

CONDITIONAL FLOW MATCHING FOR CONFORMAL REGRESSION

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper introduces Conditional flow Matching for conformal Regression (CMR), a novel framework that synergizes simulation-free conditional flow matching with conformal prediction to generate reliable and efficient prediction intervals. Unlike traditional methods that rely on quantile regression or fixed histograms, CMR leverages Continuous Normalizing Flows (CNFs) trained via Conditional Flow Matching (CFM) to accurately model complex, multimodal conditional distributions. To ensure finite-sample coverage guarantees, we introduce a novel nonconformity score defined as the minimum number of generated samples required for the shortest interval to encompass the true outcome. This mechanism allows CMR to dynamically adjust interval widths based on the learned probability density. Extensive experiments on simulated and real-world datasets demonstrate that CMR consistently produces narrower prediction intervals while maintaining the required marginal coverage and achieving superior tail coverage compared to state-of-the-art methods.

1 INTRODUCTION

In recent years, machine learning and deep learning models (Magdon-Ismael & Atiya, 1998; Meinshausen & Ridgeway, 2006) have achieved remarkable success in regression tasks and beyond. However, while point predictions from these models continue to advance, they often fall short of industrial standards due to unresolved challenges in model trustworthiness and robustness. Conformal prediction (Vovk et al., 2005; Shafer & Vovk, 2008) emerges as a critical framework to bridge this gap: by transforming point estimates into prediction intervals with guaranteed coverage probabilities, it provides statistically valid uncertainty quantification without relying on specific model architectures. Rooted in the principle of data exchangeability, this approach constructs confidence intervals aligned with predefined confidence levels, offering a model-agnostic solution to enhance reliability in real-world deployments.

A key challenge in conformal regression lies in minimizing prediction set sizes while maintaining rigorous coverage guarantees. Most existing conformal methods generally fall into two primary categories: (1) those that derive prediction intervals by leveraging model-based outputs (e.g., decision trees, random forests, or neural networks) to directly estimate interval bounds (Chipman et al., 2010; Papadopoulos et al., 2011; Kivaranovic et al., 2020; Moon et al., 2021); and (2) those that construct prediction sets by inverting conditional density estimates derived from the data (Izbicki et al., 2019; Diamant et al., 2024; Plassier et al., 2024; Zheng & Zhu, 2024). To further enhance performance, numerous approaches incorporate residual-based calibration (Chen et al., 2018b; Lei et al., 2018; Barber et al., 2021) or localized adaptation strategies (Guan, 2023; Colombo, 2024; Cheung et al., 2024; Gil et al., 2024; Hore & Barber, 2025). However, these methods often struggle with complex data distributions (e.g., bimodal distributions), producing excessively large prediction sets that introduce significant uncertainty. Furthermore, certain techniques involve overly intricate procedures, posing practical challenges for implementation. To address these issues, a subset of research (Izbicki et al., 2022; Luo & Zhou, 2025; Gao et al., 2025) has focused on handling data from multiple distributions. Nevertheless, the intervals generated by these

054 methods are typically not continuous (i.e., unions of disjoint intervals), which significantly
 055 limits their applicability in real-world scenarios.

056
 057 To address the challenges posed by complex distributions, such as excessively wide predic-
 058 tion intervals and operational complexities, we propose a novel and versatile framework:
 059 Conditional Flow Matching for Conformal Regression (CMR). By leveraging Conditional
 060 Flow Matching, which is theoretically grounded, we learn and generate the target distri-
 061 bution. We then combine this with a new nonconformity score to produce the smallest
 062 continuous prediction intervals. Our primary contributions are:

- 063 • We introduce a framework that explicitly models the conditional distribution using
 064 conditional flow matching and leverages this model to construct prediction intervals
 065 with guaranteed coverage.
- 066 • We develop a new nonconformity score based on identifying the shortest interval
 067 containing the true value among sampled predictions, significantly reducing interval
 068 sizes compared to existing methods.
- 069 • We demonstrate through rigorous theoretical analysis that our approach provides
 070 valid coverage guarantees regardless of model accuracy.
- 071 • We empirically validate our method across various simulated and real-world
 072 datasets, showing that CMR consistently achieves smaller prediction intervals while
 073 maintaining coverage requirements and delivering superior tail coverage perfor-
 074 mance.
- 075 • We decoupled CMR. It outperforms other methods by generating smaller prediction
 076 sets. Using Conditional Flow Matching as the base and combining it with other
 077 Conformal Prediction approaches yields even smaller sets, with stable conditional
 078 coverage and good tail coverage.

080 2 RELATED WORK

082 2.1 CONFORMAL PREDICTION

083 Conformal prediction (Vovk et al., 2005; Kivaranovic et al., 2020; Romano et al., 2019;
 084 Sesia & Candès, 2020; Sesia & Romano, 2021; Wang et al., 2023; Kiyani et al., 2024; Luo
 085 & Zhou, 2025; Liu et al., 2025) is a framework for constructing prediction intervals with
 086 finite-sample coverage guarantees. Given a desired miscoverage level $\alpha \in (0, 1)$, conformal
 087 prediction produces prediction regions that contain the true outcome with probability at
 088 least $1 - \alpha$.

089 Let $\mathcal{D} = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ be a dataset of i.i.d. examples from an unknown distri-
 090 bution P_{XY} on $\mathcal{X} \times \mathcal{Y}$. We begin by randomly partitioning the dataset \mathcal{D} into three sets:
 091 $\mathcal{I}_{\text{train}}$, \mathcal{I}_{cal} , and $\mathcal{I}_{\text{test}}$. A base predictive model $\hat{f}(x_i)$ is trained on $\mathcal{I}_{\text{train}}$. Subsequently, we
 092 compute a non-conformity score S_i (e.g., absolute residual $|y_i - \hat{f}(x_i)|$) on the calibration set
 093 \mathcal{I}_{cal} . Following the Split Conformal Regression framework (Papadopoulos et al., 2002; Vovk
 094 et al., 2005), the $1 - \alpha$ prediction interval for a novel test sample $x_{n+1} \in \mathcal{I}_{\text{test}}$ is constructed
 095 as:

$$096 \mathcal{C}_{n,\alpha}(x_{n+1}) = \{y \in \mathbb{R} : s(x_{n+1}, y) \leq \hat{q}_{1-\alpha}\}, \quad (1)$$

097 where $\hat{q}_{1-\alpha}$ is the $\lceil (1 - \alpha)(|\mathcal{I}_{\text{cal}}| + 1) \rceil / |\mathcal{I}_{\text{cal}}|$ -th empirical quantile of the calibration scores.
 098 Under the exchangeability assumption, this guarantees:

$$099 \mathbb{P}(Y_{n+1} \in \mathcal{C}_{n,\alpha}(x_{n+1})) \geq 1 - \alpha. \quad (2)$$

100 *Conformal Histogram Regression* (Sesia & Romano, 2021) constructs prediction intervals
 101 by estimating the full conditional density $\hat{f}_{Y|X}$ using histograms and finding the shortest
 102 interval (a, b) such that:

$$103 \mathcal{C}^{\text{CHR}}(x_{n+1}) = \arg \min_{a < b} (b - a), \quad (3)$$

$$104 \text{s.t. } \int_a^b \hat{f}_{Y|X}(y|x_{n+1}) dy \geq 1 - \alpha. \quad (4)$$

2.2 CONTINUOUS NORMALIZING FLOWS

Continuous Normalizing Flows (CNFs) (Chen et al., 2018a; Grathwohl et al., 2018) are a class of generative models that define a probability path from a simple source distribution p_0 to a complex target distribution p_1 via an Ordinary Differential Equation (ODE). The transformation is defined by:

$$\frac{dy}{dt} = v_\theta(t, y, x), \quad y(0) = y_0 \sim p_0, \quad y(1) = y_1 \sim p_1, \quad (5)$$

where $v_\theta(t, y, x)$ is a time-dependent vector field parameterized by a neural network.

2.3 CONDITIONAL FLOW MATCHING

Conditional Flow Matching (CFM) (Lipman et al., 2023; Tong et al., 2023) is a simulation-free training objective that regresses the neural vector field v_θ directly to a target vector field u_t .

Consider a conditional probability path $p_t(y|y_0, y_1)$ that interpolates between a source sample $y_0 \sim p_0$ and a target sample $y_1 \sim p_1$. If $u_t(y|y_0, y_1)$ is the vector field generating this path, the CFM objective is defined as:

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, p_0(y_0), p_1(y_1), p_t(y|y_0, y_1)} [\|v_\theta(t, y, x) - u_t(y|y_0, y_1)\|^2]. \quad (6)$$

By minimizing this objective, v_θ approximates the marginal vector field that transports p_0 to p_1 .

3 CONDITIONAL FLOW MATCHING FOR CONFORMAL REGRESSION

In this section, we elaborate on our proposed methodology, **Conditional flow Matching for conformal Regression (CMR)**. Our approach consists of three steps: (1) training a conditional generative model using CFM to estimate $p(y|x)$; (2) calibrating prediction intervals using a novel “**shortest-interval**” non-conformity score; and (3) constructing valid prediction intervals for new test points.

3.1 LEARNING THE CONDITIONAL DISTRIBUTION VIA CFM

We employ a CNF to model the conditional distribution $p(y|x)$. To ensure efficient training, we adopt the **Conditional Flow Matching (CFM)** formulation.

We define the probability path p_t as a linear interpolation between the source noise $y_0 \sim p_0(y_0) = \mathcal{N}(0, 1)$ and the data $y_1 \sim p(y|x)$, smoothed by a small noise σ :

$$p_t(y|y_0, y_1) = \mathcal{N}(y | ty_1 + (1-t)y_0, \sigma^2). \quad (7)$$

The unique vector field u_t that generates this Gaussian probability path is the constant velocity field pointing from y_0 to y_1 :

$$u_t(y|y_0, y_1) = y_1 - y_0. \quad (8)$$

Consequently, the specific regression loss for our model becomes:

$$\mathcal{L}_{\text{CMR}}(\theta) = \mathbb{E}_{\substack{t \sim \mathcal{U}(0,1), \\ x \sim p(x), y_1 \sim p(y|x), \\ y_0 \sim p_0(y_0), \\ y_t \sim p_t(y|y_0, y_1)}} [\|v_\theta(t, y_t, x) - (y_1 - y_0)\|^2]. \quad (9)$$

This objective allows us to learn the conditional distribution $p(y|x)$ without ODE simulation.

3.2 SAMPLING FROM THE LEARNED MODEL

Once trained, the network v_θ approximates the vector field. To sample from the learned conditional distribution $\hat{p}(y|x)$, we sample $z \sim \mathcal{N}(0, 1)$ and solve the ODE from $t = 0$ to $t = 1$:

$$\hat{y} = \text{ODESolve}(v_\theta, z, t \in [0, 1], x). \quad (10)$$

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

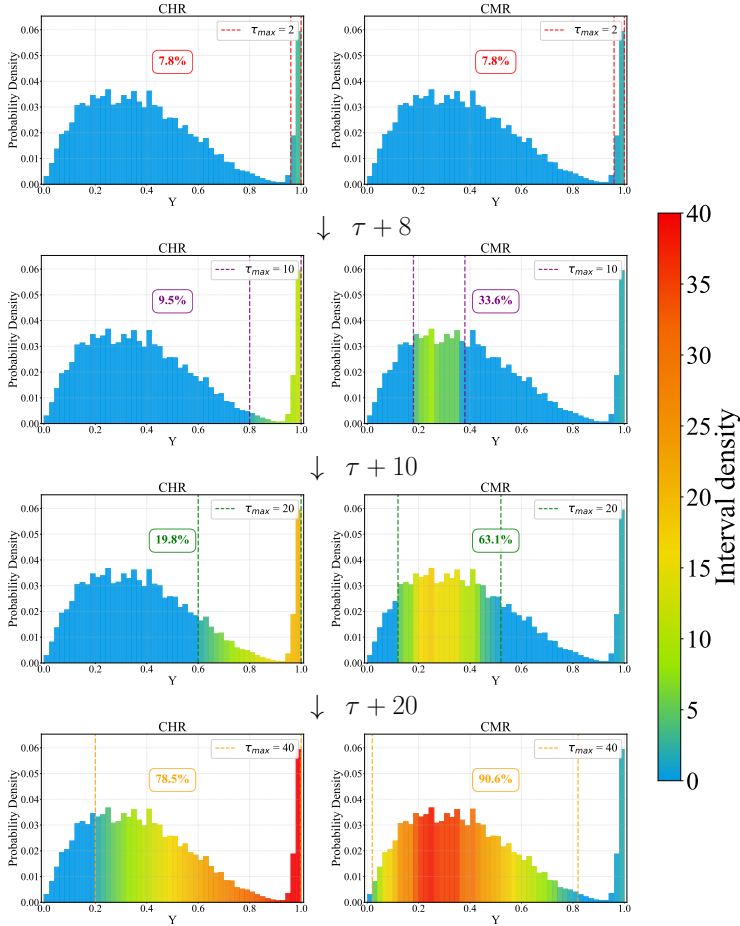


Figure 1: To demonstrate the advantages of our method in handling complex distributions, we visualize the performance of CMR versus CHR in a bimodal setting. The horizontal axis represents the response space, and the vertical-axis denotes the probability density. CMR dynamically selects the contiguous interval that maximizes the total probability density. In contrast, CHR employs an expansion-based strategy, creating a contiguous interval by extending outwards from the highest-density bin. (As shown in the bimodal distribution, if the initial starting bin does not align with the densest aggregate region, CHR requires a significantly larger number of bins to achieve the equivalent probability mass). Consequently, for a fixed number of bins, CMR captures a higher proportion of the total probability mass.

