

A PITFALL IN CONFORMAL PREDICTION: WHEN SHORTER INTERVALS ARE NOT BETTER

Anonymous authors

Paper under double-blind review

ABSTRACT

Conformal prediction has become a cornerstone of distribution-free uncertainty quantification, conventionally evaluated by its coverage and interval length. This work critically examines the sufficiency of these standard metrics. We demonstrate that *the interval length might be deceptively improved through a counter-intuitive approach* termed Prejudicial Trick (PT), while the coverage remains valid. Specifically, for any given test sample, PT probabilistically returns an interval, which is either null or constructed using an adjusted confidence level, thereby preserving marginal coverage. While PT potentially yields a deceptively lower interval length, it introduces practical vulnerabilities: the same input can yield completely different prediction intervals across repeated runs of the algorithm. We formally derive the conditions under which PT achieves these misleading improvements and provide extensive empirical evidence across various regression and classification tasks. Furthermore, we introduce a new metric *interval stability* which helps detect whether a new conformal prediction method implicitly improves the length based on such PT-like techniques.

1 INTRODUCTION

Machine learning is rapidly evolving and has been successfully applied in numerous fields (Voulodimos et al., 2018; Brown et al., 2020). However, machine learning models, particularly deep learning models, often suffer from overconfidence issues (Guo et al., 2017; Minderer et al., 2021), making them unreliable for deployment in high-stakes areas such as medicine and finance (De Prado, 2018). Therefore, it is crucial to develop techniques for uncertainty quantification and calibrate the original model to enhance the reliability of predictions (Sullivan, 2015; Minderer et al., 2021; Smith, 2024).

Among all the uncertainty quantification methods, conformal prediction stands out due to its simplicity and distribution-free characteristics (Vovk et al., 2005; Shafer & Vovk, 2008; Angelopoulos & Bates, 2021). Conformal prediction is a post hoc approach for constructing prediction intervals, based on a non-conformity score calculated on a hold-out calibration set (Algorithm 2). Conformal prediction and its variants have demonstrated promising performances in numerous applications (Lei & Candès, 2021; Angelopoulos et al., 2022).

Generally, researchers evaluate the intervals returned by conformal prediction via two criteria: *coverage* and *interval length*. Firstly, a valid coverage ensures that the actual response value has a high probability of falling within the interval. Secondly, the interval is encouraged to be as short as possible, as a shorter interval provides more precise information about prediction uncertainties. These two evaluation metrics are commonly used in the literature (Tibshirani et al., 2019; Teng et al., 2022; Angelopoulos et al., 2023; He & Lam, 2024) and a branch of works improves the length with several different meaningful approaches (Romano et al., 2019; Izbicki et al., 2020; Teng et al., 2022; Guan, 2023; Stutz et al., 2022). This raises a crucial question about the potential pitfalls of focusing narrowly on these standard metrics: If a new algorithm outperforms existing ones on coverage and length, should it automatically be considered superior for practical deployment? Specifically,

The key question:

Can a conformal prediction method maintain valid coverage and *deceptively improve interval length metrics* through counter-intuitive constructions, while *introducing practical risks*?

Consider the following Example 1:

Algorithm 1 Prejudicial Trick (PT)

```

1: Input: conformal prediction algorithm (base)  $\mathcal{A}_{1-\alpha}(\cdot; \hat{\mu})$ , test point  $\mathbf{x}'$ , probability  $p$ .
2: Generate a uniform random variable  $U \sim \text{Unif}([0, 1])$ ;
3: if  $U > p$ : then
4:   Interval  $\mathcal{C}_{1-\alpha}(\mathbf{x}') = \hat{\mu}(\mathbf{x}')$  (regression tasks) or  $\mathcal{C}_{1-\alpha}(\mathbf{x}') = \emptyset$  (classification tasks);
5: else
6:   Calculate the adjusted miscoverage rate  $\alpha' = 1 - \frac{1-\alpha}{p}$ ;
7:   Interval  $\mathcal{C}_{1-\alpha}(\mathbf{x}') = \mathcal{A}_{1-\alpha'}(\mathbf{x}'; \hat{\mu})$ ;
8: end if
9: Output: Interval  $\mathcal{C}_{1-\alpha}(\mathbf{x}')$ .

```

Example 1 (The Pitfalls of Length.). Two doctors, Alice and Bob, are estimating recovery time for patients after treatment. Conformal prediction with historical data reveals that 60% of patients recover within 4 years, and 80% within 5 years. When a new patient asks for an estimated recovery time, Alice and Bob adopt distinct strategies:

- Alice: Assign recovery time interval $[0, 4]$ years consistently;
- Bob: Assign recovery time interval $[0, 5]$ with probability 0.75, while $[0, 0]$ with probability 0.25.

For both strategies in Example 1, 60% of patients fall in the estimated interval in expectation, thus satisfying the criteria for valid marginal coverage. Besides, Bob’s approach yields a shorter average interval length $5 \times 75\% = 3.75 \ll 4$. Overall, Bob achieves a shorter interval while achieving the same coverage as Alice. However, Bob’s strategy is flawed in its practical application, since (a) from the micro-level, Bob provides different intervals for the same patient if queried multiple times, and (b) from the macro-level, Bob randomly informs 25% of patients that they will recover immediately after treatment regardless of their actual condition. The example is illustrated in Figure 3.

In this paper, inspired by the motivating example (Example 1), the *Prejudicial Trick* emerges as a practically invalid method that artificially shortens prediction intervals in conformal prediction (see Algorithm 1 and Figure 4). Instead of providing consistent intervals, PT assigns null intervals with a fixed probability to any test sample, and assigns confidence intervals with lower miscoverage rates in other cases to maintain the marginal coverage. While PT preserves marginal coverage and potentially reduces the average interval length, its rationale is less sound compared to standard methods like Vanilla Conformal Prediction (VCP). Specifically, PT suffers from two limitations:

- *Instability issues:* Repeated runs of PT produce different intervals for the same input;
- *Unfairness issues:* PT provides informative predictions for only a subset of test samples, while assigning uninformative null intervals to the rest¹.

From the theoretical perspective, we offer several theoretical results to provide deep understandings of PT regarding both coverage and length, and informally summarize them in Theorem 2.

Theorem 2 (Theorem Summary). *We term base as the base conformal prediction algorithm, term PT as the base algorithm with PT, and omit mild assumptions for clarity. For coverage, it holds that:*

- *PT satisfies marginal coverage guarantees under exchangeability assumptions (Theorem 4);*
- *PT guarantees conditional coverage if its base guarantees conditional coverage (Theorem 5);*
- *PT provably outperforms its base regarding the conditional coverage, under some conditions, even if the base does not satisfy the conditional coverage guarantees (Theorem 6).*

For length, it holds that:

- *PT achieves shorter average intervals than its base under some general conditions (Theorem 7);*
- *We provide sufficient conditions under which PT reduces the average interval length for both differentiable (Theorem 8) and non-differentiable (Theorem 11) length functions. Notably, these conditions are often satisfied in the common scenario of model misspecification (Remark 4).*
- *These results lead to corollaries for specific cases, such as when the length function is locally concave (Corollary 9) or when the base algorithm is VCP (Corollary 10).*
- *We also provide a failure case where PT cannot decrease the average length (Example 12).*

¹In this paper, unfairness stems from the fact that a portion of samples are assigned null intervals in a run, even though each has an equal probability of being prejudiced. This differs from the unfairness concept grounded in conditional coverage (Zhao et al., 2020), where individuals are prejudiced based on their features.

From the experimental perspective, we verify our findings on various real-world datasets, regarding marginal and conditional coverage (Figure 1), and interval length (Table 2). Besides, we validate our findings under different settings, including different tasks (classification regimes in Table 4) and other conformal prediction algorithms (Conformalized Quantile Regression in Table 5).

However, the improvement on length is vacuous, as discussed in Section 3.5. To detect PT, We further introduce *Interval Stability*, a new metric to quantify the variation of the prediction for the same input over multiple runs. This metric is practically meaningful since it serves to identify and alert methods that implicitly or explicitly deploy invalid techniques like PT in the future (Remark 6).

Remark 1. Notably, we present PT not as a practical solution, despite its theoretical advantages on the coverage-length metric. Instead, it acts as a cautionary example that raises issues like instability and unfairness during deployment. PT represents the most direct way to *hack* the coverage-length metric. The fact that such a simple trick can succeed exposes a fundamental blind spot in the current evaluation paradigm. If standard metrics can be fooled this easily, they are likely vulnerable to more complex and subtle manipulations from increasingly sophisticated models.

2 RELATED WORK

Conformal prediction (Vovk et al., 2005) is mainly evaluated by the coverage-length metric. For coverage, conformal prediction provides finite sample guarantees under exchangeability assumptions (Vovk et al., 2005; Tibshirani et al., 2019; Barber et al., 2023), ensuring that prediction sets achieve the expected marginal coverage. Another related metric is conditional coverage, which is unachievable in finite sample settings without further assumptions (Vovk, 2012). Therefore, recent work focuses on various relaxations of conditional coverage (Barber et al., 2020; Gibbs et al., 2024).

Another metric is the average *interval length*, as shorter intervals are generally more informative (Lei et al., 2018; Sadinle et al., 2018). Numerous methods aim to construct adaptive intervals that reduce length while maintaining valid coverage. One line of work designs alternative non-conformity score functions: for example, Romano et al. (2019) integrate quantile regression, while Guan (2023) propose a localized method that adapts to test-time information. Other score functions have been proposed in (Feldman et al., 2021; Alaa et al., 2023; Han et al., 2022; Teng et al., 2022). Of particular interest is the method of Izbicki et al. (2020), which estimates the asymptotic conditional distribution of the non-conformity scores and constructs prediction intervals based on high-density regions. Another line of work involves a training procedure for interval-length optimization, such as CPL (Kiyani et al., 2024), CP-Gen (Bai et al., 2022), ConfTr (Stutz et al., 2022), and BoostedCP (Xie et al., 2024). Different from existing methods that provide principled and meaningful advances in conformal prediction, our proposed prejudicial trick achieves efficiency gains through an illusory construction, making the length improvement non-substantive.

Besides coverage and length, several auxiliary metrics have been introduced. These include *excess* and *deficit* (Seedat et al., 2023), which measures the extent to which the prediction intervals are unnecessarily wide or insufficiently narrow; *false positive rate* (Fisch et al., 2022), which improves precision by limiting the number of incorrect labels in classification settings; and *conditional weighted coverage* (Jensen et al., 2024)—a hybrid metric that takes both coverage and length into regard. Among them, a particularly important evaluation criterion is *group coverage* (Cauchois et al., 2021), which assesses the coverage and interval length across population subgroups defined by features or response magnitudes. However, in practice, people still tend to prioritize the coverage and length metrics (Lei et al., 2017; Cresswell et al., 2024; Zhang et al., 2024; Xu et al., 2025). Notably, this paper introduces a new metric distinct from existing approaches. Unlike existing studies that design metrics mainly to showcase favorable performance but often struggle with length, we show that PT potentially surpasses its base model in length while performing poorly under the new metric. We provide additional related works on conformal prediction and interval regression in Appendix H.

3 PREJUDICIAL TRICK WITH DECEPTIVE IMPROVEMENT

In this section, we challenge the coverage-length metric in conformal prediction by constructing a trick in Section 3.2. Specifically, this trick potentially improves the interval length while maintaining the coverage, yet it introduces instability and unfairness issues. We further investigate how this trick influences the coverage (Section 3.3) and the length (Section 3.4) theoretically and empirically. Finally, we discuss more details on the deceptive improvements of this trick in Section 3.5.

Table 1: Results for the synthetic datasets (motivating example in Section 3.2). Comparison between VCP and PT-VCP under different α levels, regarding length and coverage.

α	p	VCP		PT-VCP	
		Coverage	Length	Coverage	Length
0.10	0.96	0.906 \pm 0.004	22.894 \pm 0.138	0.909 \pm 0.005	22.614 \pm 0.254
	0.98	0.906 \pm 0.004	22.894 \pm 0.138	0.904 \pm 0.004	22.714 \pm 0.165
0.20	0.96	0.792 \pm 0.011	21.886 \pm 0.125	0.799 \pm 0.011	21.255 \pm 0.136
	0.98	0.792 \pm 0.011	21.886 \pm 0.125	0.796 \pm 0.010	21.589 \pm 0.149

3.1 PRELIMINARY

Conformal Prediction. Conformal prediction creates statistically rigorous uncertainty sets for any predictive model. Given \mathbf{X} as the input and $\alpha \in (0, 1)$ as the miscoverage rate, conformal prediction returns an uncertainty set² $\mathcal{C}_{1-\alpha}(\mathbf{X})$ that satisfies

$$\mathbb{P}(y \in \mathcal{C}_{1-\alpha}(\mathbf{X})) \geq 1 - \alpha, \quad (1)$$

where y denotes the true response of feature \mathbf{X} . We omit the detailed discussion in Appendix C.

