

CONFORMAL MIRROR STATISTICS FOR MODEL ALIGNMENT: UNCERTAINTY QUANTIFICATION WITH FDR CONTROL

Anonymous authors

Paper under double-blind review

ABSTRACT

Foundation models are increasingly adopted across diverse domains, but their safe deployment requires outputs that align with human interpretation, especially in high-stakes applications. This motivates the need for rigorous uncertainty quantification (UQ) methods to assess alignment reliability. Most existing methods rely on large labeled datasets, limiting their applicability in real-world settings where labeled data is scarce or expensive. In this paper, we introduce Conformal Mirror Statistics (CMS), a novel framework for UQ in model alignment, selecting aligned outputs for unlabeled data with the false discovery rate (FDR) under control. Unlike conventional conformal methods based on p -value calibration, CMS generalizes to broader settings without restrictive calibration size requirements. We further establish theoretical guarantees by proving FDR control under weaker data assumptions than existing methods. Empirical results on simulations and a large sepsis cohort from MIMIC-III demonstrate that CMS consistently outperforms conventional methods while reliably identifying aligned outputs.

1 INTRODUCTION

With the increasing adoption of deep learning and foundation models across diverse domains, these models are progressively taking over decision-making tasks, reducing human workload. However, due to inherent biases, distribution shifts, and other limitations, these models can produce misleading outputs (Gallegos et al., 2024; Oh et al., 2025; Ranjan et al., 2024). Unchecked reliance on these outputs can lead to significant consequences, particularly in high-stakes applications such as healthcare diagnostics, legal decision-making, and financial forecasting (Hager et al., 2024; Dahl et al., 2024; Chen et al., 2025). Ensuring that foundation model outputs align with human understanding has thus become a critical concern. This underscores the necessity of robust uncertainty quantification (UQ) techniques to assess and enhance the reliability of model alignment, mitigating potential risks associated with their deployment.

While classical UQ approaches such as confidence intervals and machine learning methods provide theoretical guarantees, they rely on strong modeling assumptions and do not extend to black-box predictors (Efron, 1981; Casella & Berger, 2002; Gelman et al., 2013). Conformal inference addresses these limitations by offering distribution-free guarantees for prediction-level reliability, and has been applied to multiple testing, outlier detection, and model alignment (Jin & Candès, 2023b; Gui et al., 2024). However, existing conformal approaches often require large labeled calibration sets, which are costly and time-consuming to obtain in practice. This limitation is particularly acute in domains such as medicine, where expert annotations are scarce and expensive, or in finance and law, where ground-truth labels may be ambiguous or delayed (Olatunji et al., 2019; Hendrycks et al., 2021; Choi et al., 2025). Moreover, their validity and power can degrade under complex data structures which are common in modern large-scale dataapplications.

To address these limitations, we propose *Conformal Mirror Statistics (CMS)*, a novel framework for uncertainty-aware model alignment with minimal labeled data. Figure 1 illustrates the workflow: given an unlabeled dataset, the CMS algorithm evaluates whether the outputs of a model (e.g., transformer) are reliable. The construction of CMS leverages both magnitude and rank information to achieve asymptotic symmetry for false discovery rate (FDR) control. [This enables thresholds to be](#)

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

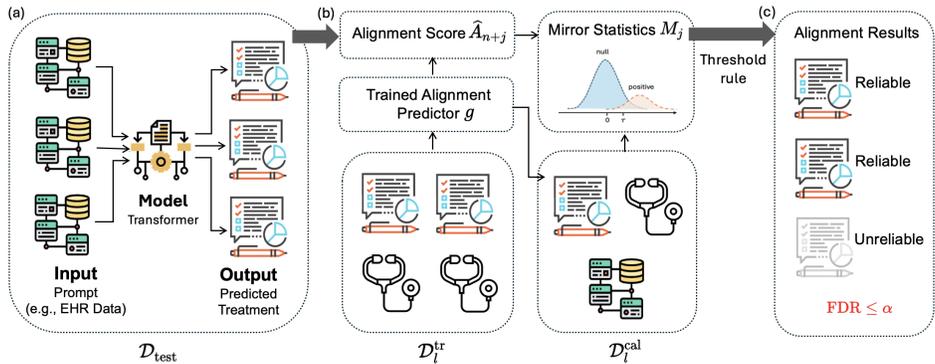


Figure 1: Workflow of the Conformal Mirror Statistics (CMS) algorithm. (a) Unlabeled test dataset consisting of model outputs awaiting evaluation. (b) CMS procedure: a small labeled dataset is used to train an alignment predictor and construct mirror statistics. (c) Alignment results on the test outputs, indicating which predictions are deemed unreliable.

