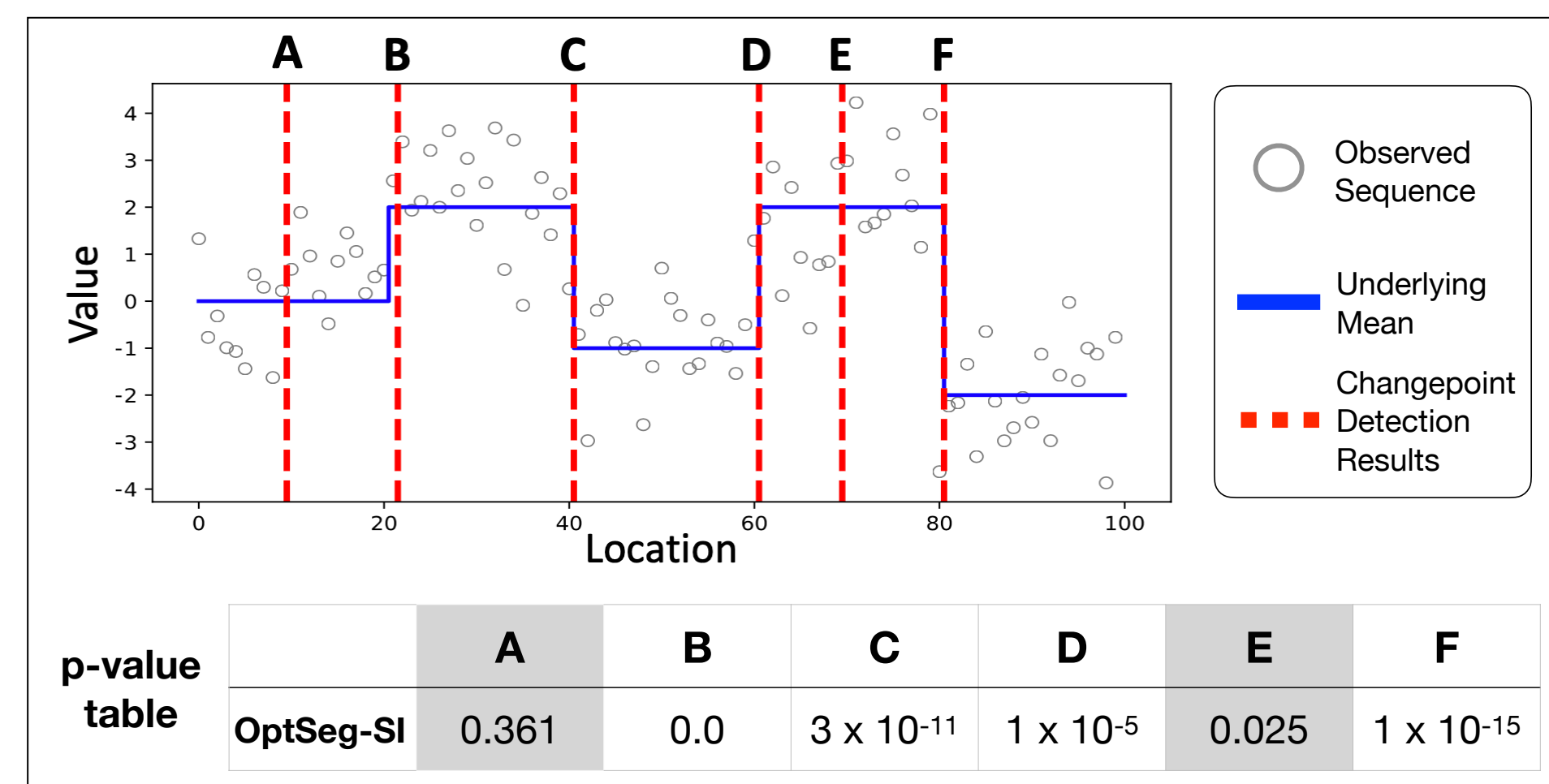# Computing Valid $p$-value for Optimal Changepoint by Selective Inference using Dynamic Programming