Notations. Let $\{Z_i\}_{i=1}^n$ denote n i.i.d. samples drawn from the distribution \mathcal{P}_Z , where $Z \in \mathbb{R}$. Denote $\{Z_{(i)}\}_{i=1}^n$ as the order statistics of $\{Z_i\}_{i=1}^n$ arranged in decreasing order, *i.e.*, $Z_{(1)} \geq Z_{(2)} \geq \dots \geq Z_{(n)}$. The empirical τ -th quantile with n samples is defined as $\hat{Q}_\tau(\{Z_i\}_{i=1}^n) := Z_{(\lceil (n+1)(1-\tau) \rceil)}$. Let \emptyset denote the empty set, and $\mathbb{I}(\cdot)$ denote the indicator function. For a given set \mathcal{C} , let $|\mathcal{C}|$ denote the measure of the set. In this paper, we denote a trained conformal prediction algorithm that directly outputs the $1 - \alpha$ confidence interval given a test point as $\mathcal{A}_{1-\alpha}(\cdot; \hat{\mu})$ for simplicity, where $\hat{\mu}$ is the machine learning algorithm used in the conformal prediction algorithm.

3.2 PREJUDICIAL TRICK

In this section, we propose a trick (Algorithm 1) used in conformal prediction with the intuition from Example 1 and deploy this trick in a motivating example on a synthetic dataset. The experiment results in Table 1 demonstrate that this trick improves the interval length while maintaining the marginal coverage compared to its base. We begin with the construction process of this trick:

Construction Process. For each test point, assign a null set³ with probability $1 - p$, and assign the interval with an adjusted miscoverage rate $\alpha' = 1 - \frac{1-\alpha}{p}$ for the remaining p portion of test points, where $p \in (1 - \alpha, 1)$. Overall, for a new test point \mathbf{x}' , the interval is constructed as Equation (2):

$$\mathcal{C}_{1-\alpha}^{\text{PT}}(\mathbf{x}') = \begin{cases} \text{null set} & \text{with the probability } 1 - p, \\ \mathcal{C}_{1-\alpha'}^{\text{CP}}(\mathbf{x}') & \text{with the probability } p. \end{cases} \quad (2)$$

We call this process Prejudicial Trick (PT, see in Algorithm 1). Note that PT can be directly applied to any base conformal prediction algorithm. To illustrate that PT improves length without sacrificing marginal coverage, we empirically consider a motivating example in Example 3

Example 3 (Synthetic Dataset). Consider a regression setting where the true underlying model is a linear model with the Gaussian mixture noise, given by $Y = \mathbf{X}^\top \beta + \epsilon$, with $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, I_2)$. The noise term ϵ follows $\mathcal{N}(\mu, 1)$ with probability 0.5, and $\mathcal{N}(-\mu, 1)$ with probability 0.5. The training fold, calibration fold, and test fold are generated based on this underlying distribution. We deploy VCP (Algorithm 2) and PT-VCP under such regimes. We refer to Appenix G.1.1 for more details.

Results and Discussions of Example 3. Tabel 1 illustrates the results of Example 3, where *PT-VCP improves the length while maintaining the marginal coverage compared to VCP*. The results validate that the length-coverage metric could be hacked by invalid tricks like PT. The insights are as follows:

²Conformal prediction either returns a set measured by its size (for classification tasks), or returns an interval measured by its length (for regression tasks). We do not distinguish between these terms throughout the paper.

³The null set represents a set of measure zero. It can be an empty set, or a single-point set in regression.

the construction of ϵ guarantees $C_{1-\alpha'}^{\text{CP}}$ close to $C_{1-\alpha}^{\text{CP}}$ regarding the length when choosing proper α and p . Therefore, PT potentially improves the average length when averaging with those null sets.

Notably, Barber et al. (2020) propose a method similar to PT. Unlike PT which emphasizes the potential length improvement, Barber et al. (2020) mainly center on conditional coverage metrics. Moreover, we extend the scope of PT in Remark 2 by relaxing the notion of the null set.

Remark 2 (The extension of PT). PT in Algorithm 1 heavily relies on the notion of null sets. Fortunately, one can extend this null set with an interval returned by conformal prediction with a small coverage rate. For example, PT can be constructed as follows:

$$C_{1-\alpha}^{\text{PT}}(\mathbf{x}') = \begin{cases} C_{1-\alpha_1'}^{\text{CP}}(\mathbf{x}') & \text{with the probability } 1-p, \\ C_{1-\alpha_2'}^{\text{CP}}(\mathbf{x}') & \text{with the probability } p, \end{cases} \quad (3)$$

where $(1-p)\alpha_1' + p\alpha_2' = \alpha$ which guarantees the marginal coverage and α_1' is sufficiently small. We mainly consider the null set in this paper to simplify the related discussions.

3.3 COVERAGE

This section investigates the coverage guarantee of conformal prediction with PT. We prove that PT maintains the marginal coverage guarantees (Theorem 4) and conditional coverage guarantees (Theorem 5, Theorem 6). The empirical validation in Figure 1 supports the theoretical findings.

Marginal Coverage. Theorem 4 proves that PT maintains the valid marginal coverage guarantees.

Theorem 4 (Marginal Coverage Guarantee). *Assume that the exchangeability assumption holds (see Proposition 14 for more details). Then the interval returned by Algorithm 1 with the adjusted miscoverage rate $\alpha' = 1 - \frac{1-\alpha}{p}$ guarantees that*

$$\mathbb{P}(y' \in C_{1-\alpha}^{\text{PT}}(\mathbf{X}')) \geq 1 - \alpha, \quad (4)$$

where (\mathbf{X}', y') denotes a new test point.

The intuition behind Theorem 4: the null set (with probability $1-p$) and the enlarged interval set with miscoverage α' (with probability p) reach the marginal coverage guarantees $p(1-\alpha') = 1-\alpha$.

Conditional Coverage. Theorem 5 further investigates the effect of PT on the condition coverage. Specifically, the interval returned by PT $C_{1-\alpha}^{\text{PT}}(\mathbf{X}')$ keeps the conditional coverage guarantees if the interval returned by its base algorithm $C_{1-\alpha}^{\text{CP}}(\mathbf{X}')$ satisfies the conditional guarantees.

Theorem 5 (Conditional Coverage Guarantee). *Let $C_{1-\alpha}^{\text{CP}}(\mathbf{X}')$ and $C_{1-\alpha}^{\text{PT}}(\mathbf{X}')$ denote the returned interval of \mathbf{X}' respectively. If for any α , $\mathbb{P}(y \in C_{1-\alpha}^{\text{CP}}(\mathbf{X}') \mid \mathbf{X}') \geq 1 - \alpha$ holds for \mathbf{X}' almost surely, then $\mathbb{P}(y \in C_{1-\alpha}^{\text{PT}}(\mathbf{X}') \mid \mathbf{X}') \geq 1 - \alpha$ holds for any α and for \mathbf{X}' almost surely as well.*

The key insight behind Theorem 5: The randomness within PT is independent of the specific input. Therefore, such randomness is averaged out given a specific input, thus keeping the conditional coverage unchanged. Specifically, the conditional coverage is calculated as $p(1-\alpha') = 1-\alpha$.

Remark 3 (Comparison to the Tradeoffs between Conditional Coverage and Interval Length). Existing works on conformal prediction have analyzed the potential tradeoffs between conditional coverage and interval length (Barber et al., 2020; Gibbs et al., 2024). Our work differs from this line, since PT does not operate by creating such a trade-off. Specifically, Theorem 5 validates that PT does not violate the conditional coverage. We will provide a further discussion in Section 3.5.

Theorem 5 requires that the base algorithm satisfies the conditional coverage guarantees. However, this requirement does not always hold in practice. We next prove in Theorem 6 that PT still exhibits the potential to outperform the base algorithm even when the requirement does not hold in practice.

Theorem 6 (Sufficient Condition of Conditional Coverage Guarantees). *Let $f_{\mathcal{A}}(\alpha)$ denote the true conditional miscoverage rate within the subset \mathcal{A} , namely, $\mathbb{P}(y \in C_{1-\alpha}^{\text{CP}}(\mathbf{X}) \mid \mathbf{X} \in \mathcal{A}) = 1 - f_{\mathcal{A}}(\alpha)$ where $C_{1-\alpha}^{\text{CP}}(\mathbf{X})$ denotes the prediction set returned by the base algorithm. Define $\mathcal{F}(p) = pf_{\mathcal{A}}(1 - (1-\alpha)/p)$ where $p \in (1-\alpha, 1)$ represents the parameter in Algorithm 1. If $\mathcal{F}(\cdot)$ satisfies that $\mathcal{F}(1) - \mathcal{F}(p) \geq 1-p$, then it holds that*

$$\mathbb{P}(y \in C_{1-\alpha}^{\text{PT}}(\mathbf{X}) \mid \mathbf{X} \in \mathcal{A}) \geq \mathbb{P}(y \in C_{1-\alpha}^{\text{CP}}(\mathbf{X}) \mid \mathbf{X} \in \mathcal{A}), \quad (5)$$

where $C_{1-\alpha}^{\text{CP}}(\mathbf{X})$ denotes the prediction set returned by PT.

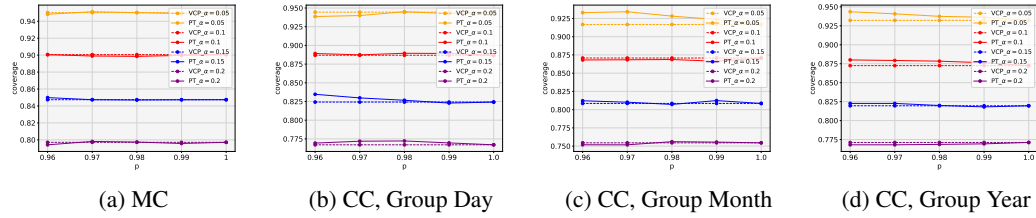


Figure 1: Comparing the (a) marginal coverage and (b, c, d) conditional coverage between VCP with and without PT. Results demonstrate that (1) PT would not significantly change the marginal coverage; (2) PT has better conditional coverage compared to the base algorithm (Theorem 6).

The key insight behind Theorem 6: The conditional miscoverage rate of PT arises from two probabilistic components. For the null set component, the miscoverage rate exceeds that of the base algorithm; for the other component, the miscoverage rate is lower. Therefore, by averaging these two components, PT potentially outperforms its base algorithm in terms of conditional coverage under certain conditions on the miscoverage function $f_{\mathcal{A}}(\alpha)$. Notably, Theorem 6 contains results on both conditional coverages (for one-point set \mathcal{A}) and group coverages (for regular set \mathcal{A}).

Experiments. We compare marginal and conditional coverage rates returned with and without PT on BIKE dataset (Fanaee-T, 2013) using VCP (Algorithm 2) as the base algorithm. We omit the implementation details here and refer to Appendix G.2.2 for details. The results in Figure 1 demonstrate that (a) PT preserves the marginal coverage (Figure 1a); (b) When VCP fails to guarantee the group coverage, PT-VCP fails as well. Experimental results demonstrate that PT achieves comparable group coverage with its base models (Figure 1b, Figure 1c, Figure 1d). We provide more experiments on group coverage with different real-world datasets in Appendix F.1.

3.4 LENGTH

In this section, we investigate the sufficient conditions under which PT improves interval length while keeping the coverage unchanged, and further conduct experiments to validate the theoretical findings. We first propose Theorem 7 as a weak sufficient condition. We then derive a more informative condition under both differentiable regimes (Theorem 8) and non-differentiable regimes (Theorem 11). We further discuss the special cases on the local concave assumption (Corollary 9) and VCP regimes (Corollary 10), and find that misspecification usually satisfies the sufficient conditions (Remark 4). We finally present a failure case in Example 12 when PT cannot outperform its base regarding the interval length.

Experiment results on various datasets in Table 2 align closely with the theoretical results. Besides, we conduct experiments on classification tasks (Table 4), different base algorithms (Table 5), and conduct ablation studies on different hyperparameters (Figure 6-Figure 15).

Additional Notations. We introduce the following notations to facilitate the discussions in this section. Let α denote the miscoverage rate and $s(\mathbf{x}, y; \hat{\mu})$ denote the score function, where $\hat{\mu}(\cdot)$ denotes the learned model. Let $\mathcal{C}_{1-\alpha}^{\text{CP}}(\mathbf{x})$ denote the interval returned by conformal prediction at point \mathbf{x} , and $\mathcal{C}_{1-\alpha}^{\text{PT}}(\mathbf{x})$ denote the interval returned by its PT-variant (Algorithm 1). Let $\mathcal{L}(\mathbf{x}, 1-\alpha; s)$ denote the length of the returned interval at point \mathbf{x} with miscoverage α , i.e., $|\mathcal{C}_{1-\alpha}^{\text{CP}}(\mathbf{x})|$.

We next prove a series of sufficient conditions under which PT improves the interval length, starting from Theorem 7 which provides a simple and straightforward sufficient condition.

Theorem 7 (General Sufficient Condition). *If \mathcal{L} satisfies the following condition:*

$$\exists p \in (1 - \tilde{\alpha}, 1) \quad \text{s.t.} \quad p\mathbb{E}\left(\mathcal{L}\left(\mathbf{x}, \frac{1 - \tilde{\alpha}}{p}; s\right)\right) < \mathbb{E}(\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s)), \quad (6)$$

where $\tilde{\alpha}$ denotes the miscoverage rate and the expectation is taken over \mathbf{x} . Then the interval length returned by PT (with parameter p) outperforms that of its base algorithm, namely

$$\mathbb{E}|\mathcal{C}_{1-\tilde{\alpha}}^{\text{PT}}(\mathbf{X}')| < \mathbb{E}|\mathcal{C}_{1-\tilde{\alpha}}^{\text{CP}}(\mathbf{X}')|, \quad (7)$$

where the expectation is taken over the testing point \mathbf{X}' .

