

ĐẠI HỌC QUỐC GIA TP.HCM  
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN



**STATISTICAL LEARNING**  
**Selection Inference for CoRT**

TTNT2025

---

*Sinh viên:*  
Phạm Bách Hiệp  
Nguyễn Ngọc Nhớ

*Giảng viên:*  
Võ Nguyễn Lê Duy

# Tóm Tắt

Trong báo cáo này, ta sẽ



# Mục Lục

<b>1 SI For CoRT</b>	<b>3</b>
1.1 CoRT Algorithm . . . . .	3
1.2 Bài toán con 1: Tính ổn định của các mô hình Huấn luyện . .	4
1.2.1 Selection event . . . . .	4
1.2.2 Biểu diễn Selection event . . . . .	4
1.3 Bài toán con 2: Tính ổn định của Bỏ phiếu Kiểm định . . .	5
1.3.1 Selection event . . . . .	5
1.3.2 Biểu diễn Selection event . . . . .	5
1.4 Bài toán con 3: Tính ổn định của Mô hình CoRT . . . . .	6
1.4.1 Selection event . . . . .	6
1.4.2 Biểu diễn Selection event . . . . .	6

# Chương 1

## SI For CoRT

### 1.1 CoRT Algorithm

---

**Thuật toán 1** Adaptive Knowledge Transfer

---

**Input:**  $\mathcal{D}^{(k)}$  với  $0 \leq k \leq K$

- 1: **Phân chia dữ liệu:** Ngẫu nhiên chia dữ liệu mục tiêu  $\mathcal{D}^{(0)}$  thành T (là số lẻ) phân bằng nhau và kí hiệu chúng là  $D_{[1]}^{(0)}, \dots, D_{[T]}^{(0)}$ .
  - 2: **for**  $k = 1$  to  $K$  **do**
  - 3:     **for**  $t = 1$  to  $T$  **do**
  - 4:         Tìm  $\hat{\beta}_{[t]}^{(0)}$  bằng cách chạy Lasso trên  $(\mathcal{D}^{(0)} \setminus \mathcal{D}_{[t]}^{(0)})$ .
  - 5:         Tìm  $\hat{\beta}_{[t]}^{(0k)}$  bằng cách chạy Lasso trên  $(\mathcal{D}^{(0)} \setminus \mathcal{D}_{[t]}^{(0)}) \cup \mathcal{D}^{(k)}$ .
  - 6:         Lần lượt tính toán hàm lỗi cho  $\ell(\hat{\beta}_{[t]}^{(0)}, \mathcal{D}_{[t]}^{(0)})$  và  $\ell(\hat{\beta}_{[t]}^{(0k)}, \mathcal{D}_{[t]}^{(0)})$ ,
  - 7:     **end for**
  - 8:     Lập hàm đếm  $C(\boldsymbol{\beta}^{(0k)})$  với  $\sum_{t=1}^T \mathbf{1} \left\{ \ell(\hat{\beta}_{[t]}^{(0k)}, \mathcal{D}_{[t]}^{(0)}) = \min_{\boldsymbol{\beta} \in \{\hat{\beta}_{[t]}^{(0k)}, \hat{\beta}_{[t]}^{(0)}\}} \ell(\boldsymbol{\beta}, \mathcal{D}_{[t]}^{(0)}) \right\}$ .
  - 9:     **end for**
  - 10: Thu được ước lượng của  $S$  với  $\hat{S} = \{k : C(\boldsymbol{\beta}^{(0k)}) \geq (T+1)/2\}$ .
  - 11: Thu được  $\hat{\beta}^{(0)}$  bằng cách chạy CoRT với  $\mathcal{D}^{(k)}, k \in \hat{S} \cup \{0\}$ .
-

## 1.2 Bài toán con 1: Tính ổn định của các mô hình Huấn luyện

### 1.2.1 Selection event

Theo mô tả của Thuật toán 1 (Adaptive Knowledge Transfer), quá trình lựa chọn nguồn phụ thuộc trực tiếp vào các mô hình trung gian được huấn luyện trong vòng lặp kiểm chứng chéo (Cross-Validation). Cụ thể:

- **Dòng 4:** Tìm  $\hat{\beta}_{[t]}^{(0)}$  bằng cách chạy Lasso trên tập huấn luyện đích  $(D^{(0)} \setminus D_{[t]}^{(0)})$ .
- **Dòng 5:** Tìm  $\hat{\beta}_{[t]}^{(0k)}$  bằng cách chạy Lasso trên tập kết hợp  $(D^{(0)} \setminus D_{[t]}^{(0)}) \cup D^{(k)}$ .

