

# GOING BEYOND STATIC: UNDERSTANDING SHIFTS WITH TIME-SERIES ATTRIBUTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Distribution shifts in time-series data are complex due to temporal dependencies, multivariable interactions, and trend changes. However, robust methods often rely on structural assumptions that lack thorough empirical validation, limiting their practical applicability. In order to support an empirically grounded inductive approach to research, we introduce our **Time-Series Shift Attribution (TSSA)** framework, which analyzes **problem-specific** patterns of distribution shifts. Our framework attributes performance degradation from various types of shifts to each *temporal data property* in a detailed manner, supported by theoretical analysis of unbiasedness and asymptotic properties. Empirical studies in real-world healthcare applications highlight how the TSSA framework enhances the understanding of time-series shifts, facilitating reliable model deployment and driving targeted improvements from both algorithmic and data-centric perspectives.

## 1 INTRODUCTION

Machine learning models are increasingly deployed in high-stakes settings such as healthcare with a variety of applications such as early disease screening (Cohn et al., 2003; Soriano et al., 2009) or patient risk assessment (Naghavi et al., 2003; Twetman & Fontana, 2009). While model reliability is vital in such stakes settings, ML model performance often degrades when faced with distribution shifts. This challenge becomes particularly pronounced in time-series data, where unlike static data, the non-stationary properties of time-series such as trends, seasonality and inherent temporal dynamics add additional complexity to the nature of these distribution shifts. In particular, we refer to a unique challenge manifested in time series namely “**shifts in non-stationary properties**”, where the patterns themselves evolve over time.

In critical fields such as healthcare, ignoring these temporal distribution shifts can have life-threatening consequences, posing risks to patient safety and care quality. For instance, predictive models trained on patient data prior to a major public health crisis, such as the COVID-19 pandemic, may exhibit significant performance drops when deployed in post-pandemic settings (Roland et al., 2022). For example, during the COVID-19 pandemic, not only did static features like the proportion of high-risk patients change (Ngiam et al., 2023; Singh et al., 2023), but also temporal relationships between vital signs like heart rate and respiratory rate or trends in blood oxygen levels also changed in complex ways over time. Consequently, understanding the changes beyond just static shifts is crucial not just to maintain model performance, but also to ensure patient safety and maintain trust in AI-assisted decision-making systems.

Addressing this problem of understanding distribution shifts in time series remains *underexplored* and particularly challenging. Existing approaches (Lu et al., 2023; Liu et al., 2024b) draw on methods similar to those used for general static out-of-distribution (OOD) generalization, such as distributionally robust optimization (DRO) (Sagawa et al.; Duchi & Namkoong, 2021; Liu et al., 2022) and causal invariant learning (Peters et al., 2016; Kuang et al., 2018; Arjovsky et al., 2019). These methods, although theoretically compelling, typically overlook the inherent temporal dynamics of time-series data. Moreover, they frequently rely on structural assumptions about distribution shifts without rigorous empirical validation, potentially limiting their practical utility (Gulrajani & Lopez-Paz; Yang et al., 2023; Gagnon-Audet et al., 2023; Liu et al., 2023a).

In response to these limitations, we emphasize the importance of adopting an *inductive* approach—one that is grounded in *understanding time-series-specific patterns of distribution shifts*—to effectively

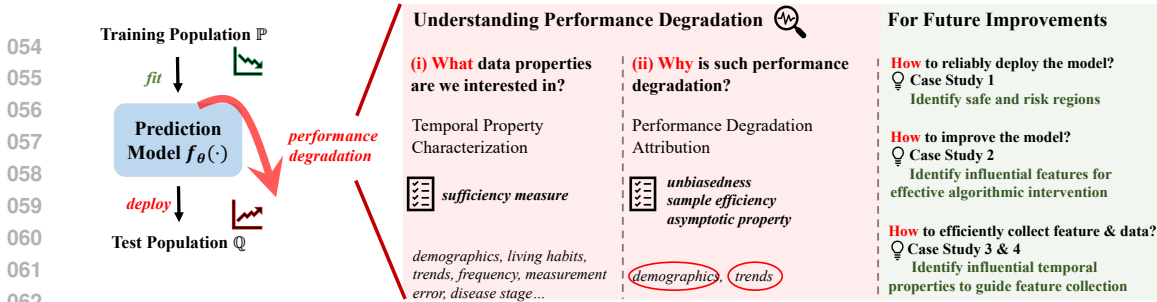


Figure 1: Overview of our TSSA framework.

address real-world distribution shifts. Rather than arbitrarily applying robust methods or fine-tuning models when performance drops on new test data, we argue for first identifying the specific temporal factors driving the decline. This targeted understanding can guide more effective model adjustments or data collection (Liu et al., 2023a) and enable more judicious model deployment. Recent works (Cai et al., 2023; Feng et al., 2024) propose decomposition approaches to understand and attribute performance drop to different factors. However, they focus on “static” settings and cannot deal with the “shifts in temporal dynamics” problem for time-series data, which necessitates modeling the changes in temporal data properties for time series. Hence, the key question of this work is: How to attribute performance degradation between populations to temporal data properties?

To tackle this problem, we propose Time-Series Shift Attribution (TSSA) framework, which analyzes distribution shifts for time-series data from a data-centric perspective. As shown in Figure 1, our framework involves two parts: (i) Define metrics that characterize the behavioral properties of time-series data, along with a *sufficiency measure* that reflects the *optimal predictive power* in estimating the deployed model’s prediction error for each sample using these metrics; (ii) Attribute performance degradation to each data property in a *detailed* manner. We derive an unbiased doubly robust estimator, incorporating a sample-efficient model architecture for attribution.

**Contributions:** TSSA is a framework which to the best of our knowledge is the first to address the issue of time-series distribution shift understanding and attribution bringing several contributions: ① **Technically:** we explicitly model and decompose time-series into interpretable properties and link them to shifts using the TSSA framework, providing a doubly robust estimator for a nuanced attribution and understanding of why and how model performance changes across distributions for time-series predictors. ② **Theoretically:** we characterize the unbiasedness and asymptotic properties for the attribution estimator. ③ **Empirically:** we demonstrate through 4 case studies in real-world health-care applications how our TSSA framework enhances the understanding of time-series shifts and informs both reliable model deployment and targeted improvements in ML models as well as data collections.

## 2 PRELIMINARIES

In the context of time-series classification, let  $X = [U, V_{1..t}]^T \in \mathcal{X}$  represent the input covariates, where  $U \in \mathbb{R}^{d_u}$  denotes static features and  $V_{1..t} \in \mathbb{R}^{d_v \times t}$  denotes multi-variate time-series data up to time step  $t$ . Let  $Y \in \mathcal{Y}$  represent the target prediction outcome. Consider a time-series prediction model  $f$ , which aims to predict the outcome  $Y$  based on the covariates  $X$ . Suppose  $f$  is trained on data pairs  $(X, Y)$  sampled from a training distribution  $\mathbb{P}$ . Let  $\ell(f(x), y)$  be a loss function quantifying the prediction error, such as cross-entropy loss or 0-1 loss.

In this work, we consider a situation where we observe a performance degradation from  $\mathbb{P}$  to a target population  $\mathbb{Q}$ , i.e.  $\mathbb{E}_{\mathbb{Q}}[\ell(f(X), Y)] > \mathbb{E}_{\mathbb{P}}[\ell(f(X), Y)]$ . In time-series data, performance degradation can arise from several factors, including changes in the distribution of input variables such as  $U$  and  $V_{1..t}$ , alterations in the temporal properties of the time series  $V_{1..t}$ , and the presence of missing variables. Understanding and attributing these shifts is crucial for effectively improving model robustness and ensuring accurate predictions in new settings.

Consider a time-series model designed to predict patient mortality risk based on demographic features and various health metrics monitored in a hospital. The model may have been trained on data collected from the patient population prior to a major public health event but deployed during or after the event,

where there could be a significant shift in patient demographics or the emergence of new factors influencing mortality. Additionally, the model may need to be deployed across diverse settings, such as different countries or hospitals. In such cases, we may have sufficient labeled data during the event to evaluate the model, but not enough to train a new one, as investigated in Section 5. Given that collecting new data is extremely costly and could delay patient care, we aim to understand and attribute the performance drop in order to determine in a principled manner the best way to adapt our model or collect the targeted data. Before presenting our methodology, we first discuss the challenges associated with time-series data, and then review the literature on performance drop attribution.

**Challenges in time-series data** Distribution shifts in time-series data are complex, involving both temporal shifts within distributions (e.g., trends, seasonality) and changes between distributions, a challenge we call the “shifts in shifts” problem. Unlike static data, time-series requires addressing these evolving patterns. For instance, pre-pandemic blood oxygen levels in healthy individuals were stable, but during COVID-19, they dropped significantly and continued fluctuating in recovered patients. This highlights the need to consider temporal shift patterns alongside standard shifts.

**Attributing performance drops** Cai et al. (2023) propose decomposing performance drops into two general sources: covariate shifts ( $X$ ) and concept shifts ( $Y|X$ ). Building on this, Feng et al. (2024) employ Shapley values for a more detailed decomposition. However, designed for static data, they cannot be directly applied to time-series data, where, as discussed above, the distribution shifts are much more complicated. Furthermore, since we are only concerned with the contribution of a single data property to the model’s performance degradation while all other properties (or features) remain constant, a comprehensive calculation of Shapley values may be unnecessary and less appropriate. Additionally, there are studies that quantify the importance of shifts in individual features within a (partially) known causal or problem structure (Wu et al., 2021; Thams et al., 2022; Zhang et al., 2023). However, these methods rely heavily on expert knowledge, and any misspecification could introduce significant risks in practice.

To the best of our knowledge, our work is the first to analyze detailed shift patterns in time-series data. While extensive research on time-series forecasting addresses temporal shifts *within a single distribution* (Kim et al., 2021; Fan et al., 2023), our work focuses on the more general problem of *shifts in temporal dynamics*. This work primarily aims to understand performance drops in time-series *classification* models. In the following sections, we will introduce metrics to quantify these temporal shifts (Section 3.1) and develop a framework to attribute performance degradation to specific shift patterns (Section 3.2), to address the aforementioned challenges and limitations. **More related works on distribution shifts, Shapley Value, and time-series anomaly detection can be found in Appendix A.**

### 3 METHOD

In this section, we introduce our **Time-Series Shift Attribution (TSSA)** framework, designed to provide a detailed analysis of model performance degradation in time-series data. The key challenges are two-fold: capturing the temporal properties of time-series data and attributing the overall performance degradation to individual features or properties. Corresponding with these two challenges, as shown in Figure 1, our framework contains two parts: (1) **Temporal Property Characterization**: We define and extract temporal properties, such as global, local and structural properties of the time series, that influence predictive performance. The sufficiency of these metrics is assessed through a novel *Sufficiency Measure*, designed to quantify the extent to which each temporal feature contributes to model performance degradation. (2) **Performance Attribution**: We attribute model performance deterioration to individual temporal properties. This attribution is conducted with a doubly robust estimator, ensuring both *unbiasedness* and *asymptotic consistency* under mild assumptions. We also leverage a shared representation space to efficiently estimate risk models and propensity scores, particularly in scenarios with limited target population data, thereby improving estimation reliability.

#### 3.1 DESIDERATA FOR TIME-SERIES METRICS

To systematically capture the complex nature of distribution shifts in time series data, we propose the following desiderata for metrics to attribute performance shifts, which cover various aspects of time series behavior subjected to distribution shift. We categorize these desiderata into four main groups:

**Global Characteristics** Capture the overall, long-term behavior of the time series. These could encompass: (1) **Overall Statistics**: The overall statistics of a sequence of data, like

the average value and standard deviation. (2) Trends: Metrics to capture long-term directional movements in the data. (3) Frequency: Metrics to identify and quantify cyclical patterns or periodic components. (4) Noise level: Metrics to quantify the signal-to-noise ratio. In this work, we introduce **8** metrics to capture the global characteristics, including the *average, max, min, standard deviation values, standardized trend, smoothed trend, maximum frequency*, and *signal-to-noise ratio*.

**Local Dynamics** Assess short-term behaviors and local patterns within the time series. These could encompass: (1) Variability Assessment: Metrics to measure the degree and nature of fluctuations in the time series. (2) Local Non-stationarity: Metrics to detect short-term shifts. (3) Outlier and Anomaly Detection: Metrics to identify outlier or unusual values. In this work, **6** metrics are included to capture the local dynamics, namely *short-term variability, high-frequency energy, normalized Jitter index, relative strength index, KPSS non-stationary test*, and *the number of breakout points*.

**Structural Changes** Identify significant shifts or changes in the underlying data-generating process. These could encompass (1) Metrics to detect abrupt shifts in the statistical properties of the time series. (2) Metrics to capture the shifts in local trends associated with abrupt shifts. We introduce **2** metrics to capture the structural changes, including *the number of change points*, and *the trend variability*.

**Inter-series Relationships** Address the interactions between multiple time series. These could encompass metrics to capture relationships and dependencies between different time series. We introduce the *covariance variability* among multiple time series to capture this.