The intuition behind Theorem 7 is pretty simple: PT assigns null sets with a fixed probability whose measure is zero, thus potentially reducing the average length. Although the sufficient condition in Theorem 7 is general, the absence of additional assumptions makes it uninformative in practice. To obtain more insights, we next introduce a differentiable assumption in Theorem 8.

Theorem 8 (First-order Condition). *Assume \mathcal{L} is first-order differentiable and satisfies*

$$\mathbb{E} \left(\frac{\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s)}{1 - \tilde{\alpha}} \right) > \mathbb{E} \left(\frac{\partial}{\partial \alpha} \mathcal{L}(\mathbf{x}, \alpha; s) \Big|_{\alpha=1-\tilde{\alpha}} \right), \quad (8)$$

where $\tilde{\alpha}$ denotes the miscoverage rate and the expectation is taken over \mathbf{x} . Then there exists a parameter p in Algorithm 1, such that the interval length returned by PT outperforms that of its base algorithm, namely

$$\mathbb{E} |C_{1-\tilde{\alpha}}^{PT}(\mathbf{X}')| < \mathbb{E} |C_{1-\tilde{\alpha}}^{CP}(\mathbf{X}')|, \quad (9)$$

where the expectation is taken over the testing point \mathbf{X}' .

Theorem 8 follows the insights of Theorem 7, and further utilizes Equation 8 as the sufficient condition, which characterizes the local behavior of the interval length function. Theorem 8 provides more insights on when and how PT outperforms its base algorithm regarding the length metrics. We next derive a localized concave condition in Corollary 9 based on Theorem 8.

Corollary 9 (Localized Concave Conditions). *Under the settings in Theorem 8, the sufficient condition in Equation 8 holds if $\mathbb{E}(\mathcal{L}(\mathbf{x}, \alpha; s))$ is locally concave on $\alpha \in [0, 1 - \tilde{\alpha}]$.*

Corollary 9 provides a condition under which PT outperforms its base algorithm regarding length. Consider a regression problem with additive noise $y = f^*(x) + \epsilon$. If the noise distribution exhibits local concavity and the base model approximates the true function f^* well, then the length function $\mathbb{E}(\mathcal{L}(\mathbf{x}, \alpha; s))$ generally satisfies the localized concavity property. Consequently, the performance improvement of PT is guaranteed by Corollary 9. This inspires the construction of Example 3.

Besides, Corollary 10 focuses on the settings of deploying VCP. Since VCP returns the same length for each individual, the expectation operator in Theorem 8 degenerates.

Corollary 10 (Deterministic Case). *Under the settings in Theorem 8, if applying VCP (Algorithm 2) as the base algorithm, the sufficient condition in Equation 8 holds if*

$$\frac{\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s)}{1 - \tilde{\alpha}} > \frac{\partial}{\partial \alpha} \mathcal{L}(\mathbf{x}, \alpha; s) \Big|_{\alpha=1-\tilde{\alpha}}, \quad (10)$$

where the expectation operator degenerates due to the characteristics of VCP.

Unfortunately, real-world applications may not satisfy the differentiability assumption in Theorem 8. Therefore, we relax this assumption and derive a secant sufficient condition in Theorem 11.

Theorem 11 (Secant Sufficient Condition). *If there exists $u \in (1 - \tilde{\alpha}, 1)$, such that*

$$\frac{\mathbb{E}(\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s))}{1 - \tilde{\alpha}} > \frac{\mathbb{E}(\mathcal{L}(\mathbf{x}, u; s)) - \mathbb{E}(\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s))}{u - (1 - \tilde{\alpha})}, \quad (11)$$

where $\tilde{\alpha}$ denotes the miscoverage rate and the expectation is taken over \mathbf{x} . Then there exists a parameter $p = (1 - \tilde{\alpha})/u$ in Algorithm 1, such that the interval length returned by PT outperforms its base algorithm, namely,

$$\mathbb{E} |C_{1-\tilde{\alpha}}^{PT}(\mathbf{X}')| < \mathbb{E} |C_{1-\tilde{\alpha}}^{CP}(\mathbf{X}')|, \quad (12)$$

where the expectation is taken over the testing point \mathbf{X}' .

Theorem 11 shares similar intuitions with Theorem 8, and further relaxes the differentiability assumption by comparing the secant slopes. Informally, PT achieves smaller average lengths than its base algorithm when the length function does not grow extremely fast within the region $(1 - \tilde{\alpha}, 1)$.

Remark 4 (Relationship Between Misspecification and Sufficient Condition). Model misspecification is a common practical scenario that aligns with our theoretical analysis, as it typically satisfies the sufficient conditions in Theorem 8 and Theorem 11 (Wang & Blei, 2020; Huang et al., 2023). Specifically, misspecification leads to a residual with a non-zero mean, resulting in a non-convex length function. This outcome is closely related to the local concavity condition in Corollary 9. For this reason, we employ misspecification regimes in most of our experiments.

Table 2: Comparison of performance between VCP and PT-VCP in regression tasks across different datasets ($\alpha = 0.1$).

METHOD	BIAS	VCP		PT-VCP	
		COVERAGE	LENGTH	COVERAGE	LENGTH
MEPS-19	20	0.90 \pm 0.000	42.34 \pm 0.228	0.90 \pm 0.000	41.92 \pm 0.389
MEPS-20	20	0.90 \pm 0.000	41.98 \pm 0.116	0.90 \pm 0.000	41.41 \pm 0.241
MEPS-21	20	0.90 \pm 0.004	42.28 \pm 0.112	0.90 \pm 0.000	41.90 \pm 0.300
BIKE	10	0.90 \pm 0.000	20.46 \pm 0.018	0.90 \pm 0.004	19.59 \pm 0.018
BLOG-DATA	20	0.90 \pm 0.004	41.67 \pm 0.336	0.90 \pm 0.000	41.13 \pm 0.416
BIO	10	0.90 \pm 0.004	21.13 \pm 0.336	0.90 \pm 0.000	20.44 \pm 0.031
FACEBOOK-1	10	0.90 \pm 0.000	20.81 \pm 0.036	0.90 \pm 0.000	20.80 \pm 0.179
FACEBOOK-2	10	0.90 \pm 0.000	20.97 \pm 0.067	0.90 \pm 0.000	21.01 \pm 0.179
CONCRETE	5	0.90 \pm 0.013	10.32 \pm 0.009	0.89 \pm 0.009	9.87 \pm 0.031
STAR	5	0.91 \pm 0.004	10.14 \pm 0.004	0.91 \pm 0.004	9.63 \pm 0.027

Remark 5 (Extensions of PT Beyond Conformal Prediction). While Theorem 7 can be extended to other interval estimation tasks, the relevance between PT and conformal prediction lies in the model misspecification (Remark 4). Specifically, model misspecification serves as a potential sufficient condition where PT works, and it generally appears in real-world applications of conformal prediction. However, the existence of model misspecification does not always hold for other interval estimation tasks, and therefore limits the extensions of PT beyond conformal prediction.

However, the aforementioned sufficient conditions are not always satisfied. We present a failure case in Example 12 and illustrate the empirical validation in Figure 5.

Example 12 (Failure Case). *If the values of non-conformity score in VCP follows a Gaussian distribution over randomness in x , then for all $\alpha \in (0, 1)$, and all $p \in (1 - \alpha, 1)$, it holds that*

$$\mathbb{E}|C_{1-\alpha}^{PT}(\mathbf{X}')| > \mathbb{E}|C_{1-\alpha}^{CP}(\mathbf{X}')|. \quad (13)$$

Example 12 demonstrates that PT does not always outperform its base algorithm regarding the length-coverage metric. However, our goal is not to present PT as a universally applicable method, but to demonstrate a fundamental flaw in the coverage-length evaluation. To achieve this, the existence of any realistic scenarios where PT can create deceptively shorter intervals is sufficient.

Experiment. We conduct several experiments comparing the length returned with and without PT using VCP (Algorithm 2) on MEPS19-21 (Cohen et al., 2009), BIKE (Fanaee-T, 2013), BLOG-DATA (Buza, 2014), BIO (Rana, 2013), FACEBOOK1-2 (Singh, 2015), CONCRETE (Yeh, 1998), STAR (Achilles et al., 2008). To simulate model misspecification, we manually add bias to the label (as shown in the bias column in Table 2). We refer to Appendix G.2 for setting details. The results in Table 2 demonstrate that PT generally achieves smaller average lengths compared to its base algorithm in most cases (9 out of 10). Besides, we conduct more experiments and ablations:

- We compare the length of VCP and PT-VCP under classification tasks in Table 4;
- We use CQR (Romano et al., 2019) as the base algorithm on regression tasks in Table 5;
- We conduct ablation studies on hyperparameter p and misspecification level μ in Appendix F.2.

3.5 DECEPTIVE IMPROVEMENT

We prove that PT preserves (conditional) coverage guarantees in Section 3.3 and achieves shorter prediction intervals under certain conditions in Section 3.4. Despite these theoretical benefits, PT is poorly suited for practical deployment. The primary issue is that PT introduces randomness, causing prediction intervals to vary across different runs. This inherent instability undermines the method’s reliability. Besides, the problem becomes more dramatic in the scenario of Remark 2, where the individuals are grouped with different misscoverage rates. This randomness makes it impossible for a user to identify their assigned group, making the confidence interval meaningless. Therefore, while PT may appear superior based on the traditional length-coverage metric, its practical instability makes it unsuitable for real-world deployment. This discrepancy challenges the sufficiency of the length-coverage metric itself, suggesting it is not a complete measure of a method’s practical utility.

Table 3: Comparison between VCP and VCP with PT regarding interval stability.

Method	meps-19	meps-20	meps-21	bike	blog-data	bio	facebook-1	facebook-2	concrete	star
VCP	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000	0.00 ±0.000
PT-VCP	1.26 ±0.015	1.24 ±0.008	1.26 ±0.011	0.58 ±0.002	1.23 ±0.014	0.61 ±0.001	0.62 ±0.005	1.19 ±0.006	1.14 ±0.012	1.14 ±0.007

4 INTERVAL STABILITY

In this section, we propose *interval stability* (Definition 1) which measures the randomness in each run of conformal prediction. We begin with the definition of interval stability in Definition 1.

Definition 1 (Interval Stability). *Let X denote a data point with returned confidence interval $C_{1-\alpha}(X)$, and let $|\cdot|$ denote a certain measure of the interval (e.g., its length). Let \mathcal{A} denote the conformal prediction algorithm, and \mathcal{D}_{ca} the calibration dataset. The interval stability is defined as*

$$\mathbb{E}_X [\text{Var}_{\mathcal{A}|X, \mathcal{D}_{ca}} (|C_{1-\alpha}(X)|)]. \quad (14)$$

The interval stability captures the expected variability of the interval size conditional on the test point and calibration randomness. Intuitively, it captures the inconsistency of the returned intervals when the algorithm is run multiple times on the same test point and calibration dataset.

Due to the stochastic nature of PT, it tends to produce a large interval stability, implying the practical instability issues. We prove in Proposition 13 that PT indeed introduces a non-zero interval stability.

Proposition 13. *Following the notations in Section 3.4, it holds that*

$$\mathbb{E}_X [\text{Var}_{\mathcal{A}|X, \mathcal{D}_{ca}} (|C_{1-\tilde{\alpha}}^{PT}(X)|)] = p(1-p) (\mathbb{E}(\mathcal{L}(\mathbf{x}, (1-\tilde{\alpha})/p; s)))^2 > 0. \quad (15)$$

Notably, the *interval stability* metric is not just for detecting our specific PT construction, but serves as a safeguard to ensure that future advancements in conformal prediction are genuine and reliable, rather than arising from the unprincipled randomness. As the community pushes for ever-shorter prediction intervals, there is a risk that increasingly complex methods might implicitly introduce forms of randomness that offer deceptive gains. We refer to Remark 6 for a detailed discussion.

Remark 6 (Why Interval Stability). Interval stability is still meaningful even if existing approaches in conformal prediction do not always rely on the randomness. In the existing literature, numerous approaches claim a superior performance through a smaller interval length, under the traditional coverage-length metric. In this paper, we show that randomness may break this metric through PT due to practical issues. This raises concerns that as methods become increasingly complex, they may implicitly utilize similar randomness to improve the length. Such effects may be unintended but hard to be recognized. To address this issue, interval stability serves as a tool for detecting such issues and highlighting the risks inherent in the current reliance on the length-coverage metric alone.

Experiment. We follow the same setting as the experiment in Section 3.4, but use interval stability as the metric. Table 3 demonstrates that interval stability successfully detects the vacuous randomness in PT. We defer to Appendix F.1 for more experiments with CQR and classification regimes.

We acknowledge that for deterministic methods, interval stability will indeed be zero. This is by design. The metric is not intended to replace coverage and length, but to complement them, acting as a specific check against the kind of vacuous randomness we identify. A value of zero is a *pass* on this specific test, confirming the method’s deterministic nature for a given input.