Vo Nguyen Le Duy, Hiroki Toda, Ryota Sugiyama, Ichiro Takeuchi
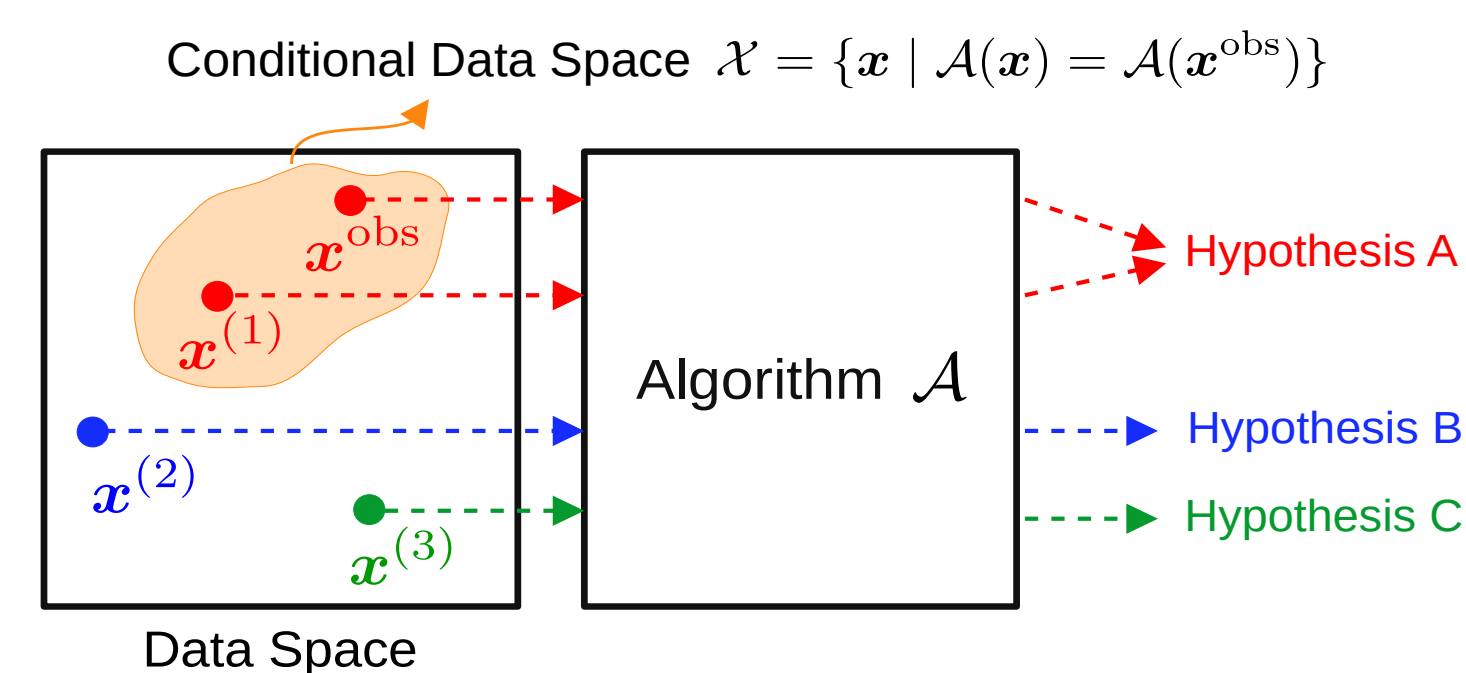
Nagoya Institute of Technology

RIKEN

NEURAL INFORMATION PROCESSING SYSTEMS

## Introduction and Motivation

❖ Changepoint (CP) detection: find changes in the underlying mechanism of the observed sequential data.

❖ CP detection is usually formulated as the problem of minimizing the segmentation cost where Dynamic Programming (DP) is commonly used.

❖ There are several CP detection methods. However, less attention has been paid to quantify the reliability of the detected CPs.



| p-value table | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| OptSeg-SI | 0.361 | 0.0 | $3 \times 10^{-11}$ | $1 \times 10^{-5}$ | 0.025 | $1 \times 10^{-15}$ |

❖ **A** and **E** are falsely detected CPs
➡ Results from CP detection algorithms are unreliable
➡ Harmful for high-stake decision making such as medical diagnosis