These desiderata are motivated by the multifaceted nature of distribution shifts in time series data. Global characteristics can reveal shifts in overall patterns, while local dynamics can capture more subtle changes that might be masked in aggregate measures. Structural changes are crucial for identifying significant alterations in the underlying process, and inter-series relationships are essential for understanding shifts in complex, multivariate time series. **While not exhaustive, as shown in Table 1, these metrics are drawn from various fields, and represent a comprehensive attempt to instantiate each desideratum, and more metrics can be added in future developments.** The detailed definitions can be found in Appendix F.

Table 1: Various metrics to quantify temporal properties from different aspects.

Temporal Property	Name	Metric	Domain
Global Characteristics	Overall Statistics	Average, Standard Deviation, Max, Min values	Statistics
	Standardized Trend	Equation (16)	Statistics
	Smoothed Trend	Savitzky-Golay Filter (Savitzky & Golay, 1964)	Analytical Chemistry
	Maximum Frequency	Dominant frequency by FFT	Signal Processing
	Signal-to-Noise Ratio	Equation (17)	Signal Processing
Local Dynamics	Short-Term Variability	Equation (19)	Signal Processing
	High-Frequency Energy	Equation (20)	Signal Processing
	Normalized Jitter Index	Equation (21)	Signal Processing
	Relative Strength Index	Equation (22)	Finance
	KPSS Non-Stationary Test	$p$ -Value from KPSS Test	Economics
Structural Changes	Breakout Points	Equation (18) (Bollinger Bands (Bollinger, 1992))	Finance
	Change Points	PELT (Killick et al., 2012)	Statistics
Multivariate Interaction	Trend Variability	Standard deviation of local trends	Statistics
	Covariance Variability	Equation (23)	Finance

With these metrics, we combine the static features  $U$  with the metrics of all time-series features  $V_{1\dots t}$ , collectively referred to as  $\tilde{X}$  in the following sections of this paper.

**Sufficiency Measure** Before moving on to the attribution, one natural question is whether the designed metrics are good enough to capture the temporal properties of time series. Since the ultimate goal is to understand the predictive performance drop, the primary requirement for these metrics is that they should relate to predictive performance. Therefore, we propose a *sufficiency measure* to evaluate the *optimal predictive power* of the metrics, defined as:

$$\text{Suff.}(\tilde{X}) := \min_{g \in \mathcal{G}} \mathbb{E}[\text{Loss}(g(\tilde{X}), \ell(f(X), Y))], \quad (1)$$

where  $\text{Loss}(\cdot, \cdot)$  denotes some loss functions (e.g., mean squared error, 0-1 loss) to measure the gap between the *predicted error* and the *real error* of the deployed (and fixed) model  $f$ ,  $\mathcal{G}$  can be chosen as any model classes (e.g., neural networks, XGBoost, etc.). This metric measures the optimal power in predicting the deployed model’s prediction error for each sample using the data property metrics. The smaller it is, the more predictive the metrics become. Therefore, it can serve as a

216 guideline to measure the quality of the metrics. In Proposition 2, we demonstrate how this sufficiency  
 217 measure affects the asymptotic properties of our attribution estimation (introduced in Section 3.2).  
 218 Moreover, in our first case study (Section 5.1), we demonstrate how temporal property metrics can  
 219 help to effectively identify, in a highly interpretable manner, both the “safe region”, where the model  
 220 performs reliably, and the “risk region”, where it is less reliable. In Appendix D, we discuss more  
 221 about the utility of our sufficiency measure.

### 223 3.2 PERFORMANCE DROP ATTRIBUTION

224  
 225 Based on the collection  $\tilde{X}$  of static features  $U$  and temporal data properties in Table 1, we propose to  
 226 attribute the performance drop, i.e.,  $\mathbb{E}_{\mathbb{Q}}[\ell(f(X), Y)] - \mathbb{E}_{\mathbb{P}}[\ell(f(X), Y)]$  to each component in  $\tilde{X}$ . To  
 227 assess how a specific data property contributes to the performance drop, inspired by treatment effect  
 228 estimation, we control for the distribution of all other features except the data property of interest. For  
 229 instance, if we are interested in understanding the effect of blood pressure on the performance decline,  
 230 we could first *control* for the distribution of all other features like demographics, blood oxygen level,  
 231 heart rate, to be the same across both populations. The remaining performance drop, after controlling  
 232 for other factors, could then be attributed to the effect of blood pressure. Additional demonstrations  
 233 on the relationship between our TSSA approach and average treatment effect (ATE) estimation can  
 234 be found at Appendix C.

235 **Objective of Attribution** For a specific data property  $S$  of interest, let  $\tilde{X}_{\setminus\{S\}}$  denote all other  
 236 properties/features in  $\tilde{X}$ , abbreviated as  $\tilde{X}_{-S}$  in the remainder of this work. To “control” for the  
 237 effects of  $\tilde{X}_{-S}$ , we define the conditional risks under distributions  $\mathbb{P}$  and  $\mathbb{Q}$  as follows:

$$238 R_{\mathbb{P}}(\tilde{X}_{-S}) := \mathbb{E}_{\mathbb{P}}[\ell(f(X), Y) \mid \tilde{X}_{-S}], \quad R_{\mathbb{Q}}(\tilde{X}_{-S}) := \mathbb{E}_{\mathbb{Q}}[\ell(f(X), Y) \mid \tilde{X}_{-S}], \quad (2)$$

240 where  $\ell(f(X), Y)$  denotes the deployed model’s *original* prediction error on the sample  $(X, Y)$ .  
 241 Note that when measuring prediction error  $\ell(f(X), Y)$ , we use the *original samples* rather than the  
 242 extracted temporal properties. Building on this, we define the attribution score of feature  $S$  as:

$$243 \text{Attr.}(S) := \mathbb{E}[R_{\mathbb{Q}}(\tilde{X}_{-S}) - R_{\mathbb{P}}(\tilde{X}_{-S})], \quad (3)$$

244 which quantifies the performance gap between  $\mathbb{Q}$  and  $\mathbb{P}$ , while holding the marginal distribution of all  
 245 features except the feature of interest,  $\tilde{X}_{-S}$ , constant. From a distribution shift perspective,  $\text{Attr.}(S)$   
 246 quantifies the performance drop introduced by  $(Y, S) \mid \tilde{X}_{\setminus\{S\}}$ -shifts between  $\mathbb{P}$  and  $\mathbb{Q}$ .

247 Furthermore, in cases where we do not isolate a specific feature (i.e. feature of interest is set to the  
 248 empty set), i.e.,  $\text{Attr.}(\emptyset)$ , the attribution score captures the “systematic” (and unavoidable) difference  
 249 between  $\mathbb{P}$  and  $\mathbb{Q}$ , potentially caused by missing information. For example, in predicting patient risk,  
 250 individuals from two populations may have similar health indices yet experience vastly different  
 251 outcomes due to missing information, such as differing living habits across regions or the absence of  
 252 key health metrics in the records. Note that this term reduces to the concept of “ $Y \mid X$ -shift” introduced  
 253 by Cai et al. (2023). To effectively estimate Equation (3), we propose a doubly robust estimator for  
 254 our objective Equation (3) as it is resilient to misspecification, ensuring reliable attribution even when  
 255 some assumptions might be violated.

256 **Doubly Robust Estimation for Attribution** Consider the original data  $(U^i, V_{1..t}^i, Y^i)_{i=1}^{n_{\mathbb{P}}} \sim \mathbb{P}$  and  
 257  $(U^j, V_{1..t}^j, Y^j)_{j=1}^{n_{\mathbb{Q}}} \sim \mathbb{Q}$ , we first calculate the temporal properties in Table 1 for  $V_{1..t}$ , and convert  
 258 the data into  $(\tilde{X}^i, Y^i)_{i=1}^{n_{\mathbb{P}}} \sim \mathbb{P}$  and  $(\tilde{X}^j, Y^j)_{j=1}^{n_{\mathbb{Q}}} \sim \mathbb{Q}$ . To estimate the attribution, i.e. Equation (3),  
 259 we first learn two predictors  $\hat{\mu}_{\mathbb{P}}(\cdot)$ ,  $\hat{\mu}_{\mathbb{Q}}(\cdot)$  to approximate the conditional risk function  $R_{\mathbb{P}}(\cdot)$ ,  $R_{\mathbb{Q}}(\cdot)$   
 260 in Equation (2) respectively from the observation data. Then we fit a domain classifier  $\hat{\pi}(\cdot)$  as  
 261 (propensity score estimator):

$$262 \hat{\pi}(x_{-S}) \approx \left\{ \pi(x_{-S}) := \Pr(x_{-S} \text{ from } \mathbb{Q} \mid \tilde{X}_{-S} = x_{-S}) \right\}. \quad (4)$$

263 We use non-parametric models for  $R_{\mathbb{P}}(\cdot)$ ,  $R_{\mathbb{Q}}(\cdot)$ ,  $\pi(\cdot)$ , ensuring flexibility in learning the relationships  
 264 between variables without relying on strong parametric assumptions. Throughout the theoretical  
 265 analysis in this paper, we consider generic nonparametric regression estimators, and all strategies  
 266 could be used directly with different ML models, e.g. tree-ensembles. Then we formulate the



270 Augmented IPW (AIPW) (Robins et al., 1994) estimator as:

$$271 \widehat{\text{Attr.}}(S) = \frac{1}{n_P + n_Q} \left( \sum_{i=1}^{n_P + n_Q} \left( \hat{\mu}_Q(\tilde{X}_{-S}^i) - \hat{\mu}_P(\tilde{X}_{-S}^i) \right) + \right. \\ 272 \left. \sum_{j=1}^{n_Q} \frac{R_Q(\tilde{X}_{-S}^j) - \hat{\mu}_Q(\tilde{X}_{-S}^j)}{\pi(\tilde{X}_{-S}^j)} - \sum_{i=1}^{n_P} \frac{R_P(\tilde{X}_{-S}^i) - \hat{\mu}_P(\tilde{X}_{-S}^i)}{1 - \pi(\tilde{X}_{-S}^i)} \right), \quad (5)$$

273 where  $R_Q(\tilde{X}_{-S}^j)$  denotes the ground-truth value on sample  $\tilde{X}_{-S}^j$  that can be calculated from our  
274 observation data (the same for  $R_P(\tilde{X}_{-S}^i)$ ). **We also demonstrate the compatibility of our attribution  
275 with Shapley Value in Appendix B.**

276 This aforementioned formulation ensures double robustness, meaning that our attribution remains  
277 consistent and unbiased as long as either the conditional risk predictors or the propensity score  
278 estimator is correctly specified. With the propensity score estimator  $\hat{\pi}(\cdot)$  in use, we rely on the  
279 overlap assumption between  $\mathbb{P}$  and  $\mathbb{Q}$ , a common requirement in causal inference (Wager, 2020).  
280 There are several approaches to address potential non-overlap in practice (Cai et al., 2023). In  
281 this work, we adopt a simpler strategy to mitigate the effects of non-overlap by excluding samples  
282 with propensity scores that are extremely close to 0 or 1. Furthermore, in Section 4, we prove the  
283 *unconfoundedness, unbiasedness*, and the *asymptotic properties* of our estimator in Equation (5).

284 **Sample-Efficient Estimation of  $\hat{\mu}_P(\cdot)$ ,  $\hat{\mu}_Q(\cdot)$ , and  $\hat{\pi}(\cdot)$  with Neural Networks** In non-parametric  
285 estimation, we generally have sufficient samples for  $\hat{\mu}_P(\cdot)$  (from the training population  $\mathbb{P}$ ). However,  
286 in real-world applications, data from the target pop-  
287 ulation  $\mathbb{Q}$  is often scarce, making it challenging to  
288 estimate  $\hat{\mu}_Q(\cdot)$  and  $\hat{\pi}(\cdot)$  accurately. To mitigate  
289 this, given the functions share the same input  $\tilde{X}_{-S}$ ,  
290 we propose to learn a shared representation space,  
291 thereby sharing information between the two popu-  
292 lations, and improving the sample efficiency of our  
293 estimates. Motivated by Shi et al. (2019), we adopt  
294 the model architecture as shown in Figure 2. Note  
295 that more advanced architectures (Curth & Van der  
296 Schaar, 2021) can be adopted here, and this is not  
297 the focus of this work. Through this model, we  
298 share information between samples from  $\mathbb{P}$  and  $\mathbb{Q}$   
299 to learn the representation space, which helps for  
300 the estimation of  $\hat{\mu}_Q(\cdot)$  and  $\hat{\pi}(\cdot)$ . Specifically, the  
301 loss function is:

$$302 \min_{\phi, \theta_1, \theta_2, \theta_3} \underbrace{\sum_{i=1}^{n_P} (w_{\theta_1}(h_\phi(\tilde{X}_{-S}^i)) - Y^i)^2}_{\text{for } \hat{\mu}_P(\cdot)} + \underbrace{\sum_{j=1}^{n_Q} (w_{\theta_3}(h_\phi(\tilde{X}_{-S}^j)) - Y^j)^2}_{\text{for } \hat{\mu}_Q(\cdot)} + \underbrace{\sum_{i=1}^{n_P + n_Q} \ell_{\text{CE}}(w_{\theta_2}(h_\phi(\tilde{X}_{-S}^i)), Y^i)}_{\text{for } \hat{\pi}(\cdot)},$$

303 where  $\ell_{\text{CE}}(\cdot, \cdot)$  denotes the cross-entropy loss.