5 CONCLUSION

In this paper, we demonstrate a pitfall in how conformal prediction methods are evaluated. We introduce PT, a technique that hacks the conventional coverage-length metric by producing deceptively shorter intervals while preserving coverage guarantees. However, PT relies on the randomness that leads to instability: the algorithm can produce different prediction sets for a given input on different runs. This creates practical issues in real-world, high-stakes scenarios. Our theoretical and empirical results confirm that while PT appears superior, its foundation is flawed. This discrepancy challenges the completeness of the coverage-length metric. Consequently, we propose Interval Stability as a diagnostic tool, which helps flag the potential vacuous randomness for a newly proposed method.

486 ETHICS STATEMENT
487

488 This research focuses on the pitfall in conformal prediction. Our work does not involve human
489 subjects. Therefore no Institutional Review Board (IRB) approval was required. All experiments
490 were conducted on standard, publicly available benchmarks, which are widely used in the machine
491 learning community. Our research does not involve the collection of new data, nor does it process
492 personally identifiable or sensitive information, thus mitigating concerns related to data privacy and
493 security.

494
495 REPRODUCIBILITY STATEMENT
496

497 To ensure the reproducibility of our work, we provide detailed descriptions of our theoretical results
498 and experimental setup. The theoretical results presented in Section 3.3 and Section 3.4 are accom-
499 panied by complete mathematical proofs in Appendix D. Our full experimental setup is described in
500 Appendix G. The source code is provided in the supplementary material and will be made publicly
501 available upon publication.

502
503 REFERENCES
504

- 505 CM Achilles, Helen Pate Bain, Fred Bellott, Jayne Boyd-Zaharias, Jeremy Finn, John Folger, John
506 Johnston, and Elizabeth Word. Tennessee’s student teacher achievement ratio (star) project. *Har-*
507 *vard Dataverse*, 1:2008, 2008.
- 508 Niccolò Ajroldi, Jacopo Diquigiovanni, Matteo Fontana, and Simone Vantini. Conformal prediction
509 bands for two-dimensional functional time series. *Computational Statistics & Data Analysis*, 187:
510 107821, 2023.
- 511 Ahmed M Alaa, Zeshan Hussain, and David Sontag. Conformalized unconditional quantile re-
512 gression. In *International conference on artificial intelligence and statistics*, pp. 10690–10702.
513 PMLR, 2023.
- 514 Anastasios Angelopoulos, Stephen Bates, Jitendra Malik, and Michael I Jordan. Uncertainty sets
515 for image classifiers using conformal prediction. *arXiv preprint arXiv:2009.14193*, 2020.
- 516 Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and
517 distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- 518 Anastasios N Angelopoulos, Amit Pal Kohli, Stephen Bates, Michael Jordan, Jitendra Malik,
519 Thayer Alshaabi, Srigoikul Upadhyayula, and Yaniv Romano. Image-to-image regression with
520 distribution-free uncertainty quantification and applications in imaging. In *International Confer-*
521 *ence on Machine Learning*, pp. 717–730. PMLR, 2022.
- 522 Anastasios N Angelopoulos, Stephen Bates, et al. Conformal prediction: A gentle introduction.
523 *Foundations and Trends® in Machine Learning*, 16(4):494–591, 2023.
- 524 Yu Bai, Song Mei, Huan Wang, Yingbo Zhou, and Caiming Xiong. Efficient and differentiable
525 conformal prediction with general function classes, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2202.11091)
526 [2202.11091](https://arxiv.org/abs/2202.11091).
- 527 Rina Foygel Barber, Emmanuel J. Candès, Aaditya Ramdas, and Ryan J. Tibshirani. The limits of
528 distribution-free conditional predictive inference, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1903.04684)
529 [1903.04684](https://arxiv.org/abs/1903.04684).
- 530 Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. Conformal
531 prediction beyond exchangeability. *The Annals of Statistics*, 51(2):816–845, 2023.
- 532 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
533 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
534 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

- 540 Krisztian Buza. Feedback prediction for blogs. In *Data analysis, machine learning and knowledge*
541 *discovery*, pp. 145–152. Springer, 2014.
- 542
- 543 Emmanuel Candès, Lihua Lei, and Zhimei Ren. Conformalized survival analysis. *Journal of the*
544 *Royal Statistical Society Series B: Statistical Methodology*, 85(1):24–45, 2023.
- 545 Maxime Cauchois, Suyash Gupta, and John C Duchi. Knowing what you know: valid and validated
546 confidence sets in multiclass and multilabel prediction. *Journal of machine learning research*, 22
547 (81):1–42, 2021.
- 548
- 549 Joel W. Cohen, Steven B. Cohen, and Jessica S. Banthin. The medical expenditure panel survey: A
550 national information resource to support healthcare cost research and inform policy and practice.
551 *Medical Care*, 47:S44–S50, 2009.
- 552 Jesse C. Cresswell, Yi Sui, Bhargava Kumar, and Noël Vouitsis. Conformal prediction sets improve
553 human decision making, 2024. URL <https://arxiv.org/abs/2401.13744>.
- 554
- 555 Lahav Dabah and Tom Tirer. On calibration and conformal prediction of deep classifiers. *arXiv*
556 *preprint arXiv:2402.05806*, 2024.
- 557 Marcos Lopez De Prado. *Advances in financial machine learning*. John Wiley & Sons, 2018.
- 558
- 559 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-scale
560 hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition*,
561 pp. 248–255, 2009.
- 562 Hadi Fanaee-T. Bike Sharing. UCI Machine Learning Repository, 2013. DOI:
563 <https://doi.org/10.24432/C5W894>.
- 564
- 565 Shai Feldman, Stephen Bates, and Yaniv Romano. Improving conditional coverage via orthogonal
566 quantile regression. *Advances in neural information processing systems*, 34:2060–2071, 2021.
- 567
- 568 Adam Fisch, Tal Schuster, Tommi Jaakkola, and Dr.Regina Barzilay. Conformal prediction sets with
569 limited false positives. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,
570 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine*
571 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 6514–6532. PMLR,
17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/fisch22a.html>.
- 572
- 573 Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. The limits of
574 distribution-free conditional predictive inference. *Information and Inference: A Journal of the*
575 *IMA*, 10(2):455–482, 2021.
- 576
- 577 Isaac Gibbs, John J. Cherian, and Emmanuel J. Candès. Conformal prediction with conditional
578 guarantees, 2024. URL <https://arxiv.org/abs/2305.12616>.
- 579
- 580 Laying Guan. Localized conformal prediction: A generalized inference framework for conformal
581 prediction. *Biometrika*, 110(1):33–50, 2023.
- 582
- 583 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
584 networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- 585
- 586 Xing Han, Ziyang Tang, Joydeep Ghosh, and Qiang Liu. Split localized conformal prediction. *arXiv*
587 *preprint arXiv:2206.13092*, 2022.
- 588
- 589 Shengyi He and Henry Lam. Statistically optimal uncertainty quantification for expensive black-box
590 models. *arXiv preprint arXiv:2408.05887*, 2024.
- 591
- 592 Daolang Huang, Ayush Bharti, Amauri Souza, Luigi Acerbi, and Samuel Kaski. Learning robust
593 statistics for simulation-based inference under model misspecification, 2023. URL <https://arxiv.org/abs/2305.15871>.
- 594
- 595 Rafael Izbicki, Gilson Shimizu, and Rafael Stern. Flexible distribution-free conditional predictive
596 bands using density estimators. In *International Conference on Artificial Intelligence and Statis-*
597 *tics*, pp. 3068–3077. PMLR, 2020.

- 594 Vilde Jensen, Filippo Maria Bianchi, and Stian Normann Anfinssen. Ensemble conformalized quan-
595 tile regression for probabilistic time series forecasting. *IEEE Transactions on Neural Networks*
596 *and Learning Systems*, 35(7):9014–9025, July 2024. ISSN 2162-2388. doi: 10.1109/tnnls.2022.
597 3217694. URL <http://dx.doi.org/10.1109/TNNLS.2022.3217694>.
- 598
- 599 Ying Jin, Zhimei Ren, and Emmanuel J Candès. Sensitivity analysis of individual treatment effects:
600 A robust conformal inference approach. *Proceedings of the National Academy of Sciences*, 120
601 (6):e2214889120, 2023.
- 602 Shayan Kiyani, George Pappas, and Hamed Hassani. Length optimization in conformal prediction,
603 2024. URL <https://arxiv.org/abs/2406.18814>.
- 604
- 605 Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning
606 using calibrated regression. In *International conference on machine learning*, pp. 2796–2804.
607 PMLR, 2018.
- 608 Jing Lei, Alessandro Rinaldo, and Larry Wasserman. A conformal prediction approach to explore
609 functional data. *Annals of Mathematics and Artificial Intelligence*, 74:29–43, 2015.
- 610
- 611 Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-
612 free predictive inference for regression, 2017. URL [https://arxiv.org/abs/1604.](https://arxiv.org/abs/1604.04173)
613 04173.
- 614
- 615 Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-
616 free predictive inference for regression. *Journal of the American Statistical Association*, 113
617 (523):1094–1111, 2018.
- 618 Lihua Lei and Emmanuel J Candès. Conformal inference of counterfactuals and individual treatment
619 effects. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 83(5):911–938,
620 2021.
- 621
- 622 Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby,
623 Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. *Advances*
624 *in Neural Information Processing Systems*, 34:15682–15694, 2021.
- 625 Jiri Navratil, Matthew Arnold, and Benjamin Elder. Uncertainty prediction for deep sequential
626 regression using meta models. *arXiv preprint arXiv:2007.01350*, 2020.
- 627
- 628 Harris Papadopoulos, Alex Gammerman, and Volodya Vovk. Normalized nonconformity measures
629 for regression conformal prediction. In *Proceedings of the IASTED International Conference on*
630 *Artificial Intelligence and Applications (AIA 2008)*, pp. 64–69, 2008.
- 631
- 632 Harris Papadopoulos, Vladimir Vovk, and Alex Gammerman. Regression conformal prediction with
633 nearest neighbours. *Journal of Artificial Intelligence Research*, 40:815–840, 2011.
- 634
- 635 Prashant Rana. Physicochemical Properties of Protein Tertiary Structure. UCI Machine Learning
636 Repository, 2013. DOI: <https://doi.org/10.24432/C5QW3H>.
- 637
- 638 Yaniv Romano, Evan Patterson, and Emmanuel Candès. Conformalized quantile regression. *Ad-*
639 *vances in neural information processing systems*, 32, 2019.
- 640
- 641 Mauricio Sadinle, Jing Lei, and Larry Wasserman. Least ambiguous set-valued classifiers with
642 bounded error levels. *Journal of the American Statistical Association*, 114(525):223–234, June
643 2018. ISSN 1537-274X. doi: 10.1080/01621459.2017.1395341. URL [http://dx.doi.org/](http://dx.doi.org/10.1080/01621459.2017.1395341)
644 10.1080/01621459.2017.1395341.
- 645
- 646 Yuya Sasaki, Takuya Ura, and Yichong Zhang. Unconditional quantile regression with high-
647 dimensional data. *Quantitative Economics*, 13(3):955–978, 2022.
- 648
- 649 Nabeel Seedat, Alan Jeffares, Fergus Imrie, and Mihaela van der Schaar. Improving adaptive con-
650 formal prediction using self-supervised learning, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2302.12238)
651 2302.12238.

- 648 Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning*
649 *Research*, 9(3), 2008.
- 650
- 651 Kamaljit Singh. Facebook Comment Volume. UCI Machine Learning Repository, 2015. DOI:
652 <https://doi.org/10.24432/C5Q886>.
- 653
- 654 Ralph C Smith. *Uncertainty quantification: theory, implementation, and applications*. SIAM, 2024.
- 655
- 656 Kamile Stankeviciute, Ahmed M Alaa, and Mihaela van der Schaar. Conformal time-series fore-
657 casting. *Advances in neural information processing systems*, 34:6216–6228, 2021.
- 658
- 659 David Stutz, Krishnamurthy, Dvijotham, Ali Taylan Cemgil, and Arnaud Doucet. Learning optimal
660 conformal classifiers, 2022. URL <https://arxiv.org/abs/2110.09192>.
- 661
- 662 Timothy John Sullivan. *Introduction to uncertainty quantification*, volume 63. Springer, 2015.
- 663
- 664 Jiaye Teng, Zeren Tan, and Yang Yuan. T-sci: A two-stage conformal inference algorithm with
665 guaranteed coverage for cox-mlp. In *International conference on machine learning*, pp. 10203–
666 10213. PMLR, 2021.
- 667
- 668 Jiaye Teng, Chuan Wen, Dinghuai Zhang, Yoshua Bengio, Yang Gao, and Yang Yuan. Predictive
669 inference with feature conformal prediction. *arXiv preprint arXiv:2210.00173*, 2022.
- 670
- 671 Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal pre-
672 diction under covariate shift. *Advances in neural information processing systems*, 32, 2019.
- 673
- 674 Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, and Eftychios Protopapadakis.
675 Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience*,
676 2018(1):7068349, 2018.
- 677
- 678 Vladimir Vovk. Conditional validity of inductive conformal predictors. In Steven C. H. Hoi and
679 Wray Buntine (eds.), *Proceedings of the Asian Conference on Machine Learning*, volume 25 of
680 *Proceedings of Machine Learning Research*, pp. 475–490, Singapore Management University,
681 Singapore, 04–06 Nov 2012. PMLR. URL [https://proceedings.mlr.press/v25/
vovk12.html](https://proceedings.mlr.press/v25/vovk12.html).
- 682
- 683 Vladimir Vovk, Alexander Gammernan, and Glenn Shafer. *Algorithmic learning in a random world*,
684 volume 29. Springer, 2005.
- 685
- 686 Kang Wang and Subhashis Ghosal. Coverage of credible intervals in bayesian multivariate isotonic
687 regression. *The Annals of Statistics*, 51(3):1376–1400, 2023.
- 688
- 689 Yixin Wang and David M. Blei. Variational bayes under model misspecification, 2020. URL
690 <https://arxiv.org/abs/1905.10859>.
- 691
- 692 Ran Xie, Rina Foygel Barber, and Emmanuel J. Candès. Boosted conformal prediction intervals,
693 2024. URL <https://arxiv.org/abs/2406.07449>.
- 694
- 695 Chen Xu and Yao Xie. Conformal prediction interval for dynamic time-series. In *International*
696 *Conference on Machine Learning*, pp. 11559–11569. PMLR, 2021.
- 697
- 698 Yunpeng Xu, Mufang Ying, Wenge Guo, and Zhi Wei. Two-stage risk control with application to
699 ranked retrieval, 2025. URL <https://arxiv.org/abs/2404.17769>.
- 700
- 701 I-Cheng Yeh. Concrete Compressive Strength. UCI Machine Learning Repository, 1998. DOI:
<https://doi.org/10.24432/C5PK67>.
- 702
- 703 Soroush H Zargarbashi and Aleksandar Bojchevski. Conformal inductive graph neural networks.
704 *arXiv preprint arXiv:2407.09173*, 2024.
- 705
- 706 Soroush H Zargarbashi, Simone Antonelli, and Aleksandar Bojchevski. Conformal prediction sets
707 for graph neural networks. In *International Conference on Machine Learning*, pp. 12292–12318.
708 PMLR, 2023.