3.3 CALIBRATION PROCEDURE

To construct valid prediction intervals, we introduce a non-conformity score based on the *tightness* of the generated distribution. We use a calibration set $\mathcal{I}_{\text{cal}} = \{(x_i, y_i)\}_{i=1}^{|\mathcal{I}_{\text{cal}}|}$. For each calibration point (x_i, y_i) :

1. **Generate Samples:** Generate m samples $\{y_i^1, \dots, y_i^m\}$ from the trained model conditioned on x_i via Eq. (10).
2. **Sort:** Sort these samples to obtain the order statistics $\{y_i^{(1)} \leq y_i^{(2)} \leq \dots \leq y_i^{(m)}\}$.
3. **Find Shortest Intervals:** For each possible number of included samples $s \in \{2, \dots, m\}$, we identify the shortest continuous interval that contains exactly s generated samples. Let $W_i(s)$ be this minimum width interval:

$$W_i(s) = [y_i^{(j^*)}, y_i^{(j^*+s-1)}], \quad \text{where } j^* = \arg \min_{1 \leq j \leq m-s+1} (y_i^{(j+s-1)} - y_i^{(j)}). \quad (11)$$

- 216 4. **Compute Score:** We define the non-conformity score L_i as the **minimum num-**
 217 **ber of samples s** required for the corresponding shortest interval $W_i(s)$ to cover
 218 the true outcome y_i^{true} :
 219

$$220 L_i = \min\{s \in \{2, \dots, m\} \mid y_i^{true} \in W_i(s)\}. \quad (12)$$

221 Finally, we compute the quantile threshold $\hat{q}_{1-\alpha}$ from the calibration scores $\{L_1, \dots, L_{|\mathcal{I}_{cal}}|\}$:
 222

$$223 \hat{q}_{1-\alpha} = \text{Quantile}_{\lceil(1-\alpha)(|\mathcal{I}_{cal}|+1)\rceil/|\mathcal{I}_{cal}|}(\{L_1, \dots, L_{|\mathcal{I}_{cal}}|\}). \quad (13)$$

224 This threshold represents the minimum ‘‘sample count’’ required to ensure coverage.
 225

226 **Theorem 1** (Marginal Coverage Guarantee). *Let (x_i, y_i) be calibration data and*
 227 *(x_{n+1}, y_{n+1}) be a test point, i.i.d. drawn from P . Let L_i be the score defined above. The*
 228 *prediction set:*

$$229 \mathcal{C}^{\text{CMR}}(x_{n+1}) = \{y \in \mathbb{R} \mid L(y; x_{n+1}) \leq \hat{q}_{1-\alpha}\}$$

230 satisfies $\mathbb{P}(y_{n+1} \in \mathcal{C}^{\text{CMR}}(x_{n+1})) \geq 1 - \alpha$.
 231

232 3.4 PREDICTION INTERVAL CONSTRUCTION

233 For a test input x_{n+1} , we generate and sort m samples. The prediction interval is constructed
 234 by selecting the shortest interval that contains $\hat{q}_{1-\alpha}$ samples:
 235

$$236 \mathcal{C}^{\text{CMR}}(x_{n+1}) = [y_{n+1}^{(j^*)}, y_{n+1}^{(j^* + \hat{q}_{1-\alpha} - 1)}], \quad (14)$$

237 where j^* minimizes the interval width.
 238

239 **Algorithm 1** Conditional Flow Matching for Conformal Regression (CMR)

240 **Require:** Dataset \mathcal{D} , noise σ , samples m , confidence α .

241 **Ensure:** Prediction interval $\mathcal{C}^{\text{CMR}}(x_{n+1})$.

242 1: Split \mathcal{D} into $\mathcal{I}_{\text{train}}$, \mathcal{I}_{cal} , $\mathcal{I}_{\text{test}}$.
 243 \triangleright *Train Conditional Flow Matching Model*

244 2: Initialize θ for $v_\theta(t, y, x)$.
 245 3: **for** each training iteration **do**
 246 4: Sample batch $\{(x_i, y_i)\}_{i=1}^B$ from $\mathcal{I}_{\text{train}}$.
 247 5: Sample $y_0^i \sim \mathcal{N}(0, 1)$ and $t^i \sim \mathcal{U}(0, 1)$.
 248 6: Compute path: $y_t^i = t^i y_i + (1 - t^i) y_0^i + \sigma \epsilon$.
 249 7: Compute target vector: $u_t^i = y_i - y_0^i$.
 250 8: Loss $\mathcal{L} = \frac{1}{B} \sum_{i=1}^B \|v_\theta(t^i, y_t^i, x_i) - u_t^i\|^2$.
 251 9: Update θ via gradient descent.
 252 10: **end for**

253 \triangleright *Calibration*
 254 11: Initialize scores $L = []$.
 255 12: **for** (x_i, y_i) in \mathcal{I}_{cal} **do**
 256 13: Generate m samples $\{y_i^j\}_{j=1}^m$ via ODE solve.
 257 14: Sort samples to get $\{y_i^{(1)}, \dots, y_i^{(m)}\}$.
 258 15: $L_i \leftarrow \text{FINDMINSAMPLECOUNT}(y_i^{(1..m)}, y_i^{true})$.
 259 16: **end for**
 260 17: $\hat{q}_{1-\alpha} \leftarrow \text{Quantile}_{\lceil(1-\alpha)(|\mathcal{I}_{cal}|+1)\rceil/|\mathcal{I}_{cal}|}(L)$.

261 \triangleright *Inference*
 262 18: **for** x_{n+1} in $\mathcal{I}_{\text{test}}$ **do**
 263 19: Generate m samples $\{y_{n+1}^{(1..m)}\}$ via ODE solve.
 264 20: $\mathcal{C}^{\text{CMR}} \leftarrow \text{FINDSHORTESTINTERVAL}(y_{n+1}^{(1..m)}, \text{count} = \hat{q}_{1-\alpha})$.
 265 21: **end for**

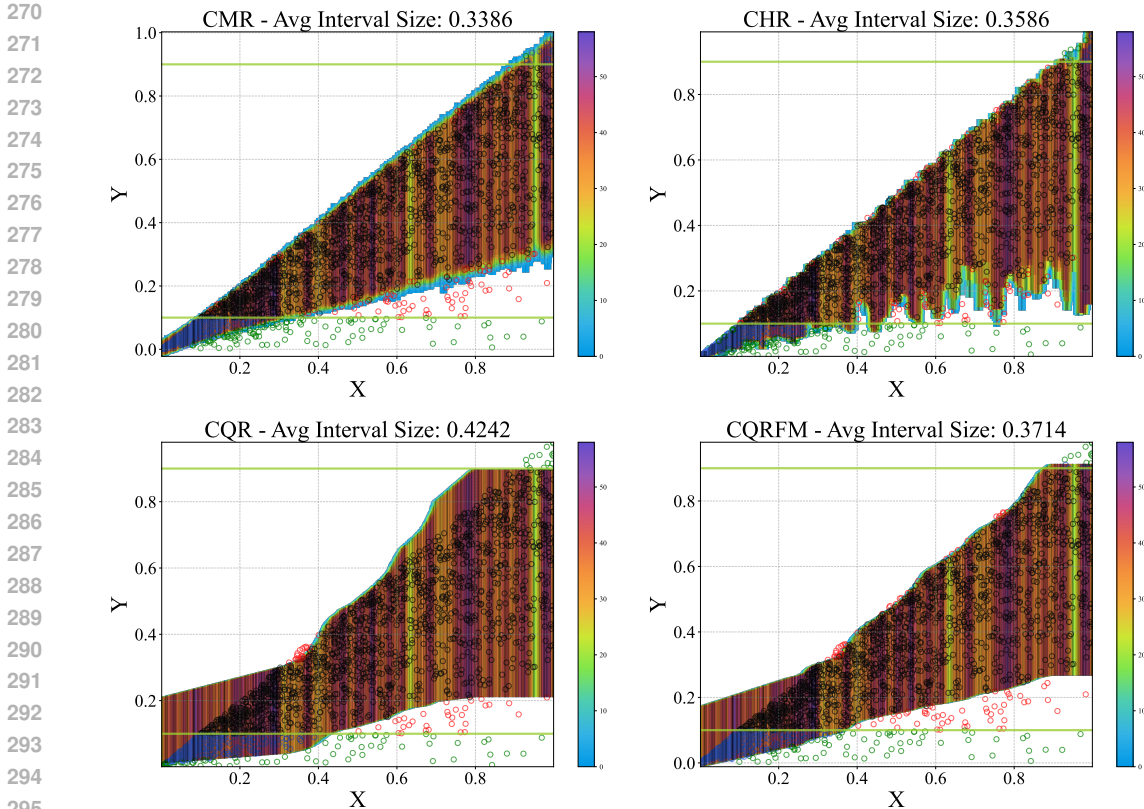


Figure 2: In the scatter plot of features (x) vs. true labels (y), points are color-coded: **green** (uncovered tail points), **blue** (covered tail points), **red** (uncovered points), and **black** (covered points). All conformal regression methods output prediction intervals visualized as vertical bars (fixed width: 0.01) with heights reflecting interval widths. Taller bars correspond to wider intervals, while overlapping regions (color-coded by density) highlight prediction set sizes: denser overlaps imply larger sets. This visualization contrasts interval characteristics and coverage patterns across methods.

4 EXPERIMENTS

In this section, we evaluated our proposed method through simulation studies ($1 - D/2 - D$ scenarios) and experiments on 12 real-world datasets. Performance was measured using Coverage, Size, Worse-Slab Coverage (WSC) (Cauchois et al., 2021) and Tail coverage rate (TCR) (Lin et al., 2021) as key metrics. Comparative analysis included six baseline methods: CHR (Sesia & Romano, 2021), CQR (Romano et al., 2019), CQRR (Sesia & Candès, 2020), CQRM (Kivaranovic et al., 2020), CQRFM (Kivaranovic et al., 2020), and Split Conformal methods (details in Appendix A.4). Experiments used fixed hyperparameters ($m = 1000$, $\sigma = 0.01$, $\alpha = 0.1$) with sensitivity analysis presented in Appendix A.8. All results were validated through 50 randomized experimental partitions to ensure statistical reliability.

4.1 EVALUATION METRICS (WSC AND TCR)

In addition to the standard coverage and size metrics, our evaluation framework incorporates Worst-Slab Coverage (WSC) and Tail Coverage Rate (TCR) to provide a more nuanced understanding of prediction interval performance, particularly concerning conditional validity and robustness across data distributions.

4.1.1 WORST-SLAB COVERAGE (WSC)

The Worst-Slab Coverage (WSC) metric, introduced by Cauchois et al. (2021), assesses the robustness of coverage guarantees across various subsets of the feature space. It identifies the minimum coverage achieved within "slabs" of the data, thereby highlighting potential regions where a model might underperform in its uncertainty quantification.

A slab $S_{v,a,b}$ is defined as a region in the feature space:

$$S_{v,a,b} := \{x \in \mathbb{R}^d \mid a \leq v^T x \leq b\}$$

where $v \in \mathbb{R}^d$ is a direction vector and $a < b \in \mathbb{R}$ define the bounds of the slab.

The Worst-Slab Coverage $WSC_n(C, v, \delta)$ for a given confidence set C , direction v , and minimum mass δ is then formally defined as:

$$WSC_n(C, v, \delta) := \inf_{a < b} \{P_n(Y \in C(X) \mid a \leq v^T X \leq b) \quad \text{s.t.} \quad P_n(a \leq v^T X \leq b) \geq \delta\} \quad (15)$$

Here, P_n denotes the empirical distribution on the observed data. This metric quantifies the lowest coverage rate observed in any slab that contains at least a δ fraction of the data, providing insight into the conditional performance of the prediction intervals.

4.1.2 TAIL COVERAGE RATE (TCR)

The Tail Coverage Rate (TCR) is a crucial metric for evaluating how well prediction intervals perform on extreme values of the response variable. As described by Lin et al. (2021), it specifically measures the coverage rate for data points whose true labels fall within the top and bottom 10% of the response distribution.