❖ We propose **OptSeg-SI** method to provide **valid** $p$-value, which is used as a criterion to quantify the reliability of the detected CPs, based on the concept of **Selective Inference (SI)**.
➡ Large $p$-value indicates false detection (**A** and **E**) and small $p$-value indicates true detection (**B**, **C**, **D** and **F**)
➡ OptSeg-SI can identify both false and true positive detections
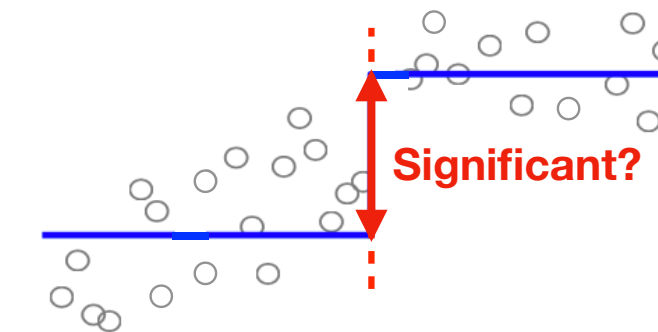
## Concept of Selective Inference (SI)

Conditional Data Space $\mathcal{X} = \{\boldsymbol{x} \mid \mathcal{A}(\boldsymbol{x}) = \mathcal{A}(\boldsymbol{x}^{\text{obs}})\}$



**Conditional inference:** $\Pr(T(\boldsymbol{x}) \mid \mathcal{A}(\boldsymbol{x}) = \mathcal{A}(\boldsymbol{x}^{\text{obs}}))$, where $T(\boldsymbol{x})$ is the test statistic.

## Problem Setting

❖ We consider the following statistical test

$$H_0 : \mu_{\text{left}} = \mu_{\text{right}}$$
vs.
$$H_1 : \mu_{\text{left}} \neq \mu_{\text{right}}$$

where $\mu$ is population mean.

Significant?

❖ The conditional $p$-value (selective $p$-value) is defined as

$$p_{\text{selective}} = \mathbb{P}_{H_0}(|\Delta| \geq |\Delta^{\text{obs}}| \mid \mathcal{X})$$

- $\Delta^{\text{obs}}$ is the difference in sample mean between the left segment and right segment in the **observed** sequence
- $\Delta$ is the mean difference in **any random sequence**
- $\mathcal{X}$ is the conditional data space defined as

$$\mathcal{X} = \{\boldsymbol{x} : \{\text{left}, \text{right}\} \leftarrow \text{DP algorithm } \mathcal{A}(\boldsymbol{x})\}$$
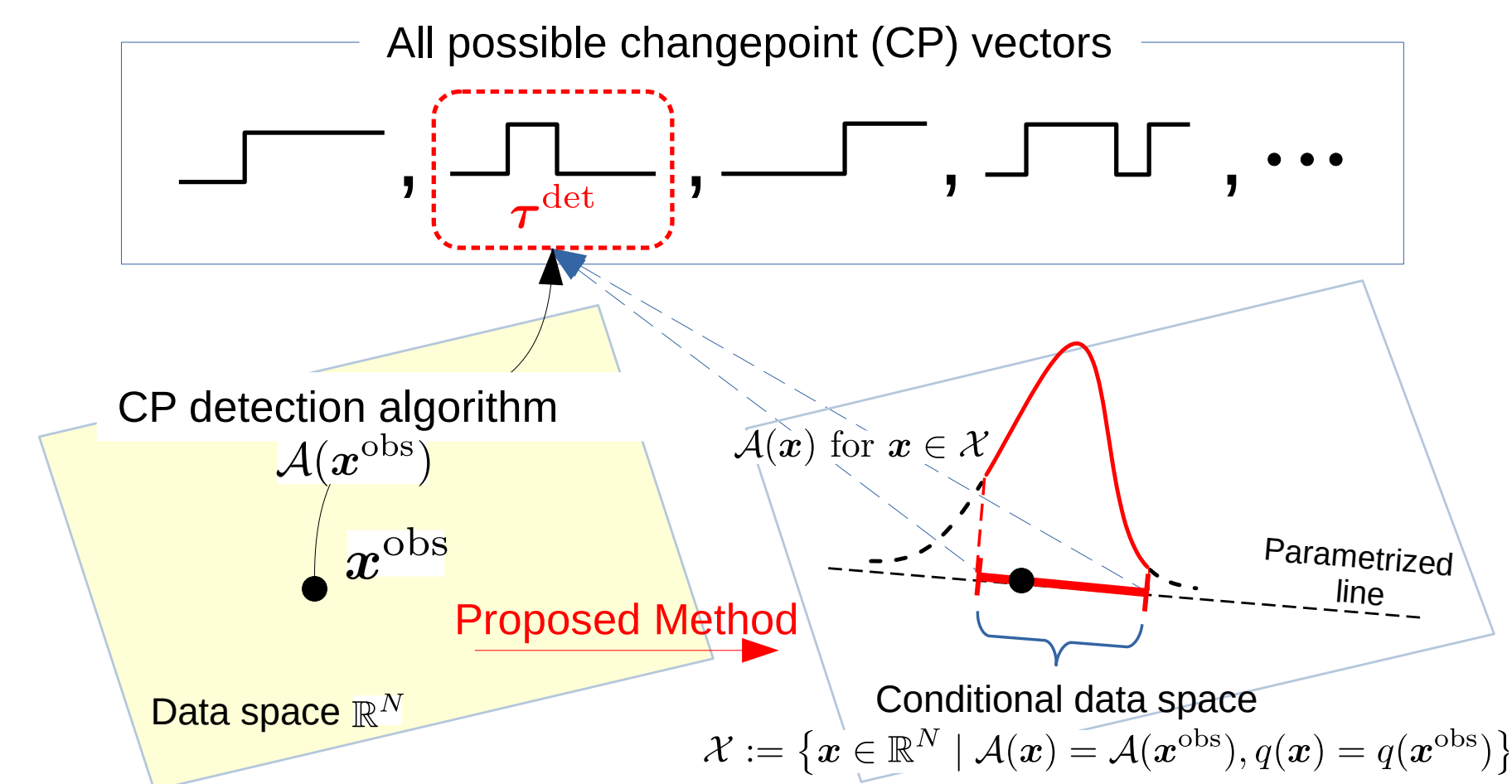
❖ In other words, $\mathcal{X}$ is the data space whose data has the same detected CP as the observed sequence.