702 Dongping Zhang, Angelos Chatzimparmpas, Negar Kamali, and Jessica Hullman. Evaluating
703 the utility of conformal prediction sets for ai-advised image labeling, 2024. URL <https://arxiv.org/abs/2401.08876>.
704
705

706 Shengjia Zhao, Tengyu Ma, and Stefano Ermon. Individual calibration with randomized forecasting,
707 2020. URL <https://arxiv.org/abs/2006.10288>.
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

Appendix

We firstly restate our contributions in Appendix A. Then we present more discussions in Appendix B. In Appendix C, we introduce the conformal prediction and we provide missing proofs in Appendix D. In Appendix E, we exhibit the omitted illustrations in our paper. In Appendix F, we illustrate the omitted experimental results. In Appendix G, we present implementation details of our experiments. In Appendix H, we exhibit more related works. In Appendix I, we clarify the use of large language models in our paper.

A CONTRIBUTIONS RESTATEMENT

We summarize our contributions as follows:

- We observe that the traditional coverage-length criteria in conformal prediction might be hacked using a counter-intuitive method PT, (Algorithm 1), since PT might deceptively improve the length while maintaining the coverage but raises fairness issues;
- We theoretically derive in Theorem 7 the conditions under which PT deceptively improves length, while keeping valid marginal coverage (Theorem 4) and conditional coverage (Theorem 5, Theorem 6). We further derive several sufficient conditions under which PT improves length with first-order differentiability assumption (Theorem 8) or without first-order differentiability assumption (Theorem 11);
- We propose a new metric in Section 4, termed interval stability. Interval stability measures the variance of the prediction interval over the input introduced by the conformal prediction algorithms, helping to mitigate the adverse impacts of PT.

B MORE DISCUSSIONS

Similarity between PT and method discussed in Barber et al. (2020). In Section 3.2, we mention the similarity between our PT method and the randomness discussed in Barber et al. (2020). While the mechanism in our work bears a structural resemblance to that in Barber et al. (2020), our motivation and conclusion are fundamentally different. Barber et al. (2020) investigate the inherent trade-offs required to achieve conditional coverage, using randomization as a tool to explore theoretical limits. In contrast, our work focuses on the evaluation paradigm itself. We use PT not to achieve a desirable property (like conditional coverage), but to demonstrate a failure mode of a widely-used metric (average interval length). Our primary contribution is to highlight this pitfall and propose a remedy (Interval Stability), a direction not explored by Barber et al. (2020).

C OMITTED PRELIMINARY

Interval Prediction. Interval prediction aims to construct a confidence interval that contains the true response value with a user-specified probability. Compared to traditional point estimation, interval prediction provides more comprehensive statistical information by quantifying the uncertainty using the interval length, which is often a more challenging goal. Definition 2 presents the formal definition.

Definition 2 (Interval Prediction). *Let (\mathbf{X}, Y) denote a feature-response pair. Given a miscoverage rate α , interval prediction aims to construct a confidence interval $\mathcal{C}_{1-\alpha}(\mathbf{X})$, such that*

$$\mathbb{P}(Y \in \mathcal{C}_{1-\alpha}(\mathbf{X})) \geq 1 - \alpha. \quad (16)$$

Given the coverage in Equation (16), a smaller confidence interval indicates a more precise estimate.

Conformal Prediction. To construct an interval prediction, we introduce a widely used approach called vanilla conformal prediction. The VCP method is typically divided into four stages: dataset splitting, training, calibration, and construction. The whole procedure is presented in Algorithm 2.

Dataset Splitting. Let $\mathcal{D} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}\}$ denote the i.i.d. samples from a distribution $\mathcal{P}_{\mathbf{X}Y}$ over the covariate $\mathbf{X} \in \mathbb{R}^d$ and the response $Y \in \mathbb{R}$. The VCP first randomly splits the dataset \mathcal{D} into

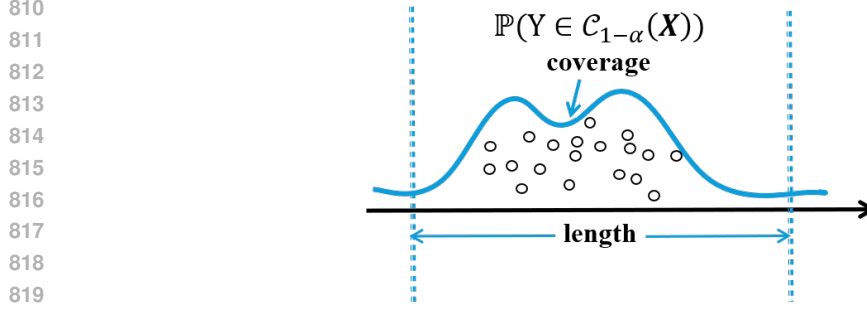


Figure 2: Illustration of coverage and interval length.

two folds: a training fold $\mathcal{D}_{\text{tr}} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}_{\text{tr}}\}$ and a calibration fold $\mathcal{D}_{\text{ca}} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}_{\text{ca}}\}$, where $\mathcal{I}_{\text{tr}} \cup \mathcal{I}_{\text{ca}} = \mathcal{I}$ and $\mathcal{I}_{\text{tr}} \cap \mathcal{I}_{\text{ca}} = \emptyset$.

Training Process. We train a model denoted by $\hat{\mu}(\cdot)$ (e.g., a neural network) via the training fold \mathcal{D}_{tr} .

Calibration Process. Given the trained model $\hat{\mu}(\cdot)$, VCP calculates the non-conformity score on the calibration fold \mathcal{D}_{ca} , denoted by $\mathcal{V} = \{s(\mathbf{x}_i, y_i; \hat{\mu}) : i \in \mathcal{I}_{\text{ca}}\}$. The non-conformity score $s(\cdot)$ measures how well the model $\hat{\mu}(\cdot)$ fits the ground truth. A commonly used non-conformity score in regression tasks is the absolute residual, defined as $s(\mathbf{x}_i, y_i; \hat{\mu}) = |y_i - \hat{\mu}(\mathbf{x}_i)|$.

Construction Process. Finally, for a given miscoverage rate α , we then compute a $(1 - \tilde{\alpha})$ -th quantile $\hat{Q}_{1-\tilde{\alpha}}(\mathcal{V})$ of the empirical distribution of the non-conformity score set \mathcal{V} calculated on the calibration set, where $1 - \tilde{\alpha} = (1 - \alpha)(1 + 1/|\mathcal{V}|)$. The prediction interval at a new point \mathbf{x}' is then given by

$$\mathcal{C}_{1-\alpha}(\mathbf{x}') = \{y : s(\mathbf{x}', y; \hat{\mu}) \leq \hat{Q}_{1-\tilde{\alpha}}(\mathcal{V})\}. \quad (17)$$

Coverage and Length. To evaluate the performance of interval prediction, two commonly used metrics: *coverage* and *length* are defined in Definition 3, as further illustrated in Figure 2.

Definition 3 (Coverage and Length). *Let (\mathbf{X}, Y) denote a feature-response pair from a joint distribution $\mathcal{P}_{\mathbf{X}Y}$, and let $\mathcal{C}_{1-\alpha}(\mathbf{X})$ denote the confidence interval to be evaluated and let $|\cdot|$ denote a certain measure of $\mathcal{C}_{1-\alpha}(\mathbf{X})$. The coverage and length of $\mathcal{C}_{1-\alpha}(\mathbf{X})$ is given by:*

$$\begin{aligned} \text{Coverage} &:= \mathbb{E}[\mathbb{I}(Y \in \mathcal{C}_{1-\alpha}(\mathbf{X}))], \\ \text{Length} &:= \mathbb{E}|\mathcal{C}_{1-\alpha}(\mathbf{X})|. \end{aligned} \quad (18)$$

For example, the length of the prediction interval given by VCP in Equation (17) is:

$$\text{Length} = \mathbb{E} \left[2\hat{Q}_{1-\tilde{\alpha}}(\mathcal{V}) \right]. \quad (19)$$

Notably, the two metrics in Definition 3 evaluate the quality of prediction intervals from different perspectives. Figure 2 illustrates the coverage and length given a distribution. Firstly, high coverage ensures that the true value falls within the interval with high probability. A valid confidence interval should guarantee that the coverage exceeds $1 - \alpha$, as suggested in Equation 16. However, setting a sufficiently large interval always guarantees Equation (16), which is impractical and meaningless. Therefore, the length metric is required to ensure the interval’s precision. Based on the above discussion, the gold standard in conformal prediction is *making the length as small as possible, given that the coverage is larger than $1 - \alpha$* .

Following the gold standard, VCP ensures the coverage guarantee under mild exchangeability assumption (Proposition 14), but pays less attention to the length. As a result, numerous works on improving the length of VCP from different perspectives (Papadopoulos et al., 2011; Romano et al., 2019) use intuitively valid approaches.

Proposition 14 (Coverage Guarantee). *The terms \mathcal{U}_i are exchangeable if arbitrary permutation leads to the same distribution, i.e., $(\mathcal{U}_1, \dots, \mathcal{U}_{|\mathcal{I}_{\text{ca}}|+1}) \stackrel{d}{=} (\mathcal{U}_{\pi(1)}, \dots, \mathcal{U}_{\pi(|\mathcal{I}_{\text{ca}}|+1)})$ with arbitrary permutation π over $1, \dots, |\mathcal{I}_{\text{ca}}| + 1$, where $\stackrel{d}{=}$ denotes equivalence in distribution. Suppose that the data pair (\mathbf{x}_i, y_i) , $i \in \mathcal{I}_{\text{ca}}$ and the test point (\mathbf{x}', y') are exchangeable, then the confidence interval $\mathcal{C}_{1-\alpha}(\mathbf{x}')$ returned by Algorithm 2 satisfies*

$$\mathbb{P}(y' \in \mathcal{C}_{1-\alpha}(x')) \geq 1 - \alpha.$$

D PROOFS FOR THEOREMS AND COROLLARIES

D.1 PROOF OF THEOREM 4

The returned interval could be in two folds: the vacuous fold and the meaningful fold. Therefore, by the definition of PT (Algorithm 1), we have

$$\mathbb{P}(y' \in \mathcal{C}_{1-\alpha}^{PT}(\mathbf{X}')) = \mathbb{P}(y' \in \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \mid \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \text{ is in the meaningful fold})p \quad (20)$$

$$+ \mathbb{P}(y' \in \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \mid \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \text{ is in the vacuous fold})(1-p) \quad (21)$$

$$\geq (1-\alpha')p = 1 - \alpha. \quad (22)$$

D.2 PROOF OF THEOREM 5

Given that for all possible values of \mathbf{X} , there holds

$$\mathbb{P}(y \in \mathcal{C}_{1-\alpha}^{CP}(\mathbf{X}) \mid \mathbf{X}) \geq 1 - \alpha \quad (23)$$

holds almost surely, we have

$$\mathbb{P}(y' \in \mathcal{C}_{1-\alpha}^{PT}(\mathbf{X}') \mid \mathbf{X}) = \mathbb{P}(y' \in \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \mid \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \text{ is in the meaningful fold, } \mathbf{X})p \quad (24)$$

$$+ \mathbb{P}(y' \in \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \mid \mathcal{C}_{1-\alpha'}^{CP}(\mathbf{X}') \text{ is in the vacuous fold, } \mathbf{X})(1-p) \quad (25)$$

$$\geq (1-\alpha')p = 1 - \alpha \quad (26)$$

holds almost surely.