Let $Y_{\text{test}} = \{Y_1, \dots, Y_N\}$ be the set of true labels in the test set, and $C(X_i)$ be the prediction interval constructed for the input X_i . To calculate the TCR, we first identify the thresholds for the extreme tails of the Y_{test} distribution. Let $Q_{0.1}$ be the 10th percentile and $Q_{0.9}$ be the 90th percentile of Y_{test} .

The set of tail examples I_{tail} is then defined as:

$$I_{\text{tail}} = \{i \mid Y_i \leq Q_{0.1} \quad \text{or} \quad Y_i \geq Q_{0.9}\}$$

The Tail Coverage Rate (TCR) is subsequently calculated as the proportion of these tail examples whose true labels are contained within their respective prediction intervals:

$$\text{TCR} = \frac{1}{|I_{\text{tail}}|} \sum_{i \in I_{\text{tail}}} \mathbf{1}_{Y_i \in C(X_i)} \quad (16)$$

This metric is particularly important in applications where accurate uncertainty quantification for extreme outcomes is critical.

4.2 SIMULATION STUDY

We conducted experiments on both 1D and 2D simulated data. Detailed parameter settings can be found in Appendix A.5 and A.6. We have visualized the method with the smallest size (CHR(RF), CQRFM) and a classic approach (CQR). As noted in the caption of Figure 2, this visualization comprehensively reflects four key metrics - coverage, size, WSC, and TCR - across different conformal regression methods (Additional information: For details about the WSC and TCR indicators, please refer to Appendix 4.1). It is particularly worth emphasizing that the size dimension can be interpreted through both density variations and the spatial extent of coverage areas. The visualization clearly demonstrates that CMR exhibits the narrowest coverage region, followed by CHR, while CQR and CHRFM show substantial wasted coverage in the lower-left quadrant. We have also conducted visual analysis on the 2D data, with detailed results provided in Appendix A.6.

4.3 REAL DATA

We conducted experiments on 12 real-world datasets (dat, h;i;j;c;b;a;d; Achilles et al., 2008; dat, e;f;g). Following the methodology outlined in Sesia & Candès (2020), we rescale the response Y by the mean absolute value. We randomly allocate 20% of the samples for testing, and from the remaining data, we utilize 70% for training the quantile regression model and 30% for calibration.

| Dataset | Metric | CMR(CFM) | CHR(RF) | CHR(NN) | CQR | CQRM | CQRR | CQRFM | Split |
|------------------------|----------|----------------------|----------------------|----------------------|---------------|----------------------|----------------------|----------------------|----------------------|
| synthetic ¹ | Coverage | 0.901 (0.007) | 0.897 (0.007) | 0.900 (0.009) | 0.900 (0.010) | 0.901 (0.010) | 0.899 (0.009) | 0.899 (0.011) | 0.900 (0.010) |
| | Size | 0.350 (0.007) | 0.361 (0.008) | 0.369 (0.014) | 0.442 (0.020) | 0.402 (0.020) | 0.434 (0.018) | 0.403 (0.018) | 0.475 (0.020) |
| | TCR | 0.865 (0.020) | 0.850 (0.021) | 0.833 (0.043) | 0.722 (0.072) | 0.804 (0.060) | 0.744 (0.068) | 0.776 (0.080) | 0.764 (0.070) |
| synthetic ² | Coverage | 0.897 (0.008) | 0.893 (0.026) | 0.899 (0.009) | 0.901 (0.009) | 0.901 (0.009) | 0.900 (0.010) | 0.902 (0.009) | 0.900 (0.009) |
| | Size | 1.023 (0.019) | 1.051 (0.016) | 1.063 (0.018) | 1.157 (0.042) | 1.113 (0.041) | 1.161 (0.028) | 1.119 (0.035) | 1.190 (0.041) |
| | TCR | 0.728 (0.025) | 0.748 (0.026) | 0.713 (0.029) | 0.677 (0.051) | 0.637 (0.057) | 0.665 (0.050) | 0.633 (0.046) | 0.544 (0.053) |
| bike | Coverage | 0.902 (0.009) | 0.902 (0.008) | 0.900 (0.009) | 0.904 (0.010) | 0.903 (0.008) | 0.905 (0.009) | 0.901 (0.008) | 0.900 (0.010) |
| | Size | 0.521 (0.016) | 1.128 (0.041) | 0.757 (0.043) | 1.564 (0.086) | 1.289 (0.089) | 1.534 (0.078) | 1.313 (0.093) | 1.355 (0.094) |
| | WSC | 0.883 (0.035) | 0.879 (0.041) | 0.885 (0.030) | 0.797 (0.065) | 0.831 (0.049) | 0.813 (0.057) | 0.832 (0.053) | 0.821 (0.049) |
| bio | Coverage | 0.899 (0.004) | 0.899 (0.004) | 0.898 (0.005) | 0.899 (0.005) | 0.900 (0.004) | 0.899 (0.005) | 0.899 (0.004) | 0.899 (0.004) |
| | Size | 1.162 (0.019) | 1.456 (0.019) | 1.578 (0.022) | 2.007 (0.024) | 1.984 (0.034) | 2.009 (0.029) | 1.980 (0.041) | 1.982 (0.053) |
| | TCR | 0.842 (0.010) | 0.749 (0.010) | 0.724 (0.012) | 0.558 (0.026) | 0.574 (0.028) | 0.556 (0.025) | 0.570 (0.027) | 0.731 (0.033) |
| blog | Coverage | 0.900 (0.004) | 0.901 (0.004) | 0.901 (0.005) | 0.942 (0.008) | 0.908 (0.011) | 0.942 (0.008) | 0.909 (0.018) | 0.909 (0.005) |
| | Size | 1.302 (0.048) | 1.596 (0.125) | 1.771 (0.171) | 3.340 (0.359) | 1.943 (0.360) | 3.303 (0.360) | 2.053 (0.648) | 1.429 (0.105) |
| | TCR | 0.699 (0.007) | 0.266 (0.017) | 0.249 (0.015) | 0.251 (0.019) | 0.751 (0.039) | 0.866 (0.035) | 0.766 (0.050) | 0.676 (0.025) |
| community | Coverage | 0.898 (0.023) | 0.899 (0.022) | 0.896 (0.019) | 0.905 (0.023) | 0.900 (0.021) | 0.903 (0.022) | 0.899 (0.023) | 0.899 (0.021) |
| | Size | 1.566 (0.080) | 1.636 (0.116) | 1.574 (0.104) | 1.758 (0.133) | 1.661 (0.116) | 1.764 (0.116) | 1.674 (0.104) | 2.183 (0.224) |
| | TCR | 0.813 (0.048) | 0.713 (0.053) | 0.771 (0.041) | 0.698 (0.094) | 0.750 (0.064) | 0.658 (0.102) | 0.771 (0.066) | 0.670 (0.062) |
| concrete | Coverage | 0.898 (0.028) | 0.903 (0.031) | 0.907 (0.024) | 0.900 (0.026) | 0.902 (0.024) | 0.904 (0.025) | 0.903 (0.022) | 0.905 (0.026) |
| | Size | 0.483 (0.043) | 0.942 (0.065) | 0.493 (0.036) | 0.700 (0.055) | 0.628 (0.050) | 0.699 (0.053) | 0.629 (0.056) | 0.607 (0.051) |
| | TCR | 0.890 (0.051) | 0.725 (0.098) | 0.888 (0.058) | 0.883 (0.094) | 0.882 (0.158) | 0.884 (0.155) | 0.918 (0.091) | 0.859 (0.119) |
| facebook1 | Coverage | 0.900 (0.004) | 0.901 (0.004) | 0.901 (0.005) | 0.943 (0.013) | 0.913 (0.022) | 0.944 (0.013) | 0.918 (0.027) | 0.902 (0.005) |
| | Size | 1.176 (0.045) | 1.559 (0.103) | 1.389 (0.081) | 2.693 (0.466) | 2.148 (0.979) | 2.696 (0.345) | 2.367 (1.565) | 2.156 (0.188) |
| | TCR | 0.789 (0.012) | 0.356 (0.013) | 0.339 (0.012) | 0.418 (0.024) | 0.835 (0.119) | 0.458 (0.013) | 0.880 (0.038) | 0.540 (0.014) |
| facebook2 | Coverage | 0.899 (0.003) | 0.899 (0.003) | 0.899 (0.003) | 0.949 (0.008) | 0.910 (0.020) | 0.935 (0.076) | 0.909 (0.016) | 0.902 (0.003) |
| | Size | 1.201 (0.050) | 1.547 (0.059) | 1.413 (0.068) | 2.778 (0.424) | 1.910 (0.414) | 2.747 (0.597) | 1.847 (0.211) | 2.140 (0.177) |
| | TCR | 0.793 (0.014) | 0.345 (0.009) | 0.344 (0.008) | 0.421 (0.020) | 0.825 (0.125) | 0.456 (0.129) | 0.871 (0.021) | 0.532 (0.013) |
| homes | Coverage | 0.901 (0.006) | 0.900 (0.006) | 0.895 (0.005) | 0.900 (0.006) | 0.901 (0.007) | 0.899 (0.005) | 0.899 (0.007) | 0.897 (0.006) |
| | Size | 0.521 (0.014) | 0.684 (0.017) | 0.538 (0.010) | 0.847 (0.055) | 0.728 (0.034) | 0.826 (0.047) | 0.734 (0.037) | 0.829 (0.077) |
| | TCR | 0.884 (0.025) | 0.834 (0.038) | 0.883 (0.018) | 0.746 (0.067) | 0.805 (0.049) | 0.758 (0.062) | 0.799 (0.053) | 0.597 (0.051) |
| meps19 | Coverage | 0.900 (0.008) | 0.902 (0.007) | 0.900 (0.007) | 0.928 (0.009) | 0.916 (0.010) | 0.927 (0.010) | 0.918 (0.009) | 0.902 (0.008) |
| | Size | 2.288 (0.146) | 2.333 (0.163) | 2.547 (0.151) | 2.898 (0.177) | 2.599 (0.168) | 2.908 (0.190) | 2.622 (0.146) | 3.061 (0.274) |
| | TCR | 0.729 (0.024) | 0.252 (0.020) | 0.296 (0.018) | 0.237 (0.028) | 0.280 (0.147) | 0.278 (0.130) | 0.708 (0.022) | 0.575 (0.031) |
| meps20 | Coverage | 0.901 (0.007) | 0.902 (0.006) | 0.902 (0.006) | 0.927 (0.008) | 0.920 (0.009) | 0.928 (0.008) | 0.918 (0.009) | 0.902 (0.006) |
| | Size | 2.146 (0.097) | 2.357 (0.148) | 2.515 (0.118) | 2.879 (0.136) | 2.745 (0.164) | 2.942 (0.147) | 2.718 (0.127) | 3.052 (0.289) |
| | TCR | 0.722 (0.018) | 0.237 (0.017) | 0.282 (0.015) | 0.227 (0.027) | 0.277 (0.134) | 0.238 (0.072) | 0.696 (0.021) | 0.573 (0.035) |
| meps21 | Coverage | 0.902 (0.008) | 0.902 (0.007) | 0.902 (0.008) | 0.930 (0.007) | 0.922 (0.009) | 0.929 (0.007) | 0.920 (0.008) | 0.903 (0.007) |
| | Size | 2.116 (0.112) | 2.474 (0.156) | 2.667 (0.179) | 2.927 (0.149) | 2.714 (0.134) | 2.934 (0.196) | 2.689 (0.164) | 2.948 (0.245) |
| | TCR | 0.718 (0.019) | 0.234 (0.016) | 0.289 (0.018) | 0.211 (0.019) | 0.239 (0.113) | 0.219 (0.071) | 0.685 (0.021) | 0.567 (0.031) |
| star | Coverage | 0.900 (0.019) | 0.900 (0.019) | 0.902 (0.019) | 0.905 (0.020) | 0.904 (0.018) | 0.901 (0.017) | 0.900 (0.023) | 0.901 (0.022) |
| | Size | 0.179 (0.005) | 0.179 (0.006) | 0.207 (0.007) | 0.181 (0.006) | 0.179 (0.006) | 0.181 (0.006) | 0.178 (0.007) | 0.177 (0.006) |
| | TCR | 0.638 (0.070) | 0.607 (0.071) | 0.724 (0.059) | 0.578 (0.079) | 0.584 (0.077) | 0.571 (0.070) | 0.573 (0.090) | 0.603 (0.090) |

Table 1: The coverage, size, WSC and TCR results for various methods are presented in the table. For various indicators, the top-ranked method is highlighted in bold blue text, while the second-ranked method is displayed in bold green.

As shown in Table 1, CMR achieves the smallest size across all comparative methods in both simulation datasets and 11 real-world datasets. Furthermore, its TCR consistently ranks first or second across most datasets. Notably, CMR demonstrates substantial superiority in tail coverage performance on the “meps19” series and Blog datasets compared to baseline methods. However, on the “star” dataset ($n = 2161$) with limited sample size, CMR exhibits a slightly larger interval size (0.02 higher than the optimal benchmark), suggesting potential challenges in fully capturing data distribution patterns under extreme sample scarcity. These results underscore the critical role of base model selection in conformal prediction performance, with CMR exhibiting remarkable competitive advantages. Boxplots of key evaluation metrics are provided in Appendix A.7 for further visualization.

