

Exact Statistical Inference for the Wasserstein Distance by Selective Inference

Vo Nguyen Le Duy

Nagoya Institute of Technology and RIKEN

duy.mllab.nit@gmail.com

Ichiro Takeuchi

Nagoya Institute of Technology and RIKEN

takeuchi.ichiro@nitech.ac.jp

January 21, 2022

Abstract

In this paper, we study statistical inference for the Wasserstein distance, which has attracted much attention and has been applied to various machine learning tasks. Several studies have been proposed in the literature, but almost all of them are based on *asymptotic* approximation and do *not* have finite-sample validity. In this study, we propose an *exact (non-asymptotic)* inference method for the Wasserstein distance inspired by the concept of conditional Selective Inference (SI). To our knowledge, this is the first method that can provide a valid confidence interval (CI) for the Wasserstein distance with finite-sample coverage guarantee, which can be applied not only to one-dimensional problems but also to multi-dimensional problems. We evaluate the performance of the proposed method on both synthetic and real-world datasets.

1 Introduction

The Wasserstein distance, which is a metric used to compare the probability distributions, has been attracted significant attention and being used more and more in Statistics and Machine Learning (ML) [Kolouri et al., 2017]. This distance measures the cost to couple one distribution with another, which arises from the notion of *optimal transport* [Villani, 2009]. The utilization of the Wasserstein distance benefits to several applications such as supervised learning [Frogner et al., 2015], generative modelling [Arjovsky et al., 2017], biology [Evans and Matsen, 2012], and computer vision [Ni et al., 2009].

When the Wasserstein distance calculated from noisy data is used for various decision-making problems, it is necessary to quantify its statistical reliability, e.g., in the form of confidence intervals (CIs). However, there is no satisfactory statistical inference method for the Wasserstein distance. The main reason is that the Wasserstein distance is defined based on the optimal solution of a linear program (LP), and it is difficult to analyze how the uncertainty in the data is transmitted to the uncertainty in the Wasserstein distance. Several studies have been proposed in literature [Bernton et al., 2017, Del Barrio et al., 1999, 2019, Ramdas et al., 2017, Imaizumi et al., 2019]. However, they all rely on asymptotic approximation to mitigate the difficulty stemming from the fact that the Wasserstein distance depends on the optimization problem, and thus most of them can only be applied to one-dimensional problems (the details will be discussed in the related work section).

When an optimization problem such as LP is applied to random data, it is intrinsically difficult to derive the *exact (non-asymptotic)* sampling distribution of the optimal solution. In this paper, to resolve the difficulty, we introduce the idea of conditional Selective Inference (SI). It is well-known that the optimal solution of an LP depends only on a subset of variables, called *basic variables*. Therefore, the LP algorithm can be interpreted as first identifying the basic variables, and then solving the linear system of equations about the basic variables.

Our basic idea is based on the fact that, in the LP problem for the Wasserstein distance, the identification of the basic variables corresponds to the process of determining the *coupling* between the source and destination of the probability mass. Since the optimal coupling is determined (selected) based on the data, the *selection bias* must be properly considered. Therefore, to address the selection bias, we propose an exact statistical inference method for the Wasserstein distance by conditioning on the basic variables of the LP, i.e., optimal coupling.

The main advantage of the proposed method is that it can provide exact (non-asymptotic) CI for the Wasserstein distance, unlike the asymptotic approximations in the existing studies. Moreover, while the existing methods are restricted to one-dimensional problems, the proposed method can be directly extended to multi-dimensional problems because the proposed CI is derived based on the fact that the Wasserstein distance depends on the optimal solution of the LP.

1.1 Contributions

Regarding the high-level conceptual contribution, we introduce a novel approach to explicitly characterize the selection bias of the data-driven coupling problem inspired by the concept of conditional Selective Inference (SI). In regard to the technical contribution, we provide an exact (non-asymptotic) inference method for the Wasserstein distance. To our knowledge, this is the first method that can provide *valid* CI, called *selective CI*, for the Wasserstein distance that guarantees the coverage property in finite sample size. Another practically important advantage of this study is that the proposed method is valid when the Wasserstein distance is computed in multi-dimensional problem, which is impossible for almost all the existing asymptotic methods since the limit distribution of the Wasserstein distance is only applicable for univariate data. We conduct experiments on both synthetic and real-world datasets to evaluate the performance of the proposed method.

1.2 Related works

In traditional statistics, reliability evaluation with Wasserstein distance has been based on asymptotic theory, i.e., sample size $\rightarrow \infty$. In the univariate case, instead of solving the optimization problem, the Wasserstein can be described by using an inverse of the distribution function. For example, let F^{-1} be the quantile function of the data and F_n^{-1} be the empirical quantile function of the generated data, the Wasserstein distance with ℓ_2 distance is computed by $\int_0^1 (F^{-1}(t) - F_n^{-1}(t))^2 dt$. Based on the quantile function, several studies [Del Barrio et al., 1999, 2019, Ramdas et al., 2017] derived the asymptotic distribution of the Wasserstein distance. Obviously, these methods can not guarantee the validity in finite sample size. Moreover, since the quantile function is only available in univariate case, these methods can not be extended to multivariate cases which are practically important.

Recently, Imaizumi et al. [2019] has proposed an approach on multidimensional problems. However, it is important to clarify that this study does *not* provide statistical inference for the “original” Wasserstein distance. Instead, the authors consider an *approximation* of the Wasserstein distance, which does not require solving a LP. Besides, this method also relies on asymptotic distribution of the test statistic which is approximated by the Gaussian multiplier bootstrap. Therefore, to our knowledge, statistical inference method for the Wasserstein distance in multi-dimensional problems is still a challenging open problem.

Conditional SI is a statistical inference framework for correcting selection bias. Traditionally, there are mainly two types of approaches for selection bias correction. The first approach is family-wise error rate (FWER) control, which includes standard multiple comparison methods such as traditional Bonferroni correction. However, the FWER control is too conservative and it is difficult to utilize it for selection bias correction in complex adaptive data analysis. Another approach is false discovery rate (FDR) control, in which the target is to control the expected proportion of discoveries that are false at a given significance level. The FDR control is less conservative than FWER control, and it is used in many high-dimensional statistical inference problems.

Conditional SI is the third approach for selection bias correction. The basic idea of conditional SI was known before, but it becomes popular by the recent seminal work proposed by Lee et al. [2016]. In that study, exact statistical inference on the selected features by Lasso was considered. Their basic idea is to employ the sampling distributions of the selected parameter coefficients *conditional on* the selection event that a subset of features is selected by the Lasso. By using the conditional distribution in statistical inference, the selection bias of the Lasso (i.e., the fact that the Lasso selected the features based on the data) can be corrected. Their contribution was to show that, even in complex data analysis methods such as Lasso, the exact sampling distribution can be characterized if appropriate selection events are considered.

After the seminal work [Lee et al., 2016], many conditional SI approaches for various feature selection methods were proposed in the literature [Loftus and Taylor, 2015, Yang et al., 2016, Tibshirani et al., 2016, Suzumura et al., 2017]. Furthermore, theoretical analyses and new computational methods for conditional SI are still being actively studied [Fithian et al., 2014, Le Duy and Takeuchi, 2021, Duy and Takeuchi, 2021, Sugiyama et al., 2021a]. However, most of conditional SI studies are focused on feature selection problems. Although there have been applications to several problems such as change point detection [Hyun et al., 2018, Duy et al., 2020b, Sugiyama et al., 2021b], outlier detection [Chen and Bien, 2019, Tsukurimichi et al., 2021], and image segmentation [Tanizaki et al., 2020, Duy et al., 2020a], these problems can also be interpreted as feature selection in a broad sense. Our novelty in this study is to first introduce conditional SI framework for statistical inference on the Wasserstein distance, which is a data-dependent adaptive distance measure. Our basic idea is based on the facts that the Wasserstein distance is formulated as the solution of a linear program (LP), and the optimal solution of an LP is characterized by the *selected* basic variables. In this study, we consider the sampling distribution of the Wasserstein distance conditional on the selected basic variables, which can be interpreted as considering the selection event on the optimal coupling between the two distributions.

2 Problem Statement

To formulate the problem, we consider two vectors corrupted with Gaussian noise as

$$\mathbf{X} = (x_1, \dots, x_n)^\top = \boldsymbol{\mu}_X + \boldsymbol{\varepsilon}_X, \quad \boldsymbol{\varepsilon}_X \sim \mathbb{N}(\mathbf{0}, \Sigma_X), \quad (1)$$

$$\mathbf{Y} = (y_1, \dots, y_m)^\top = \boldsymbol{\mu}_Y + \boldsymbol{\varepsilon}_Y, \quad \boldsymbol{\varepsilon}_Y \sim \mathbb{N}(\mathbf{0}, \Sigma_Y), \quad (2)$$

where n and m are the number of instances in each vector, $\boldsymbol{\mu}_X$ and $\boldsymbol{\mu}_Y$ are unknown mean vectors, $\boldsymbol{\varepsilon}_X$ and $\boldsymbol{\varepsilon}_Y$ are Gaussian noise vectors with covariances matrices Σ_X and Σ_Y assumed to be known or estimable from independent data. We denote by P_n and Q_m the corresponding empirical measures on \mathbf{X} and \mathbf{Y} .

2.1 Cost matrix

We define the cost matrix $C(\mathbf{X}, \mathbf{Y})$ of pairwise distances (ℓ_1 distance) between elements of \mathbf{X} and \mathbf{Y} as

$$C(\mathbf{X}, \mathbf{Y}) = [|x_i - y_j|]_{ij} \in \mathbb{R}^{n \times m}. \quad (3)$$

We can vectorize $C(\mathbf{X}, \mathbf{Y})$ in the form of

$$\begin{aligned} \mathbf{c}(\mathbf{X}, \mathbf{Y}) &= \text{vec}(C(\mathbf{X}, \mathbf{Y})) \in \mathbb{R}^{nm} \\ &= \Theta(\mathbf{X} \ \mathbf{Y})^\top, \end{aligned} \quad (4)$$

where $\text{vec}(\cdot)$ is an operator that transforms a matrix into a vector with concatenated rows. The matrix Θ is defined as

$$\begin{aligned} \Theta &= \mathcal{S}(\mathbf{X}, \mathbf{Y}) \circ \Omega \in \mathbb{R}^{nm \times (n+m)}, \\ \mathcal{S}(\mathbf{X}, \mathbf{Y}) &= \text{sign}(\Omega(\mathbf{X} \ \mathbf{Y})^\top) \in \mathbb{R}^{nm}, \\ \Omega &= \begin{pmatrix} \mathbf{1}_m & \mathbf{0}_m & \cdots & \mathbf{0}_m & -I_m \\ \mathbf{0}_m & \mathbf{1}_m & \cdots & \mathbf{0}_m & -I_m \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0}_m & \mathbf{0}_m & \cdots & \mathbf{1}_m & -I_m \end{pmatrix} \in \mathbb{R}^{nm \times (n+m)}, \end{aligned} \quad (5)$$

where the operator \circ is element-wise product, $\text{sign}(\cdot)$ is the operator that returns an element-wise indication of the sign of a number, $\mathbf{1}_m \in \mathbb{R}^m$ is the vector of ones, $\mathbf{0}_m \in \mathbb{R}^m$ is the vector of zeros, and $I_m \in \mathbb{R}^{m \times m}$ is the identity matrix.