determined directly from statistic magnitudes, rather than through the BH procedure on conformal p -value orderings, thereby preserving statistical power even in the presence of limited labeled data. Additionally, the method accommodates complex data structures, including non-exchangeability and heteroskedasticity. Unlike the original mirror statistics in Dai et al. (2023b), which are constructed from regression coefficients and require additional restrictions on both the coefficients and the statistics, our approach exploits the properties of alignment scores to build a new form of mirror statistics. As a result, our method achieves FDR control under weaker assumptions.

Empirically, we demonstrate that CMS achieves reliable alignment selection on heteroskedastic simulation settings and a large-scale sepsis cohort from MIMIC-III, outperforming conventional conformal methods in terms of FDR control and stability. More broadly, the framework is applicable to certifying outputs across a wide range of real world tasks, making it well suited for the growing landscape of foundation model applications in high-stakes domains.

1.1 MAIN RESULTS

We summarize our main contributions and findings as follows:

- **Remain power on large unlabeled datasets:** Our statistics mitigates the loss of power in existing methods that occurs when the labeled calibration set is small relative to the unlabeled test set size.
- **Greater tolerance to complex data structures:** The proposed statistic maintains FDR control under more complex data structures such as non-exchangeable data than previous conformal approaches, as supported by theory and simulation.
- **Novel statistic with weaker assumptions:** We construct a new form of mirror statistics for conformal inference and prove FDR control under weaker conditions than those required in prior mirror statistics work (Dai et al., 2023b). This broadens the theoretical foundation of methods for conformalized selection.
- **Validated on practical tasks:** We validate CMS on simulations and a large-scale MIMIC-III sepsis cohort, where it outperforms baseline methods in FDR control. More broadly, CMS provides a general framework for certifying reliable outputs across diverse machine learning applications.

2 RELATED LITERATURE

Conformal prediction and conformal inference. Conformal methods provide distribution-free guarantees and have become a powerful tool for predictive uncertainty quantification. Conformal prediction (Shafer & Vovk, 2008; Lei et al., 2018; Angelopoulos & Bates, 2021) generates prediction sets with guaranteed coverage. Conformal inference instead focuses on selecting reliable outputs using conformal p -values (Jin & Candès, 2023b;a; Liang et al., 2024; Gui et al., 2024). While

effective, these methods require large labeled calibration sets, which remain costly and limit practical adoption.

FDR control and mirror statistics. The Benjamini–Hochberg procedure (Benjamini & Hochberg, 1995) established FDR control for p -values, and knockoffs (Barber & Candès, 2015) extended this idea to more flexible settings. Recently, mirror statistics (Dai et al., 2023b) have emerged as a robust tool for selective inference, achieving exact FDR control via data splitting. Subsequent work has further generalized mirror statistics across diverse tasks (Tong et al., 2023; Dai et al., 2023a), highlighting their robustness and versatility.

3 METHODOLOGY

3.1 PROBLEM SETTINGS AND PRELIMINARY

Suppose we have a fixed pre-trained foundation model $f : \mathcal{X} \rightarrow \mathcal{Y}$ that maps an input prompt to an output. We further assume access to a labeled dataset $\mathcal{D}_l = \{(X_i, L_i)\}_{i=1}^n$, where $X_i \in \mathcal{X}$ denotes the input prompt and $L_i \in \mathcal{L}$ is a reference label used to evaluate alignment. For example, consider a pretrained model that generates treatments from electronic health record (EHR) data: here X_i corresponds to a patient’s EHR, $f(X_i)$ is the model-generated treatment plan, and L_i is the expert-provided diagnosis.