D.3 PROOF OF THEOREM 6

Given that it holds

$$\mathbb{P}(y \in \mathcal{C}_{1-\alpha}^{CP}(\mathbf{X}) \mid \mathbf{X} \in \mathcal{A}) = 1 - f_{\mathcal{A}}(\alpha), \quad (27)$$

we have

$$\mathbb{P}(y \in \mathcal{C}_{1-\alpha}^{PT}(\mathbf{X}) \mid \mathbf{X} \in \mathcal{A}) = p(1 - f_{\mathcal{A}}(\alpha')), \quad (28)$$

where $\alpha' = 1 - (1 - \alpha)/p$. And we have

$$\mathcal{F}(1) - \mathcal{F}(p) = f_{\mathcal{A}}(\alpha) - pf_{\mathcal{A}}\left(1 - \frac{1-\alpha}{p}\right) \geq 1 - p \quad (29)$$

$$\Rightarrow p\left(1 - f_{\mathcal{A}}\left(1 - \frac{1-\alpha}{p}\right)\right) \geq 1 - f_{\mathcal{A}}(\alpha) \quad (30)$$

$$\Rightarrow p(1 - \mathcal{A}(\alpha')) \geq 1 - f_{\mathcal{A}}(\alpha). \quad (31)$$

Therefore, we have

$$\mathbb{P}(y \in \mathcal{C}_{1-\alpha}^{PT}(\mathbf{X}) \mid \mathbf{X} \in \mathcal{A}) \geq \mathbb{P}(y \in \mathcal{C}_{1-\alpha}^{CP}(\mathbf{X}) \mid \mathbf{X} \in \mathcal{A}). \quad (32)$$

D.4 PROOF OF THEOREM 7

By the definition of \mathcal{L} and PT in Algorithm 1, when the returned set belongs to the meaningful fold, the length is $\mathcal{L}(x, (1 - \tilde{\alpha})/p; s)$, given the non-conformity score function and miscoverage rate $\tilde{\alpha}$. And the length of the returned set belongs to the vacuous fold is 0. Therefore, the expected length of the set returned by PT is

$$p\mathbb{E}\left(\mathcal{L}\left(x, \frac{1 - \tilde{\alpha}}{p}; s\right)\right), \quad (33)$$

where the expectation is taken over x . Then we get the general sufficient condition is

$$\exists p \in (1 - \tilde{\alpha}, 1) \quad s.t. \quad p\mathbb{E}\left(\mathcal{L}\left(x, \frac{1 - \tilde{\alpha}}{p}; s\right)\right) < \mathbb{E}(\mathcal{L}(x, 1 - \tilde{\alpha}; s)). \quad (34)$$

918 D.5 PROOF OF THEOREM 8
919

920 Let $\mathcal{G}(c) = \mathbb{E}(\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s))$, the sufficient condition in Theorem 7 is

$$921 \quad \exists p \in (c, 1), \text{ s.t. } p\mathcal{G}(c/p) < \mathcal{G}(c). \quad (35)$$

922 Note that when $p = 1$, $p\mathcal{G}(c/p) = \mathcal{G}(c)$, the sufficient condition of Eq (35) is

$$923 \quad \mathcal{F}'(1) > 0, \text{ where } \mathcal{F}(p) = p\mathcal{G}(c/p), \quad (36)$$

924 which is equivalent to

$$925 \quad \frac{\mathcal{G}(c)}{c} > \mathcal{G}'(c) \Rightarrow \mathbb{E} \left(\frac{\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s)}{1 - \tilde{\alpha}} \right) > \mathbb{E} \left(\frac{\partial}{\partial \alpha} \mathcal{L}(\mathbf{x}, \alpha; s) \Big|_{\alpha=1-\tilde{\alpha}} \right). \quad (37)$$

926 D.6 PROOF OF THEOREM 11
927

928 Use the notation in Appendix D.5, the general sufficient condition could be written as

$$929 \quad \exists p \in (c, 1), \text{ s.t. } p\mathcal{G}(c/p) < \mathcal{G}(c). \quad (38)$$

930 If there exists $u \in (1 - \tilde{\alpha})$ satisfies

$$931 \quad \frac{\mathcal{G}(c)}{c} > \frac{\mathcal{G}(u) - \mathcal{G}(c)}{u - c}, \quad (39)$$

932 we have

$$933 \quad \frac{c}{u} \mathcal{G}(u) < \mathcal{G}(c). \quad (40)$$

934 Let $p = c/u$, we have

$$935 \quad p\mathcal{G}(c/p) < \mathcal{G}(c). \quad (41)$$

936 Eq (39) is actually the condition

$$937 \quad \frac{\mathbb{E}(\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s))}{1 - \tilde{\alpha}} > \frac{\mathbb{E}(\mathcal{L}(\mathbf{x}, u; s)) - \mathbb{E}(\mathcal{L}(\mathbf{x}, 1 - \tilde{\alpha}; s))}{u - (1 - \tilde{\alpha})}. \quad (42)$$

938 D.7 PROOF OF EXAMPLE 12
939

940 When the non-conformity score follows a Gaussian distribution, the analytical solutions of the interval length returned by VCP and PT-VCP are

$$941 \quad |\mathcal{C}_{1-\tilde{\alpha}}^{VCP}(\mathbf{X}')| = 2\Phi^{-1} \left(1 - \frac{\tilde{\alpha}}{2} \right), \quad |\mathcal{C}_{1-\tilde{\alpha}}^{PT}(\mathbf{X}')| = 2p\Phi^{-1} \left(1 - \frac{1}{2} \left(1 - \frac{1 - \tilde{\alpha}}{p} \right) \right) \quad (43)$$

942 where $\Phi(\cdot)$ is the cumulative distribution function of the Gaussian distribution. Therefore, PT fails since Lemma 1.

943 **Lemma 1.** $\forall \alpha \in (0, 1), p \in (1 - \alpha, 1)$, there holds

$$944 \quad \Phi^{-1} \left(1 - \frac{\tilde{\alpha}}{2} \right) < p\Phi^{-1} \left(1 - \frac{1}{2} \left(1 - \frac{1 - \tilde{\alpha}}{p} \right) \right) \quad (44)$$

945 *Proof.* Using the symmetry identity $\Phi^{-1}(1 - u) = -\Phi^{-1}(u)$, the desired inequality is equivalent to

$$946 \quad p\Phi^{-1} \left(\frac{1}{2} \left(1 - \frac{1 - \alpha}{p} \right) \right) < \Phi^{-1}(\alpha/2). \quad (45)$$

947 Define

$$948 \quad u_0 := \alpha/2 \in (0, 1/2), \quad u(p) := \frac{1}{2} \left(1 - \frac{1 - \alpha}{p} \right) \in (0, 1/2). \quad (46)$$

949 Let $g(u) := \Phi^{-1}(u)$ on $(0, 1/2)$. Since

$$950 \quad g'(u) = \frac{1}{\phi(\Phi^{-1}(u))} > 0, \quad g''(u) = \frac{\Phi^{-1}(u)}{\phi(\Phi^{-1}(u))^2} < 0, \quad (47)$$

where $\phi(\cdot)$ denotes the p.d.f. of the standard normal distribution, the function g is increasing and strictly concave. Hence, for any $u \leq u_0$,

$$g(u) \leq g(u_0) + g'(u_0)(u - u_0). \quad (48)$$

Plugging $u = u(p)$ into Eq (48) and multiplying both sides by $p \in (0, 1)$ gives

$$p g(u(p)) \leq p g(u_0) + p g'(u_0)(u(p) - u_0). \quad (49)$$

A direct calculation yields

$$u(p) - u_0 = \frac{(1 - \alpha)(p - 1)}{2p} < 0. \quad (50)$$

Subtracting $g(u_0)$ from both sides of Eq (49) and using Eq (50), we obtain

$$p g(u(p)) - g(u_0) \leq (1 - p) \left[-g(u_0) + \frac{1 - \alpha}{2\phi(g(u_0))} \right]. \quad (51)$$

Here $1 - p > 0$. Moreover, since $g(u_0) = \Phi^{-1}(\alpha/2) < 0$ and $\phi(g(u_0)) > 0$, the bracket in Eq (51) is nonnegative. Thus the right-hand side of Eq (51) is nonpositive, implying

$$p g(u(p)) - g(u_0) < 0, \quad (52)$$

which is exactly inequality Eq (45). The inequality is strict because $u(p) \neq u_0$ and g is strictly concave.

Reverting to the upper-tail form via $\Phi^{-1}(1 - u) = -\Phi^{-1}(u)$ completes the proof:

$$\Phi^{-1}(1 - \alpha/2) < p \Phi^{-1}\left(\frac{1}{2}\left(1 + \frac{1 - \alpha}{p}\right)\right). \quad (53)$$

□

E MISSING ILLUSTRATION

In this section, we present the missing illustration of Example 1 in Section 1, the illustration of PT (Figure 4) and VCP algorithm mentioned in Section 3.1.

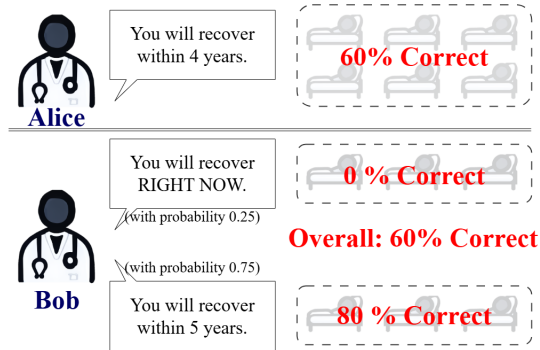


Figure 3: Illustration of Example 1. Doctor Alice and Bob both achieve 60% accuracy. Bob is more precise regarding length, but the corresponding strategy is not practically valid.

F OMITTED EXPERIMENTS

In this section, we present all the omitted experiments. In Appendix F.1, we demonstrate the missing experimental results in Section 3.3 and Section 3.4. In Appendix F.2, we exhibit the ablation study results.

Algorithm 2 Vanilla Conformal Prediction (VCP)

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

- 1: **Input:** miscoverage rate α , dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}\}$, test point \mathbf{x}' , non-conformity score function $s(\mathbf{x}_i, y_i; \hat{\mu})$.
- 2: Randomly split \mathcal{D} into a training fold $\mathcal{D}_{\text{tr}} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}_{\text{tr}}\}$ and a calibration fold $\mathcal{D}_{\text{ca}} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}_{\text{ca}}\}$;
- 3: Train a model $\hat{\mu}$ based on the training fold \mathcal{D}_{tr} ;
- 4: Calculate the non-conformity score on the calibration fold \mathcal{D}_{ca} , denoted by $\mathcal{V} = \{s(\mathbf{x}_i, y_i, \hat{\mu}) : i \in \mathcal{I}_{\text{ca}}\}$;
- 5: Compute the $(1 - \tilde{\alpha})$ -th quantile $\hat{Q}_{1-\tilde{\alpha}}(\mathcal{V})$ of the empirical distribution of the non-conformity score set \mathcal{V} calculated on the calibration set \mathcal{D}_{ca} , where $1 - \tilde{\alpha} = (1 - \alpha)(1 + 1/|\mathcal{V}|)$;
- 6: **Output:** Interval $\mathcal{C}_{1-\alpha}(\mathbf{x}') = \{y : s(\mathbf{x}', y; \hat{\mu}) \leq \hat{Q}_{1-\tilde{\alpha}}(\mathcal{V})\}$.

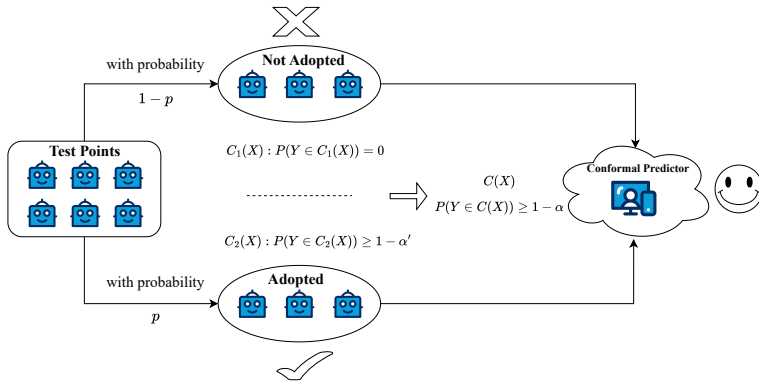


Figure 4: The illustration of Prejudicial Trick (PT). To obtain a $1 - \alpha$ confidence interval, PT first assigns empty sets for a $1 - p$ subset of the test points, and assigns $1 - \alpha'$ confidence interval for the remaining test points where $\alpha' < \alpha$. The returned confidence interval still satisfies $\mathbb{P}(Y \in \mathcal{C}(X)) \geq 1 - \alpha$ by setting a proper α' .

F.1 OMITTED EXPERIMENTAL RESULTS

Classification Tasks. We extend PT to classification tasks. We apply PT to the real-world IMAGENET-VAL dataset (Deng et al., 2009) with several pre-trained models, a similar setting with Angelopoulos et al. (2020). To simulate model misspecification, a bias is introduced to the logits of several classes before the softmax operation. The magnitude of the bias is determined based on the scale of the outputs. The experimental results on classification tasks in Table 4 perform similarly to regression tasks. Specifically, PT-VCP attains valid coverage across different models (Theorem 4) while improving the length compared to VCP (Theorem 11).