2.2 The Wasserstein distance

To compare two empirical measures P_n and Q_m with uniform weight vectors $\mathbf{1}_n/n$ and $\mathbf{1}_m/m$, we consider the following Wasserstein distance, which is defined as the solution of a linear program (LP),

$$\begin{aligned} W(P_n, Q_m) &= \min_{T \in \mathbb{R}^{n \times m}} \langle T, C(\mathbf{X}, \mathbf{Y}) \rangle \\ &\text{s.t. } T\mathbf{1}_m = \mathbf{1}_n/n, \\ &\quad T^\top \mathbf{1}_n = \mathbf{1}_m/m, \\ &\quad T \geq 0. \end{aligned} \quad (6)$$

Given \mathbf{X}^{obs} and \mathbf{Y}^{obs} respectively sampled from models (1) and (2) ¹, the Wasserstein distance in (6) on the observed data can be re-written as

$$\begin{aligned} W(P_n, Q_m) &= \min_{\mathbf{t} \in \mathbb{R}^{nm}} \mathbf{t}^\top \mathbf{c}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\ &\text{s.t. } S\mathbf{t} = \mathbf{h}, \mathbf{t} \geq \mathbf{0}, \end{aligned} \quad (7)$$

¹To make a distinction between random variables and observed variables, we use superscript ^{obs}, e.g., \mathbf{X} is a random vector and \mathbf{X}^{obs} is the observed data vector.

where $\mathbf{t} = \text{vec}(T) \in \mathbb{R}^{nm}$, $\mathbf{c}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \in \mathbb{R}^{nm}$ is defined in (4), $S = (M_r \ M_c)^\top \in \mathbb{R}^{(n+m) \times nm}$ in which

$$M_r = \begin{bmatrix} 1 & \dots & 1 & 0 & \dots & 0 & \dots & 0 & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 & \dots & 0 & \dots & 0 \\ \dots & & & & \dots & & \dots & & & \\ 0 & \dots & 0 & 0 & \dots & 0 & \dots & 1 & \dots & 1 \end{bmatrix} \in \mathbb{R}^{n \times nm}$$

that performs the sum over the rows of T and

$$M_c = \begin{bmatrix} I_m & I_m & \dots & I_m \end{bmatrix} \in \mathbb{R}^{m \times nm}$$

that performs the sum over the columns of T , and $\mathbf{h} = (\mathbf{1}_n/n \ \mathbf{1}_m/m)^\top \in \mathbb{R}^{n+m}$ ².

2.3 Optimal solution and closed-form expression of the distance

Let us denote the set of basis variables (the definition of basis variable can be found in the literature of LP, e.g., Murty [1983]) obtained when applying the LP in (7) on \mathbf{X}^{obs} and \mathbf{Y}^{obs} as

$$\mathcal{M}_{\text{obs}} = \mathcal{M}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}). \quad (8)$$

We would like to note that the identification of the basic variables can be interpreted as the process of determining the optimal coupling between the elements of \mathbf{X}^{obs} and \mathbf{Y}^{obs} in the optimal transport problem for calculating the Wasserstein distance. Therefore, \mathcal{M}_{obs} in (8) can be interpreted as the observed optimal coupling obtained after solving LP in (7) on the observed data³. An illustration of this interpretation is shown in Figure 1. We also denote by $\mathcal{M}_{\text{obs}}^c$ a set of *non-basis variables*. Then, the optimal solution of (7) can be written as

$$\hat{\mathbf{t}} \in \mathbb{R}^{nm}, \quad \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}} = S_{:, \mathcal{M}_{\text{obs}}}^{-1} \mathbf{h}, \quad \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}^c} = \mathbf{0}_{|\mathcal{M}_{\text{obs}}^c|},$$

where $S_{:, \mathcal{M}_{\text{obs}}}$ is a sub-matrix of S made up of all rows and columns in the set \mathcal{M}_{obs} .

Example 1. *Given a matrix*

$$S = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \end{bmatrix} \quad \text{and} \quad \mathcal{M}_{\text{obs}} = \{1, 2, 4\},$$

²We note that there always exists exactly one redundant equality constraint in linear equality constraint system in (7). This is due to the fact that sum of all the masses on \mathbf{X}^{obs} is always equal to sum of all the masses on \mathbf{Y}^{obs} (i.e., they are all equal to 1). Therefore, any equality constraint can be expressed as a linear combination of the others, and hence any one constraint can be dropped. In this paper, we always drop the last equality constraint (i.e., the last row of matrix S and the last element of vector \mathbf{h}) before solving (7).

³We suppose that the LP is non-degenerate. A careful discussion might be needed in the presence of degeneracy.

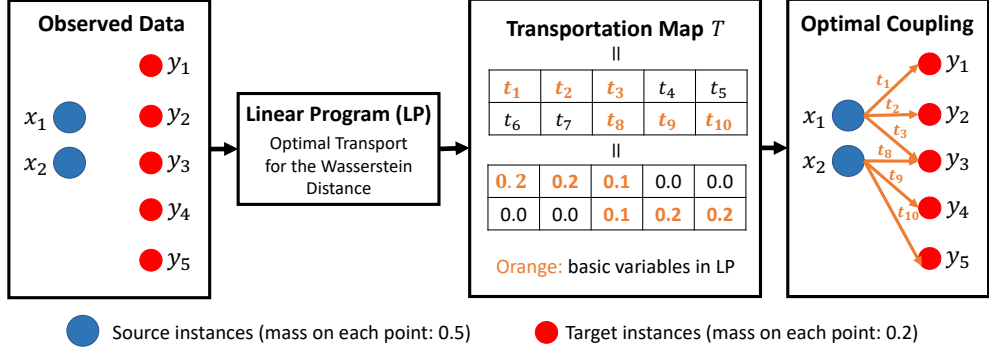


Figure 1: Illustration of the correspondence between *basic variables* in the LP and the optimal coupling. After inputting the data to the LP, we obtain the transportation matrix. The elements $t_1, t_2, t_3, t_8, t_9, t_{10}$ (whose values are non-zero) in the matrix are called basic variables in the LP, and the identification of the basic variables corresponds to the process of determining the coupling between the source and target instances in the optimal transport problem for the Wasserstein distance.

then $S_{:, \mathcal{M}_{\text{obs}}}$ is constructed by extracting the first, second and fourth columns of the matrix S , i.e.,

$$S_{:, \mathcal{M}_{\text{obs}}} = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}.$$

In the literature of LP, the matrix $S_{:, \mathcal{M}_{\text{obs}}}$ is also referred to as a *basis*. After obtaining $\hat{\mathbf{t}}$, the Wasserstein distance can be re-written as

$$\begin{aligned}
 W(P_n, Q_m) &= \hat{\mathbf{t}}^\top \mathbf{c}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\
 &= \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \mathbf{c}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\
 &= \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \underbrace{\Theta_{\mathcal{M}_{\text{obs}}, :} (\mathbf{X}^{\text{obs}} \mathbf{Y}^{\text{obs}})^\top}_{\mathbf{c}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})},
 \end{aligned} \tag{9}$$

where Θ is defined in (5), and $\Theta_{\mathcal{M}_{\text{obs}}, :}$ is a sub-matrix of Θ made up of rows in the set \mathcal{M}_{obs} and all columns.

2.4 Statistical inference (confidence interval)

The goal is to provide a confidence interval (CI) for the Wasserstein distance. The $W(P_n, Q_m)$ in (9) can be written as

$$W(P_n, Q_m) = \boldsymbol{\eta}^\top (\mathbf{X}^{\text{obs}} \mathbf{Y}^{\text{obs}})^\top, \tag{10}$$

where $\boldsymbol{\eta} = \Theta_{\mathcal{M}_{\text{obs}}, :}^\top \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}$ is the test-statistic direction. It is important to note that $\boldsymbol{\eta}$ is a *random* vector because \mathcal{M}_{obs} is selected based on the data. For the purpose of explanation, let us suppose, for now, that

the test-statistic direction $\boldsymbol{\eta}$ in (10) is fixed before observing the data; that is, non-random. Let us define

$$\tilde{\Sigma} = \begin{pmatrix} \Sigma_{\mathbf{X}} & 0 \\ 0 & \Sigma_{\mathbf{Y}} \end{pmatrix}. \quad (11)$$

Thus, we have $\boldsymbol{\eta}^\top (\mathbf{X} \ \mathbf{Y})^\top \sim \mathbb{N} \left(\boldsymbol{\eta}^\top (\boldsymbol{\mu}_{\mathbf{X}} \ \boldsymbol{\mu}_{\mathbf{Y}})^\top, \boldsymbol{\eta}^\top \tilde{\Sigma} \boldsymbol{\eta} \right)$. Given a significance level $\alpha \in [0, 1]$ (e.g., 0.05) for the inference, the *naive* (classical) CI for

$$W^* = W^*(P_m, Q_m) = \boldsymbol{\eta}^\top (\boldsymbol{\mu}_{\mathbf{X}} \ \boldsymbol{\mu}_{\mathbf{Y}})^\top$$

with $(1 - \alpha)$ -coverage is obtained by

$$C_{\text{naive}} = \left\{ w \in \mathbb{R} : \frac{\alpha}{2} \leq F_{w, \sigma^2} \left(\boldsymbol{\eta}^\top \begin{pmatrix} \mathbf{X}^{\text{obs}} \\ \mathbf{Y}^{\text{obs}} \end{pmatrix} \right) \leq 1 - \frac{\alpha}{2} \right\}, \quad (12)$$

where $\sigma^2 = \boldsymbol{\eta}^\top \tilde{\Sigma} \boldsymbol{\eta}$ and $F_{w, \sigma^2}(\cdot)$ is the c.d.f of the normal distribution with a mean w and variance σ^2 . With the assumption that $\boldsymbol{\eta}$ in (10) is fixed in advance, the naive CI is valid in the sense that

$$\mathbb{P}(W^* \in C_{\text{naive}}) = 1 - \alpha. \quad (13)$$

However, in reality, because the test-statistic direction $\boldsymbol{\eta}$ is actually not fixed in advance, the naive CI in (12) is *unreliable*. In other words, the naive CI is *invalid* in the sense that the $(1 - \alpha)$ -coverage guarantee in (13) is no longer satisfied. This is because the construction of $\boldsymbol{\eta}$ depends on the data and selection bias exists.

In the next section, we introduce an approach to correct the selection bias which is inspired by the conditional SI framework, and propose a valid CI called selective CI (C_{sel}) for the Wasserstein distance which guarantees the $(1 - \alpha)$ -coverage property in finite sample (i.e., non-asymptotic).

3 Proposed Method

We present the technical details of the proposed method in this section. We first introduce the selective CI for the Wasserstein distance in §3.1. To compute the proposed selective CI, we need to consider the sampling distribution of the Wasserstein distance conditional on the selection event. Thereafter, the selection event is explicitly characterized in a conditional data space whose identification is presented in §3.2. Finally, we end the section with the detailed algorithm.

3.1 Selective Confidence Interval for Wasserstein Distance

In this section, we propose an exact (non-asymptotic) selective CI for the Wasserstein distance by conditional SI. We interpret the computation of the Wasserstein distance in (7) as a two-step procedure:

- **Step 1:** Compute the cost matrix in (3) with the ℓ_1 distance and obtain the vectorized form of the cost matrix in (4).

- **Step 2:** Solving the LP in (7) to obtain \mathcal{M}_{obs} which is subsequently used to construct the test-statistic direction $\boldsymbol{\eta}$ in (10) and calculate the distance.

Since the distance is computed in data-driven fashion, for constructing a valid CI, it is necessary to remove the information that has been used for the initial calculation process, i.e., steps 1 and 2 in the above procedure, to correct the selection bias. This can be achieved by considering the sampling distribution of the test-statistic $\boldsymbol{\eta}^\top (\mathbf{X} \ \mathbf{Y})^\top$ conditional on the selection event; that is,

$$\boldsymbol{\eta}^\top \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \Big| \left\{ \mathcal{S}(\mathbf{X}, \mathbf{Y}) = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}), \mathcal{M}(\mathbf{X}, \mathbf{Y}) = \mathcal{M}_{\text{obs}} \right\}, \quad (14)$$

where $\mathcal{S}(\mathbf{X}, \mathbf{Y})$ is the sign vector explained in the construction of Θ in (5).