## 304 4 THEORETICAL ANALYSIS

305 In this section, we characterize the unbiasedness and asymptotic properties for our estimation.

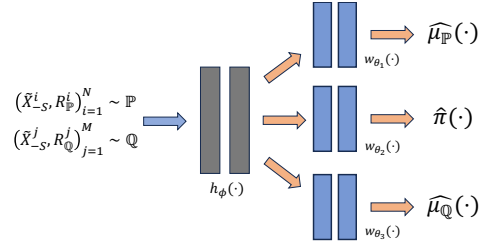
306 **Proposition 1** (Unconfoundedness & Unbiasedness). *Denote  $T \in \{0, 1\}$  as an indicator variable:*

$$307 T = 0, \quad \text{if } \tilde{X}_{-S} \text{ is from } \mathbb{P}; \quad T = 1, \quad \text{if } \tilde{X}_{-S} \text{ is from } \mathbb{Q}, \quad (6)$$

308 which can be likened to a treatment variable. For any attribute  $S \in \tilde{X}$ , we have:

$$309 \left\{ R_P(\tilde{X}_{-S}), R_Q(\tilde{X}_{-S}) \right\} \perp\!\!\!\perp T \mid \tilde{X}_{-S}. \quad (7)$$

310 Based on this, assume that for all  $\tilde{X}_{-S}$ , the overlap assumption holds, i.e.,  $0 < \pi(\tilde{X}_{-S}) < 1$ , then the  
311 estimator in Equation (3) is consistent if either the  $\hat{\mu}_P(\cdot)$ ,  $\hat{\mu}_Q(\cdot)$  are consistent or  $\hat{\pi}(\cdot)$  is consistent.



312 Figure 2: Sample-efficient estimation with  
313 neural networks.  $R_P^i$  denotes the conditional  
314 risk associated with  $i$ -th sample drawn from  $\mathbb{P}$   
315 (similar for  $R_Q^j$ ).  $h_\phi(\cdot)$  represents the shared  
316 representation learner, and  $w_{\theta_1}$ ,  $w_{\theta_2}$ ,  $w_{\theta_3}$  de-  
317 note three separate predictors.

Proof can be found in Appendix H. In addition to the unconfoundedness and unbiasedness, based on Wager (2020), we quantify the consistency between our estimator and the oracle one as follows.

**Proposition 2** (Consistency). *Consider the oracle estimator that uses the true  $\mu_{\mathbb{P}}(\cdot) = R_{\mathbb{P}}(\cdot)$ ,  $\mu_{\mathbb{Q}}(\cdot) = R_{\mathbb{Q}}(\cdot)$ , and  $\pi(\cdot)$  rather than the estimates thereof:*

$$\widehat{Attr}^*(S) = \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P + n_Q} \left( R_{\mathbb{Q}}(\tilde{X}_{-S}^i) - R_{\mathbb{P}}(\tilde{X}_{-S}^i) \right), \quad (8)$$

Assume that **all samples are independent**, and for all  $\tilde{X}_{-S}$ , the overlap assumption holds, then:

$$|\widehat{Attr}(S) - \widehat{Attr}^*(S)| = \mathcal{O}_P \left( \underbrace{\max_{\tau \in \{\mathbb{P}, \mathbb{Q}\}} \mathbb{E} \left[ (\hat{\mu}_{\tau}(\tilde{X}_{-S}) - R_{\tau}(\tilde{X}_{-S}))^2 \right]^{\frac{1}{2}}}_{\text{sufficiency measure in Equation (1)}} \mathbb{E} \left[ (\hat{\pi}(\tilde{X}_{-S}) - \pi(\tilde{X}_{-S}))^2 \right]^{\frac{1}{2}} \right).$$

Proof can be found in Appendix H.

**Remark 1.** Proposition 2 quantifies how closely our estimator approximates the oracle. We can make the following observations: (i) The first term,  $\max_{\tau \in \{\mathbb{P}, \mathbb{Q}\}} \mathbb{E} \left[ (\hat{\mu}_{\tau}(\tilde{X}_{-S}) - R_{\tau}(\tilde{X}_{-S}))^2 \right]$ , characterizes the prediction error of  $\hat{\mu}_{\mathbb{P}}(\cdot)$  and  $\hat{\mu}_{\mathbb{Q}}(\cdot)$ . This term directly corresponds to the sufficiency measure in Equation (1), which represents the objective of our designed temporal data properties, thereby linking the two stages of our framework. (ii) For any  $\mathcal{O}_P(n^{-1/4})$ -consistent machine learning methods to estimate  $\hat{\mu}_{\mathbb{P}}$ ,  $\hat{\mu}_{\mathbb{Q}}$ ,  $\hat{\pi}$ , i.e.

$$\max_{\tau \in \{\mathbb{P}, \mathbb{Q}\}} \mathbb{E} \left[ (\hat{\mu}_{\tau}(\tilde{X}_{-S}) - R_{\tau}(\tilde{X}_{-S}))^2 \right]^{\frac{1}{2}} \mathbb{E} \left[ (\hat{\pi}(\tilde{X}_{-S}) - \pi(\tilde{X}_{-S}))^2 \right]^{\frac{1}{2}} \ll \frac{1}{\sqrt{n}}, \quad \text{then we have} \quad (9)$$

$$\sqrt{n} \left( \widehat{Attr}(S) - \widehat{Attr}^*(S) \right) \rightarrow_p 0. \quad (10)$$

## 5 EXPERIMENTS

We design 4 case studies to comprehensively demonstrate the effectiveness and usage of our TSSA framework, including different real-world shift patterns and prediction tasks. Through our experiments, we use the Medical Information Mart for Intensive Care (MIMIC) (Johnson et al., 2016) dataset. It is representative of *complex real-world* medical time series and contains 23,100 patients from which 9 static demographic features (such as age, gender, admission type etc) and 53 time-series health indexes (such as blood pressure, Braden mobility, temperature etc) have been measured.

### 5.1 CASE STUDY 1: TEMPORAL PROPERTIES GUIDING RELIABLE MODEL DEPLOYMENT

In the first case study, we demonstrate the importance of our time-series data properties (Appendix F) through the interpretable guidance for model safe deployment. Consider a typical intensive-care scenario where the classifier,  $f_{\theta}(\cdot)$ , is trained on historical data but is deployed across different patients and at various stages of their care. Since different patients are likely to exhibit varying feature patterns, it is challenging for a single model to perform consistently well across all incoming patients. Therefore, in high-stakes scenarios like this, it is crucial to *identify in advance* the types of data where the model performs reliably and clinicians can trust its predictions—referred to as the *safe region*. Similarly, it is important to recognize the data patterns where the model performs poorly and should not be relied upon, referred to as the *risk region*.

**Experiment Setup** The task is to predict patient mortality based on 24-hour recordings. We follow the standard design outlined by Jarrett et al., randomly splitting the patients in the MIMIC-III dataset into a training set (18,490 patients,  $\mathbb{P}$ ) and a test set (4,610 patients,  $\mathbb{Q}$ ), ensuring no patient overlap between the two sets. For the validation set, we use the same patients as in the training set but select different time segments for their time-series features, denoted as  $\mathbb{P}_{\text{val}}$ . We train a Transformer model,  $f_{\theta}(\cdot)$ , on the training set, perform region analysis based on the validation data, and use the test set to verify the effectiveness of the identified regions.

**Methodology in Region Analysis** Given the high-stakes nature of the task, these regions must be highly interpretable, rather than relying on opaque, non-interpretable parametric models. Inspired by Lim et al. (2021); Liu et al. (2023a), we fit a decision tree model,  $h(\tilde{X})$  on the validation set ( $\mathbb{P}_{\text{val}}$ ), to predict the prediction error  $\ell(f_{\theta}(X), Y)$  of the trained model  $f_{\theta}(\cdot)$ , using the extracted data

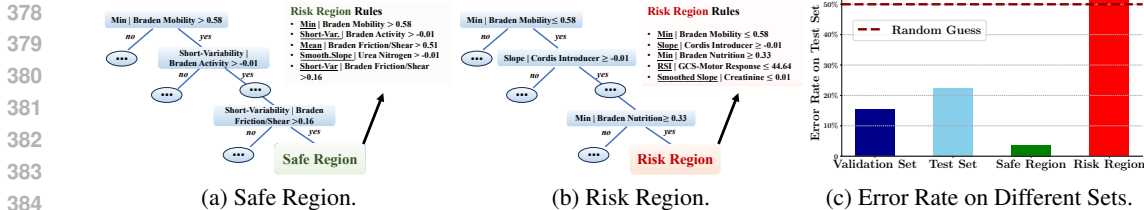


Figure 3: Region Analysis. (a)-(b): Visualizations of the safe and risk regions, defined by interpretable decision rules based on the extracted temporal properties. (c) The error rates on the validation set, the entire test set, the safe region, and the risk region. Error rates in the risk region exceed overall test error, while error in the safe region is much lower than overall test error, thereby guiding safe model deployment by identifying and avoiding usage for patients in the risk region. Note that the regions are learned on validation set  $\mathbb{P}_{\text{val}}$  with no access of the test set in advance.

properties and static demographic features  $\tilde{X}$  (Appendix F). The decision tree partitions the data into distinct regions, each defined by interpretable decision rules and corresponding to a unique leaf node. By analyzing the samples within each leaf, we can identify regions with the highest and lowest risk, as the label assigned by the decision tree represents the original prediction error.

**Analysis** In Figure 3a and Figure 3b, we visualize the safe and risk regions identified based on temporal properties, defined by interpretable and easy-to-understand decision rules. To demonstrate how this guides safe model deployment, we calculate the error rates for the entire test set, the safe region, and the risk region separately. From Figure 3c, while the overall error rate on the test set is relatively reasonable (22.4%), the error rate in the risk region is much higher (52.0%), nearing the level of random guessing. In contrast, the error rate in the safe region is significantly lower (3.7%). Our region analysis offers not only interpretability by design, but also a tool for reliable model deployment. Clinicians can confidently apply the model on patients in the safe region while avoiding use for patients in the risk region, thus ensuring reliable model deployment, which is important from the perspective of trustworthy ML in high stakes settings like healthcare. Further results in Appendix G.3 reinforce the superiority of our temporal-property-based region analysis.

**Take-away 1:** Temporal properties enable accurate and interpretable identification of safe and risk regions, ensuring the reliable and smart deployment of ML models in critical care.

## 5.2 CASE STUDY 2: AGE SHIFTS IN MORTALITY RISK PREDICTION

In real-world healthcare settings, shifts in population distribution are commonly observed. For instance, significant variations in patient age distributions often arise across different hospitals and regions. To rigorously assess the effectiveness of our attribution method, we design scenarios reflecting these variations in patient age on MIMIC-III, and verify that our framework attributes the observed performance degradation to the appropriate features or properties.

**Experiment Setup: Shifts in Patient Age** The task is to predict patient mortality based on 24-hour recordings. We consider a data collection process that oversamples patients under the age of 65, where the average age of the patients in training is 57, while in the test set, the average age is 77. The training set contains 11,476 patients (training distribution  $\mathbb{P}$ ), and the test set contains 6,408 new (but older) patients (target distribution  $\mathbb{Q}$ ). For the validation data, similar with Case Study 1 (Section 5.1), we choose the same patients as in the training set but select different time segments for their time-series features. We train a Transformer model  $f_{\theta}(\cdot)$  on  $\mathbb{P}$  and evaluate it on  $\mathbb{Q}$ , where we observe a performance drop of 13.9pp on accuracy (from 87.7% to 73.8%), and 8.8pp on Macro-F1 score (from 72.5% to 63.7%). In the subsequent analysis, we keep the model  $f_{\theta}(\cdot)$  fixed during evaluation and apply our TSSA framework to attribute its performance drop to various features. **More details can be found in Appendix G.6.**

**Analysis** We begin by presenting the feature attribution results in Figure 4a, which show the average attribution scores along with standard deviation errors from 10 random runs. The findings reveal that the top features identified by our framework are predominantly demographic variables (yellow bars), with the "Age" feature correctly attributed as the most influential to the performance drop. Other key features, such as Admission Location (e.g., emergency room, referral, transfer from other hospitals) and Previous ICU Stay Duration, have a relatively strong correlation with the Age feature. This aligns well with our problem setting, where we perform oversampling based on the "Age" feature. Furthermore, can we use the attribution to guide actions to remedy model performance.



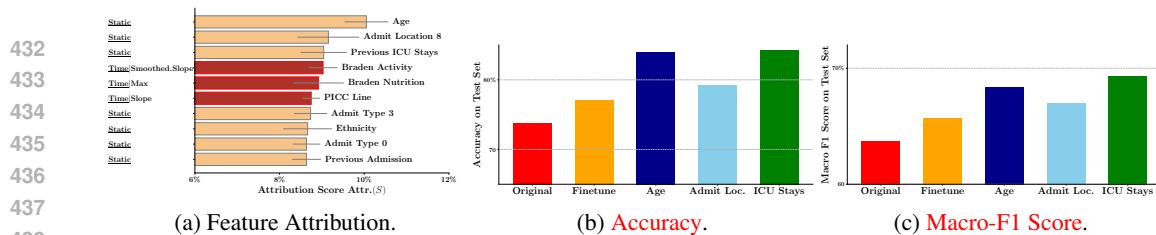


Figure 4: Results for Case Study 2: (a) We visualize the top 10 most important features, which are mainly demographic features (yellow bars), with age being the most prominent. This corresponds with our “prior” knowledge and verify that our framework attributes to the appropriate features. (b)-(c) We do some balancing the top three features respectively during the training phase and plot the accuracy and Macro-F1 score for each model. The results show a significant improvement in performance over original model and simple finetuning when we intervene on these top features.