4.4 ABLATION EXPERIMENT

| Dataset | Metric | CMR(CFM) | CMR(RF) | CMR(NN) | CHR(CFM) | CHR(RF) | CHR(NN) | CQR(CFM) | CQR |
|------------------------|----------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| synthetic ¹ | Coverage | 0.901 (0.007) | 0.898 (0.007) | 0.902 (0.013) | 0.905 (0.010) | 0.897 (0.007) | 0.900 (0.009) | 0.904 (0.009) | 0.900 (0.010) |
| | Size | 0.350 (0.007) | 0.356 (0.004) | 0.355 (0.009) | 0.356 (0.027) | 0.361 (0.008) | 0.369 (0.014) | 0.383 (0.012) | 0.442 (0.020) |
| | TCR | 0.865 (0.020) | 0.832 (0.015) | 0.852 (0.026) | 0.880 (0.031) | 0.850 (0.021) | 0.833 (0.043) | 0.819 (0.035) | 0.722 (0.072) |
| synthetic ² | Coverage | 0.897 (0.008) | 0.896 (0.005) | 0.899 (0.012) | 0.904 (0.008) | 0.893 (0.026) | 0.899 (0.009) | 0.898 (0.008) | 0.901 (0.009) |
| | Size | 1.023 (0.019) | 1.031 (0.014) | 1.024 (0.029) | 1.073 (0.026) | 1.051 (0.016) | 1.063 (0.018) | 1.134 (0.047) | 1.157 (0.042) |
| | TCR | 0.728 (0.025) | 0.734 (0.039) | 0.691 (0.043) | 0.711 (0.035) | 0.748 (0.026) | 0.713 (0.029) | 0.733 (0.081) | 0.677 (0.051) |
| bike | Coverage | 0.902 (0.009) | 0.899 (0.007) | 0.902 (0.005) | 0.901 (0.010) | 0.902 (0.008) | 0.900 (0.009) | 0.900 (0.012) | 0.904 (0.010) |
| | Size | 0.521 (0.016) | 1.024 (0.040) | 0.745 (0.041) | 0.618 (0.056) | 1.128 (0.041) | 0.757 (0.043) | 0.595 (0.028) | 1.564 (0.086) |
| | TCR | 0.883 (0.035) | 0.867 (0.027) | 0.875 (0.049) | 0.892 (0.011) | 0.879 (0.041) | 0.885 (0.030) | 0.891 (0.013) | 0.797 (0.065) |
| bio | Coverage | 0.899 (0.004) | 0.891 (0.005) | 0.898 (0.005) | 0.899 (0.003) | 0.899 (0.004) | 0.898 (0.005) | 0.901 (0.004) | 0.899 (0.005) |
| | Size | 1.162 (0.019) | 1.368 (0.020) | 1.578 (0.016) | 1.469 (0.011) | 1.456 (0.019) | 1.578 (0.022) | 1.606 (0.059) | 2.007 (0.024) |
| | TCR | 0.842 (0.010) | 0.752 (0.009) | 0.758 (0.008) | 0.788 (0.018) | 0.883 (0.020) | 0.877 (0.015) | 0.895 (0.005) | 0.805 (0.025) |
| blog | Coverage | 0.900 (0.004) | 0.903 (0.004) | 0.901 (0.002) | 0.900 (0.003) | 0.901 (0.004) | 0.901 (0.005) | 0.900 (0.002) | 0.942 (0.008) |
| | Size | 1.362 (0.048) | 1.712 (0.088) | 1.997 (0.182) | 1.571 (0.189) | 1.596 (0.125) | 1.771 (0.171) | 2.646 (0.494) | 3.340 (0.359) |
| | TCR | 0.770 (0.018) | 0.886 (0.017) | 0.878 (0.028) | 0.904 (0.004) | 0.862 (0.022) | 0.859 (0.022) | 0.905 (0.003) | 0.751 (0.039) |
| community | Coverage | 0.898 (0.023) | 0.909 (0.019) | 0.897 (0.020) | 0.904 (0.025) | 0.899 (0.022) | 0.896 (0.019) | 0.901 (0.024) | 0.905 (0.023) |
| | Size | 1.566 (0.080) | 1.545 (0.093) | 1.565 (0.097) | 1.463 (0.103) | 1.636 (0.116) | 1.574 (0.104) | 1.619 (0.091) | 1.758 (0.133) |
| | TCR | 0.880 (0.080) | 0.887 (0.048) | 0.896 (0.038) | 0.895 (0.021) | 0.892 (0.057) | 0.895 (0.057) | 0.897 (0.019) | 0.889 (0.068) |
| concrete | Coverage | 0.898 (0.028) | 0.909 (0.016) | 0.905 (0.015) | 0.903 (0.028) | 0.903 (0.031) | 0.907 (0.024) | 0.905 (0.032) | 0.900 (0.026) |
| | Size | 0.483 (0.043) | 0.900 (0.060) | 0.470 (0.054) | 0.539 (0.019) | 0.942 (0.065) | 0.493 (0.036) | 0.538 (0.030) | 0.700 (0.055) |
| | TCR | 0.848 (0.157) | 0.875 (0.092) | 0.870 (0.102) | 0.890 (0.024) | 0.879 (0.107) | 0.888 (0.088) | 0.885 (0.035) | 0.883 (0.094) |
| facebook1 | Coverage | 0.900 (0.004) | 0.903 (0.005) | 0.899 (0.006) | 0.902 (0.005) | 0.901 (0.004) | 0.901 (0.005) | 0.900 (0.004) | 0.943 (0.013) |
| | Size | 1.176 (0.045) | 1.600 (0.055) | 1.572 (0.093) | 1.285 (0.013) | 1.559 (0.103) | 1.389 (0.081) | 2.093 (0.157) | 2.693 (0.466) |
| | TCR | 0.746 (0.032) | 0.879 (0.022) | 0.884 (0.014) | 0.901 (0.004) | 0.838 (0.025) | 0.838 (0.030) | 0.897 (0.005) | 0.878 (0.033) |
| facebook2 | Coverage | 0.899 (0.003) | 0.900 (0.002) | 0.900 (0.003) | 0.899 (0.002) | 0.899 (0.003) | 0.899 (0.003) | 0.898 (0.003) | 0.949 (0.008) |
| | Size | 1.201 (0.050) | 1.521 (0.056) | 1.562 (0.051) | 1.412 (0.052) | 1.547 (0.059) | 1.413 (0.068) | 2.146 (0.109) | 2.778 (0.424) |
| | TCR | 0.793 (0.014) | 0.361 (0.007) | 0.876 (0.007) | 0.945 (0.004) | 0.345 (0.009) | 0.344 (0.008) | 0.932 (0.011) | 0.421 (0.020) |
| homes | Coverage | 0.901 (0.006) | 0.903 (0.005) | 0.899 (0.009) | 0.900 (0.007) | 0.900 (0.006) | 0.895 (0.005) | 0.899 (0.005) | 0.900 (0.006) |
| | Size | 0.521 (0.014) | 0.682 (0.016) | 0.527 (0.009) | 0.691 (0.146) | 0.684 (0.017) | 0.538 (0.010) | 0.759 (0.270) | 0.847 (0.055) |
| | TCR | 0.884 (0.025) | 0.830 (0.030) | 0.884 (0.021) | 0.896 (0.008) | 0.834 (0.038) | 0.883 (0.018) | 0.895 (0.006) | 0.746 (0.067) |
| meps19 | Coverage | 0.900 (0.008) | 0.902 (0.007) | 0.902 (0.004) | 0.895 (0.005) | 0.902 (0.007) | 0.900 (0.007) | 0.899 (0.004) | 0.928 (0.009) |
| | Size | 2.288 (0.146) | 2.380 (0.160) | 2.602 (0.160) | 2.280 (0.311) | 2.333 (0.163) | 2.547 (0.151) | 2.853 (0.246) | 2.898 (0.177) |
| | TCR | 0.855 (0.042) | 0.897 (0.026) | 0.865 (0.056) | 0.891 (0.005) | 0.890 (0.025) | 0.897 (0.021) | 0.895 (0.004) | 0.863 (0.034) |
| meps20 | Coverage | 0.901 (0.007) | 0.901 (0.007) | 0.901 (0.005) | 0.905 (0.005) | 0.902 (0.006) | 0.902 (0.006) | 0.902 (0.005) | 0.927 (0.008) |
| | Size | 2.146 (0.097) | 2.321 (0.152) | 2.586 (0.105) | 2.288 (0.181) | 2.357 (0.148) | 2.515 (0.118) | 3.046 (0.224) | 2.879 (0.136) |
| | TCR | 0.847 (0.029) | 0.896 (0.022) | 0.901 (0.016) | 0.900 (0.004) | 0.893 (0.025) | 0.896 (0.017) | 0.896 (0.005) | 0.856 (0.038) |
| meps21 | Coverage | 0.902 (0.008) | 0.906 (0.005) | 0.898 (0.006) | 0.903 (0.008) | 0.902 (0.007) | 0.902 (0.008) | 0.902 (0.003) | 0.930 (0.007) |
| | Size | 2.116 (0.112) | 2.340 (0.143) | 2.580 (0.107) | 2.116 (0.113) | 2.474 (0.156) | 2.667 (0.179) | 2.763 (0.240) | 2.927 (0.149) |
| | TCR | 0.854 (0.044) | 0.895 (0.026) | 0.887 (0.020) | 0.899 (0.005) | 0.896 (0.022) | 0.895 (0.021) | 0.897 (0.003) | 0.854 (0.029) |
| star | Coverage | 0.900 (0.019) | 0.899 (0.017) | 0.899 (0.020) | 0.904 (0.018) | 0.900 (0.019) | 0.902 (0.019) | 0.905 (0.023) | 0.905 (0.020) |
| | Size | 0.179 (0.005) | 0.179 (0.006) | 0.201 (0.008) | 0.201 (0.011) | 0.179 (0.006) | 0.207 (0.007) | 0.211 (0.015) | 0.181 (0.006) |
| | TCR | 0.914 (0.052) | 0.913 (0.052) | 0.892 (0.045) | 0.888 (0.021) | 0.892 (0.050) | 0.901 (0.057) | 0.889 (0.025) | 0.894 (0.056) |

Table 2: The coverage, size, WSC and TCR results for various methods are presented in the table. An orange background indicates the best metric within the dataset. A blue background denotes the better metric between CMR and CHR when using the same underlying predictive model. Additionally, bold text represents the better metric among different underlying models under the same non-conformity score condition (CMR, CHR, or CQR).

In this section, we will decouple CMR to demonstrate how the non-conformity scores of CFM and CMR jointly enhance the efficiency of prediction sets. Specifically, we conducted comparative experiments where we tested neural networks and random forests as base predictive models respectively, and also tested CFM as the base predictive model, comparing their performances across three non-conformity score functions: CMR, CHR, and CQR. Each dataset was randomly split 10 times.

As shown in Table 2, we use an orange background to indicate the best metric within the dataset, and a blue background to denote the better metric between CMR and CHR under the same base predictive model. Bold text is used to highlight the better metric among different base models under the same non-conformity score condition (CMR, CHR, CQR). It can be observed that the optimal metrics are either achieved by using CMR as the non-conformity score or by using CFM as the base predictive model. Secondly, the prediction set sizes generated by CMR (with CFM as the base model) are optimal in 11 out of 14 datasets. Using CFM as the base model ensures stable conditional coverage and TCR. Compared to CHR, CMR generally outperforms CHR across three different base predictive models. Furthermore, we have also analyzed the resource consumption of CFM under different calibration methods, with details provided in Appendix A.9.

5 CONCLUSION

We have introduced Conditional Flow Matching for Conformal Regression (CMR), a novel framework that effectively combines conditional flow matching with conformal prediction to generate statistically valid prediction intervals for regression tasks. Our approach employs a simulation-free training method to learn the conditional distribution of the target variable and then applies a calibration procedure to construct prediction intervals with reduced widths. Our method offers a brand-new perspective for one-dimensional regression tasks, generating continuous prediction intervals that are not only confidence-guaranteed but also the most efficient.