Interpretation of the selection event in (14). The first and second conditions in (14) respectively represent the selection event for steps 1 and 2 in the procedure described in the beginning of this section. The first condition $\mathcal{S}(\mathbf{X}, \mathbf{Y}) = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})$ indicates that $\text{sign}(x_i - y_j) = \text{sign}(x_i^{\text{obs}} - y_j^{\text{obs}})$ for $i \in [n], j \in [m]$. In other words, for $i \in [n], j \in [m]$, we condition on the event whether x_i^{obs} is greater than y_j^{obs} or not. The second condition $\mathcal{M}(\mathbf{X}, \mathbf{Y}) = \mathcal{M}_{\text{obs}}$ indicates that the set of selected basic variables for random vectors \mathbf{X} and \mathbf{Y} is the same as that for \mathbf{X}^{obs} and \mathbf{Y}^{obs} . This condition can be interpreted as conditioning on the observed optimal coupling \mathcal{M}_{obs} between the elements of \mathbf{X}^{obs} and \mathbf{Y}^{obs} , which is obtained after solving the LP in (7) on the observed data (see Figure 1).

Selective CI. To derive exact statistical inference for the Wasserstein distance, we introduce so-called selective CI for $W^* = W^*(P_m, Q_m) = \boldsymbol{\eta}^\top (\boldsymbol{\mu}_{\mathbf{X}} \ \boldsymbol{\mu}_{\mathbf{Y}})^\top$ that satisfies the following $(1 - \alpha)$ -coverage property conditional on the selection event:

$$\mathbb{P} \left(W^* \in C_{\text{sel}} \mid \mathcal{S}(\mathbf{X}, \mathbf{Y}) = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}), \mathcal{M}(\mathbf{X}, \mathbf{Y}) = \mathcal{M}_{\text{obs}} \right) = 1 - \alpha, \quad (15)$$

for any $\alpha \in [0, 1]$. The selective CI is defined as

$$C_{\text{sel}} = \left\{ w \in \mathbb{R} : \frac{\alpha}{2} \leq F_{w, \sigma^2}^{\mathcal{Z}} \left(\boldsymbol{\eta}^\top \begin{pmatrix} \mathbf{X}^{\text{obs}} \\ \mathbf{Y}^{\text{obs}} \end{pmatrix} \right) \leq 1 - \frac{\alpha}{2} \right\}. \quad (16)$$

where the pivotal quantity

$$F_{w, \sigma^2}^{\mathcal{Z}} \left(\boldsymbol{\eta}^\top \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \right) \Big| \left\{ \mathcal{S}(\mathbf{X}, \mathbf{Y}) = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}), \mathcal{M}(\mathbf{X}, \mathbf{Y}) = \mathcal{M}_{\text{obs}}, \mathbf{q}(\mathbf{X}, \mathbf{Y}) = \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \right\} \quad (17)$$

is the c.d.f of the *truncated* normal distribution with a mean $w \in \mathbb{R}$, variance $\sigma^2 = \boldsymbol{\eta}^\top \tilde{\Sigma} \boldsymbol{\eta}$, and truncation region \mathcal{Z} (the detailed construction of \mathcal{Z} will be discussed later in §3.2) which is calculated based on the selection event in (17). The $\mathbf{q}(\mathbf{X}, \mathbf{Y})$ in the additional third condition is the sufficient statistic of nuisance parameter defined as

$$\mathbf{q}(\mathbf{X}, \mathbf{Y}) = \left(I_{n+m} - \mathbf{c} \boldsymbol{\eta}^\top \right) (\mathbf{X} \ \mathbf{Y})^\top$$

in which $\mathbf{c} = \tilde{\Sigma} \boldsymbol{\eta} (\boldsymbol{\eta}^\top \tilde{\Sigma} \boldsymbol{\eta})^{-1}$ with $\tilde{\Sigma}$ is defined in (11). Here, we note that the selective CI depends on $\mathbf{q}(\mathbf{X}, \mathbf{Y})$ because the pivotal quantity in (17) depends on this component, but the sampling property in (15)

is kept satisfied without this additional condition because we can marginalize over all values of $\mathbf{q}(\mathbf{X}, \mathbf{Y})$. The $\mathbf{q}(\mathbf{X}, \mathbf{Y})$ corresponds to the component \mathbf{z} in the seminal paper of Lee et al. [2016] (see Section 5, Eq. 5.2 and Theorem 5.2). We note that additionally conditioning on $\mathbf{q}(\mathbf{X}, \mathbf{Y})$ is a standard approach in the SI literature and it is used in almost all the SI-related works that we cited in this paper.

To obtain the selective CI in (16), we need to compute the quantity in (17) which depends on the truncation region \mathcal{Z} . Therefore, the remaining task is to identify \mathcal{Z} whose characterization will be introduced in the next section.

3.2 Conditional Data Space Characterization

We define the set of $(\mathbf{X} \ \mathbf{Y})^\top \in \mathbb{R}^{n+m}$ that satisfies the conditions in Equation (17) as

$$\mathcal{D} = \left\{ \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \in \mathbb{R}^{n+m} \mid \mathcal{S}(\mathbf{X}, \mathbf{Y}) = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}), \mathcal{M}(\mathbf{X}, \mathbf{Y}) = \mathcal{M}_{\text{obs}}, \mathbf{q}(\mathbf{X}, \mathbf{Y}) = \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \right\}. \quad (18)$$

According to the third condition $\mathbf{q}(\mathbf{X}, \mathbf{Y}) = \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})$, the data in \mathcal{D} is restricted to a line in \mathbb{R}^{n+m} as stated in the following Lemma.

Lemma 1. *Let us define*

$$\mathbf{a} = \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \quad \text{and} \quad \mathbf{b} = \tilde{\Sigma} \boldsymbol{\eta} (\boldsymbol{\eta}^\top \tilde{\Sigma} \boldsymbol{\eta})^{-1}, \quad (19)$$

where $\tilde{\Sigma}$ is defined in (11). Then, the set \mathcal{D} in (18) can be rewritten using the scalar parameter $z \in \mathbb{R}$ as follows:

$$\mathcal{D} = \left\{ (\mathbf{X} \ \mathbf{Y})^\top = \mathbf{a} + \mathbf{b}z \mid z \in \mathcal{Z} \right\}, \quad (20)$$

where

$$\mathcal{Z} = \left\{ z \in \mathbb{R} \mid \begin{array}{l} \mathcal{S}(\mathbf{a} + \mathbf{b}z) = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}), \\ \mathcal{M}(\mathbf{a} + \mathbf{b}z) = \mathcal{M}_{\text{obs}} \end{array} \right\}. \quad (21)$$

Here, with a slight abuse of notation, $\mathcal{S}(\mathbf{a} + \mathbf{b}z) = \mathcal{S}((\mathbf{X} \ \mathbf{Y})^\top)$ is equivalent to $\mathcal{S}(\mathbf{X}, \mathbf{Y})$. This similarly applies to $\mathcal{M}(\mathbf{a} + \mathbf{b}z)$.

Proof. According to the third condition in (18), we have

$$\begin{aligned} \mathbf{q}(\mathbf{X}, \mathbf{Y}) &= \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\ \Leftrightarrow (I_{n+m} - \mathbf{c}\boldsymbol{\eta}^\top) \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} &= \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\ \Leftrightarrow \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} &= \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) + \frac{\tilde{\Sigma} \boldsymbol{\eta}}{\boldsymbol{\eta}^\top \tilde{\Sigma} \boldsymbol{\eta}} \boldsymbol{\eta}^\top \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}. \end{aligned}$$

By defining $\mathbf{a} = \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})$, $\mathbf{b} = \tilde{\Sigma} \boldsymbol{\eta} (\boldsymbol{\eta}^\top \tilde{\Sigma} \boldsymbol{\eta})^{-1}$, $z = \boldsymbol{\eta}^\top (\mathbf{X} \ \mathbf{Y})^\top$, and incorporating the first and second conditions in (18), we obtain the results in Lemma 1. We note that the fact of restricting the data to the line has been already implicitly exploited in the seminal conditional SI work of Lee et al. [2016], but explicitly discussed for the first time in Section 6 of Liu et al. [2018]. ■

Lemma 1 indicates that we do not have to consider the $(n + m)$ -dimensional data space. Instead, we only to consider the *one-dimensional projected* data space \mathcal{Z} in (21), which is the truncation region that is important for computing the pivotal quantity in (17) and constructing the selective CI C_{sel} in (16).

Characterization of truncation region \mathcal{Z} . We can decompose \mathcal{Z} into two separate sets as $\mathcal{Z} = \mathcal{Z}_1 \cap \mathcal{Z}_2$, where

$$\begin{aligned} \mathcal{Z}_1 &= \{z \in \mathbb{R} \mid \mathcal{S}(\mathbf{a} + \mathbf{b}z) = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})\} \\ \text{and } \mathcal{Z}_2 &= \{z \in \mathbb{R} \mid \mathcal{M}(\mathbf{a} + \mathbf{b}z) = \mathcal{M}_{\text{obs}}\}. \end{aligned}$$

We first present the construction of \mathcal{Z}_1 in the following Lemma.

Lemma 2. *For notational simplicity, we denote $\mathbf{s}_{\text{obs}} = \mathcal{S}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})$. Then, the \mathcal{Z}_1 is an interval defined as:*

$$\mathcal{Z}_1 = \left\{ z \in \mathbb{R} \mid \max_{j: \nu_j^{(2)} > 0} \frac{-\nu_j^{(1)}}{\nu_j^{(2)}} \leq z \leq \min_{j: \nu_j^{(2)} < 0} \frac{-\nu_j^{(1)}}{\nu_j^{(2)}} \right\},$$

where $\boldsymbol{\nu}^{(1)} = \mathbf{s}_{\text{obs}} \circ \Omega \mathbf{a}$ and $\boldsymbol{\nu}^{(2)} = \mathbf{s}_{\text{obs}} \circ \Omega \mathbf{b}$.

Proof. Let us first remind that $\mathcal{S}(\mathbf{X}, \mathbf{Y}) = \text{sign}(\Omega(\mathbf{X} \ \mathbf{Y})^\top)$ where Ω is defined in (4). Then, the \mathcal{Z}_1 can be re-written as follows:

$$\begin{aligned} \mathcal{Z}_1 &= \{z \in \mathbb{R} \mid \mathcal{S}(\mathbf{a} + \mathbf{b}z) = \mathbf{s}_{\text{obs}}\} \\ &= \left\{ z \in \mathbb{R} \mid \text{sign}(\Omega(\mathbf{a} + \mathbf{b}z)) = \mathbf{s}_{\text{obs}} \right\} \\ &= \{z \in \mathbb{R} \mid \mathbf{s}_{\text{obs}} \circ \Omega(\mathbf{a} + \mathbf{b}z) \geq \mathbf{0}\}. \end{aligned}$$

By defining $\boldsymbol{\nu}^{(1)} = \mathbf{s}_{\text{obs}} \circ \Omega \mathbf{a}$ and $\boldsymbol{\nu}^{(2)} = \mathbf{s}_{\text{obs}} \circ \Omega \mathbf{b}$, the result of Lemma 2 is straightforward by solving the above system of linear inequalities. ■

Next, we present the construction of \mathcal{Z}_2 . Here, \mathcal{Z}_2 can be interpreted as the set of values of z in which we obtain the same set of the selected basic variables \mathcal{M}_{obs} when applying the LP in (7) on the parametrized data $\mathbf{a} + \mathbf{b}z$.

Lemma 3. *The set \mathcal{Z}_2 is an interval defined as:*

$$\mathcal{Z}_2 = \left\{ z \in \mathbb{R} \mid \max_{j \in \mathcal{M}_{\text{obs}}^c: \tilde{v}_j > 0} \frac{-\tilde{u}_j}{\tilde{v}_j} \leq z \leq \min_{j \in \mathcal{M}_{\text{obs}}^c: \tilde{v}_j < 0} \frac{-\tilde{u}_j}{\tilde{v}_j} \right\},$$

where

$$\begin{aligned} \tilde{\mathbf{u}} &= \left(\mathbf{u}_{\mathcal{M}^c}^\top - \mathbf{u}_{\mathcal{M}}^\top S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c} \right)^\top, \\ \tilde{\mathbf{v}} &= \left(\mathbf{v}_{\mathcal{M}^c}^\top - \mathbf{v}_{\mathcal{M}}^\top S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c} \right)^\top, \end{aligned}$$

$\mathbf{u} = \Theta \mathbf{a}$, $\mathbf{v} = \Theta \mathbf{b}$, Θ is defined as in (5) with observed data.