Conformalized Quantile Regression. We deploy PT into other variants of conformal prediction. Specifically, we choose CQR as a baseline (Romano et al., 2019). CQR inherits the advantages of both conformal prediction and classical quantile regression. We use the same datasets and evaluation metrics as in Section 3.4. To mimic the model misspecification, we add bias directly to the lower and upper quantiles obtained by the quantile regression. The experimental results on the real-world CQR tasks are exhibited in Table 5. It illustrates that *PT achieves shorter interval length while maintaining valid coverage on CQR.*

Group Coverage. In Section 3.3, we conduct experiments to evaluate the different performance of VCP and PT-VCP regarding group coverage. The experimental results shown in Table 6 demonstrate that PT not only achieves shorter confidence intervals while maintaining overall coverage, *but also improves the group coverage in regression tasks*⁴.

⁴Group coverage is defined as the lowest coverage rate among all the groups.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

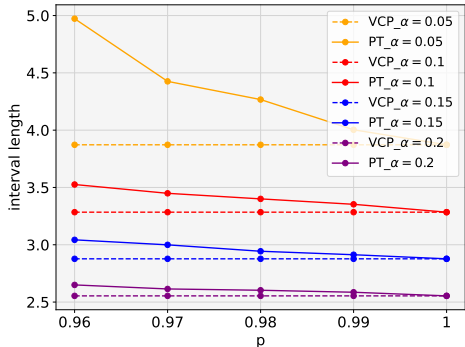


Figure 5: PT fails to improve length when the conditions on the distribution of non-conformity score are not satisfied.

Table 4: Comparison between VCP and PT-VCP in classification tasks across different models. Experiments on RAPS ($\alpha = 0.1, p = 0.95$), with index range chosen as 300.

METHOD	BIAS	VCP		PT-VCP	
		COVERAGE	LENGTH	COVERAGE	LENGTH
RESNET18	40	0.90 ±0.000	304.02 ±0.004	0.90 ±0.000	295.60 ±0.233
RESNET50	40	0.90 ±0.000	302.09 ±0.027	0.90 ±0.000	290.29 ±0.224
RESNET101	40	0.90 ±0.000	302.01 ±0.004	0.90 ±0.000	289.56 ±0.174
RESNET152	40	0.89 ±0.004	301.53 ±0.054	0.90 ±0.000	288.98 ±0.165
RESNEXT101	40	0.90 ±0.000	301.48 ±0.013	0.90 ±0.000	288.49 ±0.201
VGG16	40	0.90 ±0.004	303.39 ±0.143	0.90 ±0.000	293.34 ±0.304
SHUFFLENET	40	0.90 ±0.000	304.05 ±0.040	0.90 ±0.000	295.79 ±0.282
INCEPTION	40	0.90 ±0.000	304.10 ±0.013	0.90 ±0.000	297.25 ±0.228
DENSENET161	40	0.90 ±0.000	302.03 ±0.009	0.90 ±0.000	289.29 ±0.197

Interval Stability. In Section 4, we introduce a new evaluation criterion, termed *interval stability* and conduct several empirical evaluations using the datasets described in Section 3.4, the results of which are listed in Table 3. We further investigate the performance of the interval stability metric using CQR as the base algorithm in Table 7. The results are similar to the results in Table 3. We also evaluate the interval stability metric on classification tasks. As shown in Table 8, interval stability successfully identify the vacuous randomness in PT.

F.2 ABLATION STUDIES

This section exhibits the ablation studies on the probability hyperparameter p in PT and the bias parameter μ on different base algorithms (Figure 6-Figure 15). All the experiments are conducted based on various miscoverage rates α . The experiment results demonstrate that, although not all the probability hyperparameters p outperform the base algorithm, our goal is to show that *there exist multiple (at least one) probability hyperparameters such that PT-VCP outperforms VCP, which suffices to challenge the coverage-length gold standard.* Furthermore, we find that the bias parameter actually matters here, implying that PT-VCP performs better than VCP under misspecification, which validates Theorem 11.

G EXPERIMENT DETAILS

In this section, we provide implementation details of the experiments in this paper, including experiments on synthetic datasets in Appendix G.1 and experiments on real-world datasets in Appendix G.2.

Table 5: Comparison between CQR and PT-CQR in quantile regression task across different datasets.

Method	Bias	CQR		PT-CQR	
Dataset		Coverage	Length	Coverage	Length
meps-19	1	0.91 ± 0.000	4.60 ± 0.148	0.91 ± 0.246	4.44 ± 0.143
meps-20	1	0.91 ± 0.000	4.58 ± 0.192	0.91 ± 0.179	4.41 ± 0.188
meps-21	1	0.91 ± 0.000	4.65 ± 0.080	0.91 ± 0.161	4.52 ± 0.107
bike	1	0.91 ± 0.000	2.61 ± 0.013	0.90 ± 0.268	2.51 ± 0.009
blog-data	1	0.91 ± 0.000	3.80 ± 0.107	0.93 ± 0.116	3.61 ± 0.098
bio	1	0.91 ± 0.000	3.45 ± 0.009	0.90 ± 0.112	3.32 ± 0.009
facebook-1	1	0.91 ± 0.000	3.38 ± 0.022	0.92 ± 0.125	3.22 ± 0.027
facebook-2	1	0.91 ± 0.000	3.57 ± 0.027	0.92 ± 0.085	3.39 ± 0.027
concrete	2	0.91 ± 0.000	4.39 ± 0.022	0.88 ± 0.648	4.23 ± 0.018
star	2	0.91 ± 0.000	4.15 ± 0.004	0.90 ± 0.349	3.96 ± 0.009

Table 6: Comparison of group coverage between VCP and PT-VCP on regression tasks across different datasets ($\alpha = 0.1$).

Dataset	Group	VCP	PT-VCP
bike	Day	0.878 ± 0.007	0.884 ± 0.010
	Month	0.826 ± 0.010	0.857 ± 0.011
	Year	0.851 ± 0.005	0.871 ± 0.004
star	Gender	0.905 ± 0.008	0.905 ± 0.002
	Stark	0.890 ± 0.005	0.895 ± 0.007
	School1	0.902 ± 0.008	0.899 ± 0.022
meps-19	SEX=1	0.883 ± 0.004	0.895 ± 0.001
	MARRY=1	0.901 ± 0.004	0.901 ± 0.003
	REGION=1	0.862 ± 0.005	0.877 ± 0.006
meps-20	FTSTU=1	0.893 ± 0.004	0.900 ± 0.002
	ACTDTY=1	0.897 ± 0.003	0.902 ± 0.002
	HONRDC=1	0.792 ± 0.010	0.846 ± 0.008
meps-21	RTHLTH=1	0.864 ± 0.004	0.877 ± 0.004
	MNHLTH=1	0.856 ± 0.004	0.873 ± 0.004
	HIBPDX=1	0.755 ± 0.013	0.818 ± 0.009

G.1 SYNTHETIC DATASETS

This section presents experiment details about the motivating example (Section 3.2) in Appendix G.1.1 and failure case (Section 3.4) in Appendix G.1.2.

G.1.1 MOTIVATING EXAMPLE

In our motivating example, we consider a simple data-generating process where the true underlying model is linear with Gaussian mixture noise:

$$Y = \mathbf{X}^\top \boldsymbol{\beta} + \epsilon, \quad \mathbf{X} \sim \mathcal{N}(\mathbf{0}, I_2).$$

The noise term ϵ follows $\mathcal{N}(\mu, 1)$ with probability 0.5 and $\mathcal{N}(-\mu, 1)$ with probability 0.5. The training, calibration, and test folds are all generated from this distribution. To emulate model misspecification, we fit the training fold using a linear model with Gaussian noise. Throughout the experiments, we set $\mu = 20$, $\alpha \in \{0.1, 0.2\}$, and $p \in \{0.96, 0.98\}$. We further average results over 5 random seeds and report the corresponding standard errors. Both VCP (Algorithm 2) and PT-VCP are evaluated under this setting.

Table 7: Comparison between CQR and PT-CQR regarding interval stability.

Dataset	CQR	PT-CQR
meps-19	0.00 ± 0.000	0.13 ± 0.005
meps-20	0.00 ± 0.000	0.13 ± 0.006
meps-21	0.00 ± 0.000	0.14 ± 0.003
bike	0.00 ± 0.000	0.08 ± 0.001
blog-data	0.00 ± 0.000	0.11 ± 0.003
bio	0.00 ± 0.000	0.10 ± 0.000
facebook-1	0.00 ± 0.000	0.10 ± 0.001
facebook-2	0.00 ± 0.000	0.10 ± 0.001
concrete	0.00 ± 0.000	0.13 ± 0.002
star	0.00 ± 0.000	0.12 ± 0.001

Table 8: Comparison between VCP and PT-VCP in classification tasks regarding interval stability. Experiments on RAPS ($\alpha = 0.1, p = 0.95$), with index range chosen as 300.

MODEL	BIAS	VCP	PT-VCP
RESNET18	40	0.02 ± 0.004	12.08 ± 0.060
RESNET50	40	0.07 ± 0.017	11.86 ± 0.057
RESNET101	40	0.01 ± 0.004	11.93 ± 0.089
RESNET152	40	0.19 ± 0.002	11.91 ± 0.058
RESNEXT101	40	0.20 ± 0.000	11.89 ± 0.083
VGG16	40	0.11 ± 0.030	12.00 ± 0.060
SHUFFLENET	40	0.03 ± 0.025	12.08 ± 0.055
INCEPTION	40	0.07 ± 0.010	12.19 ± 0.080
DENSENET161	40	0.03 ± 0.006	11.90 ± 0.057

G.1.2 FAILURE CASE

In Figure 5, we present a failure case where VCP outperforms PT-VCP. Here, we modify the data-generating process to a linear model with Gaussian noise and fit it with the same linear model, so that no model misspecification arises in contrast to our motivating example. In this setting, the distribution of the score function fails to satisfy the sufficient condition in Theorem 11, which explains why VCP outperforms PT-VCP.

G.2 REAL WORLD DATASETS

In this section, we firstly introduce the model structure in our experiments in Appendix G.2.1. Then we present the experiment details, including experiments on marginal and group coverage (Appendix G.2.2, Appendix G.2.6), regression tasks (Appendix G.2.3), classification tasks (Appendix G.2.4), ablation studies (Appendix G.2.5) and interval stability (Appendix G.2.7).

G.2.1 MODEL STRUCTURE

In this section, we present the details of the structure of our model on real world datasets. Specifically, our model shares the same structure with (Romano et al., 2019).

Neural Net. Our neural network design includes three fully connected layers, with ReLU activation functions applied between each layer. The initial layer accepts an input feature vector X of n dimensions and produces 64 hidden units. The second layer mirrors this structure, generating another set of 64 hidden units. The final layer is a linear output layer that provides a pointwise prediction for the response variable Y . The network’s parameters are optimized by minimizing a quadratic loss function. We used the Adam optimization algorithm with a constant learning rate of 5×10^{-4} , minibatch size of 64, and a weight decay coefficient of 10^{-6} . In addition, regularization of dropout is implemented, with a retention probability of 0.1 for hidden units. To avoid overfitting, early stop is used and the number of training epochs is determined by cross-validation, with a maximum cap of 1000 epochs.

CQR Neural Net. We utilize neural networks to implement CQR for quantile regression. The net-

1242
1243
1244
1245
1246
1247
1248
1249

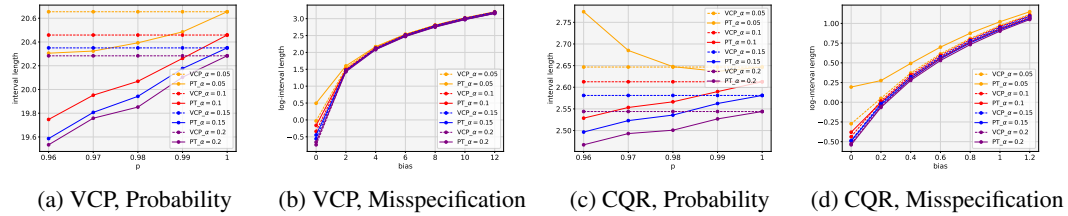


Figure 6: Ablation studies of dataset BIKE on different misspecification levels (b, d) and probability hyperparameters (a, c), including comparisons with VCP (a–b) and CQR (c–d).

1250
1251
1252
1253
1254
1255
1256
1257
1258
1259

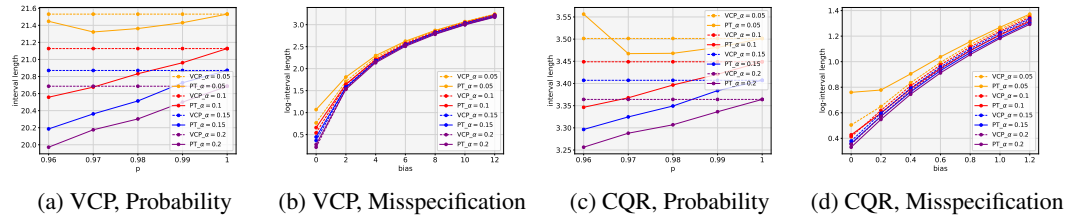


Figure 7: Ablation studies of dataset BIO on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271

work structure is consistent with the one described above, with the sole difference being that the output of the quantile regression network is a two-dimensional vector, which indicates the lower and upper conditional quantiles. Additionally, the training process remains the same, except that the pinball loss function in equation is employed instead of the quadratic loss.