To evaluate whether the foundation model output $f(X_i)$ aligns with the true label L_i , we introduce an alignment function $\varphi : \mathcal{Y} \times \mathcal{L} \rightarrow \mathbb{R}$. This function compares the generated output $f(X_i)$ against the reference label L_i and returns an alignment score $A_i = \varphi(f(X_i), L_i)$. In the treatment plan generation example, A_i could represent a similarity measure between the LLM-generated plan and the expert-authored plan. Throughout this paper, given a suitable choice of alignment function, we treat A_i as the true alignment score. If $A_i \leq c$, the output $f(X_i)$ is regarded as misaligned with L_i ; if $A_i > c$, it is considered aligned, where $c \in \mathbb{R}$ is a pre-specified threshold.

In addition to the labeled data, we also have an unlabeled test dataset $\mathcal{D}_{\text{test}} = \{X_{n+j}\}_{j=1}^m$. Note that for a unit data X_{n+j} without reference label information, the true alignment score of its generated output is unknown, and we do not seek to evaluate it, which often requires expert annotation or human judgment. Instead, the goal of this paper is to select a subset $\hat{S} \subseteq [m] := \{1, \dots, m\}$ such that most of their (unobserved) true alignment scores are above the pre-fixed threshold $c \in \mathbb{R}$, with the FDR under control. Let $S = \{j \in [m] : A_{n+j} > c\}$ be the true alignment set. Formally, we measure reliability by the false discovery proportion (FDP) and its expectation, the FDR:

$$\text{FDP} = \frac{\sum_{j \in [m]} I\{A_{n+j} \leq c, j \in \hat{S}\}}{\max(|\hat{S}|, 1)} \quad \text{with} \quad \text{FDR} = \mathbb{E}(\text{FDP}),$$

where $I(\cdot)$ denotes the indicator function, and $|\cdot|$ represents the cardinality of a set. In particular, we aim to enforce that $\text{FDR} \leq \alpha$ for a pre-specified level $\alpha \in (0, 1)$. The FDR represents the average proportion of selected units that do not meet alignment criteria, offering a direct measure of the risk incurred when deploying outputs in \hat{S} . Apart from controlling the FDR, it is also desirable to select as many aligned units as possible, which corresponds to maximizing the power:

$$\text{Power} = \mathbb{E} \left[\frac{\sum_{j \in [m]} I\{A_{n+j} > c, j \in \hat{S}\}}{\max\left(\sum_{j \in [m]} I\{A_{n+j} > c\}, 1\right)} \right].$$

3.2 CONFORMAL MIRROR STATISTICS FOR CONFORMAL ALIGNMENT

Following the conformalized selection framework (Jin & Candès, 2023a;b), we formalize alignment detection as a multiple testing problem. For each test unit $j \in [m]$, this induces the unit-level hypotheses

$$H_{0,j} : A_{n+j} \leq c, \quad H_{1,j} : A_{n+j} > c.$$

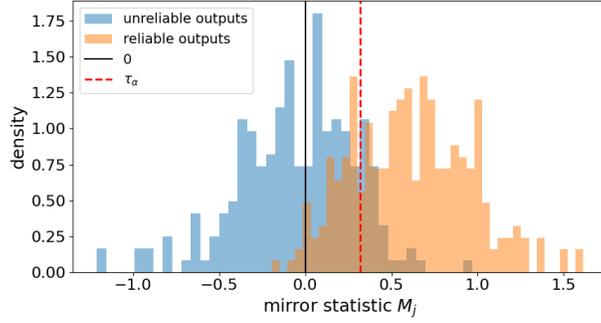


Figure 2: Empirical density of CMS. Outputs with mirror statistics above the cutoff τ_α are selected.

We define the (true) null index set and alternative index set as

$$\mathcal{H}_0 := \{j \in [m] : A_{n+j} \leq c\}, \quad \mathcal{H}_1 := \{j \in [m] : A_{n+j} > c\}.$$

Rejecting $H_{0,j}$ provides evidence that the alignment score of unit j exceeds the threshold c . Thus aligned-unit selection reduces to simultaneously testing the hypotheses $\{H_{0,j}\}_{j=1}^m$.