Proof. We consider the LP in (7) with the parametrized data $\mathbf{a} + \mathbf{b}z$ as follows:

$$\min_{\mathbf{t} \in \mathbb{R}^{nm}} \mathbf{t}^\top \Theta(\mathbf{a} + \mathbf{b}z) \quad \text{s.t.} \quad S\mathbf{t} = \mathbf{h}, \mathbf{t} \geq \mathbf{0}. \quad (22)$$

Here, we remind that $\Theta(\mathbf{a} + \mathbf{b}z)$ is the vectorized version of the cost matrix defined in (4). The optimization problem in (22) is well-known as the *parametric cost problem* in LP literature (e.g., see Section 8.2 in Murty [1983]). Let us denote $\mathbf{u} = \Theta\mathbf{a}$ and $\mathbf{v} = \Theta\mathbf{b}$, the LP in (22) can be re-written as

$$\min_{\mathbf{t} \in \mathbb{R}^{nm}} (\mathbf{u} + \mathbf{v}z)^\top \mathbf{t} \quad \text{s.t.} \quad S\mathbf{t} = \mathbf{h}, \mathbf{t} \geq \mathbf{0}. \quad (23)$$

Given a fixed value z_0 , let \mathcal{M} be an optimal basic index set of the LP in (23) at $z = z_0$ and \mathcal{M}^c be its complement. Then by partitioning S , \mathbf{t} , \mathbf{u} , and \mathbf{v} as

$$\begin{aligned} S &= [S_{:, \mathcal{M}}, S_{:, \mathcal{M}^c}] \\ \mathbf{t} &= (\mathbf{t}_{\mathcal{M}}, \mathbf{t}_{\mathcal{M}^c}), \quad \mathbf{u} = (\mathbf{u}_{\mathcal{M}}, \mathbf{u}_{\mathcal{M}^c}), \quad \mathbf{v} = (\mathbf{v}_{\mathcal{M}}, \mathbf{v}_{\mathcal{M}^c}), \end{aligned}$$

the LP in (23) becomes

$$\begin{aligned} \min_{\mathbf{t}_{\mathcal{M}}, \mathbf{t}_{\mathcal{M}^c}} & (\mathbf{u}_{\mathcal{M}} + \mathbf{v}_{\mathcal{M}}z)^\top \mathbf{t}_{\mathcal{M}} + (\mathbf{u}_{\mathcal{M}^c} + \mathbf{v}_{\mathcal{M}^c}z)^\top \mathbf{t}_{\mathcal{M}^c} \\ \text{s.t.} & S_{:, \mathcal{M}}\mathbf{t}_{\mathcal{M}} + S_{:, \mathcal{M}^c}\mathbf{t}_{\mathcal{M}^c} = \mathbf{h}, \\ & \mathbf{t}_{\mathcal{M}} \geq \mathbf{0}, \mathbf{t}_{\mathcal{M}^c} \geq \mathbf{0}. \end{aligned} \quad (24)$$

The value of $\mathbf{t}_{\mathcal{M}}$ can be computed as

$$\mathbf{t}_{\mathcal{M}} = S_{:, \mathcal{M}}^{-1} \mathbf{h} - S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c} \mathbf{t}_{\mathcal{M}^c},$$

and this general expression when substituted in the objective (cost) function of (23) yields

$$\begin{aligned} f &= (\mathbf{u}_{\mathcal{M}} + \mathbf{v}_{\mathcal{M}}z)^\top (S_{:, \mathcal{M}}^{-1} \mathbf{h} - S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c} \mathbf{t}_{\mathcal{M}^c}) + (\mathbf{u}_{\mathcal{M}^c} + \mathbf{v}_{\mathcal{M}^c}z)^\top \mathbf{t}_{\mathcal{M}^c} \\ &= (\mathbf{u}_{\mathcal{M}} + \mathbf{v}_{\mathcal{M}}z)^\top S_{:, \mathcal{M}}^{-1} \mathbf{h} + \left[(\mathbf{u}_{\mathcal{M}^c}^\top - \mathbf{u}_{\mathcal{M}}^\top S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c}) + (\mathbf{v}_{\mathcal{M}^c}^\top - \mathbf{v}_{\mathcal{M}}^\top S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c}) \times z \right] \mathbf{t}_{\mathcal{M}^c}, \end{aligned}$$

which expresses the cost of (24) in terms of $\mathbf{t}_{\mathcal{M}^c}$. Let us denote

$$\tilde{\mathbf{u}} = \left(\mathbf{u}_{\mathcal{M}^c}^\top - \mathbf{u}_{\mathcal{M}}^\top S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c} \right)^\top \quad \text{and} \quad \tilde{\mathbf{v}} = \left(\mathbf{v}_{\mathcal{M}^c}^\top - \mathbf{v}_{\mathcal{M}}^\top S_{:, \mathcal{M}}^{-1} S_{:, \mathcal{M}^c} \right)^\top,$$

we can write $\mathbf{r}_{\mathcal{M}^c} = \tilde{\mathbf{u}} + \tilde{\mathbf{v}}z$ which is known as the *relative cost vector* in the LP literature. Then, \mathcal{M} is an optimal basic index set of (23) for all values of the parameter z satisfying

$$\mathbf{r}_{\mathcal{M}^c} = \tilde{\mathbf{u}} + \tilde{\mathbf{v}}z \geq \mathbf{0}, \quad (25)$$

which is also explicitly discussed in Section 8.2 of Murty [1983]. Finally, the results in Lemma 3 are obtained by respectively replacing \mathcal{M} and \mathcal{M}^c by \mathcal{M}_{obs} and $\mathcal{M}_{\text{obs}}^c$, and solving the linear inequality system in (25). ■

Once \mathcal{Z}_1 and \mathcal{Z}_2 are identified, we can compute the truncation region $\mathcal{Z} = \mathcal{Z}_1 \cap \mathcal{Z}_2$. Finally, we can use \mathcal{Z} to calculate the pivotal quantity in (17) which is subsequently used to construct the proposed selective CI in (16). The details of the algorithm is presented in Algorithm 1.

Algorithm 1 Selective CI for the Wasserstein Distance

Input: $\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}$

- 1: Compute the cost matrix as in (3) and obtained its vectorized version $\mathbf{c}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})$ as in (4)
- 2: Solve LP in (7) to obtain \mathcal{M}_{obs}
- 3: Compute $\boldsymbol{\eta}$ based on $\mathcal{M}_{\text{obs}} \leftarrow$ Equation (10)
- 4: Calculate \mathbf{a} and \mathbf{b} based on $\boldsymbol{\eta} \leftarrow$ Equation (19)
- 5: Construct $\mathcal{Z}_1 \leftarrow$ Lemma 2
- 6: Identify $\mathcal{Z}_2 \leftarrow$ Lemma 3
- 7: Truncation region $\mathcal{Z} = \mathcal{Z}_1 \cap \mathcal{Z}_2$
- 8: $C_{\text{sel}} \leftarrow$ Equation (16)

Output: C_{sel}

4 Extension to Multi-Dimension

In §2 and §3, we mainly focus on the Wasserstein distance in one-dimension, i.e., $x_{i \in [n]} \in \mathbb{R}$ and $y_{j \in [m]} \in \mathbb{R}$. In this section, we generalize the problem setup and extend the proposed method for the Wasserstein distance in multi-dimension. We consider two random sets X and Y of d -dimensional vectors

$$\begin{aligned} X &= (\mathbf{x}_{1,:}, \dots, \mathbf{x}_{n,:})^\top \in \mathbb{R}^{n \times d}, \\ Y &= (\mathbf{y}_{1,:}, \dots, \mathbf{y}_{m,:})^\top \in \mathbb{R}^{m \times d}, \end{aligned} \tag{26}$$

corrupted with Gaussian noise as

$$\begin{aligned} \mathbb{R}^{nd} \ni \mathbf{X}_{\text{vec}} = \text{vec}(X) &= (\mathbf{x}_{1,:}^\top, \dots, \mathbf{x}_{n,:}^\top)^\top = \boldsymbol{\mu}_{\mathbf{X}_{\text{vec}}} + \boldsymbol{\varepsilon}_{\mathbf{X}_{\text{vec}}}, \quad \boldsymbol{\varepsilon}_{\mathbf{X}_{\text{vec}}} \sim \mathbb{N}(\mathbf{0}, \Sigma_{\mathbf{X}_{\text{vec}}}), \\ \mathbb{R}^{md} \ni \mathbf{Y}_{\text{vec}} = \text{vec}(Y) &= (\mathbf{y}_{1,:}^\top, \dots, \mathbf{y}_{m,:}^\top)^\top = \boldsymbol{\mu}_{\mathbf{Y}_{\text{vec}}} + \boldsymbol{\varepsilon}_{\mathbf{Y}_{\text{vec}}}, \quad \boldsymbol{\varepsilon}_{\mathbf{Y}_{\text{vec}}} \sim \mathbb{N}(\mathbf{0}, \Sigma_{\mathbf{Y}_{\text{vec}}}), \end{aligned}$$

where n and m are the number of instances in each set, $\boldsymbol{\mu}_{\mathbf{X}_{\text{vec}}}$ and $\boldsymbol{\mu}_{\mathbf{Y}_{\text{vec}}}$ are unknown mean vectors, $\boldsymbol{\varepsilon}_{\mathbf{X}_{\text{vec}}}$ and $\boldsymbol{\varepsilon}_{\mathbf{Y}_{\text{vec}}}$ are Gaussian noise vectors with covariances matrices $\Sigma_{\mathbf{X}_{\text{vec}}}$ and $\Sigma_{\mathbf{Y}_{\text{vec}}}$ assumed to be known or estimable from independent data.

Cost matrix. The cost matrix $C(X, Y)$ of pairwise distances (ℓ_1 distance) between elements of X and Y as

$$C(X, Y) = \left[\sum_{k=1}^d |\mathbf{x}_{i,k} - \mathbf{y}_{j,k}| \right]_{ij} \in \mathbb{R}^{n \times m}. \tag{27}$$

Then, the vectorized form of $C(X, Y)$ can be defined as

$$\begin{aligned} \mathbf{c}(X, Y) &= \text{vec}(C(X, Y)) \\ &= \Theta^{\text{mul}}(\mathbf{X}_{\text{vec}} \mathbf{Y}_{\text{vec}})^\top \in \mathbb{R}^{nm}, \end{aligned} \tag{28}$$

where

$$\Theta^{\text{mul}} = \sum_{k=1}^d \mathcal{S}_k(X, Y) \circ (\Omega \otimes \mathbf{e}_{d,k}^\top) \in \mathbb{R}^{nm \times (nd+md)},$$

$$\mathcal{S}_k(X, Y) = \text{sign} \left((\Omega \otimes \mathbf{e}_{d,k}^\top) (\mathbf{X}_{\text{vec}} \mathbf{Y}_{\text{vec}})^\top \right) \in \mathbb{R}^{nm},$$

the matrix Ω is defined in (5), the operator \otimes is Kronecker product, and $\mathbf{e}_{d,k} \in \mathbb{R}^d$ is a d -dimensional unit vector with 1 at position $k \in [d]$.