For each of the top 3 identified features (i.e., Age, Admit Location, and Previous ICU Stays), we apply *simple* balancing by reweighting the data based solely on the inverse density ratio of that specific feature to achieve a uniform distribution. As shown in Figure 4b and Figure 4c, we observe *significant* improvements in both test accuracy and Macro-F1 score. This demonstrates that the top features identified by our framework directly influence the performance drop observed between the training and test patients. Additionally, this highlights how our attribution results can inspire further algorithmic interventions. In this case study, this inductive approach shows that even simple algorithmic modifications can lead to significant improvements.

**Take-away 2:** Our attribution framework accurately attributes the performance drop to relevant demographic features and inspires straightforward yet effective algorithmic interventions, such as data balancing (as done in our intervention), and targeted data augmentation.

### 5.3 CASE STUDY 3: PREEMPTIVE DIAGNOSIS UNDER TEMPORAL SHIFTS

Beyond shifts in patient age, temporal variations in patients’ conditions represent another prevalent type of shift in healthcare data. In preemptive diagnostic scenarios, our target is for the model to detect mortality risks at the earliest possible stage (Filippin et al., 2015). This proactive approach not only enhances patient outcomes by facilitating timely interventions but also underscores the importance of adapting models to account for the dynamic nature of patient health over time. In this case study, we explore how our attribution framework can offer valuable insights for this.

**Experiment Setup: Preemptive Diagnosis under Temporal Shifts** To evaluate reliability in the early detection of mortality risk, we examine the temporal shifts between the training and test patients. Specifically, for the training set  $\mathbb{P}$ , we utilize the last 24-hour time segments for all time-series features, while for the test set  $\mathbb{Q}$ , we select the first 24-hour time segments for all features. This setup allows us to assess whether the model can effectively withstand these temporal shifts and accurately identify patients at high risk of mortality in the early stage. We train a Transformer model  $f_{\theta}(\cdot)$  on  $\mathbb{P}$ , which comprises 12,574 patients, and validate it on an additional 5,547 patients. To control for other shifts, we use the same set of patients for both the validation and test sets  $\mathbb{Q}$ ; the only difference lies in the time segments used: the last 24 hours for validation and the first 24 hours for testing. In this case, we observe a performance drop of 13.8pp in accuracy (from 90.0% to 76.2%) and 24.8pp in Macro-F1 score (from 64.1% to 39.3%), highlighting the necessity to investigate the reasons behind this significant decline. Then for the subsequent analysis, we keep the model  $f_{\theta}(\cdot)$  fixed during evaluation and apply our TSSA framework to attribute its performance drop to various features.

**Analysis** In Figure 5a, we present the average attribution scores along with standard deviation errors from 10 random runs. In contrast to Figure 4a, which displays attribution based on age shifts, the prominent features in this analysis are primarily temporal properties, including the smoothed slope, high-frequency energy, and break points. For instance, for the top feature—the smoothed slope of Braden Activity—we randomly selected ten patients and visualized the last 24-hour (blue) and the first 24-hour (red) time series for Braden Activity in Figure 5b. Here, we observe that the trends, as indicated by the smoothed slope metric (top feature), and break points (sixth feature) of this time series differ significantly. As demonstrate by Valiani et al. (2017), mobility status during hospitalization provides substantial prognostic value in hospitalized older adults. The Braden Activity score could be an efficient and valuable source of information about mobility status for targeting post-hospital care of older adults. Our attribution results provide valuable insights for clinicians.

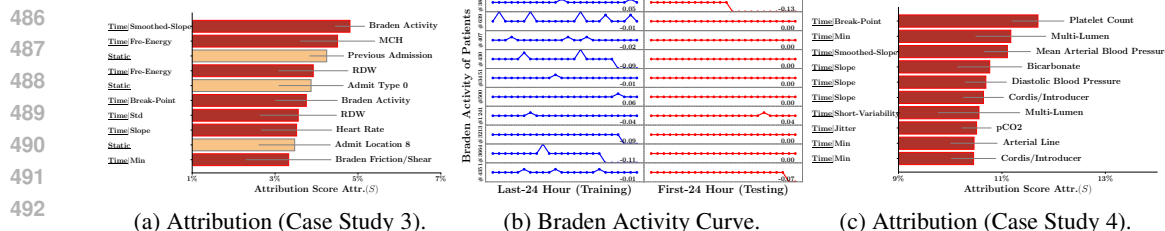


Figure 5: Results of Case Study 3&4: (a) We visualize the top 10 most important features for Case Study 3. We can see that the performance drop is primarily driven by extracted temporal properties. (b) We randomly select 10 patients to compare the curves of the top feature (Braden Activity). A significant difference in both the slope and breakpoints is observed. (c) We visualize the top 10 most important features for Case Study 4, where they are all temporal properties.

First, it is essential to monitor changes in the Braden Activity feature closely. Second, as illustrated on the right side of Figure 5b, many of these changes may go unrecorded or unmeasured during the early stages of patient care. To enable timely detection, clinicians should prioritize the assessment of this feature from the outset. Overall, these findings underscore the potential of our results to enhance clinical practices and improve feature collection.

**Take-away 3:** Our attribution framework provides valuable insights to guide ML engineers and clinicians, like implementing timely and effective feature monitoring and collection.

### 5.4 CASE STUDY 4: VENTILATOR PREDICTION

We focus on predicting the use of mechanical ventilation in intensive care—a procedure that is both invasive and uncomfortable, requiring the induction of an artificial coma and carrying a significant risk of mortality. Accurate predictions are crucial, as errors can lead to serious consequences.

**Experiment Setup** The task is to predict whether a patient requires mechanical ventilation. We consider a realistic scenario in which a model trained on historical data (12,574 patients, first 24-hour time series, denoted as  $\mathbb{P}$ ) must be deployed for new patients and future time segments (an additional 5,547 patients, *second* 24-hour time series, denoted as  $\mathbb{Q}$ ). This scenario accounts for shifts in both demographics (across different patients) and temporal factors (across various time periods, such as the first day versus the second day). We train a Transformer model, and observe a performance drop of 12.4pp in accuracy (from 93.7% to 81.3%) and 10.1pp in Macro-F1 score (from 62.4% to 52.3%), highlighting the necessity to investigate the reasons behind this significant decline.

**Analysis** We present the feature attribution results in Figure 5c. Notably, all top 10 features are temporal properties, indicating that the expected demographic shifts among patients are *unexpectedly minor*. To further explore this phenomenon, we exclude temporal shifts and build a new test set denoted as  $\mathbb{Q}'$  where the time-series features are also derived from the first 24 hours (the same as  $\mathbb{P}$ ). After evaluation, the drop from  $\mathbb{P}$  to  $\mathbb{Q}'$  is only 0.6pp in accuracy and 0.3pp in Macro-F1 score. Thus, this drop can be primarily attributed to demographic shifts among different patients. This further supports our attribution results that indicate minor shifts among patients in this scenario. Furthermore, our identified top features are closely related to the ventilator prediction. For example, Platelet Count (Top-1) is recognized as a significant prognostic marker in intensive care (Mackay et al., 2010; Ilban & Simsek, 2023). Therefore, our attribution results offer valuable insights for ensuring reliable deployment and aiding clinicians in making more informed decisions.

**Take-away 4:** Our attribution framework effectively identifies the primary sources of distributional shifts, providing actionable insights for guiding subsequent algorithmic interventions.

## 6 CONCLUSION

This paper presented the Time-Series Shift Attribution (TSSA) framework, which effectively attributes performance degradation due to distribution shifts in time-series data, with a focus on healthcare applications. Our empirical and theoretical results demonstrate its potential to enhance model reliability and inspire both algorithmic and data-centric interventions in the future.

## REFERENCES

- 540  
541  
542 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.  
543 *arXiv preprint arXiv:1907.02893*, 2019. 1, 15
- 544 Patrick Billingsley. *Probability and measure*. John Wiley & Sons, 2017. 23
- 545  
546 John Bollinger. Using bollinger bands. 1992. 4, 19
- 547  
548 Mohammad Braei and Sebastian Wagner. Anomaly detection in univariate time-series: A survey on  
549 the state-of-the-art. *arXiv preprint arXiv:2004.00433*, 2020. 15
- 550  
551 Tiffany Tianhui Cai, Hongseok Namkoong, and Steve Yeadowsky. Diagnosing model performance  
552 under distribution shift. *arXiv preprint arXiv:2303.02011*, 2023. 2, 3, 5, 6
- 553  
554 Xian Yao Chen, Zhaohua Wu, and Norden E Huang. The time-dependent intrinsic correlation based  
555 on the empirical mode decomposition. *Advances in Adaptive Data Analysis*, 2(02):233–265, 2010.  
20
- 556  
557 Jay N Cohn, Lynn Hoke, Wayne Whitwam, Paul A Sommers, Anne L Taylor, Daniel Duprez,  
558 Renee Roessler, and Natalia Florea. Screening for early detection of cardiovascular disease in  
559 asymptomatic individuals. *American heart journal*, 146(4):679–685, 2003. 1
- 560  
561 Alicia Curth and Mihaela Van der Schaar. Nonparametric estimation of heterogeneous treatment  
562 effects: From theory to learning algorithms. In *International Conference on Artificial Intelligence  
and Statistics*, pp. 1810–1818. PMLR, 2021. 6
- 563  
564 John Duchi, Tatsunori Hashimoto, and Hongseok Namkoong. Distributionally robust losses for latent  
565 covariate mixtures. *Operations Research*, 71(2):649–664, 2023. 15
- 566  
567 John C Duchi and Hongseok Namkoong. Learning models with uniform performance via distribu-  
568 tionally robust optimization. *The Annals of Statistics*, 49(3), 2021. 1, 15
- 569  
570 Wei Fan, Pengyang Wang, Dongkun Wang, Dongjie Wang, Yuanchun Zhou, and Yanjie Fu. Dish-ts:  
571 a general paradigm for alleviating distribution shift in time series forecasting. In *Proceedings of  
the AAAI conference on artificial intelligence*, volume 37, pp. 7522–7529, 2023. 3
- 572  
573 Jean Feng, Harvineet Singh, Fan Xia, Adarsh Subbaswamy, and Alexej Gossmann. A hierarchical  
574 decomposition for explaining ml performance discrepancies. *arXiv preprint arXiv:2402.14254*,  
2024. 2, 3
- 575  
576 Lidiane Isabel Filippin, Vivian Nunes de Oliveira Teixeira, Magali Pilz Monteiro da Silva, Fernanda  
577 Miraglia, and Fabiano Silva da Silva. Sarcopenia: a predictor of mortality and the need for early  
578 diagnosis and intervention. *Aging clinical and experimental research*, 27:249–254, 2015. 9
- 579  
580 Jean-Christophe Gagnon-Audet, Kartik Ahuja, Mohammad Javad Darvishi Bayazi, Pooneh Mousavi,  
581 Guillaume Dumas, and Irina Rish. Woods: Benchmarks for out-of-distribution generalization in  
time series. *Transactions on Machine Learning Research*, 2023. 1
- 582  
583 Josh Gardner, Zoran Popovic, and Ludwig Schmidt. Subgroup robustness grows on trees: An  
584 empirical baseline investigation. *Advances in Neural Information Processing Systems*, 35:9939–  
585 9954, 2022. 22
- 586  
587 Ishaan Gulrajani and David Lopez-Paz. In search of lost domain generalization. In *International  
Conference on Learning Representations*. 1
- 588  
589 Omur Ilban and Fatih Simsek. Prognostic significance of platelet and mean platelet volume kinetics  
590 for ventilator-associated pneumonia. *Journal of Anesthesia/Anestezi Dergisi (JARSS)*, 31(1), 2023.  
591 10
- 592  
593 Daniel Jarrett, Jinsung Yoon, Ioana Bica, Zhaozhi Qian, Ari Ercole, and Mihaela van der Schaar.  
Clairvoyance: A pipeline toolkit for medical time series. In *International Conference on Learning  
Representations*. 7