Given the labeled data \mathcal{D}_l with reference information, we begin by splitting \mathcal{D}_l into two disjoint parts: a training set $\mathcal{D}_l^{\text{tr}}$ and a calibration set $\mathcal{D}_l^{\text{cal}}$. A prediction model $g : \mathcal{X} \rightarrow \mathbb{R}$ is then fitted on $\mathcal{D}_l^{\text{tr}}$ to estimate the alignment score from the prompt X and label L , which may also incorporate information from f . This yields predicted alignment scores $\hat{A}_i = g(X_i)$ for $i \in [n]$ and $\hat{A}_{n+j} = g(X_{n+j})$ for $j \in [m]$. The calibration pairs $\{(A_i, \hat{A}_i)\}_{i \in \mathcal{D}_l^{\text{cal}}}$ are subsequently used to construct the selection rule.

Limitations of conventional methods. Conformal p -values (Jin & Candès, 2023b; Gui et al., 2024) take the form

$$p_j = \frac{1 + \sum_{i \in \mathcal{D}_{\text{cal}}} I\{A_i \leq c, \hat{A}_i \geq \hat{A}_{n+j}\}}{|\mathcal{D}_{\text{cal}}| + 1},$$

and apply the BH procedure (Benjamini & Hochberg, 1995) to identify reliable sets:

$$\hat{S} = \left\{ j \in [m] : p_j \leq \frac{\alpha k^*}{m} \right\}, \quad k^* = \max \left\{ k \in [m] : p_{(k)} \leq \frac{\alpha k}{m} \right\},$$

where $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$ denote the p -values sorted in ascending order. From Gui et al. (2024), conformal p -value in this setting is valid only under the exchangeability assumption. This requirement is rarely satisfied in real applications where calibration and test units typically come from heterogeneous or non-exchangeable sources. Moreover, each p_j takes one of the discrete levels $\{(r+1)/(n_{\text{cal}}+1)\}_{r=0}^{n_{\text{cal}}}$. When m is large relative to n_{cal} , these coarse p -values cannot distinguish fine differences among test scores, leading to substantial power loss in alignment problems with large test sets.

To overcome these limitations, we propose the Conformal Mirror Statistic, which incorporates both rank information and the magnitude differences between predicted scores. Specifically, for each $j \in [m]$ we define

$$M_j = \frac{\sum_{i \in \mathcal{D}_l^{\text{cal}}} I\{A_i \leq c\} (\hat{A}_{n+j} - \hat{A}_i)}{\sum_{i \in \mathcal{D}_l^{\text{cal}}} I\{A_i \leq c\}}. \quad (1)$$

Figure 2 shows the density distribution of conformal mirror statistics in simulated normal distribution data. The key intuition behind equation 1 is as follows. If a generated output for the test unit X_{n+j} is unlikely to be aligned (i.e., $A_{n+j} \leq c$), then conditional on $A_i \leq c$, the distribution of $\hat{A}_{n+j} - \hat{A}_i$ should be symmetric around zero. Conversely, if X_{n+j} is aligned (i.e., $A_{n+j} > c$), we expect $\hat{A}_{n+j} - \hat{A}_i$ to have a significantly positive magnitude given that $A_i \leq c$. This construction leads to two desirable properties of M_j :

- 216 (i) Under the null, the distribution of M_j is symmetric about zero.
 217 (ii) Larger values of M_j indicate that the test unit X_{n+j} is more likely to be aligned.
 218

219 The symmetry guaranteed by Property (i) is particularly useful: it allows us to directly control the
 220 FDP without resorting to the BH procedure based on conformal p -values ordering. Consequently,
 221 our method does not suffer from the previously mentioned restriction that the test size m must
 222 greatly exceed the calibration size n . We will rigorously establish these properties later.