### 1.2.2 Biểu diễn Selection event

Ta cần điều kiện hóa sự kiện tất cả các mô hình train này là không thay đổi. Cụ thể, tập biến hoạt động (active sets)  $\hat{M}$  và dấu của các hệ số hồi quy (signs)  $\hat{s}$  của mọi mô hình Baseline và Augmented trong suốt quá trình CV phải giữ nguyên giá trị quan sát được khi dữ liệu  $y$  thay đổi dọc theo đường  $z$ .

**Kết luận Toán học:** Áp dụng điều kiện KKT cho mỗi mô hình Lasso, ta được bất phương trình:

$$\underbrace{\begin{pmatrix} -\text{diag}(\mathbf{s}_M)\mathbf{v} \\ \mathbf{q} \\ -\mathbf{q} \end{pmatrix}}_{\psi} z \leq \underbrace{\begin{pmatrix} \text{diag}(\mathbf{s}_M)\mathbf{u} \\ 1-p \\ 1+p \end{pmatrix}}_{\gamma} \quad (1.1)$$

Từ đó ta xây dựng được ràng buộc tuyến tính cho tất cả mô hình Baseline và Augmented như sau:

$$\Psi z \leq \Gamma \quad (1.2)$$

trong đó:

$$\Psi = \begin{bmatrix} \vdots \\ \psi_{[t]}^{(0)} \\ \vdots \\ \psi_{[t]}^{(0k)} \\ \vdots \end{bmatrix}, \quad \Gamma = \begin{bmatrix} \vdots \\ \gamma_{[t]}^{(0)} \\ \vdots \\ \gamma_{[t]}^{(0k)} \\ \vdots \end{bmatrix} \quad (1.3)$$

## 1.3 Bài toán con 2: Tính ẩn định của Bỏ phiếu Kiểm định

### 1.3.1 Selection event

Giai đoạn này quyết định việc một nguồn  $k$  có được đưa vào tập hợp  $\hat{S}$  hay không dựa trên kết quả so sánh sai số dự báo trên tập kiểm định (Validation set). Theo Thuật toán 1:

- **Dòng 6:** Tính toán hàm lỗi  $\ell(\cdot)$  cho mô hình Baseline và Augmented trên tập kiểm định  $D_{[t]}^{(0)}$ .
- **Dòng 8:** So sánh sai số và thực hiện bỏ phiếu thông qua hàm chỉ thị  $1\{\ell_{aug} < \ell_{base}\}$ .
- **Dòng 10:** Thu được ước lượng của  $S$  với  $\hat{S} = \{k : C(\beta^{(0k)}) \geq \frac{T+1}{2}\}$ .

### 1.3.2 Biểu diễn Selection event

Ta cần điều kiện hóa sự kiện "Thắng/Thua" của mỗi cặp so sánh là không đổi. Vì nghiệm Lasso  $\hat{\beta}(z)$  là tuyến tính theo  $z$ , hàm mất mát (bình phương sai số) trên tập kiểm định sẽ trở thành một hàm bậc hai theo  $z$ .

Xét hiệu sai số giữa mô hình Augmented và Baseline tại fold  $t$ :

$$\Delta_{t,k}(z) = \ell(\hat{\beta}_{[t]}^{(0k)}(z)) - \ell(\hat{\beta}_{[t]}^{(0)}(z)) \quad (1.4)$$

**Kết luận Toán học:** Triển khai biểu thức sai số, ta thu được đa thức bậc hai đặc trưng cho sự chênh lệch hiệu suất:

$$\Delta_{t,k}(z) = \mathcal{A}_{t,k}z^2 + \mathcal{B}_{t,k}z + \mathcal{C}_{t,k} \quad (1.5)$$

trong đó các hệ số  $\mathcal{A}, \mathcal{B}, \mathcal{C}$  được tính toán dựa trên các tham số  $(\mathbf{u}, \mathbf{v})$  từ Bài toán con 1 và dữ liệu kiểm định.