❖ The selective $p$-value is **valid** since

$$\mathbb{P}_{H_0}(p_{\text{selective}} < \alpha) = \alpha, \quad \forall \alpha \in [0,1].$$

However, characterization of the conditional data space $\mathcal{X}$ is challenging

## Proposed Method - Schematic illustration

All possible changepoint (CP) vectors



CP detection algorithm $\mathcal{A}(\boldsymbol{x}^{\text{obs}})$

$\mathcal{A}(\boldsymbol{x})$ for $\boldsymbol{x} \in \mathcal{X}$

Proposed Method

Parametrized line

Conditional data space

$\mathcal{X} := \{\boldsymbol{x} \in \mathbb{R}^N \mid \mathcal{A}(\boldsymbol{x}) = \mathcal{A}(\boldsymbol{x}^{\text{obs}}), q(\boldsymbol{x}) = q(\boldsymbol{x}^{\text{obs}})\}$

Data space $\mathbb{R}^N$

❖ Step 1: Obtain CP results from the observed data $\boldsymbol{x}^{\text{obs}}$
❖ Step 2: By restricting data on the line, we perform DP on parametrized data and identify the sub-space whose data has the same CP results as $\boldsymbol{x}^{\text{obs}}$

## Proposed Method - Details

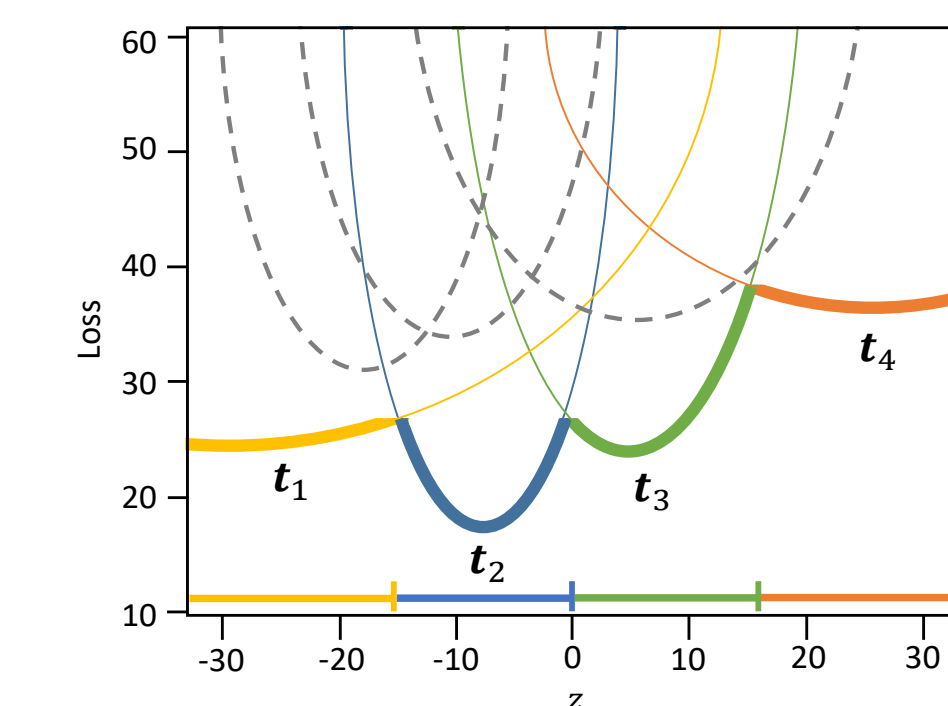❖ We first restrict the data to the line by using a scalar parameter $z \in \mathbb{R}$

$$\boldsymbol{x}(z) = \boldsymbol{a} + \boldsymbol{b}z,$$

where $\boldsymbol{a}$ and $\boldsymbol{b}$ have specific forms.

❖ The conditional data space $\mathcal{X}$ is then re-written as

$$\mathcal{X} = \{\boldsymbol{x}(z) = \boldsymbol{a} + \boldsymbol{b}z \mid z \in \mathcal{Z}\},$$
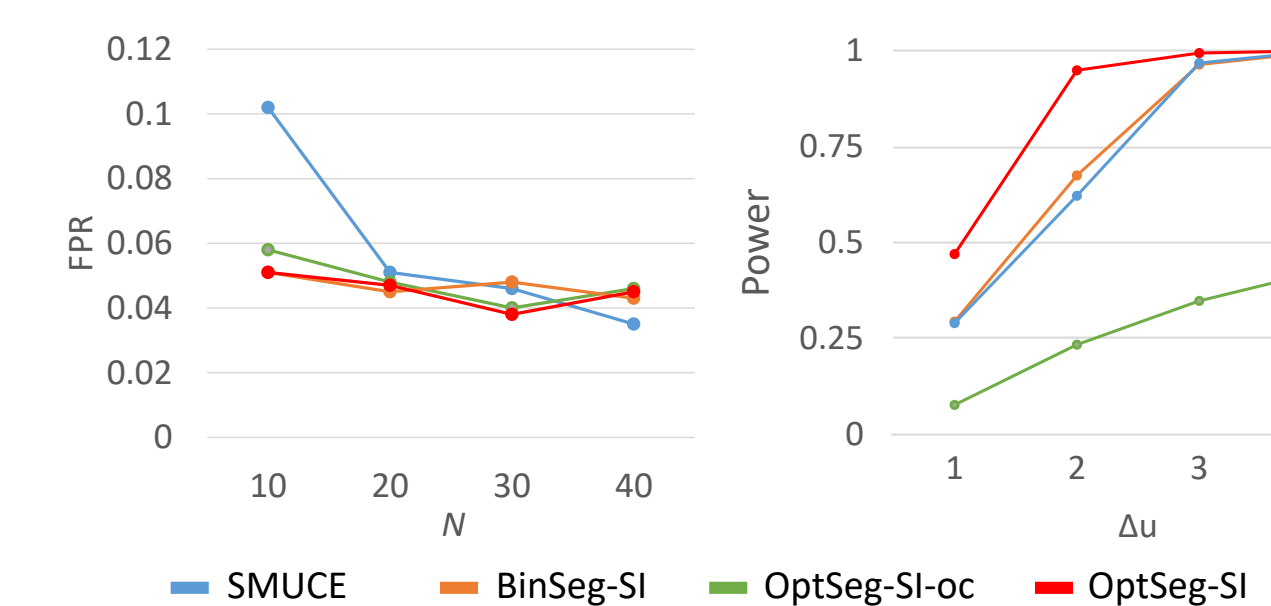
where $\mathcal{Z} = \{z \in \mathbb{R} : \{\text{left}, \text{right}\} \leftarrow \text{DP algorithm } \mathcal{A}(\boldsymbol{x}(z))\}$.