REFERENCES

- Bike sharing dataset data set. <https://archive.ics.uci.edu/ml/datasets/bike+sharing+dataset>, a. Accessed: July, 2019.
- Physicochemical properties of protein tertiary structure data set. <https://archive.ics.uci.edu/ml/datasets/Physicochemical+Properties+of+Protein+Tertiary+Structure>, b. Accessed: July, 2019.
- BlogFeedback data set. <https://archive.ics.uci.edu/ml/datasets/BlogFeedback>, c. Accessed: July, 2019.
- Communities and crime data set. <http://archive.ics.uci.edu/ml/datasets/communities+and+crime>, d. Accessed: July, 2019.
- Concrete compressive strength data set. <http://archive.ics.uci.edu/ml/datasets/concrete+compressive+strength>, e. Accessed: July, 2019.
- Facebook comment volume data set. <https://archive.ics.uci.edu/ml/datasets/Facebook+Comment+Volume+Dataset>, f. Accessed: July, 2019.
- House sales in King County, USA. <https://www.kaggle.com/harlfoxem/housesalesprediction/metadata>, g. Accessed: August, 2019.
- Medical expenditure panel survey, panel 19. https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-181, h. Accessed: July, 2019.
- Medical expenditure panel survey, panel 20. https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-181, i. Accessed: July, 2019.
- Medical expenditure panel survey, panel 21. https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-192, j. Accessed: July, 2019.
- C.M. Achilles, Helen Pate Bain, Fred Bellott, Jayne Boyd-Zaharias, Jeremy Finn, John Folger, John Johnston, and Elizabeth Word. Tennessee’s student teacher achievement ratio (STAR) project, 2008. URL <https://doi.org/10.7910/DVN/SIWH9F>. Accessed: July, 2019.
- Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. Predictive inference with the jackknife+. *The Annals of Statistics*, 49(1):486–507, 2021.
- Maxime Cauchois, Suyash Gupta, and John C Duchi. Knowing what you know: valid and validated confidence sets in multiclass and multilabel prediction. *Journal of machine learning research*, 22(81):1–42, 2021.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018a.

- 540 Wenyu Chen, Kelli-Jean Chun, and Rina Foygel Barber. Discretized conformal prediction
541 for efficient distribution-free inference. *Stat*, 7(1):e173, 2018b.
- 542
- 543 Matt Y Cheung, Tucker J Netherton, Laurence E Court, Ashok Veeraraghavan, and
544 Guha Balakrishnan. Regression conformal prediction under bias. *arXiv preprint*
545 *arXiv:2410.05263*, 2024.
- 546 Hugh A Chipman, Edward I George, and Robert E McCulloch. Bart: Bayesian additive
547 regression trees. *Annals of Applied Statistics*, 4(1):266–298, 2010.
- 548
- 549 Nicolo Colombo. Normalizing flows for conformal regression. *arXiv preprint*
550 *arXiv:2406.03346*, 2024.
- 551 Nathaniel Diamant, Ehsan Hajiramezanali, Tommaso Biancalani, and Gabriele Scalia. Con-
552 formalized deep splines for optimal and efficient prediction sets. In *International Confer-*
553 *ence on Artificial Intelligence and Statistics*, pp. 1657–1665. PMLR, 2024.
- 554
- 555 Chao Gao, Liren Shan, Vaidehi Srinivas, and Aravindan Vijayaraghavan. Volume optimality
556 in conformal prediction with structured prediction sets. In *Forty-second International*
557 *Conference on Machine Learning*, 2025.
- 558
- 559 Natalia Martinez Gil, Dhaval Patel, Chandra Reddy, Giridhar Ganapavarapu, Roman Va-
560 culin, and Jayant Kalagnanam. Identifying homogeneous and interpretable groups for
561 conformal prediction. In *Proceedings of the Fortieth Conference on Uncertainty in Arti-*
562 *ficial Intelligence*, pp. 2471–2485, 2024.
- 563 Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud.
564 Ffjord: Free-form continuous dynamics for scalable reversible generative models. In *In-*
565 *ternational Conference on Learning Representations*, 2018.
- 566
- 567 Leying Guan. Localized conformal prediction: A generalized inference framework for con-
568 formal prediction. *Biometrika*, 110(1):33–50, 2023.
- 569 Rohan Hore and Rina Foygel Barber. Conformal prediction with local weights: randomiza-
570 tion enables robust guarantees. *Journal of the Royal Statistical Society Series B: Statistical*
571 *Methodology*, 87(2):549–578, 2025.
- 572
- 573 Rafael Izbicki, Gilson T Shimizu, and Rafael B Stern. Flexible distribution-free conditional
574 predictive bands using density estimators. *arXiv preprint arXiv:1910.05575*, 2019.
- 575
- 576 Rafael Izbicki, Gilson Shimizu, and Rafael B Stern. Cd-split and hpd-split: Efficient con-
577 formal regions in high dimensions. *Journal of Machine Learning Research*, 23(87):1–32,
578 2022.
- 579 Danijel Kivaranovic, Kory D Johnson, and Hannes Leeb. Adaptive, distribution-free pre-
580 diction intervals for deep networks. In *International Conference on Artificial Intelligence*
581 *and Statistics*, pp. 4346–4356. PMLR, 2020.
- 582 Shayan Kiyani, George J Pappas, and Hamed Hassani. Length optimization in conformal
583 prediction. *Advances in Neural Information Processing Systems*, 37:99519–99563, 2024.
- 584
- 585 Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman.
586 Distribution-free predictive inference for regression. *Journal of the American Statisti-*
587 *cal Association*, 113(523):1094–1111, 2018.
- 588
- 589 Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. Locally valid and discriminative prediction
590 intervals for deep learning models. *Advances in Neural Information Processing Systems*,
591 34:8378–8391, 2021.
- 592 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow
593 matching for generative modeling. *International Conference on Learning Representations*
(*ICLR*), 2023.

- 594 Weiguang Liu, Áureo de Paula, and Elie T Tamer. Prediction sets and conformal inference
595 with censored outcomes. Technical report, cemmap working paper, 2025.
596
- 597 Rui Luo and Zhixin Zhou. Conformal thresholded intervals for efficient regression. *Proceed-*
598 *ings of the AAAI Conference on Artificial Intelligence*, 39(18):19216–19223, 2025.
- 599 Malik Magdon-Ismail and Amir Atiya. Neural networks for density estimation. *Advances*
600 *in Neural Information Processing Systems*, 11, 1998.
601
- 602 Nicolai Meinshausen and Greg Ridgeway. Quantile regression forests. *Journal of machine*
603 *learning research*, 7(6), 2006.
- 604 Sang Jun Moon, Jong-June Jeon, Jason Sang Hun Lee, and Yongdai Kim. Learning multiple
605 quantiles with neural networks. *Journal of Computational and Graphical Statistics*, 30
606 (4):1238–1248, 2021.
607
- 608 Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive
609 confidence machines for regression. In *Machine learning: ECML 2002: 13th European*
610 *conference on machine learning Helsinki, Finland, August 19–23, 2002 proceedings 13*,
611 pp. 345–356. Springer, 2002.
- 612 Harris Papadopoulos, Vladimir Vovk, and Alex Gammerman. Regression conformal pre-
613 diction with nearest neighbours. *Journal of Artificial Intelligence Research*, 40:815–840,
614 2011.
- 615 Vincent Plassier, Alexander Fishkov, Mohsen Guizani, Maxim Panov, and Eric Moulines.
616 Probabilistic conformal prediction with approximate conditional validity. *arXiv preprint*
617 *arXiv:2407.01794*, 2024.
618
- 619 Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression.
620 *Advances in neural information processing systems*, 32, 2019.
- 621 Matteo Sesia and Emmanuel J Candès. A comparison of some conformal quantile regression
622 methods. *Stat*, 9(1):e261, 2020.
623
- 624 Matteo Sesia and Yaniv Romano. Conformal prediction using conditional histograms. *Ad-*
625 *vances in Neural Information Processing Systems*, 34:6304–6315, 2021.
- 626 Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine*
627 *Learning Research*, 9(3), 2008.
628
- 629 Alexander Tong, Kilian Fatras, Nikolay Malkin, Guillaume Huguet, Yanlei Zhang, Jarrid
630 Rector-Brooks, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based
631 generative models with minibatch optimal transport. *arXiv preprint arXiv:2302.00482*,
632 2023.
- 633 Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a ran-*
634 *dom world*, volume 29. Springer, 2005.
635
- 636 Zhendong Wang, Ruijiang Gao, Mingzhang Yin, Mingyuan Zhou, and David Blei. Proba-
637 bilistic conformal prediction using conditional random samples. In *International Confer-*
638 *ence on Artificial Intelligence and Statistics*, pp. 8814–8836. PMLR, 2023.
- 639 Minxing Zheng and Shixiang Zhu. Optimizing probabilistic conformal prediction with vec-
640 torized non-conformity scores. *arXiv preprint arXiv:2410.13735*, 2024.
641
642
643
644
645
646
647

A APPENDIX

LARGE LANGUAGE MODEL (LLM) USAGE

During the preparation of this manuscript, a large language model (LLM), specifically [Gemini-2.5-Flash], was employed as a general-purpose assist tool. The LLM’s contributions were primarily in the following areas:

- **Code Optimization for Experimental Visualization:** The LLM assisted in optimizing and refining Python code snippets used for experimental visualization and data processing routines. This collaboration led to more efficient and readable implementations, particularly for generating the figures and tables presented in the paper.
- **Writing Assistance and Refinement:** The LLM was utilized for drafting and refining certain sections of the paper, including improving clarity, grammar, and stylistic coherence. This involved generating initial textual descriptions and polishing existing content to enhance its overall quality and readability.

The authors maintained full responsibility for reviewing, editing, and validating all content generated or optimized with the assistance of the LLM, ensuring its accuracy, originality, and adherence to scientific standards. The LLM was not involved in the core research ideation, experimental design, data collection, or primary analysis leading to the scientific conclusions. The scientific content, conclusions, and any potential errors remain solely the responsibility of the authors.

A.1 THEORETICAL ANALYSIS

A.2 PROOF OF THEOREM 1

Proof: Let n be the size of the calibration set \mathcal{I}_{cal} . We consider n calibration points $(x_1, y_1), \dots, (x_n, y_n)$ and one test point (x_{n+1}, y_{n+1}) . All these $n+1$ data points are assumed to be drawn independently and identically distributed (i.i.d.) from the same underlying distribution P .

For each data point (x_i, y_i) , we compute its nonconformity score L_i using the procedure described in Section 3.3. **Specifically, L_i is defined as the minimum number of samples s such that the shortest interval constructed from s generated samples contains the true outcome y_i^{true} .** Let L_{n+1} denote the nonconformity score for the test point (x_{n+1}, y_{n+1}) , i.e., $L_{n+1} = L(y_{n+1}; x_{n+1})$.

Since the data points $(x_1, y_1), \dots, (x_n, y_n), (x_{n+1}, y_{n+1})$ are i.i.d., and the nonconformity score function $L(\cdot; \cdot)$ is applied identically to each point (after the model is trained on a separate training set, if applicable), the resulting nonconformity scores L_1, \dots, L_n, L_{n+1} are also i.i.d., and therefore **exchangeable**.

Let $k_{\text{cal}} = \lceil (n+1)(1-\alpha) \rceil$. The threshold $\hat{q}_{1-\alpha}$ is defined as the k_{cal} -th order statistic of the calibration scores $\{L_1, \dots, L_n\}$. We denote this as $\hat{q} = L_{(k_{\text{cal}})}$. **Note that in our method, \hat{q} represents a discrete count of samples.**

The prediction set for x_{n+1} is $\mathcal{C}^{\text{CMR}}(x_{n+1}) = \{y' \mid L(y'; x_{n+1}) \leq \hat{q}\}$. We want to show that $\mathbb{P}(y_{n+1} \in \mathcal{C}^{\text{CMR}}(x_{n+1})) \geq 1 - \alpha$. This is equivalent to showing $\mathbb{P}(L_{n+1} \leq \hat{q}) \geq 1 - \alpha$.

The event $y_{n+1} \notin \mathcal{C}^{\text{CMR}}(x_{n+1})$ occurs if and only if its nonconformity score L_{n+1} is strictly greater than the threshold \hat{q} :

$$y_{n+1} \notin \mathcal{C}^{\text{CMR}}(x_{n+1}) \iff L_{n+1} > \hat{q}$$

The event $L_{n+1} > \hat{q} = L_{(k_{\text{cal}})}$ means that L_{n+1} is strictly larger than at least $n - k_{\text{cal}} + 1$ of the scores from the calibration set $\{L_1, \dots, L_n\}$. This implies that the rank of L_{n+1} within the full set of $n + 1$ scores $\{L_1, \dots, L_n, L_{n+1}\}$ must be at least $k_{\text{cal}} + 1$.

Therefore, the probability of non-coverage is bounded by the probability of this rank condition:

702
703
704
705
706
707
708
709
710

$$\begin{aligned}\mathbb{P}(y_{n+1} \notin \mathcal{C}^{\text{CMR}}(x_{n+1})) &= \mathbb{P}(L_{n+1} > \hat{q}) & (17) \\ &\leq \mathbb{P}(\text{rank}(L_{n+1}) \geq k_{cal} + 1) & (18)\end{aligned}$$

Since the scores L_1, \dots, L_n, L_{n+1} are exchangeable, any of the $n + 1$ scores is equally likely to take any rank position. Thus, the probability that L_{n+1} has a rank of r (from 1 to $n + 1$) is $\frac{1}{n+1}$.