223 Building on Property (i), for any threshold $\tau > 0$ the number of false positives under the null
 224 hypothesis can be controlled as

$$225 \#\{j \in \mathcal{H}_0 : M_j > \tau\} \approx \#\{j \in \mathcal{H}_0 : M_j < -\tau\} \lesssim \#\{j : M_j < -\tau\}, \quad (2)$$

227 where $\#\{\cdot\}$ denotes the cardinality of a set. Consequently, this symmetry allows us to derive an
 228 approximate upper bound for the FDP directly:

$$230 \text{FDP}(\tau) = \frac{\#\{j \in \mathcal{H}_0 : M_j > \tau\}}{\#\{j : M_j > \tau\} \vee 1} \lesssim \frac{\#\{j : M_j < -\tau\}}{\#\{j : M_j > \tau\} \vee 1}.$$

233 Together with Property (ii), this motivates using magnitude information rather than rank to select
 234 the discovery set \hat{S} at a target FDR level $\alpha \in (0, 1)$ as

$$235 \hat{S} = \{j : M_j > \tau_\alpha\},$$

237 where the threshold τ_α is determined by

$$238 \tau_\alpha = \min\left\{\tau > 0 : \widehat{\text{FDP}}(\tau) := \frac{\#\{j : M_j < -\tau\}}{\#\{j : M_j > \tau\} \vee 1} \leq \alpha\right\}. \quad (3)$$

242 The complete procedure is summarized in Algorithm 1. For Step 3, If the labeled data \mathcal{D}_l are
 243 exchangeable, we adopt a random split.

244 **Algorithm 1** Conformal mirror statistics for testing alignment with FDR control

246 **Require:** Pre-trained foundation model f ; alignment score function φ ; labeled dataset $\mathcal{D}_l =$
 247 $\{X_i, L_i\}_{i=1}^n$; test dataset $\mathcal{D}_{\text{test}} = \{X_{n+j}\}_{j=1}^m$; method for fitting prediction model g ; align-
 248 ment level c ; target FDR level α .

249 **Ensure:** The final selected units set \hat{S} .

- 250 1: Compute the alignment score $A_i = \varphi(f(X_i), L_i), \forall i \in \mathcal{D}_l$.
 - 251 2: Partition \mathcal{D}_l into two disjoint sets: the training set $\mathcal{D}_l^{\text{tr}}$ and the calibration set $\mathcal{D}_l^{\text{cal}}$.
 - 252 3: Fit the alignment score predictor g based on training dataset $\mathcal{D}_l^{\text{tr}}$.
 - 253 4: Compute the predicted alignment score: $\hat{A}_i \leftarrow g(X_i), \forall i \in \mathcal{D}_l^{\text{cal}}$, and $\hat{A}_{n+j} \leftarrow$
 254 $g(X_{n+j}), \forall j \in \mathcal{D}_{\text{test}}$.
 - 255 5: **for** $j \in [m]$ **do**
 - 256 6: Compute the conformal mirror statistics M_j according to equation 1.
 - 257 7: **end for**
 - 258 8: For a given FDR level $\alpha \in (0, 1)$, determine the threshold τ_α by equation 3.
 - 259 9: Select the units set by $\hat{S} = \{j : M_j > \tau_\alpha\}$.
-

261 Next, we provide formal justifications for Properties (i) and (ii) and rigorously establish FDR control.
 262 Property (ii) follows directly from the definition in equation 1, while Property (i) requires the
 263 following assumptions.

264 **Assumption 3.1** (Asymptotic Symmetry). Fix a threshold $c \in \mathbb{R}$ with $\mathbb{P}(A_{n+j} \leq c) > 0$ for all
 265 $j \in [m]$. For each $X_j \in \mathcal{D}_l^{\text{cal}} \cup \mathcal{D}_{\text{test}}$, let $G_{0,j}$ denote the distribution of $g(X_j) - \mu_c$ conditional on
 266 $H_{0,j}$, where $\mu_c := \mathbb{E}[g(X_j) | H_{0,j}]$, and let $G_{0,j}^-$ be the reflected distribution defined by $G_{0,j}^-(B) =$
 267 $G_{0,j}(-B)$ for any Borel set B . We assume there exists a sequence $\eta_s = \eta_s(n_{\text{cal}}) \rightarrow 0$ as $n_{\text{cal}} \rightarrow \infty$,
 268 such that