- 594 Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad  
595 Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a  
596 freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016. 7, 20
- 597  
598 Rebecca Killick, Paul Fearnhead, and Idris A Eckley. Optimal detection of changepoints with a linear  
599 computational cost. *Journal of the American Statistical Association*, 107(500):1590–1598, 2012.  
600 4, 20
- 601 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-  
602 versible instance normalization for accurate time-series forecasting against distribution shift. In  
603 *International Conference on Learning Representations*, 2021. 3
- 604 Masanori Koyama and Shoichiro Yamaguchi. When is invariance useful in an out-of-distribution  
605 generalization problem? *arXiv preprint arXiv:2008.01883*, 2020. 15
- 606  
607 Kun Kuang, Peng Cui, Susan Athey, Ruoxuan Xiong, and Bo Li. Stable prediction across unknown  
608 environments. In *proceedings of the 24th ACM SIGKDD international conference on knowledge  
609 discovery & data mining*, pp. 1617–1626, 2018. 1
- 610 Kim-Hung Le and Paolo Papotti. User-driven error detection for time series with events. In *2020  
611 IEEE 36th International Conference on Data Engineering (ICDE)*, pp. 745–757. IEEE, 2020. 15
- 612  
613 Justin Lim, Christina X Ji, Michael Oberst, Saul Blecker, Leora Horwitz, and David Sontag. Finding  
614 regions of heterogeneity in decision-making via expected conditional covariance. *Advances in  
615 Neural Information Processing Systems*, 34:15328–15343, 2021. 7
- 616 Haoxin Liu, Harshavardhan Kamarthi, Lingkai Kong, Zhiyuan Zhao, Chao Zhang, and  
617 B Aditya Prakash. Time-series forecasting for out-of-distribution generalization using invari-  
618 ant learning. *arXiv e-prints*, pp. arXiv–2406, 2024a. 15
- 619 Haoxin Liu, Harshavardhan Kamarthi, Lingkai Kong, Zhiyuan Zhao, Chao Zhang, B Aditya Prakash,  
620 et al. Time-series forecasting for out-of-distribution generalization using invariant learning. In  
621 *Forty-first International Conference on Machine Learning*, 2024b. 1
- 622  
623 Jiashuo Liu, Zheyuan Hu, Peng Cui, Bo Li, and Zheyuan Shen. Heterogeneous risk minimization. In  
624 *International Conference on Machine Learning*, pp. 6804–6814. PMLR, 2021a. 15
- 625 Jiashuo Liu, Zheyuan Shen, Yue He, Xingxuan Zhang, Renzhe Xu, Han Yu, and Peng Cui. Towards  
626 out-of-distribution generalization: A survey. *arXiv preprint arXiv:2108.13624*, 2021b. 15
- 627  
628 Jiashuo Liu, Jiayun Wu, Bo Li, and Peng Cui. Distributionally robust optimization with data geometry.  
629 *Advances in neural information processing systems*, 35:33689–33701, 2022. 1
- 630 Jiashuo Liu, Tianyu Wang, Peng Cui, and Hongseok Namkoong. On the need for a language  
631 describing distribution shifts: Illustrations on tabular datasets. *Advances in Neural Information  
632 Processing Systems*, 36, 2023a. 1, 2, 7, 15, 18, 22
- 633 Jiayi Liu, Donghua Yang, Kaiqi Zhang, Hong Gao, and Jianzhong Li. Anomaly and change point  
634 detection for time series with concept drift. *World Wide Web*, 26(5):3229–3252, 2023b. 15
- 635  
636 Wang Lu, Jindong Wang, Xinwei Sun, Yiqiang Chen, and Xing Xie. Out-of-distribution representation  
637 learning for time series classification. In *The Eleventh International Conference on Learning  
638 Representations*, 2023. 1
- 639 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon,  
640 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.),  
641 *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.  
642 15, 16
- 643 A Mackay, J Erskine, K Connor, M McCusker, and S Noble. Platelet count as a prognostic marker in  
644 intensive care. *Critical Care*, 14(Suppl 1):P365, 2010. 10
- 645  
646 Duncan McElfresh, Sujay Khandagale, Jonathan Valverde, Vishak Prasad C, Ganesh Ramakrishnan,  
647 Micah Goldblum, and Colin White. When do neural nets outperform boosted trees on tabular data?  
*Advances in Neural Information Processing Systems*, 36, 2024. 22

- 648 Morteza Naghavi, Peter Libby, Erling Falk, S Ward Casscells, Silvio Litovsky, John Rumberger,  
649 Juan Jose Badimon, Christodoulos Stefanadis, Pedro Moreno, Gerard Pasterkamp, et al. From  
650 vulnerable plaque to vulnerable patient: a call for new definitions and risk assessment strategies:  
651 Part i. *Circulation*, 108(14):1664–1672, 2003. 1
- 652  
653 Jinghao Nicholas Ngiam, Srishti Chhabra, Wilson Goh, Meng Ying Sim, Nicholas WS Chew, Ching-  
654 Hui Sia, Gail Brenda Cross, and Paul Anantharajah Tambyah. Continued demographic shifts in  
655 hospitalised patients with covid-19 from migrant workers to a vulnerable and more elderly local  
656 population at risk of severe disease. *International Journal of Infectious Diseases*, 127:77–84, 2023.  
657 1
- 658 Spiros Papadimitriou, Jimeng Sun, and S Yu Philip. Local correlation tracking in time series. In *Sixth*  
659 *International Conference on Data Mining (ICDM’06)*, pp. 456–465. IEEE, 2006. 20
- 660  
661 Jonas Peters, Peter Bühlmann, and Nicolai Meinshausen. Causal inference by using invariant  
662 prediction: identification and confidence intervals. *Journal of the Royal Statistical Society Series*  
663 *B: Statistical Methodology*, 78(5):947–1012, 2016. 1
- 664 James M Robins, Andrea Rotnitzky, and Lue Ping Zhao. Estimation of regression coefficients when  
665 some regressors are not always observed. *Journal of the American statistical Association*, 89(427):  
666 846–866, 1994. 6
- 667  
668 Theresa Roland, Carl Böck, Thomas Tschoellitsch, Alexander Maletzky, Sepp Hochreiter, Jens Meier,  
669 and Günter Klambauer. Domain shifts in machine learning based covid-19 diagnosis from blood  
670 tests. *Journal of Medical Systems*, 46(5):23, 2022. 1
- 671 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust  
672 neural networks. In *International Conference on Learning Representations*. 1, 15
- 673  
674 Abraham Savitzky and Marcel JE Golay. Smoothing and differentiation of data by simplified least  
675 squares procedures. *Analytical chemistry*, 36(8):1627–1639, 1964. 4, 19
- 676  
677 Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment  
678 effects. *Advances in neural information processing systems*, 32, 2019. 6
- 679 Sally J Singh, Molly M Baldwin, Enya Daynes, Rachael A Evans, Neil J Greening, R Gisli Jenkins,  
680 Nazir I Lone, Hamish McAuley, Puja Mehta, Joseph Newman, et al. Respiratory sequelae of  
681 covid-19: pulmonary and extrapulmonary origins, and approaches to clinical care and rehabilitation.  
682 *The Lancet Respiratory Medicine*, 11(8):709–725, 2023. 1
- 683  
684 Joan B Soriano, Jan Zielinski, and David Price. Screening for and early detection of chronic  
685 obstructive pulmonary disease. *The Lancet*, 374(9691):721–732, 2009. 1
- 686  
687 Nikolaj Thams, Michael Oberst, and David Sontag. Evaluating robustness to dataset shift via  
688 parametric robustness sets. In *Advances in Neural Information Processing Systems*, 2022. 3
- 689  
690 Svante Twetman and Margherita Fontana. Patient caries risk assessment. *Detection, assessment,*  
*diagnosis and monitoring of caries*, 21:91–101, 2009. 1
- 691  
692 Vincenzo Valiani, Zhiguo Chen, Gigi Lipori, Marco Pahor, Carlo Sabbá, and Todd M Manini.  
693 Prognostic value of braden activity subscale for mobility status in hospitalized older adults. *Journal*  
694 *of hospital medicine*, 12(6):396–401, 2017. 9
- 695  
696 Stefan Wager. Stats 361: Causal inference. Technical report, Technical report, Technical report,  
Stanford University, 2020. URL: [https . . .](https://www.stat.cmu.edu/~wager/), 2020. 6, 7, 24
- 697  
698 Eric Wu, Kevin Wu, and James Zou. Explaining medical ai performance disparities across sites with  
699 confounder shapley value analysis. *arXiv preprint arXiv:2111.08168*, 2021. 3
- 700  
701 Yuzhe Yang, Haoran Zhang, Dina Katabi, and Marzyeh Ghassemi. Change is hard: A closer look at  
subpopulation shift. In *International Conference on Machine Learning*, pp. 39584–39622. PMLR,  
2023. 1



702 Haoran Zhang, Harvineet Singh, Marzyeh Ghassemi, and Shalmali Joshi. "Why did the model fail?":  
703 Attributing model performance changes to distribution shifts. In Andreas Krause, Emma Brunskill,  
704 Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of*  
705 *the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
706 *Learning Research*, pp. 41550–41578. PMLR, 23–29 Jul 2023. 3  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## A RELATED WORK

Besides the works introduced in Section 2, we will briefly review other methods related to our topic in this section.

**Distribution Shifts** To tackle the challenge of distribution shifts (Liu et al., 2021b), several methods have been proposed, with distributionally robust optimization (DRO) (Sagawa et al.; Duchi & Namkoong, 2021; Duchi et al., 2023; Liu et al., 2023a) and invariant learning (Arjovsky et al., 2019; Koyama & Yamaguchi, 2020; Liu et al., 2021a; 2024a) being among the most prominent approaches. DRO techniques involve perturbing the data distribution within a predefined uncertainty set and optimizing for the worst-case prediction risk. However, these uncertainty sets are often selected based on theoretical considerations rather than empirical evidence, which can lead to overly pessimistic outcomes in practice (Liu et al., 2023a). Invariant learning methods, on the other hand, aim to discover representations that maintain a consistent relationship with the outcome variable across different domains. Nevertheless, these methods depend heavily on the availability of high-quality environments (Liu et al., 2021a), and the assumption of invariance may not hold in real-world settings. For instance, unobserved confounders or missing variables—common in practical scenarios—can prevent the existence of a robust invariant representation. Our work addresses the distribution shift problem through an inductive approach, first aiming to understand the shift patterns. This understanding can then inform targeted interventions from both algorithmic and data-centric perspectives.

**Time-Series Anomaly Detection** In time-series anomaly detection, several works have addressed concept drifts over time, developing both supervised and unsupervised methods (Le & Papotti, 2020; Braei & Wagner, 2020; Liu et al., 2023b). Our approach differs from these methods in several key ways: (i) our goal is to understand why a model’s performance declines in a prediction task by attributing this drop to specific properties of interest, whereas anomaly detection methods focus primarily on reliably detecting breakpoints; (ii) the shifts considered in time-series anomaly detection are typically associated with breakpoints, whereas our work covers a broader range of temporal properties (as shown in Table 1). Additionally, our attribution framework can be applied to diagnose performance drops in anomaly detection methods as well.

Furthermore, in the following Appendix B, we demonstrate in detail the related works with Shapley Value (Lundberg & Lee, 2017), as well as the compatibility of our approach with Shapley Value.

## B SHAPLEY VALUE

One typical attribution method is the SHapley Additive exPlanations (SHAP) (Lundberg & Lee, 2017), which uses cooperative game theory to compute explanations of model predictions. However, our work differs from previous ones in the following aspects:

- **Different goals:** Previous works focus on attributing *model predictions* to input features, whereas our work aims to understand the *performance degradation* of the model when transferring from one distribution to another.
- **Different philosophies:** Shapley values compute the *average* contribution of a feature across all possible subsets of features. In contrast, our approach focuses on the impact of a single data property on the model’s performance degradation, while keeping *all other properties (or features) constant*. As a result, a full calculation of Shapley values may be unnecessary and less suitable for our objectives.
- **Different settings:** Previous works typically address static settings, where models can be refit with any subset of features. Our work, however, deals with time-series data, where each time series involves multiple temporal properties, making it infeasible to “remove” one property and refit the model to estimate SHAP values.

810 Furthermore, our proposed attribution method is compatible with SHAP values, as our attribution  
 811 score can be integrated into SHAP as the “effect” function. In the following, we would like to first  
 812 demonstrate why we currently do not use Shapley Value, and then illustrate how our approach can be  
 813 integrated with Shapley Value.

814 **Why we do not use Shapley Value** One typical attribution method is the SHapley Additive  
 815 exPlanations (SHAP) (Lundberg & Lee, 2017), which uses cooperative game theory to compute  
 816 explanations of model predictions. However, our work differs from previous ones in the following  
 817 aspects:

- 818 • **Different goals:** Previous works focus on attributing *model predictions* to input features,  
 819 whereas our work aims to understand the *performance degradation* of the model when  
 820 transferring from one distribution to another.
- 821 • **Different philosophies:** Shapley values compute the *average* contribution of a feature  
 822 across all possible subsets of features. In contrast, our approach focuses on the impact of  
 823 a single data property on the model’s performance degradation, while keeping *all other*  
 824 *properties (or features) constant*. As a result, a full calculation of Shapley values may be  
 825 unnecessary and less suitable for our objectives.
- 826 • **Different settings:** Previous works typically address static settings, where models can be  
 827 refit with any subset of features. Our work, however, deals with time-series data, where  
 828 each time series involves multiple temporal properties, making it infeasible to “remove” one  
 829 property and refit the model to estimate SHAP values.

830 **Compatibility with Shapley Value** Furthermore, our proposed attribution method is compatible  
 831 with SHAP values, as our attribution score can be integrated into SHAP as the “effect” function.