711
712
713
714
715

$$\mathbb{P}(\text{rank}(L_{n+1}) \geq k_{cal} + 1) = \sum_{r=k_{cal}+1}^{n+1} \mathbb{P}(\text{rank}(L_{n+1}) = r) = \frac{(n+1) - k_{cal}}{n+1}$$

So, we have:

716
717

$$\mathbb{P}(y_{n+1} \notin \mathcal{C}^{\text{CMR}}(x_{n+1})) \leq \frac{(n+1) - k_{cal}}{n+1}$$

718
719

Now, we substitute the value of $k_{cal} = \lceil (n+1)(1-\alpha) \rceil$. By the definition of the ceiling function, for any real number x , we have $\lceil x \rceil \geq x$. Therefore,

720
721

$$k_{cal} \geq (n+1)(1-\alpha)$$

722
723

Substituting this into our probability bound:

724
725

$$\mathbb{P}(y_{n+1} \notin \mathcal{C}^{\text{CMR}}(x_{n+1})) \leq \frac{n+1 - \lceil (n+1)(1-\alpha) \rceil}{n+1} \quad (19)$$

727
728

$$\leq \frac{n+1 - (n+1)(1-\alpha)}{n+1} \quad (20)$$

729
730

$$= \frac{(n+1)\alpha}{n+1} = \alpha \quad (21)$$

731
732

Finally, the probability of coverage is:

733
734

$$\mathbb{P}(y_{n+1} \in \mathcal{C}^{\text{CMR}}(x_{n+1})) = 1 - \mathbb{P}(y_{n+1} \notin \mathcal{C}^{\text{CMR}}(x_{n+1})) \geq 1 - \alpha$$

736
737

This proof holds unconditionally over the joint distribution of all $n + 1$ samples.

738
739

A.3 CALIBRATE AND TEST THE ALGORITHM

740
741
742
743
744

In this section, we provide a detailed description of our algorithm implementation during the calibration and testing phases. By combining the sliding window approach with the binary segmentation algorithm, we have reduced the time complexity of the algorithm in the calibration phase from $O(nm)$ to $O(n \log m)$. Additionally, the time complexity of the testing phase is $O(n)$.

745

746
747

A.3.1 FIND MINIMUM SAMPLE COUNT (SCORE CALCULATION)

748
749

A.3.2 FIND SHORTEST INTERVAL (INFERENCE)

750
751

A.4 OTHER METHODS (CONFORMAL PREDICTION FOR REGRESSION)

752
753
754
755

In this section, we provide a detailed exposition of the methodological foundations and operational frameworks underlying the CHR Sesia & Romano (2021), CQR Romano et al. (2019), CQRR Sesia & Candès (2020), CQRM Kivaranovic et al. (2020), and CQRFM Kivaranovic et al. (2020) methodologies.

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

Algorithm 2 **FINDMINSAMPLECOUNT** - Find minimum sample number of interval containing a target value

Require: Sorted samples $Y = \{y^{(1)}, y^{(2)}, \dots, y^{(m)}\}$, target value y^{true} .

Ensure: Minimum samples L_{min} .

```

1: Initialize search range:  $low \leftarrow 2, high \leftarrow m$ .
2:  $result\_samples \leftarrow m$  {Default to full range}
3: while  $low \leq high$  do
4:    $mid \leftarrow \lfloor (low + high)/2 \rfloor$ . {Current sample count  $s$ }
5:    $min\_width \leftarrow \infty$ .
6:    $start^* \leftarrow 0, end^* \leftarrow 0$ 
7:   for  $start \in \{0, 1, \dots, m - mid\}$  do
8:      $end \leftarrow start + mid - 1$ .
9:      $left \leftarrow y^{(start+1)}, right \leftarrow y^{(end+1)}$ . {Adjust for 1-based index}
10:     $width \leftarrow right - left$ .
11:    if  $width < min\_width$  then
12:       $min\_width \leftarrow width$ .
13:       $start^* \leftarrow start + 1, end^* \leftarrow end + 1$ . {Store indices of shortest interval}
14:    end if
15:  end for
16:  {Check if the shortest interval covers the true value}
17:  if  $y^{(start^*)} \leq y^{true} \leq y^{(end^*)}$  then
18:     $result\_samples \leftarrow mid$ .
19:     $high \leftarrow mid - 1$ . Try smaller sample count
20:  else
21:     $low \leftarrow mid + 1$ . Need larger sample count
22:  end if
23: end while
24: Return  $result\_samples$ 

```

Algorithm 3 **FINDSHORTESTINTERVAL** - Find shortest interval with fixed sample count

Require: Sorted samples $Y = \{y^{(1)}, y^{(2)}, \dots, y^{(m)}\}$, sample count threshold $\hat{q}_{1-\alpha}$ (integer)

Ensure: Shortest interval $[y_{low}, y_{up}]$ containing $\hat{q}_{1-\alpha}$ samples

```

1:  $min\_width \leftarrow \infty$ .
2:  $best\_interval \leftarrow [y^{(1)}, y^{(m)}]$ .
3:  $s \leftarrow \hat{q}_{1-\alpha}$ .
4: for  $i \in \{1, \dots, m - s + 1\}$  do
5:    $j \leftarrow i + s - 1$ .
6:    $width \leftarrow y^{(j)} - y^{(i)}$ .
7:   if  $width < min\_width$  then
8:      $min\_width \leftarrow width$ .
9:      $best\_interval \leftarrow [y^{(i)}, y^{(j)}]$ .
10:  end if
11: end for
12: Return  $best\_interval$ .

```

A.4.1 CONFORMAL QUANTILE REGRESSION (1).

Conformal Quantile Regression constructs intervals based on quantile regression:

$$\mathcal{C}^{\text{CQR}}(x_{n+1}) = \left[\widehat{q}_{\frac{\alpha}{2}}(x_{n+1}) - t_{1-\alpha}^{\text{CQR}}, \widehat{q}_{1-\frac{\alpha}{2}}(x_{n+1}) + t_{1-\alpha}^{\text{CQR}} \right], \quad (22)$$

where $\widehat{q}_{\alpha/2}$ and $\widehat{q}_{1-\alpha/2}$ are conditional quantile estimates, and $t_{1-\alpha}^{\text{CQR}}$ is the $(1-\alpha)(1+1/|\mathcal{I}_{\text{cal}}|)$ -th empirical quantile of $\{S_i^{\text{CQR}}\}_{i \in \mathcal{I}_{\text{cal}}} \cup \{\infty\}$, with:

$$S_i^{\text{CQR}} = \max \left(\widehat{q}_{\frac{\alpha}{2}}(x_i) - y_i, y_i - \widehat{q}_{1-\frac{\alpha}{2}}(x_i) \right). \quad (23)$$

A.4.2 CONFORMAL QUANTILE REGRESSION (2).

Similar to Sesia & Candès (2020), we name the methods proposed by Kivaranovic et al. (2020) as CQR-m, respectively, which differ from the CQR proposed by Romano et al. (2019). First, let's consider CQR-m. The model defined based on quantile regression is as follows:

$$\begin{aligned} \mathcal{C}^{\text{CQR-m}}(x_{n+1}) &= \left[\widehat{q}_{\frac{\alpha}{2}}(x_{n+1}) - \widehat{\Delta}_{\alpha, \text{lo}}^{\text{CQR-m}}, \widehat{q}_{1-\frac{\alpha}{2}}(x_{n+1}) + \widehat{\Delta}_{\alpha, \text{up}}^{\text{CQR-m}} \right], \\ \widehat{\Delta}_{\alpha, \text{lo}}^{\text{CQR-m}} &= t_{1-\alpha}^{\text{CQR-m}} \left[\widehat{q}_{\frac{1}{2}}(x_i) - \widehat{q}_{\frac{\alpha}{2}}(x_i) \right], \\ \widehat{\Delta}_{\alpha, \text{up}}^{\text{CQR-m}} &= t_{1-\alpha}^{\text{CQR-m}} \left[\widehat{q}_{1-\frac{\alpha}{2}}(x_i) - \widehat{q}_{\frac{1}{2}}(x_i) \right], \end{aligned} \quad (24)$$

where $\widehat{q}_{1/2}$ indicates an estimated median regression function obtained with the same black-box algorithm as $\widehat{q}_{\alpha/2}$ and $\widehat{q}_{1-\alpha/2}$, and $t_{1-\alpha}^{\text{CQR-m}}$ is the same $(1-\alpha)(1+1/|\mathcal{I}_{\text{cal}}|)$ -th empirical quantile of $\{S_i^{\text{CQR-m}}\}_{i \in \mathcal{I}_{\text{cal}}} \cup \{\infty\}$, with:

$$S_i^{\text{CQR-m}} = \max \left(\frac{\widehat{q}_{\frac{\alpha}{2}}(x_i) - y_i}{\widehat{q}_{\frac{1}{2}}(x_i) - \widehat{q}_{\frac{\alpha}{2}}(x_i)}, \frac{y_i - \widehat{q}_{1-\frac{\alpha}{2}}(x_i)}{\widehat{q}_{1-\frac{\alpha}{2}}(x_i) - \widehat{q}_{\frac{1}{2}}(x_i)} \right). \quad (25)$$

In addition, there is an improved version of CQR-m, which does not require estimating the quantile of the regression median Sesia & Candès (2020). The prediction interval it constructs is as follows:

$$\begin{aligned} \mathcal{C}^{\text{CQR-r}}(x_{n+1}) &= \left[\widehat{q}_{\frac{\alpha}{2}}(x_{n+1}) - \widehat{\Delta}_{\alpha}^{\text{CQR-r}}, \widehat{q}_{1-\frac{\alpha}{2}}(x_{n+1}) + \widehat{\Delta}_{\alpha}^{\text{CQR-r}} \right], \\ \widehat{\Delta}_{\alpha}^{\text{CQR-r}} &= t_{1-\alpha}^{\text{CQR-r}} \left[\widehat{q}_{1-\frac{\alpha}{2}}(x_i) - \widehat{q}_{\frac{\alpha}{2}}(x_i) \right], \end{aligned} \quad (26)$$

where $t_{1-\alpha}^{\text{CQR-r}}$ is the $(1-\alpha)(1+1/|\mathcal{I}_{\text{cal}}|)$ -th empirical quantile of $\{S_i^{\text{CQR-r}}\}_{i \in \mathcal{I}_{\text{cal}}} \cup \{\infty\}$, with:

$$S_i^{\text{CQR-r}} = \max \left(\frac{\widehat{q}_{\frac{\alpha}{2}}(x_i) - y_i}{\widehat{q}_{1-\frac{\alpha}{2}}(x_i) - \widehat{q}_{\frac{\alpha}{2}}(x_i)}, \frac{y_i - \widehat{q}_{1-\frac{\alpha}{2}}(x_i)}{\widehat{q}_{1-\frac{\alpha}{2}}(x_i) - \widehat{q}_{\frac{\alpha}{2}}(x_i)} \right). \quad (27)$$

A.4.3 CONFORMAL QUANTILE REGRESSION WITH FULL MODEL.

CQRFM builds upon CQR-m by introducing a modification that allows the model to output three distinct values: the lower bound, median, and upper bound simultaneously from a single neural network. The key idea is to train a neural network $\mathcal{N} : \mathbb{R}^d \rightarrow \mathbb{R}^3$ such that $\mathcal{N}(x) = (l(x), m(x), u(x))$, where l , m , and u are functions that estimate the $\alpha/2$ -quantile, the median, and the $(1-\alpha/2)$ -quantile, respectively, with the constraint that $l(x) \leq m(x) \leq u(x)$ for all $x \in \mathbb{R}^d$.

The network is trained using a modified quantile regression loss function:

$$L_{\tau}(\mathcal{N}(x), y) = h_{\tau/2}(y - l(x)) + h_{1/2}(y - m(x)) + h_{1-\tau/2}(y - u(x)), \quad (28)$$

where $h_{\tau}(u) = (\tau - \mathbf{1}_{u \leq 0})u$ is the standard quantile regression loss function.