$$269 \sup_t |G_{0,j}(t) - G_{0,j}^-(t)| \leq \eta_s, \quad \forall j \in [m].$$

Assumption 3.2 (Uniform moment and tail regularity). *The calibration set $\{(X_i, L_i)\}_{i \in \mathcal{D}_i^{\text{cal}}}$ and the test set $\{(X_{n+j}, L_{n+j})\}_{j \in [m]}$ are independent. The calibration estimated scores $\{\widehat{A}_i\}_{i \in \mathcal{D}_i^{\text{cal}}}$ satisfy a uniform second-moment bound $\sup_{i \in \mathcal{D}_i^{\text{cal}}} \mathbb{E}[\widehat{A}_i^2] < \infty$. For each test index j , define the tail function*

$$Q_j(\delta) := \mathbb{P}(\widehat{A}_{n+j} - \mu_c > \delta).$$

We assume that each Q_j is continuously differentiable in δ and that the derivatives are uniformly bounded:

$$\sup_{j, \delta} |Q'_j(\delta)| < \infty.$$

Remark 3.3 (On the asymptotic symmetry assumption). Assumption 3.1 requires only that the null distribution $G_{0,j}$ of $g(X_j) - \mu_c$ be approximately symmetric in the sense that $\sup_t |G_{0,j}(t) - G_{0,j}^-(t)|$ vanishes with the calibration sample size. This is a mild condition. For example, if $g(X)$ is normally distributed then the symmetry holds exactly. More generally, for any distribution of $g(X)$ one can always construct a monotone transformation $F = S^{-1} \circ G$, where G is the CDF of $g(X) \mid A \leq c$ and S is any symmetric target CDF (e.g. the standard normal). Under this transformation, $F(g(X)) \mid A \leq c$ is exactly symmetric, and because F is monotone our method can be applied after replacing g by $F(g(X))$ without loss of generality.

Remark 3.4 (On distributional assumptions). The i.i.d. setting is only a special case of our Assumption 3.2. CMS requires neither identical distributions nor exchangeability, and our conditions are strictly weaker than those imposed in existing conformal p -value methods. We note that CMS can also remain valid under certain forms of weak dependence (e.g., α -mixing sequences), although for generality and clarity we impose independence in our theoretical development.

We are now ready to establish the asymptotic symmetry of the conformal mirror statistic.

Proposition 3.5 (Asymptotic symmetry of M_j). *Under Assumptions 3.1 and 3.2, define $n_{\text{cal}} := |\mathcal{D}_i^{\text{cal}}|$. Under $H_{0,j}$, the mirror statistic M_j is asymptotically symmetric about 0; that is,*

$$\lim_{n_{\text{cal}} \rightarrow \infty} \mathbb{P}(M_j \leq t \mid H_{0,j}) = 1 - \lim_{n_{\text{cal}} \rightarrow \infty} \mathbb{P}(M_j \leq -t \mid H_{0,j}), \quad \forall t \in \mathbb{R}.$$

With Proposition 3.5, we can prove that Algorithm 1 achieves FDR control:

Theorem 3.6. *Assume that the calibration size n_{cal} grows polynomially with $m_0 := |\mathcal{H}_0|$, i.e. $n_{\text{cal}} \asymp m_0^\beta$ for some $\beta > 0$. For any nominal FDR level $\alpha \in (0, 1)$, assume that there exists a constant $\tau > 0$ such that $\mathbb{P}(\text{FDP}(\tau) \leq \alpha) \rightarrow 1$ as $m \rightarrow \infty$. Then, under Assumptions 3.1 and 3.2, Algorithm 1 satisfies*

$$\text{FDP}(\tau_\alpha) \leq \alpha + o_m(1) \quad \text{and} \quad \limsup_{n_{\text{cal}}, m \rightarrow \infty} \text{FDR}(\tau_\alpha) \leq \alpha.$$

Remark 3.7 (On Asymptotic Guarantees). The result above establishes asymptotic FDR control. Compared with the finite-sample guarantee of classical conformal methods, this reflects a tradeoff: CMS can work with the much weaker non-exchangeable condition (see Remark 3.4), which allows heteroskedasticity in the data. For this reason, we view the absence of a finite-sample theorem as a modest weakness. Moreover, both our simulation studies and real-data analysis show that CMS achieves reliable finite-sample FDR control in practice, suggesting that the asymptotic guarantee captures the method’s practical behavior.

We also provide the following proposition to show that our method achieves the asymptotic power.