832 First, recall the definition of Shapley Value: for a player  $i$ , the Shapley value  $\phi_i(v)$  is calculated as:

$$833 \phi_i(v) = \sum_{\mathcal{S} \subseteq \mathcal{N} \setminus \{i\}} \frac{|\mathcal{S}|!(|\mathcal{N}| - |\mathcal{S}| - 1)!}{|\mathcal{N}|!} (v(\mathcal{S} \cup \{i\}) - v(\mathcal{S})) \quad (11)$$

834 Where  $\mathcal{N}$  is the set of all players,  $\mathcal{S}$  is a subset of players that does not include  $i$ ,  $|\mathcal{S}|$  and  $|\mathcal{N}|$  are the  
 835 sizes of sets  $\mathcal{S}$  and  $\mathcal{N}$ , respectively,  $v(\mathcal{S})$  is the value function for coalition  $\mathcal{S}$ , representing the payoff  
 836 that the players in  $\mathcal{S}$  can generate.  $v(\mathcal{S} \cup \{i\}) - v(\mathcal{S})$  is the marginal contribution of player  $i$  when  
 837 joining coalition  $\mathcal{S}$ .

838 Second, we demonstrate how our TSSA approach can be integrated with Shapley Value. Denote the  
 839 attribution results of feature  $S$  given by our original approach as  $\widehat{\text{Attr}}.(S)$ . Then we can define the  
 840 Shapley Value as:

$$841 \widehat{\text{SHAP}}.\widehat{\text{Attr}}.(S) = \sum_{\mathcal{V} \subseteq X_{-S}} \frac{|\mathcal{V}|!(d - |\mathcal{V}| - 1)!}{d!} (\widehat{\text{Attr}}.(\mathcal{V} \cup S) - \widehat{\text{Attr}}.(\mathcal{V})), \quad (12)$$

842 where  $d$  denotes the number of extracted features,  $X_{-S}$  denotes the set of all features except  $S$ .  
 843 Therefore, our TSSA approach is compatible with Shapley Value, which we leave as a promising way  
 844 of future extension of this work.

## 845 C DISCUSSION ON THE RELATIONSHIP WITH AVERAGE TREATMENT EFFECT

850 **Relationship with Average Treatment Effect (ATE)** First, we would like to clarify that our TSSA  
 851 is an attribution approach, which is not designed to estimate ATE or solve causal problems. And we  
 852 only “interpret” our objective function as a special kind of ATE. That is, denote  $T \in \{0, 1\}$  as an  
 853 indicator variable:

$$854 T = 0, \quad \text{if } \tilde{X}_{-S} \text{ is from } \mathbb{P}; \quad T = 1, \quad \text{if } \tilde{X}_{-S} \text{ is from } \mathbb{Q}, \quad (13)$$

which can be likened to a *treatment variable*. And our attribution objective can then be rewritten as:

$$\text{Attr.}(S) = \mathbb{E}[R(\tilde{X}_{-S}, T = 1) - R(\tilde{X}_{-S}, T = 0)], \quad (14)$$

where  $R(\cdot, T = 1)$  denotes  $R_{\mathbb{Q}}(\cdot)$ , and  $R(\cdot, T = 0)$  denotes  $R_{\mathbb{P}}(\cdot)$  that can be viewed as two *outcome function* (like in causal literature). Therefore, when studying the effect of one feature to the performance drop, our attribution approach controls all the other attributes (to be identical between group  $T = 1$  and group  $T = 0$ ), and then calculate its effect. Note that we are *not* studying original ATE estimation problems.

**Theoretical Analysis Beyond ATE** Unlike in causal inference literature, where verifying the unconfoundedness assumption can be challenging, our problem formulation enables us to *prove* the unconfoundedness directly in Proposition 1, highlighting a key advantage of this approach. Based on this, it follows naturally that the double robustness property, as characterized in Proposition 1, ensures the unbiasedness of our estimator.

## D UTILITY OF SUFFICIENCY MEASURE

In this section, we would like to clarify the utility of the sufficiency measure by addressing it from the following two aspects.

**Experiment Perspective** First, for case study 1 (see Section 5.1), we select top- $K$  features with  $K \in \{10, 20, 30, \text{all}\}$ , calculate the sufficiency measure as well as the worst-accuracy (among  $\mu_{\mathbb{P}}(\cdot)$ ,  $\mu_{\mathbb{Q}}(\cdot)$ , and  $\pi(\cdot)$ ) of our attribution model. The results in Table 2 show a clear trend: better sufficiency (*corresponding to lower sufficiency score, which is the MSE*) correspond to improved prediction performance and more precise attribution, thereby highlighting the practical utility of this measure. *Note that lower sufficiency score represents better sufficiency, since the measure is defined as mean-square error.*

**Theory Perspective** Second, as demonstrated in Proposition 2, the *attribution estimation error* is derived as:

$$\mathcal{O}_P \left( \underbrace{\max_{\tau \in \{\mathbb{P}, \mathbb{Q}\}} \mathbb{E} \left[ (\hat{\mu}_{\tau}(\tilde{X}_{-S}) - R_{\tau}(\tilde{X}_{-S}))^2 \right]^{\frac{1}{2}}}_{\text{sufficiency measure in Equation (1), error of } \mu_{\mathbb{P}}, \mu_{\mathbb{Q}}} \underbrace{\mathbb{E} \left[ (\hat{\pi}(\tilde{X}_{-S}) - \pi(\tilde{X}_{-S}))^2 \right]^{\frac{1}{2}}}_{\text{error of } \pi(\cdot)} \right), \quad (15)$$

which is directly controlled by the sufficiency measure. This theoretical result further underscores the importance and relevance of this metric in ensuring robust and reliable attribution.

**Why to report the worst accuracy?** The attribution estimation error in Equation (15) is the *product* of (1) the worst error of  $\mu_{\mathbb{P}}(\cdot)$ ,  $\mu_{\mathbb{Q}}(\cdot)$  and (2) the error of  $\pi(\cdot)$ . Therefore, we care about the performance of all three predictors  $\mu_{\mathbb{P}}(\cdot)$ ,  $\mu_{\mathbb{Q}}(\cdot)$  and  $\pi(\cdot)$ . Thus, we report the worst accuracy among  $\mu_{\mathbb{P}}(\cdot)$ ,  $\mu_{\mathbb{Q}}(\cdot)$  and  $\pi(\cdot)$  to reflect the sufficiency of extracted features.

Table 2: Utility of sufficiency measure, where lower sufficiency score represents better sufficiency, since the measure is defined as mean-square error.

	Top-10 Features	Top-20 Features	Top-30 Features	All Features
Sufficiency↓	0.202	0.182	0.176	0.166
Worst-Acc↑ (among $\mu_{\mathbb{P}}, \mu_{\mathbb{Q}}, \pi$ )	68.9	72.8	73.4	76.0

## E VALUE OF OUR ATTRIBUTION ON INTERVENTIONS

In this section, we would like to demonstrate how our attribution can guide further interventions.

**What can TSSA provide** Our attribution results are designed to provide a comprehensive framework to understand the distribution shifts responsible for performance drops. Specifically, the interpretability of both static features and extracted non-stationary properties allows clinicians to identify the underlying reasons behind model failures.

**How to leverage** Our results provide a basis for both model-centric and data-centric intervention, and can also inform a smart deployment.

- **Model Intervention:** When the model can be finetuned/re-trained, our TSSA can guide target model interventions. For example, our results can be incorporated with distributionally robust optimization to form some directed/targeted uncertainty set, as shown in Liu et al. (2023a). Also, our results can guide a more data-driven algorithmic design for better robustness or fairness among demographic groups (e.g., defined by the identified feature “Age”).
- **Data-Centric Intervention:** Our results can guide efficient data collection (for example, collect more data from the risk region), data balancing (as done in our case study 2 in Section 5.2), and targeted data augmentation (only perturb sensitive features).
- **Smart Deployment:** When the model cannot be changed, our safe/risk regions can guide engineers/clinicians about where to (not to) deploy the model, as demonstrated in our case study 1 in Section 5.1 and Appendix G.3.

**How to choose** In practice, these three aspects should be incorporated together so as to further improve the model. And we acknowledge that further efforts can be made on how to select better algorithms in each kind based on the specific situations one face, which to the best of our knowledge is still an untouched field.

## F DETAILS OF TEMPORAL PROPERTIES

In this section, we first clarify the guidance and transparency in our metric selection, and then introduce the metrics we used in detail.

### F.1 GUIDANCE AND TRANSPARENCY IN METRIC SELECTION

Before introducing metrics used for extracting non-stationarity, we would like to clarify the guidance and transparency in our metric selection.

**Desiderata for metrics** : To provide a principled and transparent basis for metric selection, we introduce desiderata for the types of metrics (Section 3.1). Specifically, our desiderata define the properties that metrics should capture in time-series data, grouped into four categories:

- **Global Characteristics:** Capturing long-term trends, averages, and periodicities.
- **Local Dynamics:** Measuring short-term variability and anomalies.
- **Structural Changes:** Identifying abrupt shifts in data generation processes.
- **Multivariate (Inter-Series) Relationships:** Quantifying dependencies between multiple time series.

By organizing metrics into these categories, we offer a structured framework to guide practitioners in selecting or tailoring metrics for their specific datasets. To systematically assess the quality and relevance of metrics, we employ a sufficiency measure, which quantifies the predictive power of selected metrics in explaining performance degradation. As theoretically demonstrated by Proposition 2, metrics with better sufficiency (lower sufficiency score, experiments see Appendix D) contribute enable more reliable attribution. We believe such a measure provides a data-driven measure to evaluate and refine metric selection in addition.

While we acknowledge that our current set of metrics is not exhaustive, TSSA as a framework permits users to extend it with additional metrics that align with the defined desiderata. By linking sufficiency to predictive performance, this then ensures that newly added metrics have utility for attribution.



## 972 F.2 PROPERTY DEFINITION

973  
974 To address the non-stationarity of time series, inspired from various fields, such as finance, statistics,  
975 and signal processing, we identify numerous data property metrics. Specifically, for a sequence  $V_{1\dots t}$   
976 of length  $t$ , we define metrics corresponding to different temporal properties:

### 977 1. Global Characteristics:

- 978 • Overall Statistics: We calculate the average, standard deviation, maximum, and minimum  
979 values of each time-series feature.
- 980 • Standardized Slope: Widely used in financial and climate analysis, the standardized slope is  
981 defined as:

$$982 \text{Standardized Slope} := \text{Slope}(V_{1\dots t})/\text{Std}(V_{1\dots t}), \quad (16)$$

983 which quantifies the strength of a sequence relative to its variability.

- 984 • Smoothed Trend: Drawing inspiration from analytical chemistry, we use the Savitzky-Golay  
985 filter (Savitzky & Golay, 1964) to characterize the smoothed trend. This filter smooths  
986 the data by fitting a polynomial to each segment of data points, effectively extracting the  
987 underlying trend.
- 988 • Frequency: We calculate the dominant frequency for each time series using the Fast Fourier  
989 Transform (FFT) i.e. we extract the dominant frequency (with the maximal amplitude) from  
990 positive frequencies, where the positive frequency values are real numbers, calculated as  
991  $f_k = \frac{k}{NT_s}$  ( $T_s = 1$  in our experiments).
- 992 • Signal-to-Noise Ratio: We compute the Signal-to-Noise Ratio (SNR) for each feature in a  
993 time-series dataset by using a moving average to estimate the signal and residuals as noise.  
994 The metric is defined as:

$$995 \text{SNR} = \frac{\mathbb{E}[\|\tilde{V}\|^2]}{\mathbb{E}[\|V - \tilde{V}\|^2]}, \quad (17)$$

996 where  $V$  represents the original time-series feature, and  $\tilde{V}$  is the smoothed data (estimated  
997 signal).

### 1000 2. Local Dynamics:

- 1001 • Breakout points: Inspired by the Bollinger Bands (Bollinger, 1992) widely applied in  
1002 financial analysis, we calculate the number of breakout points within the sequence  $V_{1\dots t}$  as:

$$1003 \left| \mathcal{V}(V_{1\dots t}) \right| := \left| \left\{ i : |V_i| \geq |\text{Mean}(V_{1\dots t})| + 2 \cdot \text{Std}(V_{1\dots t}) \text{ for } i = 1, \dots, t \right\} \right| \quad (18)$$

1004  
1005 which identifies the number of points that fall outside the 2-standard-deviation bands to  
1006 capture its local non-stationarity.

- 1007 • Short-term Variability: We first calculate first differences  $\Delta V_i = V_i - V_{i-1}$ , and the  
1008 short-term variability can be defined as:

$$1009 \sigma_{\Delta V} = \text{Std}(\Delta V_2, \dots, \Delta V_t), \quad (19)$$

1010 which is the standard deviation of first differences, and a larger standard deviation indicates  
1011 greater short-term fluctuation.

- 1012 • High-Frequency Energy: To capture the high-frequency components, based on Discrete  
1013 Fourier Transform, we define the high-frequency energy as:

$$1014 E_{\text{high}} := \sum_{k=\lceil \frac{t}{2} \rceil}^{t-1} \sum_{j=1}^t V_j e^{-i2\pi k j/t}, \quad (20)$$

1015 which calculates the squared magnitudes of the upper half of the frequency spectrum.