The Wasserstein distance in multi-dimension. Given X^{obs} and Y^{obs} sampled from (26), after obtaining $\mathbf{c}(X^{\text{obs}}, Y^{\text{obs}})$ as in (28), the Wasserstein distance in multi-dimension is defined as

$$W^{\text{mul}}(P_n, Q_m) = \min_{\mathbf{t} \in \mathbb{R}^{nm}} \mathbf{t}^\top \mathbf{c}(X^{\text{obs}}, Y^{\text{obs}}) \quad (29)$$

s.t. $S\mathbf{t} = \mathbf{h}, \mathbf{t} \geq \mathbf{0},$

where S and \mathbf{h} are defined in (7). By solving LP in (29), we obtain the set of selected basic variables

$$\mathcal{M}_{\text{obs}} = \mathcal{M}(X^{\text{obs}}, Y^{\text{obs}}), \quad (30)$$

Then, the Wasserstein distance can be re-written as

$$\begin{aligned} W^{\text{mul}}(P_n, Q_m) &= \hat{\mathbf{t}}^\top \mathbf{c}(X^{\text{obs}}, Y^{\text{obs}}) \\ &= \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \mathbf{c}_{\mathcal{M}_{\text{obs}}}(X^{\text{obs}}, Y^{\text{obs}}) \\ &= \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \Theta_{\mathcal{M}_{\text{obs}},:}^{\text{mul}} (\mathbf{X}_{\text{vec}}^{\text{obs}} \mathbf{Y}_{\text{vec}}^{\text{obs}})^\top \\ &= \boldsymbol{\eta}_{\text{mul}}^\top (\mathbf{X}_{\text{vec}}^{\text{obs}} \mathbf{Y}_{\text{vec}}^{\text{obs}})^\top \end{aligned} \quad (31)$$

where $\boldsymbol{\eta}_{\text{mul}} = (\hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \Theta_{\mathcal{M}_{\text{obs}},:}^{\text{mul}})^\top$ is the test-statistic direction, $\hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}} = S_{:, \mathcal{M}_{\text{obs}}}^{-1} \mathbf{h}$ is the optimal solution of (29), and the matrix Θ^{mul} is defined in (28).

Selection event and selective CI. Since we are dealing with multi-dimensional case, the selection event is slightly different from but more general than the one presented in (17) of §3. Specifically, we consider the following conditional inference

$$\boldsymbol{\eta}_{\text{mul}}^\top (\mathbf{X}_{\text{vec}} \mathbf{Y}_{\text{vec}})^\top \mid \mathcal{E}^{\text{mul}}, \quad (32)$$

where

$$\mathcal{E}^{\text{mul}} = \left\{ \bigcup_{k=1}^d \mathcal{S}_k(X, Y) = \mathcal{S}_k(X^{\text{obs}}, Y^{\text{obs}}), \mathcal{M}(X, Y) = \mathcal{M}_{\text{obs}}, \mathbf{q}(X, Y) = \mathbf{q}(X^{\text{obs}}, Y^{\text{obs}}) \right\}.$$

Once the selection event \mathcal{E}^{mul} has been identified, the pivotal quantity can be computed:

$$F_{\boldsymbol{\eta}_{\text{mul}}^\top (\boldsymbol{\mu}_{\mathbf{X}_{\text{vec}}} \boldsymbol{\mu}_{\mathbf{Y}_{\text{vec}}})^\top, \sigma_{\text{mul}}^2}^{\mathcal{Z}^{\text{mul}}} \left(\boldsymbol{\eta}_{\text{mul}}^\top (\mathbf{X}_{\text{vec}} \mathbf{Y}_{\text{vec}})^\top \mid \mathcal{E}^{\text{mul}} \right)$$

where $\sigma_{\text{mul}}^2 = \boldsymbol{\eta}_{\text{mul}}^\top \tilde{\Sigma}^{\text{mul}} \boldsymbol{\eta}_{\text{mul}}$ with $\tilde{\Sigma}^{\text{mul}} = \begin{pmatrix} \Sigma_{\mathbf{X}_{\text{vec}}} & 0 \\ 0 & \Sigma_{\mathbf{Y}_{\text{vec}}} \end{pmatrix}$, and truncation region \mathcal{Z}^{mul} is calculated based on the selection event \mathcal{E}^{mul} which we will discuss later. After \mathcal{Z}^{mul} is identified, the selective CI is defined as

$$C_{\text{sel}}^{\text{mul}} = \left\{ w \in \mathbb{R} : \frac{\alpha}{2} \leq F_{w, \sigma_{\text{mul}}^2}^{\mathcal{Z}^{\text{mul}}} \left(\boldsymbol{\eta}_{\text{mul}}^\top \begin{pmatrix} \mathbf{X}_{\text{vec}}^{\text{obs}} \\ \mathbf{Y}_{\text{vec}}^{\text{obs}} \end{pmatrix} \right) \leq 1 - \frac{\alpha}{2} \right\}. \quad (33)$$

The remaining task is to identify \mathcal{Z}^{mul} .

Characterization of truncation region \mathcal{Z}^{mul} . Similar to the discussion in §3, the data is restricted on the line due to the conditioning on the nuisance component $\mathbf{q}(X, Y)$. Then, the set of data that satisfies the condition in (32) is defined as

$$\mathcal{D}^{\text{mul}} = \left\{ (\mathbf{X}_{\text{vec}} \ \mathbf{Y}_{\text{vec}})^\top = \mathbf{a}^{\text{mul}} + \mathbf{b}^{\text{mul}} z \mid z \in \mathcal{Z}^{\text{mul}} \right\},$$

where

$$\begin{aligned} \mathbf{a}^{\text{mul}} &= \mathbf{q}(X^{\text{obs}}, Y^{\text{obs}}), \\ \mathbf{b}^{\text{mul}} &= \tilde{\Sigma}^{\text{mul}} \boldsymbol{\eta}_{\text{mul}} (\boldsymbol{\eta}_{\text{mul}}^\top \tilde{\Sigma}^{\text{mul}} \boldsymbol{\eta}_{\text{mul}})^{-1}, \\ \mathcal{Z}^{\text{mul}} &= \left\{ z \in \mathbb{R} \mid \begin{array}{l} \bigcup_{k=1}^d \mathcal{S}_k(\mathbf{a}^{\text{mul}} + \mathbf{b}^{\text{mul}} z) = \mathcal{S}_k(X^{\text{obs}}, Y^{\text{obs}}), \\ \mathcal{M}(\mathbf{a}^{\text{mul}} + \mathbf{b}^{\text{mul}} z) = \mathcal{M}_{\text{obs}} \end{array} \right\} \end{aligned}$$

with $z \in \mathbb{R}$. Next, we can decompose \mathcal{Z}^{mul} into two separate sets as $\mathcal{Z}^{\text{mul}} = \mathcal{Z}_1^{\text{mul}} \cap \mathcal{Z}_2^{\text{mul}}$, where

$$\begin{aligned} \mathcal{Z}_1^{\text{mul}} &= \left\{ z \in \mathbb{R} \mid \bigcup_{k=1}^d \mathcal{S}_k(\mathbf{a}^{\text{mul}} + \mathbf{b}^{\text{mul}} z) = \mathcal{S}_k(X^{\text{obs}}, Y^{\text{obs}}) \right\}, \\ \mathcal{Z}_2^{\text{mul}} &= \{ z \in \mathbb{R} \mid \mathcal{M}(\mathbf{a}^{\text{mul}} + \mathbf{b}^{\text{mul}} z) = \mathcal{M}_{\text{obs}} \}. \end{aligned}$$

From now, the identification of $\mathcal{Z}_1^{\text{mul}}$ and $\mathcal{Z}_2^{\text{mul}}$ is straightforward and similar to the construction of \mathcal{Z}_1 and \mathcal{Z}_2 discussed in §3. Once $\mathcal{Z}_1^{\text{mul}}$ and $\mathcal{Z}_2^{\text{mul}}$ are identified, we can compute the truncation region $\mathcal{Z}^{\text{mul}} = \mathcal{Z}_1^{\text{mul}} \cap \mathcal{Z}_2^{\text{mul}}$ and use it to compute the selective CI in (33).

5 Experiment

In this section, we demonstrate the performance of the proposed method in both univariate case and multi-dimensional case. We present the results on synthetic data in §5.1. Thereafter, the results on real data are shown in §5.2. In all the experiments, we set the significance level $\alpha = 0.05$, i.e., all the experiments were conducted with the coverage level of $1 - \alpha = 0.95$.

5.1 Numerical Experiment

In this section, we evaluate the performance of the proposed selective CI in terms of coverage guarantee, CI length and computational cost. We also show the results of comparison between our selective CI and the naive

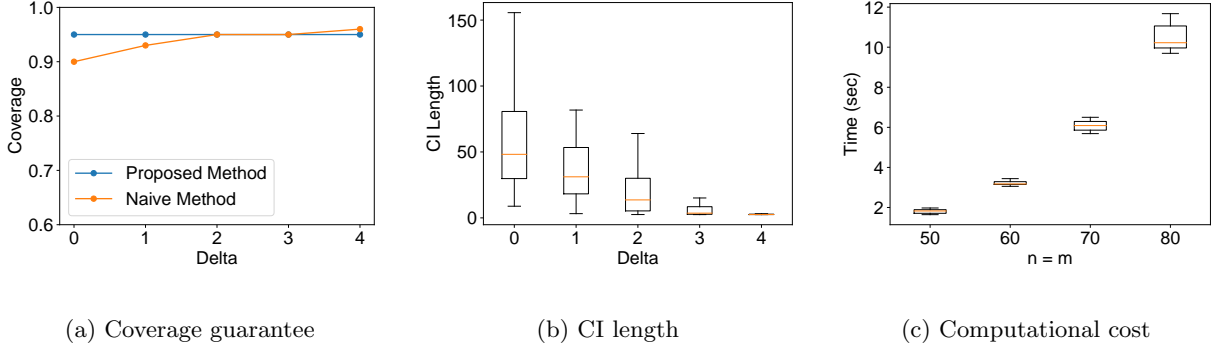


Figure 2: Coverage guarantee, CI length and computational cost in univariate case ($d = 1$).

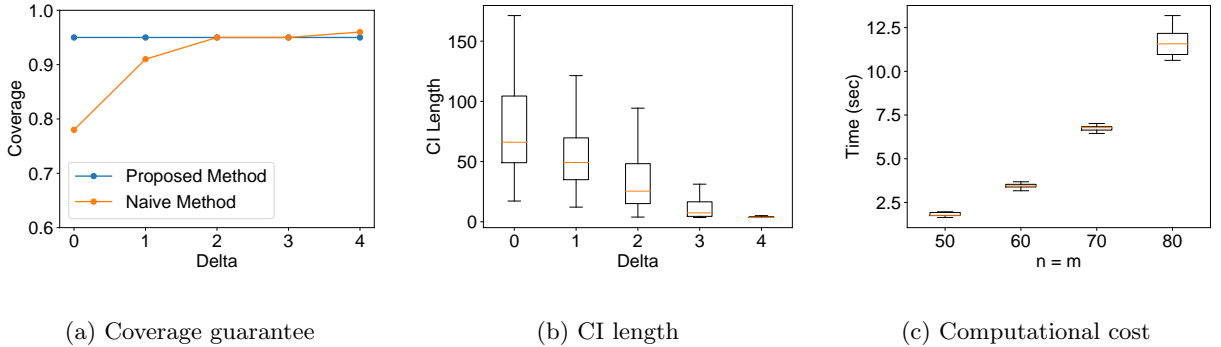


Figure 3: Coverage guarantee, CI length and computational cost in multi-dimensional case ($d = 2$).

CI in (12) in terms of coverage guarantee. We would like to note that we did not conduct the comparison in terms of CI length because the naive CI could not guarantee the coverage property. In statistical viewpoint, if the CI is unreliable, i.e., invalid or does not satisfy the coverage property, the demonstration of CI length does not make sense. Besides, we additionally compare the performance of the proposed method with the latest asymptotic study [Imaizumi et al., 2019].