Similar to CQR-m, the prediction interval is constructed as:

$$\begin{aligned} \mathcal{C}^{\text{CQRFM}}(x_{n+1}) &= \left[l(x_{n+1}) - \widehat{\Delta}_{\alpha, \text{lo}}^{\text{CQRFM}}, u(x_{n+1}) + \widehat{\Delta}_{\alpha, \text{up}}^{\text{CQRFM}} \right], \\ \widehat{\Delta}_{\alpha, \text{lo}}^{\text{CQRFM}} &= t_{1-\alpha}^{\text{CQRFM}} [m(x_i) - l(x_i)], \\ \widehat{\Delta}_{\alpha, \text{up}}^{\text{CQRFM}} &= t_{1-\alpha}^{\text{CQRFM}} [u(x_i) - m(x_i)], \end{aligned} \quad (29)$$

where $t_{1-\alpha}^{\text{CQRFM}}$ is the $(1-\alpha)(1+1/|\mathcal{I}_{\text{cal}}|)$ -th empirical quantile of $\{S_i^{\text{CQRFM}}\}_{i \in \mathcal{I}_{\text{cal}}} \cup \{\infty\}$, with:

$$S_i^{\text{CQRFM}} = \max\left(\frac{l(x_i) - y_i}{m(x_i) - l(x_i)}, \frac{y_i - u(x_i)}{u(x_i) - m(x_i)}\right). \quad (30)$$

A.5 MIXTURE DISTRIBUTION WITH 1D INPUT

Following Luo & Zhou (2025), we generate $n = 10000$ samples where predictors X_i are drawn independently from $\text{Uniform}(0, 1)$. The response variables are sampled i.i.d. according to:

$$y \sim \text{Triangular}(0, x, x), \quad (31)$$

with conditional density:

$$f(y|x) = \frac{2y}{x^2} \mathbb{1}\{y \in (0, x)\}. \quad (32)$$

A.6 MIXTURE DISTRIBUTION WITH 2D INPUT

We generate $n = 10000$ samples where 2D predictors $X_i = (X_{i1}, X_{i2})$ are drawn independently from $\text{Uniform}(-1, 1)^2$. The response variable Y_i is then sampled according to the quadrant-specific distributions:

$$Y_i|X_i \sim p(y|x_{i1}, x_{i2}), \quad (33)$$

where the conditional distribution $p(y|x_{i1}, x_{i2})$ is defined as:

$$p(y|x_{i1}, x_{i2}) = \begin{cases} \text{Uniform}(0, 1), & \text{if } x_{i1} \geq 0, x_{i2} \geq 0 \\ \text{Normal}(0, (\frac{1}{5})^2), & \text{if } x_{i1} < 0, x_{i2} \geq 0 \\ \text{Exponential}(\frac{1}{2}), & \text{if } x_{i1} < 0, x_{i2} < 0 \\ \frac{1}{2} \cdot \text{Normal}(-\frac{1}{2}, (\frac{1}{10})^2) + \frac{1}{2} \cdot \text{Normal}(\frac{1}{2}, (\frac{1}{10})^2). & \text{if } x_{i1} \geq 0, x_{i2} < 0 \end{cases} \quad (34)$$

We have also visualized the simulation experiment results (synthetic²) for multiple features. We provide two perspectives: a front view and a back view. It is evident from both views that CMR occupies a smaller proportion of the entire space. In this dataset, CHR exhibits lower tail coverage (as indicated by fewer green stars), while CMR demonstrates even lower tail coverage compared to CQR and CQRFM.

A.7 BOXPLOTS

In this section, we present box plots comparing various conformal regression methods across 14 datasets, evaluating their performance in terms of coverage, size, and tail coverage. For each dataset, each conformal regression method underwent 50 randomized partitioning tests with stratified training/test splits to ensure robustness of results.

918
 919
 920
 921
 922
 923
 924
 925
 926
 927
 928
 929
 930
 931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

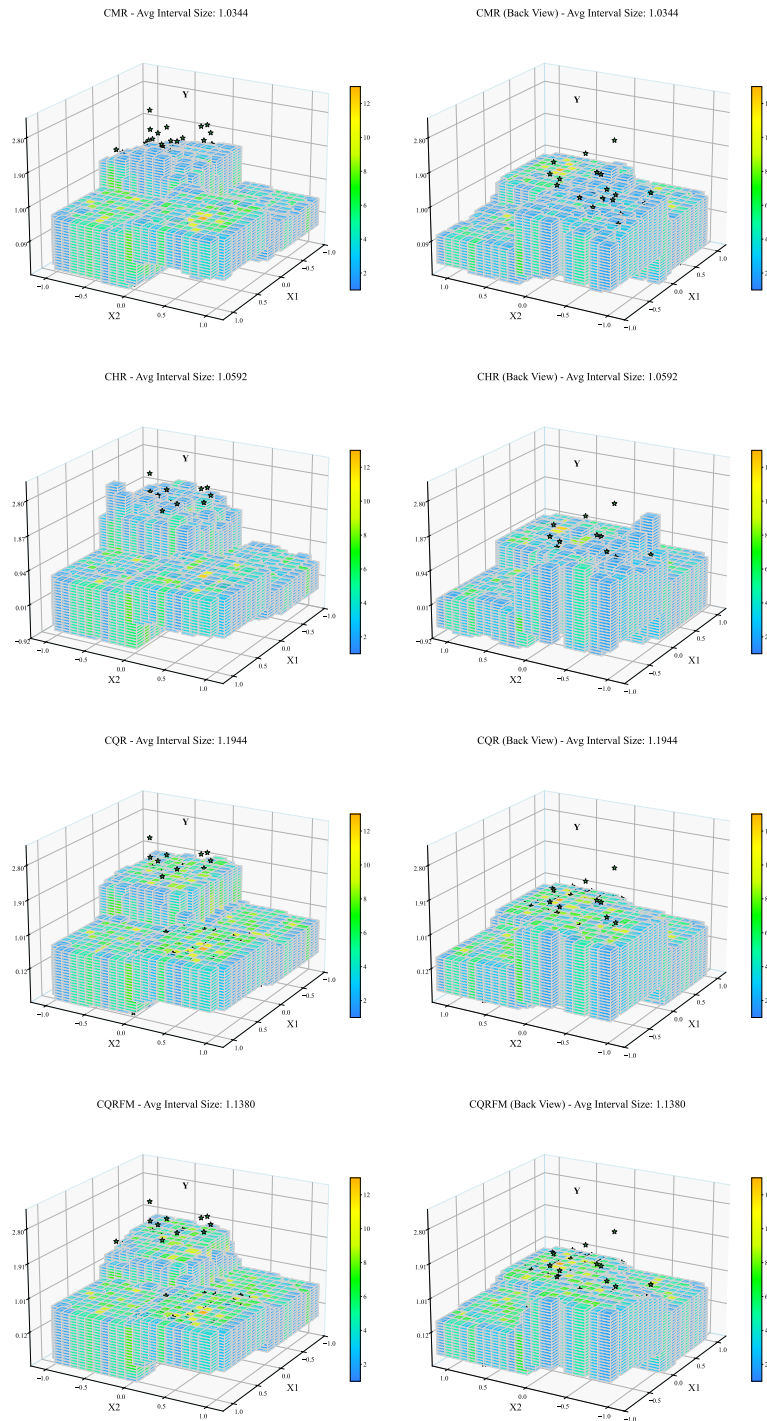


Figure 3: In this figure, we present the three-dimensional prediction intervals and coverage rates for CMR, CHR, CQR, and CQRFM. Similar to Figure 2, we use regional density to reflect the size of the intervals for different methods. Additionally, we provide the tail coverage performance of the models, where the green stars represent the points that are not covered. We offer two perspectives: a front view and a back view. From both views, it can be observed that the region occupied by CMR is smaller compared to the other methods.

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025

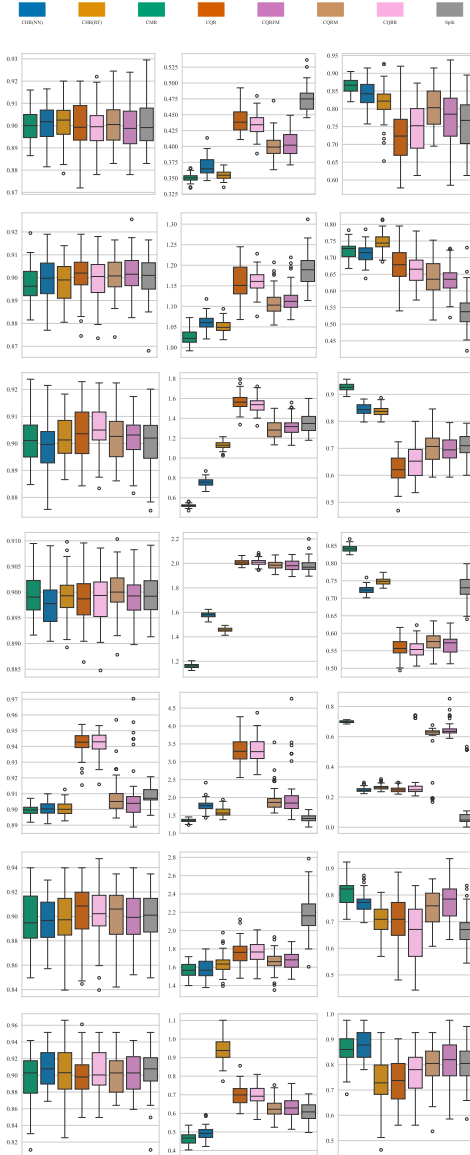


Figure 4: Box plots comparing various conformal regression methods across multiple datasets and evaluation metrics. Each experiment was conducted with 50 random dataset splits. The datasets are arranged vertically from top to bottom: "synthetic¹", "synthetic²", "bike", "bio", "blog", "community", and "concrete". The evaluation metrics, displayed horizontally from left to right, include Coverage, Size, and TCR.

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

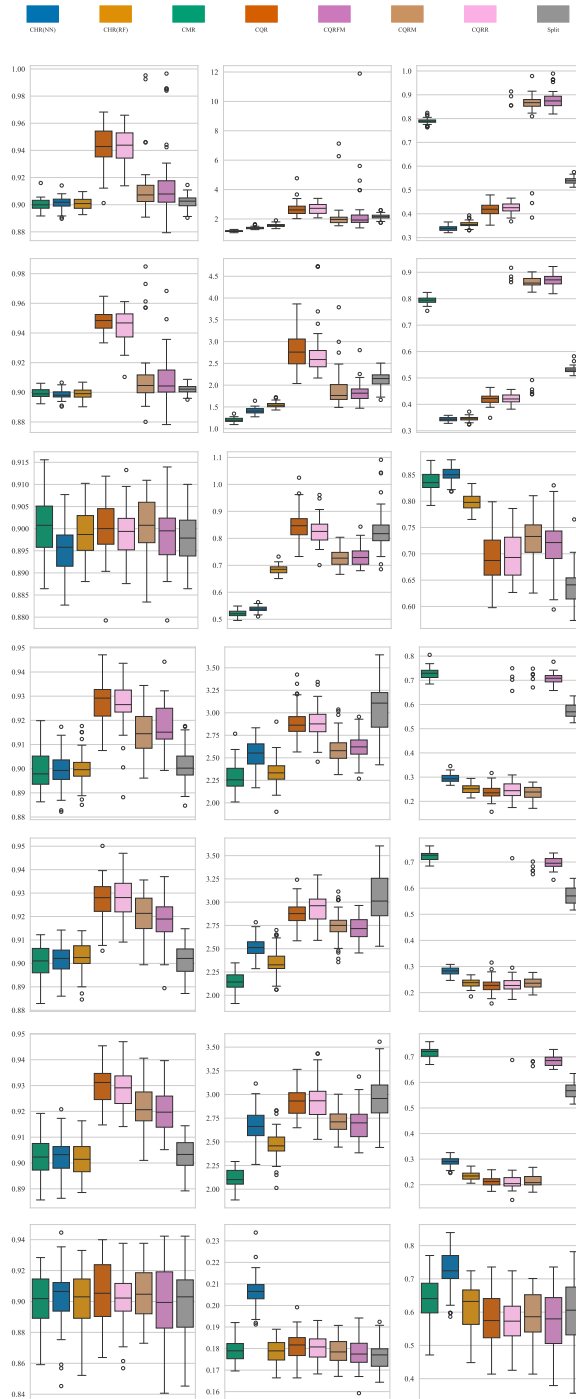


Figure 5: Box plots comparing various conformal regression methods across multiple datasets and evaluation metrics. Each experiment was conducted with 50 random dataset splits. The datasets are arranged vertically from top to bottom: "facebook1", "facebook2", "homes", "meps19", "meps20", "meps21", and "star". The evaluation metrics, displayed horizontally from left to right, include Coverage, Size, and TCR.

A.8 SENSITIVITY ANALYSIS

For synthetic¹ experiments, each parameter configuration was randomly partitioned and evaluated across 10 independent trials, with results averaged to ensure statistical robustness. The hyperparameter grid included $m = [100, 200, 500, 1000, 2000, 5000]$ and $\sigma = [0.001, 0.01, 0.1, 0.5]$. Dataset partitions followed a 5600-sample training set, 2400-sample calibration set, and 2000-sample test set configuration. All experimental protocols (including supplementary experiments) were executed using an NVIDIA GeForce RTX 4060 GPU with 16GB memory.

We denote the best metric in bold for cases with the same m but different σ values, and underline the best metric for cases with the same σ but different m values. As can be observed from the table, the model performs optimally when σ is set to 0.01. As m increases, which corresponds to an increase in the number of generated predictive samples, the size of the prediction set decreases while the tail coverage rate increases. However, this also incurs a certain computational cost in terms of calibration time. Additionally, even when m is set to 100, the model can still maintain sufficient coverage and a relatively small prediction set size.