Khi đó, miền giá trị  $z$  thỏa mãn điều kiện ổn định bỏ phiếu được xác định bởi tập hợp các bất phương trình bậc hai:

$$\begin{cases} \mathcal{A}_{t,k}z^2 + \mathcal{B}_{t,k}z + \mathcal{C}_{t,k} < 0 & \text{nếu } \Delta_{t,k}(z_{obs}) < 0 \text{ (Augmented thắng)} \\ \mathcal{A}_{t,k}z^2 + \mathcal{B}_{t,k}z + \mathcal{C}_{t,k} \geq 0 & \text{nếu } \Delta_{t,k}(z_{obs}) \geq 0 \text{ (Baseline thắng)} \end{cases} \quad (1.6)$$

Cần lưu ý rằng dòng 10 của thuật toán ("Thu được ước lượng của  $S\dots$ ") không yêu cầu xây dựng thêm một hệ bất phương trình mới. Đây là hệ quả trực tiếp của việc ổn định bỏ phiếu ở trên.

Cụ thể, tập hợp nguồn được chọn  $\hat{S}$  được định nghĩa dựa trên tổng số phiếu bầu:

$$\hat{S} = \left\{ k : \sum_{t=1}^T 1\{\Delta_{t,k}(z) < 0\} \geq \frac{T+1}{2} \right\} \quad (1.7)$$

Do hệ bất phương trình ở Bài toán con 2 đã đảm bảo dấu của mọi  $\Delta_{t,k}(z)$  là không đổi trong khoảng  $\mathcal{Z}_{val}$ , nên giá trị của hàm chỉ thị  $1\{\cdot\}$  và tổng số phiếu cũng bất biến. Điều này suy ra tập hợp  $\hat{S}$  là cố định trong miền  $\mathcal{Z}_{val}$ .

## 1.4 Bài toán con 3: Tính ổn định của Mô hình CoRT

### 1.4.1 Selection event

Sau khi xác định được tập hợp các nguồn tương đồng  $\hat{S}$  từ giai đoạn trước, thuật toán tiến hành bước ước lượng tham số cuối cùng. Theo Dòng 11 của Thuật toán 1:

- **Dòng 11:** Thu được  $\hat{\beta}^{(0)}$  bằng cách chạy CoRT với  $\mathcal{D}^{(k)}, k \in \hat{S} \cup \{0\}$ .

### 1.4.2 Biểu diễn Selection event

Cụ thể, tập biến hoạt động  $\mathcal{M}_{CoRT}$  và vectơ dấu  $\mathbf{s}_{CoRT}$  của mô hình Lasso trên dữ liệu kết hợp phải giữ nguyên giá trị quan sát được khi  $z$  thay đổi.

**Kết luận Toán học:** Áp dụng điều kiện KKT của bài toán Lasso trên tập dữ liệu kết hợp, ta thu được hệ bất phương trình tuyến tính xác định khoảng  $\mathcal{Z}_{CoRT}$ :

$$\Omega z \leq \Upsilon \quad (1.8)$$

Trong đó, các ma trận hệ số được xác định cụ thể từ điều kiện KKT của mô hình cuối cùng (tương tự như Bài toán con 1 nhưng áp dụng trên  $X_{combined}$ ):

$$\Omega = \begin{pmatrix} -\text{diag}(\mathbf{s}_{CoRT})\mathbf{v}_{CoRT} \\ \mathbf{q}_{CoRT} \\ -\mathbf{q}_{CoRT} \end{pmatrix}, \quad \Upsilon = \begin{pmatrix} \text{diag}(\mathbf{s}_{CoRT})\mathbf{u}_{CoRT} \\ 1 - \mathbf{p}_{CoRT} \\ 1 + \mathbf{p}_{CoRT} \end{pmatrix} \quad (1.9)$$