5.1.1 Univariate case ($d = 1$)

We generated the dataset \mathbf{X} and \mathbf{Y} with $\boldsymbol{\mu}_{\mathbf{X}} = \mathbf{1}_n$, $\boldsymbol{\mu}_{\mathbf{Y}} = \mathbf{1}_m + \Delta$ (element-wise addition), $\boldsymbol{\varepsilon}_{\mathbf{X}} \sim \mathbb{N}(\mathbf{0}, I_n)$, $\boldsymbol{\varepsilon}_{\mathbf{Y}} \sim \mathbb{N}(\mathbf{0}, I_m)$. Regarding the experiments of coverage guarantee and CI length, we set $n = m = 5$ and ran 120 trials for each $\Delta \in \{0, 1, 2, 3, 4\}$. In regard to the experiments of computational cost, we set $\Delta = 2$ and ran 10 trials for each $n = m \in \{50, 60, 70, 80\}$. The results are shown in Figure 2. In the left plot, the naive CI can not properly guarantee the coverage property while the proposed selective CI does. The results in the middle plot indicate that the larger the true distance between \mathbf{X} and \mathbf{Y} , the shorter selective CI we obtain. The right plot shows that the proposed method also has reasonable computational cost.

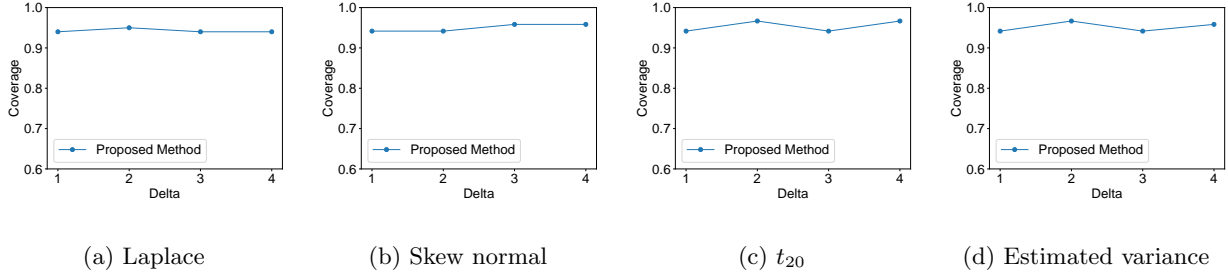


Figure 4: Robustness of the proposed selective CI in terms of coverage guarantee.

5.1.2 Multi-dimensional case ($d = 2$)

We generated the dataset $X = \{\mathbf{x}_{i,:}\}_{i \in [n]}$ with $\mathbf{x}_{i,:} \sim \mathcal{N}(\mathbf{1}_d, I_d)$ and $Y = \{\mathbf{y}_{j,:}\}_{j \in [m]}$ with $\mathbf{y}_{j,:} \sim \mathcal{N}(\mathbf{1}_d + \Delta, I_d)$ (element-wise addition). Similar to the univariate case, we set $n = m = 5$ and ran 120 trials for each $\Delta \in \{0, 1, 2, 3, 4\}$ for the experiments of coverage guarantee and CI length as well as setting $\Delta = 2$ and ran 10 trials for each $n = m \in \{50, 60, 70, 80\}$ for the experiments of computational cost. The results are shown in Figure 3. The interpretation of the results is similar and consistent with the univariate case.

5.1.3 Robustness of the proposed selective CI in terms of coverage guarantee

We additionally demonstrate the robustness of the proposed selective CI in terms of coverage guarantee by considering the following cases:

- Non-normal noise: we considered the noises ε_X and ε_Y following the Laplace distribution, skew normal distribution (skewness coefficient: 10), and t_{20} distribution.
- Unknown variance: we considered the case in which the variance of the noises was also estimated from the data.

The dataset \mathbf{X} and \mathbf{Y} were generated with $\boldsymbol{\mu}_X = \mathbf{1}_n$, $\boldsymbol{\mu}_Y = \mathbf{1}_m + \Delta$. We set $n = m = 5$ and ran 120 trials for each $\Delta \in \{1, 2, 3, 4\}$. We confirmed that our selective CI maintained good performance in terms of coverage guarantee. The results are shown in Figure 4.

5.1.4 Comparison with asymptotic method in Imaizumi et al. [2019]

The authors of Imaizumi et al. [2019] provided us their code of computing the p -value in hypothesis testing framework. Therefore, we conducted the comparisons in terms of false positive rate (FPR) control and true positive rate (TPR). Although we mainly focused on the CI in the previous sections, the corresponding hypothesis testing problem is defined as follows:

$$H_0 : \boldsymbol{\eta}^\top \begin{pmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{pmatrix} = 0 \quad \text{v.s.} \quad H_1 : \boldsymbol{\eta}^\top \begin{pmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{pmatrix} \neq 0.$$

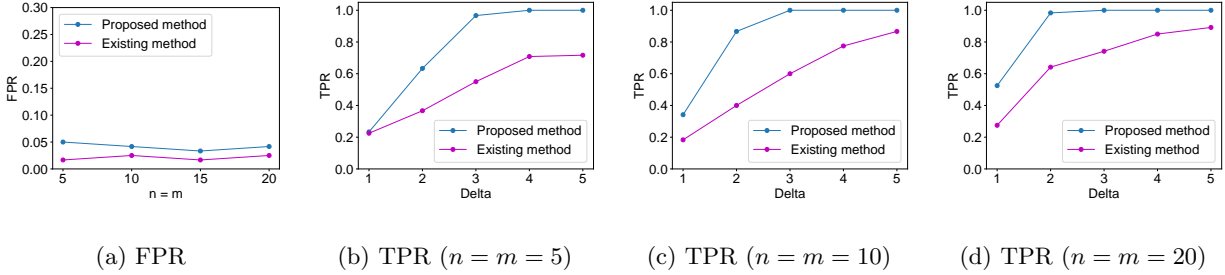


Figure 5: Comparisons with asymptotic method [Imaizumi et al., 2019] in terms of FPR and TPR. While both methods can successfully control the FPR under the significance level $\alpha = 0.05$, the proposed method has higher TPR than the existing asymptotic method in all the cases.

The details are presented in Appendix A. We generated the dataset \mathbf{X} and \mathbf{Y} with $\boldsymbol{\mu}_{\mathbf{X}} = \mathbf{1}_n$, $\boldsymbol{\mu}_{\mathbf{Y}} = \mathbf{1}_m + \Delta$ (element-wise addition), $\boldsymbol{\varepsilon}_{\mathbf{X}} \sim \mathbb{N}(\mathbf{0}, I_n)$, $\boldsymbol{\varepsilon}_{\mathbf{Y}} \sim \mathbb{N}(\mathbf{0}, I_m)$. Regarding the FPR experiments, we set $\Delta = 0$ and ran 120 trials for each $n = m \in \{5, 10, 15, 20\}$. In regard to the TPR experiments, we set $\Delta \in \{1, 2, 3, 4, 5\}$ and ran 120 trials for each $n = m \in \{5, 10, 20\}$. The results are shown in Figure 5. In terms for FPR control, both methods could successfully control the FPR under $\alpha = 0.05$. However, in terms of TPR, the proposed method outperformed the existing asymptotic one in all the cases. As observed in Figure 5 (a), the existing method is conservative in the sense that the FPR is smaller than the specified significance level $\alpha = 0.05$. As a consequence of this conservativeness, the power of their method was consistently lower than ours. Such a phenomenon is commonly observed in approximate statistical inference.

5.2 Real Data Experiment

In this section, we evaluate the proposed selective CI on four real-world datasets. We used Iris dataset, Wine dataset, Breast Cancer dataset which are available in the UCI machine learning repository, and Lung Cancer dataset ⁴ [Feltus et al., 2019]. The experiments were conducted with the following settings:

- **Setting 1:** For each pair of classes in the dataset:
 - Randomly select $n = m = 5$ instances from each class. Here, each instance is represented by a d -dimensional vector where d is the number of features.
 - Compute the selective CI.
 - Repeat the above process up to 120 times.
- **Setting 2:** Given a dataset with two classes $\mathbf{C1}$ and $\mathbf{C2}$, we either chose $n = 5$ instances from $\mathbf{C1}$ and $m = 5$ instances from $\mathbf{C2}$ (X^{obs} and Y^{obs} are from different classes); or both X^{obs} and Y^{obs} from either $\mathbf{C1}$ or $\mathbf{C2}$ (X^{obs} and Y^{obs} are from the same class). Then, we compute the selective CI. We repeated this process up to 120 times.

⁴We used dataset Lung_GSE7670 which is available at <https://sbcb.inf.ufrgs.br/cumida>.

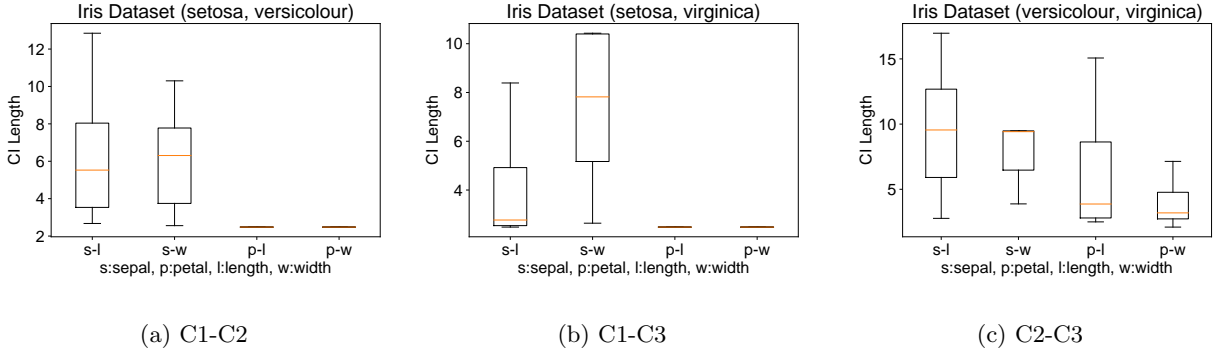


Figure 6: Results on Iris dataset in univariate case ($d = 1$).

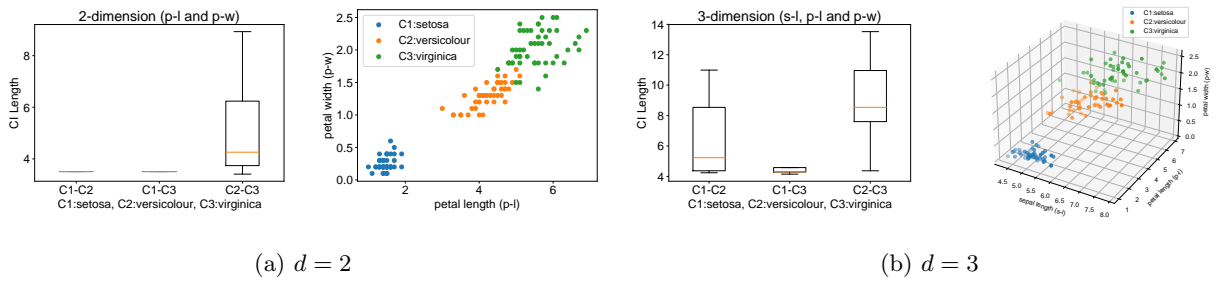


Figure 7: Results on Iris dataset in multi-dimensional case ($d \in \{2, 3\}$)

5.2.1 Univariate case ($d = 1$) with Setting 1

We conducted the experiments on Iris dataset which contains three classes: Iris Setosa (C1), Iris Versicolour (C2), and Iris Virginica (C3). This dataset also contains four features: sepal length ($\mathbf{s-l}$), sepal width ($\mathbf{s-w}$), petal length ($\mathbf{p-l}$), and petal width ($\mathbf{p-w}$). We ran the procedure described in Setting 1 on each individual feature. The results are shown in Figure 6. In all three plots of this figure, the two features $\mathbf{p-l}$ and $\mathbf{p-w}$ always have the shortest CI length among the four features which indicates that these two features are informative to discriminate between the classes. Besides, the results of Figure 6 are also consistent with the results obtained after plotting the histogram of each feature in each class. In other words, the farther the two histograms, the smaller the length of selective CI.