- 1016 • Normalized Jitter Index: To provide a comprehensive characteristic of variability, we design  
1017 the normalized Jitter index as:

$$1018 \text{Jitter Index} := \frac{\alpha \sigma_{\Delta V} + (1 - \alpha) E_{\text{high}}}{\text{mean}(|V_{1\dots t}|)}, \quad (21)$$

1019 where  $\alpha \in (0, 1)$  is the hyper-parameter to adjust the information from time and frequency  
1020 domains to provide a comprehensive measure of fluctuation in a time series.  
1021  
1022  
1023  
1024  
1025

- **Relative Strength Index:** In order to capture the speed and change of a signal, we use the relative strength index (RSI) defined as:

$$RS := \frac{\sum_{i=2}^t \mathbb{I}(V_i > V_{i-1}) \cdot (V_i - V_{i-1})}{\sum_{i=2}^t \mathbb{I}(V_i < V_{i-1}) \cdot (V_{i-1} - V_i)}, \quad RSI := 100 - \frac{100}{1 + RS}, \quad (22)$$

where RS captures the ratio of average gains to average declines of  $V_{1..t}$ . Thus, the RSI measures the momentum strength of a signal, particularly the relative magnitude of recent gains versus declines.

- **KPSS Non-Stationary Test:** We calculate the  $p$ -value from the KPSS test (Kwiatkowski-Phillips-Schmidt-Shin test), which is used to assess the stationarity of a time series.

### 3. Structural Changes

- **Change points:** From statistics, we utilize Pruned Exact Linear Time (PELT (Killick et al., 2012)) to capture the optimal change point set for sequence  $V_{1..t}$ , which identifies multiple change points in the sequence such that the statistical properties (e.g., mean, variance) remain consistent within each segment.
- **Trend Variability:** We calculate the local trend changes associated with the change points.

### 4. Multivariate Interaction

- **Covariance Variability:** To capture the varying relationships (w.r.t. time) among multiple time-series features, inspired by literatures on local covariance (Papadimitriou et al., 2006; Chen et al., 2010), we design a covariance variability for time-series features  $\mathbf{V}_{1..T} \in \mathbb{R}^{d \times T}$  as follows:

$$\text{Cov. Var}(\mathbf{V}_{1..T}) := \text{Std}(\lambda_{\max}(C(t))), \quad (23)$$

where  $C(t)$  denotes the local covariance matrix at time  $t$ .

We list the metrics in Table 1, and more details as well as some other metrics are provided in the Appendix. With these metrics, we combine the static features  $U$  with the metrics of all time-series features  $V_{1..t}$ , collectively referred to as  $\tilde{X}$  in the following sections of this paper.

## G ADDITIONAL EXPERIMENT RESULTS

In this section, we add more details of our experiments.

### G.1 DATASET DESCRIPTIONS

**Dataset Details** Through our experiments, we use the Medical Information Mart for Intensive Care (MIMIC) (Johnson et al., 2016) dataset, which is representative of *complex real-world* medical time series. The whole dataset contains 23, 100 patients, from which **9** static demographic features—including insurance status, marital status, ethnicity, gender, age, previous admission, previous ICU stay time, admission type, and admission location—and **53** time-series health indexes, including BUN, Braden activity, Braden friction/shear, Braden mobility, Braden moisture, Braden nutrition, GCS (eye opening), GCS (motor response), GCS (verbal response), HCO<sub>3</sub>, MCH, MCHC, MCV, O<sub>2</sub> fraction, O<sub>2</sub> pressure, O<sub>2</sub> saturation, PTT, RDW, anion gap, arterial line, bicarbonate, calcium, chloride, cordis/introducer, creatinine, dialysis catheter, diastolic blood pressure, glucose, heart rate, hematocrit, hemoglobin, magnesium, mean airway pressure, mean arterial blood pressure, multi lumen, norepinephrine, pCO<sub>2</sub>, pH, pO<sub>2</sub>, phenylephrine, phosphate, PICC line, platelet count, potassium, red blood cells, respiratory rate, sodium, systolic blood pressure, temperature, tidal volume, urea nitrogen, ventilator usage, and white blood cell counts. The time-series features are measured every hour, and the average length is 85.4.

In case studies 1 ~ 3, the outcome variable is mortality, and in case study 4, the outcome variable is ventilator usage (where we exclude the mortality feature).

**Advantages of the MIMIC Dataset** The MIMIC dataset offers several advantages, making it especially suitable for our study:

- **Real-World Complexity:** MIMIC represents real-world clinical settings, including data from a diverse set of patients across different demographics and medical conditions. This is crucial for developing models that generalize well to actual hospital environments, where variability is high and patient trajectories are complex.
- **Granularity and Richness of Time-Series Data:** The 53 time-series features in MIMIC cover a wide array of physiological signals, lab measurements, and treatments. This granularity allows for more comprehensive modeling of patient health. In comparison, many other datasets may focus on a narrower subset of features, limiting their ability to capture nuanced patient dynamics.
- **Diverse Outcome Variables:** The dataset supports the study of multiple clinical outcomes, such as mortality and ventilator usage, which are critical in the context of ICU care. This allows for different case study designs and enables us to explore various clinical scenarios and broaden the scope of potential applications.
- **Longitudinal Data with High Temporal Resolution:** MIMIC provides high-resolution, time-stamped data that tracks patients throughout their ICU stays. This temporal depth allows for the detailed study of health trajectories over time, which is essential for building predictive models that can anticipate patient outcomes based on continuous monitoring—a capability that many smaller or less detailed datasets lack.
- **Widely Used and Well-Validated:** MIMIC has been extensively validated in the research community and is a well-established benchmark dataset for various tasks in medical machine learning. Its widespread use ensures the reliability and comparability of our results with those of previous studies. Additionally, its consistent presence in peer-reviewed research provides confidence in its data quality and facilitates benchmarking

**Justification of the Temporal Shift Setting** Then we would like to further clarify the early diagnosis setup in case study 3. This partition reflects a practical clinical scenario where early-stage predictions are critical in healthcare to identify patients at risk. Note, we are not predicting the past but rather assessing how well a model trained on later-stage data can generalize and perform on earlier data (which is temporally shifted). This is important in healthcare since while later-stage data may eventually become available, clinicians often need to make decisions based on the first 24 hours of patient data to prevent adverse outcomes or to take actions. Hence, our setup evaluates the model’s ability to perform early predictions and correctly identify patients at risk, which is essential for preemptive and early care in clinical practice.

## G.2 MODEL TRAINING DETAILS

As for the original model (under evaluation), we use Transformer model (n\_head:4, n\_layer:3, hidden\_dim:32), learning rate is  $1e^{-3}$ , the total epoch number is 200, batch size is 256, and the early stop is used during training (according to last 10 epoch). As for the attribution model: The model architecture is shown in Figure 2, where we use two-layer MLP with hidden size selected from {16,32,64,128} for each part according to the validation results, learning rate  $1e^{-3}$ , and batch size 64.

## G.3 DIFFERENT TEMPORAL PROPERTIES MATTER FOR DIFFERENT TIME SERIES

To illustrate the necessity of incorporating various temporal properties, as discussed in Section 3.1, we compute the feature importance scores within the conditional risk predictor  $\hat{\mu}_{\mathbb{P}}(\cdot)$ . The importance score for each feature is determined using the gradient norm of that feature, given by:

$$\mathbb{E}_{\mathbb{P}}[|\partial \hat{\mu}_{\mathbb{P}}(X)/\partial \tilde{X}_j|],$$

where a higher score indicates that the feature plays a more significant role in predicting the error of the deployed model  $f(\cdot)$ . We visualize the feature importance in Figure 6. The results reveal that different temporal metrics are important for different time-series features, highlighting the intricate

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145

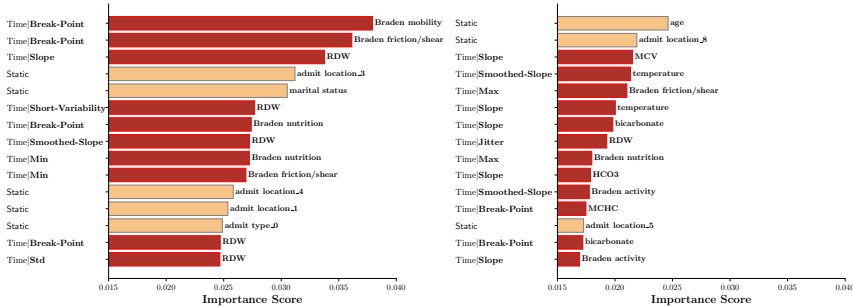


Figure 6: Feature Importance. In case study 1, we visualize the top 30 most important features for  $\hat{\mu}_{\mathbb{P}}(\cdot)$ . The results highlight that different temporal characteristics are significant for different time series, demonstrating the complex nature of time series.

1146  
1147  
1148  
1149  
1150  
1151

nature of time-series data. This not only highlights the complexity of the MIMIC dataset, but also justifies the necessity of the temporal property characterization stage in our framework.

1152  
1153

#### G.4 COMPARISON WITH “STATIC” REGION ANALYSIS

1154  
1155  
1156  
1157  
1158  
1159  
1160

In Case Study 1, a straightforward approach would be to treat all time-series features as static features and then apply the existing static region analysis method (Liu et al., 2023a). This approach offers a natural baseline for comparison, as it simplifies the temporal aspect of the data. However, this simplification may overlook important dynamic patterns inherent in the time-series data.

1161  
1162

Table 3: Comparison of Region Analysis Methods on Test Set.

	Safe Region		Risk Region	
	Error Rate ↓	# Samples ↑	Error Rate ↑	# Samples ↑
Static	8.8%	215	36.7%	30
Ours	<b>3.7%</b>	<b>241</b>	<b>52.0%</b>	<b>296</b>

1163  
1164  
1165  
1166  
1167

In Table 3, we present the error rate and the number of test samples falling within each region for both the static method and our proposed method. A lower error rate indicates a better safe region, while a higher error rate corresponds to a more effective risk region. Regarding sample size, a larger number of test samples in a region suggests that the region is more robust and reliable.

1173  
1174  
1175  
1176

From the results, it is evident that our region analysis method significantly outperforms the static feature-based approach. Additionally, the risk region in our method encompasses a much larger sample size, suggesting that our method, by incorporating temporal properties, captures more reliable and generalizable regions.

1177  
1178

#### G.5 SAMPLE EFFICIENCY OF OUR ARCHITECTURE

1179  
1180  
1181  
1182  
1183  
1184

In case study 2, we examine the performance of our proposed architecture (Figure 2) in fitting the functions  $\hat{\mu}_{\mathbb{P}}(\cdot)$ ,  $\hat{\mu}_{\mathbb{Q}}(\cdot)$ , and  $\hat{\pi}(\cdot)$  across different target sample sizes. Since the outcome variables are binary, we compute the *worst balanced accuracy* of our model, using XGBoost as a baseline for comparison. For XGBoost, we train three independent models, one for each of the three functions. Note that we mainly compare with XGBoost here since it has shown superior prediction power on tabular data, and even outperforms neural networks (Gardner et al., 2022; McElfresh et al., 2024).

1185  
1186  
1187

As shown in Figure 7, when the target sample size exceeds 30% of the training data, our proposed model consistently outperforms XGBoost, highlighting the effectiveness of the shared representation

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

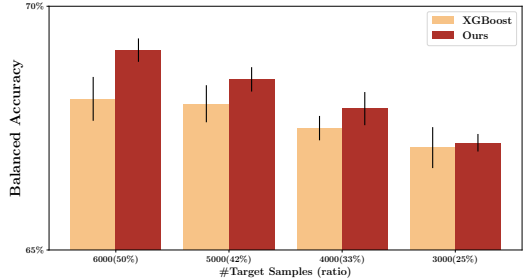


Figure 7: The worst balanced accuracy under different target sample sizes.

Table 4: Intervention results. We compare our attribution-guided simple intervention with simple finetuning.

	Original	Finetune	Age Intervention	Admit Loc. Intervention	ICU Stay Time Intervention
Accuracy	73.8	77.0	<b>83.9</b>	<b>79.2</b>	<b>84.3</b>
Macro-F1	63.7	65.7	<b>68.4</b>	<b>67.0</b>	<b>69.3</b>

space. However, as the target sample size decreases further, our model performs similarly to XGBoost. This outcome is expected, as tree-ensemble models like XGBoost are known to excel on tabular data, even in low-data regimes. Additionally, for our attribution objective, a wide range of models can be employed to estimate  $\hat{\mu}_{\mathbb{P}}(\cdot)$ ,  $\hat{\mu}_{\mathbb{Q}}(\cdot)$ , and  $\hat{\pi}(\cdot)$ , offering flexibility in model selection.

### G.6 COMPARISON WITH SIMPLE FINETUNING

For case study 2, we add the results of finetuning, where we use 50% samples (3200) from the original test set for finetuning, and test on the remaining 50% test samples. For our model intervention, we mix the original training set with the finetuning samples, and then retrain a classifier with balanced weights according to the identified features (age, ICU stay time, etc.). The results are shown in Table 4, where we can see that our intervention can significantly outperform simple finetuning.