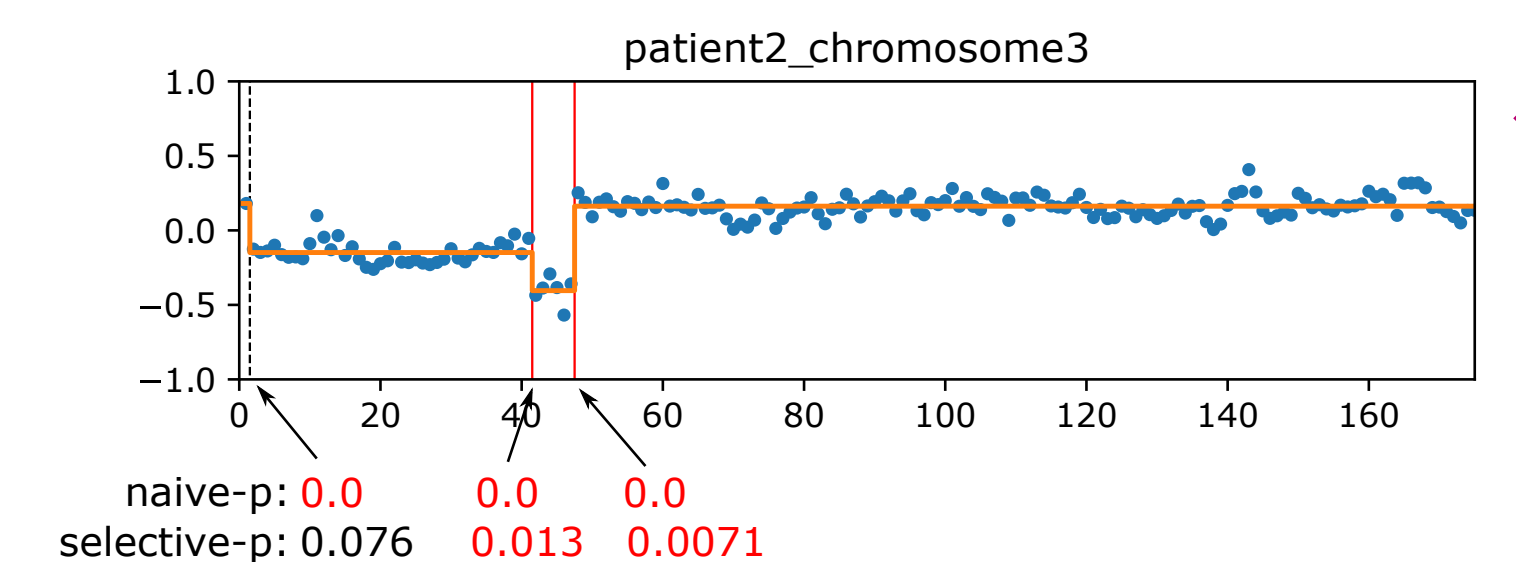⟹ **The remaining task is to identify truncation region $\mathcal{Z}$**



❖ We propose a **parametric DP approach** to compute CP results of $\boldsymbol{x}(z)$ for all $z \in \mathbb{R}$

❖ The region $\mathcal{Z}$ is then the union of intervals of $z$ on which we obtain the same CP results as the observed data

## Experimental Results



SMUCE    BinSeg-SI    OptSeg-SI-oc    OptSeg-SI

❖ The OptSeg-SI method (red) is powerful while successfully controlling the false positive rate (FPR)

patient2_chromosome3



naive-p: 0.0    0.0    0.0
selective-p: 0.076    0.013    0.0071

❖ Application to real-world bioinformatics dataset

### References

Lee et al. (2016). "Exact post-selection inference, with application to the lasso" In: The Annals of Statistics.

Duy et al. (2020). "Parametric programming approach for more powerful and general lasso selective inference". In: arXiv:2004.09749.

Duy et al. (2020). "Quantifying statistical significance of neural network representation-driven hypotheses by selective inference". In: arXiv:2010.01823.