5.2.2 Multi-dimensional case ($d \in \{2, 3\}$) with Setting 1

Regarding the experiments on Iris dataset in multi-dimensional case, we chose two features $\mathbf{p-l}$ and $\mathbf{p-w}$ when $d = 2$ and additionally include feature $\mathbf{s-l}$ when $d = 3$. The results are shown in Figure 7. In each sub-plot, we show the results of the length of selective CI and the corresponding scatter plot which is used to verify the CI length results. For example, in Figure 7a, it is obvious that the distance between C1 and C2 as well as the distance between C1 and C3 are larger than the distance between C2 and C3 by seeing the scatter plot. Therefore, the CI lengths of C1-C2 and C1-C3 tend to be smaller than that of C2-C3. Besides,

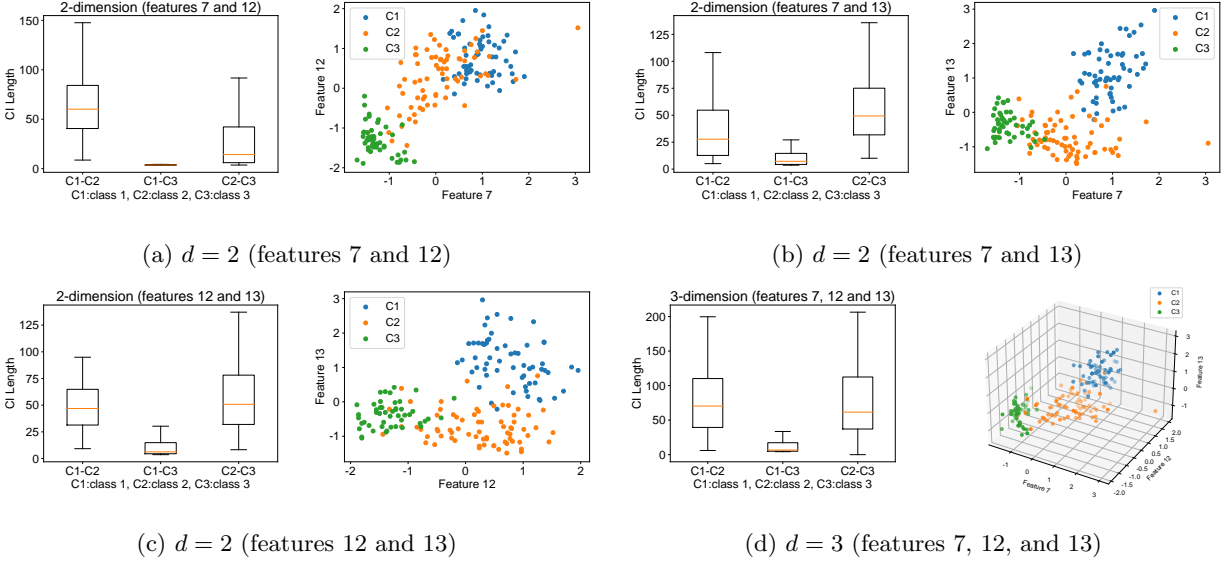


Figure 8: Results on Wine dataset in multi-dimensional case ($d \in \{2, 3\}$)

we also additionally conducted experiments on Wine dataset. This dataset contains 3 classes of wine and 13 features. In the case of $d = 2$, we conducted the experiments on each pair features in the set $\{7, 12, 13\}$ (feature 7: flavanoids, features 12: od280/od315 of diluted wines, feature 13: proline). In the case of $d = 3$, we conducted the experiments on both three features. The results are shown in Figure 8. In general, the results of CI length are consistent with the scatter plots, i.e., the farther the scatter plots between two classes, the smaller the length of selective CI.

5.2.3 Multi-dimensional case with Setting 2

We conducted experiments on Breast Cancer and Lung Cancer datasets. In Breast Cancer dataset, there are two classes (malignant and benign) and $d = 30$ features. In Lung Cancer dataset, there are two classes (normal and adenocarcinoma) and we choose $d = 1,000$ (we selected the top 1,000 genes which have the largest standard deviations as it is commonly done in the literature). The results on these datasets with Setting 2 are shown in Figure 9. The results are consistent with the intuitive expectation. When X^{obs} and Y^{obs} are from different classes, the Wasserstein distance tends to be larger than the one computed when X^{obs} and Y^{obs} are from the same class. Therefore, the CI for the Wasserstein distance in the case of different classes is shorter than the one computed in the case of same class. In other words, the larger the Wasserstein distance is, the shorter the CI becomes.

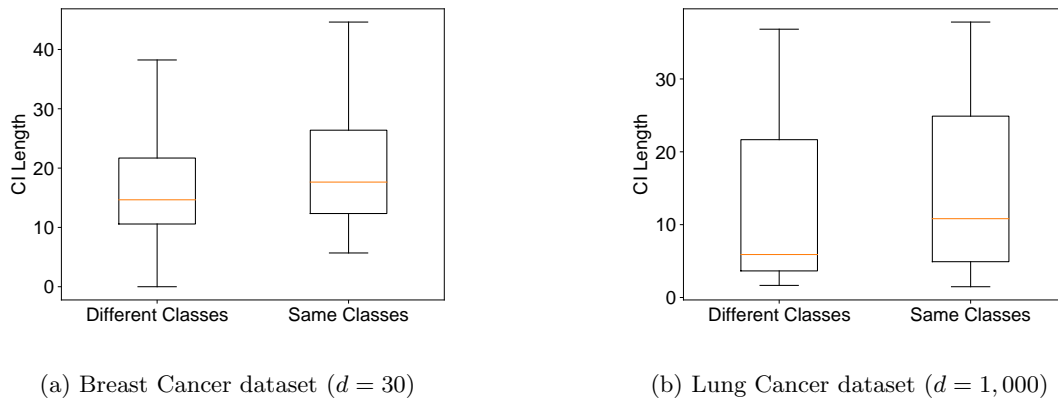


Figure 9: Results on Breast Cancer and Lung Cancer datasets in multi-dimensional case.

6 Conclusion

In this paper, we present an exact (non-asymptotic) statistical inference method for the Wasserstein distance. We first introduce the problem setup and present the proposed method for univariate case. We next provide the extension to multi-dimensional case. We finally conduct the experiments on both synthetic and real-world datasets to evaluate the performance of our method. To our knowledge, this is the first method that can provide a valid confidence interval (CI) for the Wasserstein distance with finite-sample coverage guarantee. We believe this study is an important contribution toward reliable ML, which is one of the most critical issues in the ML community.

References

- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- E. Bernton, P. E. Jacob, M. Gerber, and C. P. Robert. Inference in generative models using the wasserstein distance. *arXiv preprint arXiv:1701.05146*, 1(8):9, 2017.
- S. Chen and J. Bien. Valid inference corrected for outlier removal. *Journal of Computational and Graphical Statistics*, pages 1–12, 2019.
- E. Del Barrio, J. A. Cuesta-Albertos, C. Matrán, and J. M. Rodríguez-Rodríguez. Tests of goodness of fit based on the l2-wasserstein distance. *Annals of Statistics*, pages 1230–1239, 1999.
- E. Del Barrio, P. Gordaliza, H. Lescornel, and J.-M. Loubes. Central limit theorem and bootstrap procedure for wasserstein’s variations with an application to structural relationships between distributions. *Journal of Multivariate Analysis*, 169:341–362, 2019.

- V. N. L. Duy and I. Takeuchi. More powerful conditional selective inference for generalized lasso by parametric programming. *arXiv preprint arXiv:2105.04920*, 2021.
- V. N. L. Duy, S. Iwazaki, and I. Takeuchi. Quantifying statistical significance of neural network representation-driven hypotheses by selective inference. *arXiv preprint arXiv:2010.01823*, 2020a.
- V. N. L. Duy, H. Toda, R. Sugiyama, and I. Takeuchi. Computing valid p-value for optimal changepoint by selective inference using dynamic programming. In *Advances in Neural Information Processing Systems*, 2020b.
- S. N. Evans and F. A. Matsen. The phylogenetic kantorovich–rubinstein metric for environmental sequence samples. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(3):569–592, 2012.
- B. C. Feltes, E. B. Chandelier, B. I. Grisci, and M. Dorn. Cumida: an extensively curated microarray database for benchmarking and testing of machine learning approaches in cancer research. *Journal of Computational Biology*, 26(4):376–386, 2019.
- W. Fithian, D. Sun, and J. Taylor. Optimal inference after model selection. *arXiv preprint arXiv:1410.2597*, 2014.
- C. Frogner, C. Zhang, H. Mobahi, M. Araya-Polo, and T. Poggio. Learning with a wasserstein loss. *arXiv preprint arXiv:1506.05439*, 2015.
- S. Hyun, K. Lin, M. G’Sell, and R. J. Tibshirani. Post-selection inference for changepoint detection algorithms with application to copy number variation data. *arXiv preprint arXiv:1812.03644*, 2018.
- M. Imaizumi, H. Ota, and T. Hamaguchi. Hypothesis test and confidence analysis with wasserstein distance with general dimension. *arXiv preprint arXiv:1910.07773*, 2019.
- S. Kolouri, S. R. Park, M. Thorpe, D. Slepcev, and G. K. Rohde. Optimal mass transport: Signal processing and machine-learning applications. *IEEE signal processing magazine*, 34(4):43–59, 2017.
- V. N. Le Duy and I. Takeuchi. Parametric programming approach for more powerful and general lasso selective inference. In *International Conference on Artificial Intelligence and Statistics*, pages 901–909. PMLR, 2021.
- J. D. Lee, D. L. Sun, Y. Sun, and J. E. Taylor. Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907–927, 2016.
- K. Liu, J. Markovic, and R. Tibshirani. More powerful post-selection inference, with application to the lasso. *arXiv preprint arXiv:1801.09037*, 2018.
- J. R. Loftus and J. E. Taylor. Selective inference in regression models with groups of variables. *arXiv preprint arXiv:1511.01478*, 2015.

- K. Murty. *Linear Programming*. Wiley, 1983. ISBN 9780471097259.
- K. Ni, X. Bresson, T. Chan, and S. Esedoglu. Local histogram based segmentation using the wasserstein distance. *International journal of computer vision*, 84(1):97–111, 2009.
- A. Ramdas, N. G. Trillos, and M. Cuturi. On wasserstein two-sample testing and related families of non-parametric tests. *Entropy*, 19(2):47, 2017.
- K. Sugiyama, V. N. Le Duy, and I. Takeuchi. More powerful and general selective inference for stepwise feature selection using homotopy method. In *International Conference on Machine Learning*, pages 9891–9901. PMLR, 2021a.
- R. Sugiyama, H. Toda, V. N. L. Duy, Y. Inatsu, and I. Takeuchi. Valid and exact statistical inference for multi-dimensional multiple change-points by selective inference. *arXiv preprint arXiv:2110.08989*, 2021b.
- S. Suzumura, K. Nakagawa, Y. Umezu, K. Tsuda, and I. Takeuchi. Selective inference for sparse high-order interaction models. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3338–3347. JMLR. org, 2017.
- K. Tanizaki, N. Hashimoto, Y. Inatsu, H. Hontani, and I. Takeuchi. Computing valid p-values for image segmentation by selective inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9553–9562, 2020.
- R. J. Tibshirani, J. Taylor, R. Lockhart, and R. Tibshirani. Exact post-selection inference for sequential regression procedures. *Journal of the American Statistical Association*, 111(514):600–620, 2016.
- T. Tsukurimichi, Y. Inatsu, V. N. L. Duy, and I. Takeuchi. Conditional selective inference for robust regression and outlier detection using piecewise-linear homotopy continuation. *arXiv preprint arXiv:2104.10840*, 2021.
- C. Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.
- F. Yang, R. F. Barber, P. Jain, and J. Lafferty. Selective inference for group-sparse linear models. In *Advances in Neural Information Processing Systems*, pages 2469–2477, 2016.

A Proposed method in hypothesis testing framework

We present the proposed method in the setting of hypothesis testing and consider the case when the cost matrix is defined by using squared ℓ_2 distance.