## H PROOFS

*Proof for Proposition 1.* As for unconfoundedness, the core is to prove that  $R_{\mathbb{P}}(\tilde{X}_{-S})$  and  $R_{\mathbb{Q}}(\tilde{X}_{-S})$  are both functions relying solely on  $\tilde{X}_{-S}$ . Since  $\tilde{X}_{-S}$  is derived from the original input  $X$  via some kind of transformations, we define  $g(X) := \tilde{X}_{-S}$ . We first prove that  $R_{\mathbb{P}}(\tilde{X}_{-S}) = R_{\mathbb{P}}(g(X)) = \mathbb{E}_{\mathbb{P}}[\ell(f_{\theta}(X), Y)|g(X)]$  is a function of  $g(X)$ . From measure theory, the conditional expectation  $\mathbb{E}_{\mathbb{P}}[\ell(f_{\theta}(X), Y)|g(X)]$  is defined with respect to the  $\sigma$ -algebra  $\sigma(g(X))$  and satisfies:

$$\mathbb{E}_{\mathbb{P}}[\ell(f_{\theta}(X), Y)|g(X)] = h(g(X)), \tag{24}$$

where  $h(g(X))$  is a measurable function of  $g(X)$ . This implies that for any event  $A \in \sigma(g(X))$ , we have:

$$\mathbb{E}[\mathbb{E}_{\mathbb{P}}[\ell(f_{\theta}(X), Y)|g(X)] \cdot \mathbb{I}_A] = \mathbb{E}[\ell(f_{\theta}(X), Y) \cdot \mathbb{I}_A], \tag{25}$$

which implies that the conditional expectation  $\mathbb{E}_{\mathbb{P}}[\ell(f_{\theta}(X), Y)|g(X)]$  depends only on the value of  $g(X)$ . Note that here we use the Tower property (Billingsley, 2017), which states that  $\mathbb{E}[\mathbb{E}[X|\mathcal{G}]] = \mathbb{E}[X]$ , where  $\mathcal{G}$  is a sub- $\sigma$ -algebra. This means that when taking the conditional expectation of a random variable  $X$ , and then the result is equal to the total expectation of  $X$ . For the above equation, the inner conditional expectation  $\mathbb{E}_{\mathbb{P}}[\ell(X, Y)|g(X)]$  can be viewed as a function of  $g(X)$ , denoted as  $h(g(X))$ . And the outer expectation considers the expectation of  $h(g(X))$  over some event  $A \in \sigma(g(X))$ . Thus, by tower property, the result of LHS equals to  $\mathbb{E}[\ell(X, Y) \cdot \mathbb{I}_A]$  because  $\mathbb{I}_A$  is measurable with respect to  $\sigma(g(X))$ , the  $\sigma$ -algebra generated by  $g(X)$ . And this ensures that  $\mathbb{E}_P[\ell(X, Y)|g(X)]$  is well-defined and depends only on  $g(X)$ .



Based on this, we can re-write  $R_{\mathbb{P}}(\tilde{X}_{-S})$  and  $R_{\mathbb{Q}}(\tilde{X}_{-S})$  as  $h_1(g(X))$  and  $h_2(g(X))$ . Then we have:

$$\Pr(T, R_{\mathbb{P}}(\tilde{X}_{-S}), R_{\mathbb{Q}}(\tilde{X}_{-S}) | \tilde{X}_{-S}) \quad (26)$$

$$= \Pr(T, h_1(g(X)), h_2(g(X)) | g(X)) \quad (27)$$

$$= \Pr(T | h_1(g(X)), h_2(g(X)), g(X)) \cdot \Pr(h_1(g(X)), h_2(g(X)) | g(X)) \quad (28)$$

$$= \Pr(T | g(X)) \cdot \Pr(h_1(g(X)), h_2(g(X)) | g(X)), \quad (29)$$

which proves the conditional independence. Note that the last equation holds because  $h_1(g(X))$  and  $h_2(g(X))$  are both functions of  $g(X)$ .

Then we proceed to proving the unbiasedness (also known as the double robustness). Note that the SUTVA assumption holds naturally in our problem setting. First, if  $\hat{\mu}_{\mathbb{P}}(\cdot)$  and  $\hat{\mu}_{\mathbb{Q}}(\cdot)$  are consistent, i.e.  $\hat{\mu}_{\mathbb{P}}(\cdot) = R_{\mathbb{P}}(\cdot)$  and  $\hat{\mu}_{\mathbb{Q}}(\cdot) = R_{\mathbb{Q}}(\cdot)$ , it is easy to show the consistency of our estimator (by plugging-in  $\hat{\mu}_{\mathbb{P}}$  and  $\hat{\mu}_{\mathbb{Q}}$ ). Second, if  $\hat{\pi}(\cdot)$  is consistent, i.e.  $\hat{\pi}(\cdot) \approx \pi(\cdot)$ , denote the indicator

$$T = \begin{cases} 1, & \text{if } \tilde{X}_{-S} \text{ is from } \mathbb{Q} \\ 0, & \text{if } \tilde{X}_{-S} \text{ is from } \mathbb{P}, \end{cases} \quad (30)$$

our estimator can be re-written as:

$$\begin{aligned} \widehat{\text{Attr.}}(S) &= \underbrace{\frac{1}{n_P + n_Q} \sum_{i=1}^{n_Q + n_P} \left( \frac{T_i R_{\mathbb{Q}}(\tilde{X}_{-S}^i)}{\pi(\tilde{X}_{-S}^i)} - \frac{(1 - T_i) R_{\mathbb{P}}(\tilde{X}_{-S}^i)}{\pi(\tilde{X}_{-S}^i)} \right)}_{\text{optimal IPW estimator}} \\ &+ \underbrace{\frac{1}{n_P + n_Q} \sum_{i=1}^{n_Q + n_P} \left( \hat{\mu}_{\mathbb{Q}}(\tilde{X}_{-S}^i) \left(1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)}\right) - \hat{\mu}_{\mathbb{P}}(\tilde{X}_{-S}^i) \left(1 - \frac{1 - T_i}{1 - \pi(\tilde{X}_{-S}^i)}\right) \right)}_{\text{additional term}}. \end{aligned} \quad (31)$$

To complete the proof, we need to show (i) the optimal IPW estimator is consistent, and (ii) the additional term is equal to 0. The proof of consistency of the optimal IPW estimator is standard and one can refer to [Wager \(2020, Chapter 2\)](#). For the additional term, we have:

$$\mathbb{E}\left[1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)} | \tilde{X}_{-S}^i\right] = 0, \quad (32)$$

and therefore complete the proof of unbiasedness.  $\square$

*Proof of Proposition 2.* Our proof builds on the established techniques presented in ([Wager, 2020, Chapter 3](#)), with tailored adaptations and simplifications specific to our problem setting.

First, since in our problem setting,  $\mu_{\mathbb{P}}(\cdot) = R_{\mathbb{P}}(\cdot)$  and  $\mu_{\mathbb{Q}}(\cdot) = R_{\mathbb{Q}}(\cdot)$ , the oracle estimator is simplified to:

$$\widehat{\text{Attr.}}^*(S) = \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P + n_Q} \left( R_{\mathbb{Q}}(\tilde{X}_{-S}^i) - R_{\mathbb{P}}(\tilde{X}_{-S}^i) \right), \quad (33)$$

where we do not have the propensity score term. Then we decompose  $\widehat{\text{Attr.}}(S) - \widehat{\text{Attr.}}^*(S)$  as:

$$\widehat{\text{Attr.}}(S) - \widehat{\text{Attr.}}^*(S) = \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P + n_Q} \left( \hat{\mu}_{\mathbb{Q}}(\tilde{X}_{-S}^i) - \hat{\mu}_{\mathbb{P}}(\tilde{X}_{-S}^i) - R_{\mathbb{Q}}(\tilde{X}_{-S}^i) + R_{\mathbb{P}}(\tilde{X}_{-S}^i) \right) \quad (34)$$

$$+ \frac{R_{\mathbb{Q}}(\tilde{X}_{-S}^i) - \hat{\mu}_{\mathbb{Q}}(\tilde{X}_{-S}^i)}{\hat{\pi}(\tilde{X}_{-S}^i)} T_i - \frac{R_{\mathbb{P}}(\tilde{X}_{-S}^i) - \hat{\mu}_{\mathbb{P}}(\tilde{X}_{-S}^i)}{1 - \hat{\pi}(\tilde{X}_{-S}^i)} (1 - T_i) \quad (35)$$

$$= \Delta_{\mu_{\mathbb{Q}}} - \Delta_{\mu_{\mathbb{P}}}, \quad (36)$$

where we define

$$\Delta_{\mu_{\mathbb{Q}}} := \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P + n_Q} \left( \hat{\mu}_{\mathbb{Q}}(\tilde{X}_{-S}^i) - R_{\mathbb{Q}}(\tilde{X}_{-S}^i) + \frac{R_{\mathbb{Q}}(\tilde{X}_{-S}^i) - \hat{\mu}_{\mathbb{Q}}(\tilde{X}_{-S}^i)}{\hat{\pi}(\tilde{X}_{-S}^i)} T_i \right), \quad (37)$$

and define  $\Delta_{\mu_P}$  analogously. In order to prove Proposition 2, it suffices to show that  $\Delta_{\mu_Q}$  satisfies that conclusion. To prove this, we decompose  $\Delta_{\mu_Q}$  as follows:

$$\Delta_{\mu_Q} = \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P+n_Q} (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i)) \left(1 - \frac{T_i}{\hat{\pi}(\tilde{X}_{-S}^i)}\right) \quad (38)$$

$$= \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P+n_Q} (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i)) \left(1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)}\right) \quad (39)$$

$$- \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P+n_Q} T_i (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i)) \left(\frac{1}{\hat{\pi}(\tilde{X}_{-S}^i)} - \frac{1}{\pi(\tilde{X}_{-S}^i)}\right). \quad (40)$$

We first deal with the first term. Note that in practice, we typically use cross-fitting and therefore  $\hat{\mu}_Q(\cdot)$  can be viewed as deterministic in the following. From Equation (32), the summands used to build the first term become mean-zero. Therefore, we have:

$$\mathbb{E} \left[ \left( \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P+n_Q} (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i)) \left(1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)}\right) \right)^2 \right] \quad (41)$$

$$= \text{Var} \left( \frac{1}{n_P + n_Q} \sum_{i=1}^{n_P+n_Q} (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i)) \left(1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)}\right) \right) \quad (42)$$

$$= \frac{1}{n_P + n_Q} \text{Var} \left( (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i)) \left(1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)}\right) \right) \quad (\text{independent terms}) \quad (43)$$

$$= \frac{1}{n_P + n_Q} \mathbb{E} \left[ (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i))^2 \left(1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)}\right)^2 \right] \quad (44)$$

$$= \frac{1}{n_P + n_Q} \mathbb{E} \left[ (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i))^2 \left(\frac{1}{\pi(\tilde{X}_{-S}^i)} - 1\right)^2 \right], \quad (45)$$

where the last equality is because of:

$$\mathbb{E} \left[ \left(1 - \frac{T_i}{\pi(\tilde{X}_{-S}^i)}\right)^2 \middle| \tilde{X}_{-S}^i \right] = \mathbb{E} \left[ 1 - \frac{2T_i}{\pi(\tilde{X}_{-S}^i)} + \frac{T_i}{\pi^2(\tilde{X}_{-S}^i)} \middle| \tilde{X}_{-S}^i \right] = \left(\frac{1}{\pi(\tilde{X}_{-S}^i)} - 1\right)^2. \quad (46)$$

Then from the overlap assumption, we assume that for all  $\tilde{X}_{-S}^i$ ,  $\eta < \pi(\tilde{X}_{-S}^i) < 1 - \eta$ , which gives that

$$\frac{1}{n_P + n_Q} \mathbb{E} \left[ (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i))^2 \left(\frac{1}{\pi(\tilde{X}_{-S}^i)} - 1\right)^2 \right] \quad (47)$$

$$\leq \frac{1}{\eta(n_P + n_Q)} \mathbb{E} [(\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i))^2]. \quad (48)$$

Then for the second term, we have:

$$\frac{1}{n_P + n_Q} \sum_{i=1}^{n_P+n_Q} T_i (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i)) \left(\frac{1}{\hat{\pi}(\tilde{X}_{-S}^i)} - \frac{1}{\pi(\tilde{X}_{-S}^i)}\right) \quad (49)$$

$$\leq \sqrt{\frac{1}{n_P + n_Q} \sum_{i:T_i=1} (\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i))^2} \cdot \sqrt{\frac{1}{n_P + n_Q} \sum_{i:T_i=1} \left(\frac{1}{\hat{\pi}(\tilde{X}_{-S}^i)} - \frac{1}{\pi(\tilde{X}_{-S}^i)}\right)^2} \quad (50)$$

$$= \mathcal{O}_P \left( \mathbb{E} [(\hat{\mu}_Q(\tilde{X}_{-S}^i) - R_Q(\tilde{X}_{-S}^i))^2]^{\frac{1}{2}} \mathbb{E} [(\hat{\pi}(\tilde{X}_{-S}^i) - \pi(\tilde{X}_{-S}^i))^2]^{\frac{1}{2}} \right). \quad (51)$$

□