Table 3: Performance Comparison: CMR Model Performance Under Different m and σ Values

| Dataset | m | σ | Coverage (%) | Interval Size | Tail Coverage (%) | Train Time (s) | Calib Time (s) | Test Time (s) |
|------------------------|------|----------|----------------|----------------------|-----------------------|----------------|----------------|---------------|
| synthetic ¹ | 100 | 0.001 | 89.825 (1.070) | 0.362 (0.008) | 85.050 (1.920) | 59.704 | 83.568 | 74.288 |
| synthetic ¹ | 100 | 0.010 | 89.650 (1.282) | 0.358 (0.008) | 85.575 (1.743) | 58.474 | 82.629 | 74.781 |
| synthetic ¹ | 100 | 0.100 | 89.985 (0.938) | 0.385 (0.008) | 91.150 (1.463) | 59.371 | 82.578 | 75.118 |
| synthetic ¹ | 100 | 0.500 | 89.790 (1.005) | 0.594 (0.004) | <u>85.825 (2.264)</u> | 59.265 | 82.949 | 74.765 |
| synthetic ¹ | 200 | 0.001 | 89.650 (1.046) | 0.354 (0.005) | <u>85.600 (1.586)</u> | 59.560 | 83.027 | 74.687 |
| synthetic ¹ | 200 | 0.010 | 89.610 (0.867) | 0.352 (0.007) | 86.025 (1.429) | 59.006 | 82.983 | 77.334 |
| synthetic ¹ | 200 | 0.100 | 89.840 (0.691) | 0.378 (0.004) | 92.175 (1.401) | 74.496 | 103.748 | 96.643 |
| synthetic ¹ | 200 | 0.500 | 89.350 (0.919) | 0.545 (0.008) | 84.150 (2.762) | 73.409 | 103.822 | 94.213 |
| synthetic ¹ | 500 | 0.001 | 89.665 (1.057) | 0.350 (0.006) | 85.825 (1.323) | 73.252 | 104.748 | 94.268 |
| synthetic ¹ | 500 | 0.010 | 89.575 (1.081) | 0.347 (0.006) | 86.450 (1.308) | 73.261 | 105.462 | 94.479 |
| synthetic ¹ | 500 | 0.100 | 89.885 (0.993) | 0.374 (0.007) | 93.375 (1.398) | 73.586 | 104.831 | 94.194 |
| synthetic ¹ | 500 | 0.500 | 89.215 (1.002) | 0.503 (0.008) | 83.325 (2.620) | 72.968 | 107.867 | 94.514 |
| synthetic ¹ | 1000 | 0.001 | 89.775 (0.879) | 0.348 (0.005) | 86.175 (1.173) | 70.581 | 111.678 | 97.950 |
| synthetic ¹ | 1000 | 0.010 | 89.850 (1.000) | 0.346 (0.005) | 86.975 (1.325) | 69.100 | 111.217 | 96.935 |
| synthetic ¹ | 1000 | 0.100 | 89.930 (1.057) | 0.373 (0.006) | 94.075 (1.475) | 64.220 | 94.174 | 79.935 |
| synthetic ¹ | 1000 | 0.500 | 89.690 (1.141) | 0.485 (0.009) | 83.475 (2.916) | 66.124 | 98.262 | 85.317 |
| synthetic ¹ | 2000 | 0.001 | 89.750 (1.032) | 0.347 (0.006) | 86.700 (1.208) | 66.229 | 100.781 | 83.691 |
| synthetic ¹ | 2000 | 0.010 | 89.845 (1.009) | 0.344 (0.005) | <u>87.350 (1.361)</u> | 57.984 | 88.392 | 74.747 |
| synthetic ¹ | 2000 | 0.100 | 89.900 (1.070) | 0.371 (0.007) | 94.400 (1.184) | 49.222 | 78.070 | 65.324 |
| synthetic ¹ | 2000 | 0.500 | 89.585 (1.101) | 0.468 (0.009) | 82.125 (3.509) | 50.378 | 80.410 | 67.836 |
| synthetic ¹ | 5000 | 0.001 | 89.865 (0.903) | 0.346 (0.006) | 86.700 (0.843) | 70.462 | 116.747 | 87.807 |
| synthetic ¹ | 5000 | 0.010 | 89.820 (1.044) | 0.343 (0.005) | 87.300 (1.150) | 64.484 | 106.798 | 76.157 |
| synthetic ¹ | 5000 | 0.100 | 89.850 (0.956) | <u>0.370 (0.005)</u> | 94.775 (1.186) | 57.479 | 108.471 | 79.874 |
| synthetic ¹ | 5000 | 0.500 | 89.740 (0.981) | <u>0.458 (0.007)</u> | <u>82.025 (3.552)</u> | 60.124 | 123.556 | 81.857 |

A.9 COMPARISON OF TIME CONSUMPTION OF DIFFERENT METHODS

The inference process can be broken down into two stages:

- Base Model Sampling:** Generating m samples for a given test point using the trained CFM model, RF and NN. This cost is a fundamental overhead shared by all CFM - based conformal methods we compared (CHR, CQR, and ours).
- Conformal Interval Generation:** Utilizing these samples to compute the final prediction interval. The cost difference in this stage reflects the intrinsic complexity of each conformal algorithm.

To quantify these costs, we conducted an ablation study using CFM as the fixed base model and measured the time taken by different conformal algorithms. The results are as follows:

A.10 PART 1: CONDITIONAL FLOW MATCHING MODEL

This section defines the core generative model, which learns the conditional distribution $p(y|x)$ using a simulation-free objective.

| Method | Train Time | Calibration Sampling Time | \hat{q} Calculation Time | Test Sampling Time | Prediction Interval Generation Time | Generation time of prediction interval for a single sample |
|-----------|------------|---------------------------|----------------------------|--------------------|-------------------------------------|--|
| CHR (CFM) | 64.220s | 1.368s | 70.230s | 1.104s | 0.107s | 0.606ms |
| CQR (CFM) | 64.325s | 0.207s | 1.245s | 0.925s | 0.030s | 0.478ms |
| CMR (CFM) | 64.224s | 1.346s | 92.264s | 0.988s | 78.560s | 39.774ms |
| CMR (RF) | 10.781s | 0.073s | 0.063s | 0.059s | 0.005s | 0.003ms |
| CMR (NN) | 120.854s | 0.003s | 0.068s | 0.003s | 0.005s | 0.003ms |
| CHR (RF) | 10.792s | 0.049s | 79.102s | 0.048s | 66.006s | 33.003ms |
| CHR (NN) | 120.894s | 0.002S | 80.825s | 0.003s | 66.332s | 33.166ms |
| CQR | 12.183S | 0.002S | 0.046s | 0.003s | 5.183s | 2.592ms |

Table 4: A comparison of the time consumption of different methods. In the simulation experiment, each method was tested 10 times.

```

1146 1 import torch
1147 2 import torch.nn as nn
1148 3 import numpy as np
1149 4
1149 5 class FlowMatchingModel(nn.Module):
1150 6     """Defines the neural network and the training/sampling logic for the CFM model
1151 7     """
1151 8     def __init__(self, context_dim, hidden_dim=64, num_layers=3, sigma=0.01):
1152 9         super().__init__()
1153 10         self.sigma = sigma
1154 11
1154 12         # Neural network to approximate the vector field v_theta(y, t, x)
1155 13         layers = [nn.Linear(1 + 1 + context_dim, hidden_dim), nn.SiLU()]
1156 14         for _ in range(num_layers - 1):
1157 15             layers.extend([nn.Linear(hidden_dim, hidden_dim), nn.SiLU()])
1158 16         layers.append(nn.Linear(hidden_dim, 1))
1159 17         self.net = nn.Sequential(*layers)
1160 18
1160 19     def forward(self, y, t, context):
1161 20         net_input = torch.cat([y, t.reshape(-1, 1), context], dim=1)
1162 21         return self.net(net_input)
1163 22
1163 23     def train_model(self, train_loader, epochs=50, lr=1e-3, device='cpu'):
1164 24         self.to(device)
1165 25         optimizer = torch.optim.Adam(self.parameters(), lr=lr)
1166 26
1166 27         self.train()
1167 28         for epoch in range(epochs):
1168 29             for x_batch, y_batch in train_loader:
1169 30                 x_batch, y_batch = x_batch.to(device), y_batch.to(device)
1170 31                 optimizer.zero_grad()
1171 32
1171 33                 t = torch.rand(x_batch.size(0), 1, device=device)
1172 34                 z0 = torch.randn_like(y_batch)
1173 35
1173 36                 yt = t * y_batch + (1 - t) * z0 + self.sigma * torch.randn_like(
1174 37                 y_batch)
1175 38                 ut = y_batch - z0 # Target vector field
1176 39
1176 40                 vt_pred = self(yt, t, x_batch)
1177 41                 loss = nn.functional.mse_loss(vt_pred, ut)
1178 42
1178 43                 loss.backward()
1179 44                 optimizer.step()
1180 45
1180 46     def sample(self, x_cond, num_samples, steps=100):
1181 47         # Generate samples by solving the ODE from t=0 to t=1
1182 48         self.eval()
1183 49         with torch.no_grad():
1184 50             y_t = torch.randn(num_samples, 1, device=next(self.parameters()).device)
1185 51
1185 52             dt = 1.0 / steps
1186 53             for t_val in np.linspace(0, 1 - dt, steps):
1187 54                 t = torch.full((num_samples, 1), t_val, device=y_t.device)
1188 55                 context = x_cond.repeat(num_samples, 1)
1189 56                 dy = self(y_t, t, context) * dt
1190 57                 y_t = y_t + dy
1191 58             return y_t.cpu().numpy().flatten()

```

Listing 1: Flow Matching Model: The Generative Core

A.11 PART 2: CONFORMAL REGRESSION ALGORITHM

This section implements the conformalization procedure. It uses a trained Flow Matching Model to perform calibration and construct prediction intervals.

```

1192 1 class CMR:
1193 2     """Implements the calibration and prediction logic for Conformal Regression."""
1194 3
1195 4     def __init__(self, model):
1196 5         self.model = model
1197 6         self.q_threshold = None # Calibrated quantile
1198 7
1199 8     def calibrate(self, cal_loader, alpha, m_samples=1000):
1200 9         # Phase 2: Calibration to find the threshold q
1201 10         scores = []
1202 11         device = next(self.model.parameters()).device
1203 12         for x_cal, y_cal in cal_loader:
1204 13             for i in range(x_cal.size(0)):
1205 14                 xi, yi = x_cal[i:i+1].to(device), y_cal[i].item()
1206 15
1207 16                 y_samples = self.model.sample(xi, m_samples)
1208 17                 y_samples.sort()
1209 18
1210 19                 # Nonconformity score: min number of points in the shortest
1211 20 interval covering yi
1212 21                 score = self._find_min_points_to_cover(y_samples, yi)
1213 22                 scores.append(score)
1214 23
1215 24         n_cal = len(scores)
1216 25         q_level = np.ceil((1 - alpha) * (n_cal + 1)) / n_cal
1217 26         self.q_threshold = int(np.quantile(scores, q_level, method="higher"))
1218 27
1219 28     def predict(self, x_test, m_samples=1000):
1220 29         # Phase 3: Prediction Interval Construction
1221 30         if self.q_threshold is None:
1222 31             raise RuntimeError("CMR must be calibrated first. Call .calibrate()")
1223 32
1224 33         device = next(self.model.parameters()).device
1225 34         y_samples = self.model.sample(x_test.to(device), m_samples)
1226 35         y_samples.sort()
1227 36
1228 37         # Find the shortest interval containing q_threshold points
1229 38         interval = self._find_shortest_interval(y_samples, self.q_threshold)
1230 39         return interval
1231 40
1232 41     def _find_min_points_to_cover(self, sorted_samples, y_true):
1233 42         # Find the smallest k such that the shortest interval of k points covers
1234 43 y_true
1235 44         for k in range(2, len(sorted_samples) + 1):
1236 45             shortest_interval = self._find_shortest_interval(sorted_samples, k)
1237 46             if shortest_interval[0] <= y_true <= shortest_interval[1]:
1238 47                 return k
1239 48         return len(sorted_samples)
1240 49
1241 50     def _find_shortest_interval(self, sorted_samples, k_points):
1242 51         # Find the shortest interval containing exactly k_points using a sliding
1243 52 window
1244 53         k_points = min(k_points, len(sorted_samples))
1245 54         min_width = float('inf')
1246 55         best_interval = (sorted_samples[0], sorted_samples[-1])
1247 56
1248 57         for i in range(len(sorted_samples) - k_points + 1):
1249 58             width = sorted_samples[i + k_points - 1] - sorted_samples[i]
1250 59             if width < min_width:
1251 60                 min_width = width
1252 61                 best_interval = (sorted_samples[i], sorted_samples[i + k_points -
1253 62 1])
1254 63         return best_interval

```

Listing 2: CMR: The Conformal Prediction Algorithm