Cost matrix. We define the cost matrix $C(\mathbf{X}, \mathbf{Y})$ of pairwise distances (squared ℓ_2 distance) between elements of \mathbf{X} and \mathbf{Y} as

$$C(\mathbf{X}, \mathbf{Y}) = [(x_i - y_j)^2]_{ij} \in \mathbb{R}^{n \times m}. \quad (34)$$

Then, the vectorized form of $C(\mathbf{X}, \mathbf{Y})$ can be defined as

$$\begin{aligned} \mathbf{c}(\mathbf{X}, \mathbf{Y}) &= \text{vec}(C(\mathbf{X}, \mathbf{Y})) \in \mathbb{R}^{nm} \\ &= \begin{bmatrix} \Omega(\mathbf{X}) \\ \mathbf{Y} \end{bmatrix} \circ \begin{bmatrix} \Omega(\mathbf{X}) \\ \mathbf{Y} \end{bmatrix}, \end{aligned} \quad (35)$$

where Ω is defined as in (5) and the operation \circ is element-wise product.

The Wasserstein distance. By solving LP with the cost vector defined in (35) on the observed data \mathbf{X}^{obs} and \mathbf{Y}^{obs} , we obtain the set of selected basic variables

$$\mathcal{M}_{\text{obs}} = \mathcal{M}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}).$$

Then, the Wasserstein distance can be re-written as (we denote $W = W(P_n, Q_m)$ for notational simplicity)

$$\begin{aligned} W &= \hat{\mathbf{t}}^\top \mathbf{c}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\ &= \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \mathbf{c}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\ &= \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \left[\begin{bmatrix} \Omega_{\mathcal{M}_{\text{obs},:}}(\mathbf{X}^{\text{obs}}) \\ \mathbf{Y}^{\text{obs}} \end{bmatrix} \circ \begin{bmatrix} \Omega_{\mathcal{M}_{\text{obs},:}}(\mathbf{X}^{\text{obs}}) \\ \mathbf{Y}^{\text{obs}} \end{bmatrix} \right]. \end{aligned}$$

Hypothesis testing. Our goal is to test the following hypothesis:

$$\text{H}_0 : \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \left[\begin{bmatrix} \Omega_{\mathcal{M}_{\text{obs},:}}(\boldsymbol{\mu}_X) \\ \boldsymbol{\mu}_Y \end{bmatrix} \circ \begin{bmatrix} \Omega_{\mathcal{M}_{\text{obs},:}}(\boldsymbol{\mu}_X) \\ \boldsymbol{\mu}_Y \end{bmatrix} \right] = 0.$$

Unfortunately, it is technically difficult to directly test the above hypothesis. Therefore, we propose to test the following equivalent one:

$$\begin{aligned} \text{H}_0 &: \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}^\top \left[\Theta_{\mathcal{M}_{\text{obs},:}} \begin{bmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{bmatrix} \right] = 0 \\ \Leftrightarrow \text{H}_0 &: \boldsymbol{\eta}^\top \begin{bmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{bmatrix} = 0 \end{aligned}$$

where Θ is defined as in (5) and $\boldsymbol{\eta} = \Theta_{\mathcal{M}_{\text{obs},:}}^\top \hat{\mathbf{t}}_{\mathcal{M}_{\text{obs}}}$.

Conditional SI. To test the aforementioned hypothesis, we consider the following selective p -value:

$$p_{\text{selective}} = \mathbb{P}_{\text{H}_0} \left(\left| \boldsymbol{\eta}^\top \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \right| \geq \left| \boldsymbol{\eta}^\top \begin{bmatrix} \mathbf{X}^{\text{obs}} \\ \mathbf{Y}^{\text{obs}} \end{bmatrix} \right| \mid \mathcal{E} \right)$$

where the conditional selection event is defined as

$$\mathcal{E} = \left\{ \begin{array}{l} \mathcal{M}(\mathbf{X}, \mathbf{Y}) = \mathcal{M}_{\text{obs}}, \\ \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}, \mathbf{Y}) = \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \\ \mathbf{q}(\mathbf{X}, \mathbf{Y}) = \mathbf{q}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \end{array} \right\}.$$

Our next task is to identify the conditional data space whose data satisfies \mathcal{E} .

Characterization of the conditional data space. Similar to the discussion in §3, the data is restricted on the line due to the conditioning on the nuisance component $\mathbf{q}(\mathbf{X}, \mathbf{Y})$. Then, the conditional data space is defined as

$$\mathcal{D} = \left\{ (\mathbf{X} \ \mathbf{Y})^\top = \mathbf{a} + \mathbf{b}z \mid z \in \mathcal{Z} \right\},$$

where

$$\mathcal{Z} = \left\{ z \in \mathbb{R} \mid \begin{array}{l} \mathcal{M}(\mathbf{a} + \mathbf{b}z) = \mathcal{M}_{\text{obs}}, \\ \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{a} + \mathbf{b}z) = \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}}) \end{array} \right\}.$$

The remaining task is to construct \mathcal{Z} . We can decompose \mathcal{Z} into two separate sets as $\mathcal{Z} = \mathcal{Z}_1 \cap \mathcal{Z}_2$, where

$$\begin{aligned} \mathcal{Z}_1 &= \{z \in \mathbb{R} \mid \mathcal{M}(\mathbf{a} + \mathbf{b}z) = \mathcal{M}_{\text{obs}}\}, \\ \mathcal{Z}_2 &= \{z \in \mathbb{R} \mid \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{a} + \mathbf{b}z) = \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})\}. \end{aligned}$$

The construction of \mathcal{Z}_2 is as follows (we denote $\mathbf{s}_{\mathcal{M}_{\text{obs}}} = \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{X}^{\text{obs}}, \mathbf{Y}^{\text{obs}})$ for notational simplicity):

$$\begin{aligned} \mathcal{Z}_2 &= \{z \in \mathbb{R} \mid \mathcal{S}_{\mathcal{M}_{\text{obs}}}(\mathbf{a} + \mathbf{b}z) = \mathbf{s}_{\mathcal{M}_{\text{obs}}}\} \\ &= \left\{ z \in \mathbb{R} \mid \text{sign}\left(\Omega_{\mathcal{M}_{\text{obs}},:}(\mathbf{a} + \mathbf{b}z)\right) = \mathbf{s}_{\mathcal{M}_{\text{obs}}}\right\} \\ &= \{z \in \mathbb{R} \mid \mathbf{s}_{\mathcal{M}_{\text{obs}}} \circ \Omega_{\mathcal{M}_{\text{obs}},:}(\mathbf{a} + \mathbf{b}z) \geq \mathbf{0}\}, \end{aligned}$$

which can be obtained by solving the system of linear inequalities. Next, we present the identification of \mathcal{Z}_1 . Because we use squared ℓ_2 distance to define the cost matrix, the LP with the parametrized data $\mathbf{a} + \mathbf{b}z$ is written as follows:

$$\begin{aligned} &\min_{\mathbf{t} \in \mathbb{R}^{nm}} \mathbf{t}^\top [\Omega(\mathbf{a} + \mathbf{b}z) \circ \Omega(\mathbf{a} + \mathbf{b}z)] \quad \text{s.t.} \quad \mathbf{S}\mathbf{t} = \mathbf{h}, \mathbf{t} \geq \mathbf{0} \\ \Leftrightarrow &\min_{\mathbf{t} \in \mathbb{R}^{nm}} (\mathbf{u} + \mathbf{v}z + \mathbf{w}z^2)^\top \mathbf{t} \quad \text{s.t.} \quad \mathbf{S}\mathbf{t} = \mathbf{h}, \mathbf{t} \geq \mathbf{0}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{u} &= (\Omega\mathbf{a}) \circ (\Omega\mathbf{a}), \\ \mathbf{v} &= (\Omega\mathbf{a}) \circ (\Omega\mathbf{b}) + (\Omega\mathbf{b}) \circ (\Omega\mathbf{a}), \\ \mathbf{w} &= (\Omega\mathbf{b}) \circ (\Omega\mathbf{b}), \end{aligned}$$

and \mathbf{S} and \mathbf{h} are the same as in (7). By fixing \mathcal{M}_{obs} as the optimal basic index set, the *relative cost vector* w.r.t to the set of non-basic variables is defines as

$$\mathbf{r}_{\mathcal{M}_{\text{obs}}^c} = \tilde{\mathbf{u}} + \tilde{\mathbf{v}}z + \tilde{\mathbf{w}}z^2,$$

where

$$\begin{aligned} \tilde{\mathbf{u}} &= \left(\mathbf{u}_{\mathcal{M}_{\text{obs}}^c}^\top - \mathbf{u}_{\mathcal{M}_{\text{obs}}}^\top \mathbf{S}_{:, \mathcal{M}_{\text{obs}}}^{-1} \mathbf{S}_{:, \mathcal{M}_{\text{obs}}^c} \right)^\top, \\ \tilde{\mathbf{v}} &= \left(\mathbf{v}_{\mathcal{M}_{\text{obs}}^c}^\top - \mathbf{v}_{\mathcal{M}_{\text{obs}}}^\top \mathbf{S}_{:, \mathcal{M}_{\text{obs}}}^{-1} \mathbf{S}_{:, \mathcal{M}_{\text{obs}}^c} \right)^\top, \\ \text{and } \tilde{\mathbf{w}} &= \left(\mathbf{w}_{\mathcal{M}_{\text{obs}}^c}^\top - \mathbf{w}_{\mathcal{M}_{\text{obs}}}^\top \mathbf{S}_{:, \mathcal{M}_{\text{obs}}}^{-1} \mathbf{S}_{:, \mathcal{M}_{\text{obs}}^c} \right)^\top. \end{aligned}$$

The requirement for \mathcal{M}_{obs} to be the optimal basis index set is $\mathbf{r}_{\mathcal{M}_{\text{obs}}^c} \geq \mathbf{0}$ (i.e., the cost in minimization problem will never decrease when the non-basic variables become positive and enter the basis). Finally, the

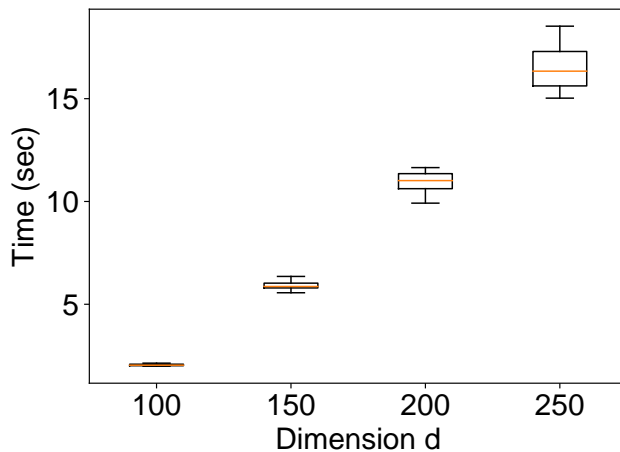


Figure 10: Computational time of the proposed method when increasing the dimension d .

set \mathcal{Z}_1 can be defined as

$$\begin{aligned} \mathcal{Z}_1 &= \{z \in \mathbb{R} \mid \mathcal{M}(\mathbf{a} + \mathbf{b}z) = \mathcal{M}_{\text{obs}}\}, \\ \mathcal{Z}_1 &= \{z \in \mathbb{R} \mid \mathbf{r}_{\mathcal{M}_{\text{obs}}^c} = \tilde{\mathbf{u}} + \tilde{\mathbf{v}}z + \tilde{\mathbf{w}}z^2 \geq \mathbf{0}\}, \end{aligned}$$

which can be obtained by solving the system of quadratic inequalities.

B Experiment on High-dimensional Data

We generated the dataset $X = \{\mathbf{x}_{i,:}\}_{i \in [n]}$ with $\mathbf{x}_{i,:} \sim \mathbb{N}(\mathbf{1}_d, I_d)$ and $Y = \{\mathbf{y}_{j,:}\}_{j \in [m]}$ with $\mathbf{y}_{j,:} \sim \mathbb{N}(\mathbf{1}_d + \Delta, I_d)$ (element-wise addition). We set $n = m = 20, \Delta = 2$, and ran 10 trials for each $d \in \{100, 150, 200, 250\}$. The results in Figure 10 show that the proposed method still has reasonable computational time.