle (X(t) - \mu(t)), \phi_j(t) \rangle$
- ▶ Methods:
  - ▶ Sparse X with possible measurement errors: PACE
  - ▶ Sparse X with functional outliers: ROB

# Sparse TFBoost: progress

- Have reviewed literature and got familiar with the software of PACE and ROB
- Most of this project remains to be done.

# Challenges

Past:

- ▶ Late completion of the comprehensive exam
- ▶ Limited computational resources (solved)
- ▶ Time management

Future:

- ▶ Writing speed
- ▶ Additional time to revise the TFBoost paper after submission

# Timeline

- ▶ 2021/08:
    - ▶ TFBoost: submit paper and package
    - ▶ Thesis: draft the RRBoost and TFBoost chapters
    - ▶ Robust TFBoost: simulation and real example
- ▶ 2021/09:
    - ▶ Thesis: draft Robust TFBoost chapter
    - ▶ Sparse TFBoost: simulation
    - ▶ Revision: RRBoost and TFBoost chapters
- ▶ 2021/10:
    - ▶ Thesis: draft Sparse TFBoost chapter
    - ▶ Sparse TFBoost: simulation and real example
    - ▶ Revision: RRBoost and TFBoost chapters
- ▶ 2021/11:
    - ▶ Thesis: draft Sparse TFBoost chapter, conclusion and future work
    - ▶ Revision: Robust TFBoost chapter
- ▶ 2021/12:
    - ▶ Thesis: draft conclusion and future work
    - ▶ Revision: Sparse TFBoost chapter

- 2022/01:
  - Revision: Conclusion and future work
  - Complete the first revised draft
- 2022/02:
  - Complete the second revised draft
  - Send the draft to committee members
- 2022/06:
  - Thesis defense

# Contents of the Thesis

# Contents of the Thesis

# Contents of the Thesis