G.2.2 MARGINAL AND GROUP COVERAGE

In Figure 1, we compare the coverage of VCP and PT-VCP on the BIKE dataset (Fanaee-T, 2013). Specifically, we evaluate several choices of α and p , and report both the marginal coverage and the group coverage (an empirical indicator of conditional coverage). The group coverage is obtained by partitioning the data according to Day, Month, and Year.

1272
1273
1274
1275
1276
1277
1278
1279

G.2.3 REGRESSION TASKS

1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

In ordinary regression tasks, we employ the neural network described in Section G.2.1 to fit several real-world datasets: MEPS19–21 (Cohen et al., 2009), BIKE (Fanaee-T, 2013), BLOG-DATA (Buza, 2014), BIO (Rana, 2013), FACEBOOK1–2 (Singh, 2015), CONCRETE (Yeh, 1998), and STAR (Achilles et al., 2008). To mimic model misspecification, we introduce a bias term that is directly added to the logits output by the neural network. The magnitude of this bias term, which varies across datasets, is reported in Table 2. Throughout these experiments, we set $\alpha = 0.1$ and $p = 0.95$, under which both VCP (Algorithm 2) and PT-VCP are evaluated. We further average results over 5 random seeds and report the corresponding standard errors.

In CQR tasks, we employ the CQR neural network described in Section G.2.1 to fit several real-world datasets: MEPS19–21 (Cohen et al., 2009), BIKE (Fanaee-T, 2013), BLOG-DATA (Buza, 2014), BIO (Rana, 2013), FACEBOOK1–2 (Singh, 2015), CONCRETE (Yeh, 1998), and STAR (Achilles et al., 2008). To mimic model misspecification, we introduce a bias term by directly adding it to both the lower and upper quantiles estimated by the quantile regression. The magnitude of this bias varies across datasets. Throughout the experiments, we set $\alpha = 1$ and $p = 0.95$, under which both CQR and PT-CQR are evaluated, as reported in Table 5. We further average results over 5 random seeds and report the corresponding standard errors.

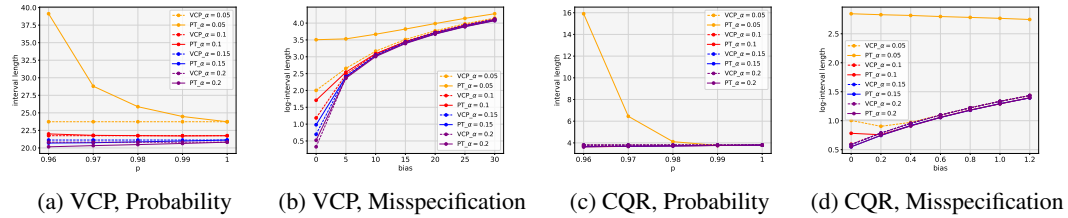


Figure 8: Ablation studies of dataset BLOGDATA on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

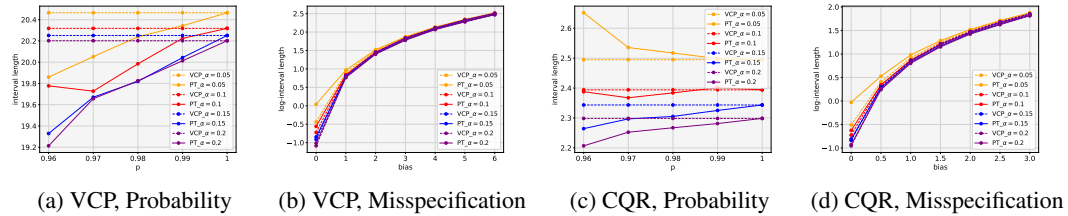


Figure 9: Ablation studies of dataset CONCRETE on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

G.2.4 CLASSIFICATION TASK

In classification tasks, we apply PT to the real-world IMAGENET-VAL dataset (Deng et al., 2009) using several pre-trained models listed in Table 4, following a setting similar to Angelopoulos et al. (2020). To simulate model misspecification, we introduce a bias to the logits of several classes before the softmax operation. The magnitude of this bias is scaled according to the outputs, and the number of biased classes is specified by an index range in our experiments. For the classification task, we adopt RAPS (Angelopoulos et al., 2020) as the score function and set $\alpha = 0.1$ and $p = 0.95$. We further average results over 5 random seeds and report the corresponding standard errors.

G.2.5 ABLATION STUDIES

The ablation studies mainly focus on regression tasks, including both ordinary regression and CQR. We conduct experiments on all datasets used in the regression setting. For these ablation studies, we set $\alpha \in \{0.05, 0.1, 0.15, 0.2\}$, $p \in \{0.96, 0.97, 0.98, 0.99, 1.00\}$, and apply dataset-specific bias magnitudes.

G.2.6 GROUP COVERAGE

To demonstrate that PT-VCP does not degrade conditional coverage compared to VCP, we conduct experiments measuring group coverage on the MEPS19–21 (Cohen et al., 2009), BIKE (Fanace-T, 2013), and STAR (Achilles et al., 2008) datasets. The grouping strategies are summarized in Table 6. In these experiments, we set $\alpha = 0.1$ and $p = 0.95$, and further average results over 5 random seeds, reporting the corresponding standard errors.

G.2.7 INTERVAL STABILITY

To compare the newly proposed *interval stability* between VCP and PT-VCP, we conduct experiments on both regression and classification tasks, using different datasets and pre-trained models, respectively. The results are reported in Table 3, Table 7, and Table 8. We measure *interval stability* by computing the variance of the interval (or set) length (or size) for repeated predictions on the same input, and then averaging this variance over all inputs X . Results are further averaged over 5 random seeds, with the corresponding standard errors reported.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

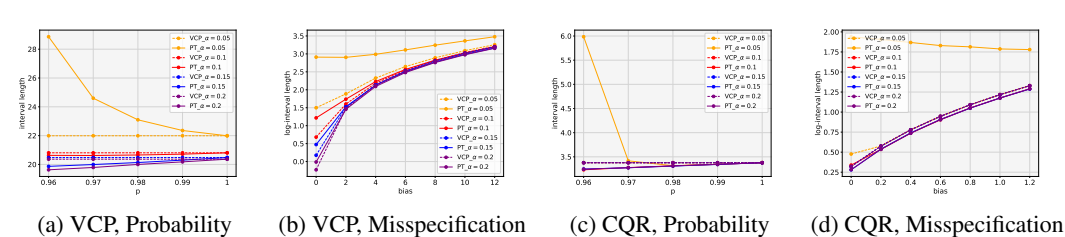


Figure 10: Ablation studies of dataset FACEBOOK1 on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

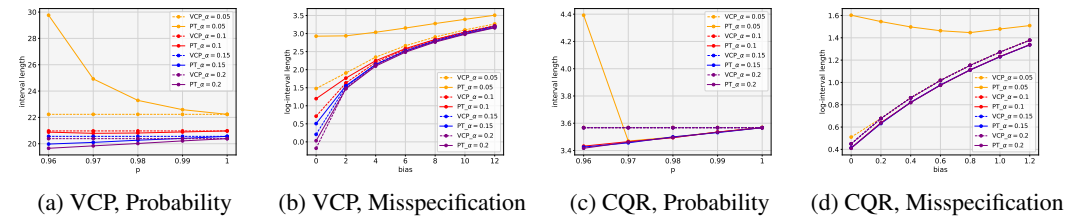


Figure 11: Ablation studies of dataset FACEBOOK2 on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

H ADDITIONAL RELATED WORKS

Conformal prediction. Conformal prediction is a post hoc calibration framework that constructs statistically rigorous uncertainty sets for predictions from machine learning models (Vovk et al., 2005; Shafer & Vovk, 2008; Lei et al., 2018; Foygel Barber et al., 2021; Angelopoulos & Bates, 2021; Papadopoulos et al., 2008). Traditionally, vanilla conformal prediction is deployed in regression tasks (Vovk et al., 2005; Shafer & Vovk, 2008; Lei et al., 2018). Later, a branch of research expands vanilla conformal prediction to diverse data structures and applications, including classification tasks (Angelopoulos et al., 2020; Dabah & Tirer, 2024), censored data in survival analysis (Teng et al., 2021; Candès et al., 2023), functional data (Lei et al., 2015; Ajroldi et al., 2023), graph-based models (Zargarbashi et al., 2023; Zargarbashi & Bojchevski, 2024), time series data (Xu & Xie, 2021; Stankeviciute et al., 2021), treatment effects (Lei & Candès, 2021; Jin et al., 2023), *etc.*

Interval regression. While coverage guarantees and interval length serve as fundamental metrics for evaluating conformal prediction (Vovk et al., 2005; Lei et al., 2018; Barber et al., 2020), these criteria are deeply entrenched in the broader paradigm of interval regression methodologies. Established approaches including quantile regression (Alaa et al., 2023; Sasaki et al., 2022) and Bayesian credible intervals (Kuleshov et al., 2018; Wang & Ghosal, 2023) similarly prioritize the dual metrics. Of particular relevance is Navratil et al. (2020) who proposes an excess and deficit metrics beyond the traditional coverage-length metric. Our paper differs from Navratil et al. (2020) in that our main contributions center on uncovering the inherent limitations of coverage-length metrics. Additionally, we contend that the proposed excess and deficit metrics cannot be directly applied to PT-VCP.

I THE USE OF LARGE LANGUAGE MODELS (LLMs)

In preparing this paper, a large language model (LLM) was employed solely for language refinement purposes, such as improving the clarity and fluency of expressions. The LLM did not contribute to research ideation, methodology, data analysis, or substantive content generation. The authors fully acknowledge responsibility for all contents of the paper, including any text polished with the assistance of the LLM.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

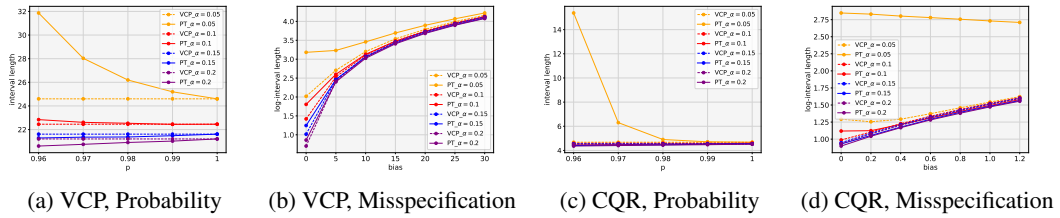


Figure 12: Ablation studies of dataset MEPS19 on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

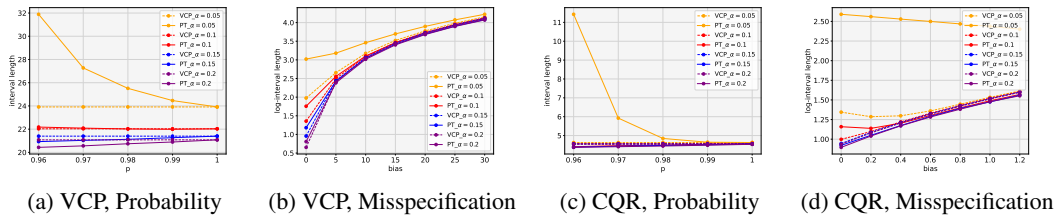


Figure 13: Ablation studies of dataset MEPS20 on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

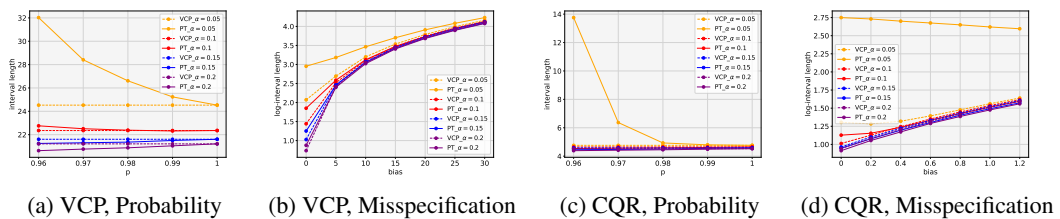


Figure 14: Ablation studies of dataset MEPS21 on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).

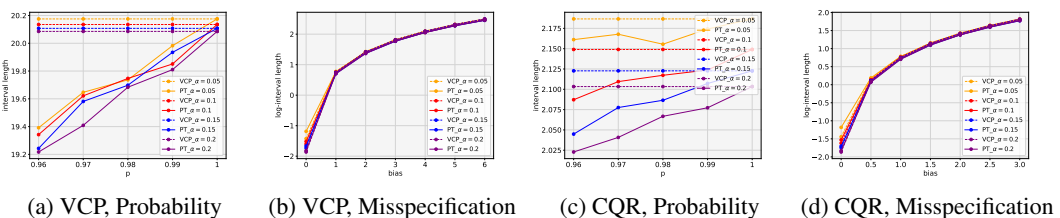


Figure 15: Ablation studies of dataset STAR on different misspecification level (a, c) and probability hyperparameter (b, d), including the comparison with VCP (a-b) and CQR (c-d).