Proposition 3.8 (Asymptotic Power of CMS). *Under Assumptions 3.1 and assume that the calibration set $\{(X_i, L_i)\}_{i \in \mathcal{D}_i^{\text{cal}}}$ and the test set $\{(X_{n+j}, L_{n+j})\}_{j \in [m]}$ are independent and identically distributed. Under the i.i.d. population model where H_0 denotes the null population ($A \leq c$) and H_1 the alternative population ($A > c$), let $M = g(X) - \mu_c$ with $\mu_c = \mathbb{E}[g(X) \mid H_0]$, and assume M is continuously distributed.*

At the population level, the false discovery rate associated with a threshold τ is defined as

$$\text{FDR}(\tau) = \frac{\mathbb{P}(M > \tau, A \leq c)}{\mathbb{P}(M > \tau)}.$$

Further assume that $\text{FDR}(\tau)$ is strictly decreasing at the oracle threshold $\tau^*(\alpha) = \inf\{\tau : \text{FDR}(\tau) \leq \alpha\}$ (see Appendix D), so that the CMS threshold τ_α converges in probability to $\tau^*(\alpha)$ as $n_{\text{cal}}, m \rightarrow \infty$. Consequently,

$$\lim_{n_{\text{cal}}, m \rightarrow \infty} \text{Power}(\tau_\alpha) = \mathbb{P}(M > \tau^*(\alpha) \mid H_1), \quad (4)$$

$$\lim_{n_{\text{cal}}, m \rightarrow \infty} \frac{1}{m} \sum_{j=1}^m I\{M_j > \tau_\alpha, A_{n+j} > c\} = \mathbb{P}(M > \tau^*(\alpha), A > c). \quad (5)$$

Remark 3.9 (On the Role of g and Power Equivalence). A more informative g assigns larger values under H_1 and smaller values under H_0 , which makes $\mathbb{P}(M > \tau, A \leq c)$ small. Consequently, the oracle threshold $\tau^*(\alpha)$ required to satisfy the FDR constraint also becomes small. Combined with a larger M under H_1 , the power $\mathbb{P}(M > \tau^*(\alpha) \mid H_1)$ becomes larger. Moreover, the power expression has the same structural form as that of conformal p -value methods, since both approaches result in comparing a single data-dependent score against a threshold. It means that CMS remains valid under non-exchangeable and heterogeneous settings, while in the i.i.d. case its power does not suffer compared to conformal p -value methods. A detailed comparison is provided in Appendix E.

Remark 3.10 (On the Role of g and Power Equivalence).

Generality of the framework. Theorem 3.6 establishes that our method achieves asymptotic FDR control under mild assumptions. Importantly, these guarantees depend only on generic properties of the alignment predictor g , rather than on the specific form of the underlying model f . This model-agnostic design enables the framework to generalize across different architectures and application domains. Moreover, because the procedure is distribution-free and requires only a small labeled dataset, it remains statistically valid even in heterogeneous real-world scenarios. In addition, the estimation accuracy of the alignment scores does not affect FDR validity, as the guarantees rely only on the symmetry property of g rather than its predictive precision. Because a single data split can introduce variance, we also construct a more robust multi-splitting version of CMS. The method and supporting results are provided in Appendix F.

4 SIMULATIONS

In this section, we will demonstrate the performance of CMS compared with baseline method under homoscedastic and heteroscedastic settings.

We generate covariates $X \in \mathbb{R}^{20}$ independently as $X_j \sim \text{Unif}[-1, 1]$. A latent continuous outcome is defined by $A = \mu(X) + \varepsilon$, $\varepsilon \sim N(0, \sigma(X)^2)$, where the nonlinear signal is

$$\mu(X) = 2(X_1 X_2 + X_3^2 + e^{X_4 - 1} - 1).$$

Heterogeneity is introduced through the noise level $\sigma(X)$. We consider two simulation settings:

- (i) a **homoscedastic** case $\sigma(X) = 1.5$,
- (ii) a **heteroscedastic** case $\sigma(X) = \frac{5.5 - |\mu(X)|}{2}$. Although the covariates X are generated i.i.d., the heteroscedastic noise structure makes the alignment scores non-exchangeable, which violates the exchangeability assumption required for conformal p -value methods.

Throughout all experiments we set the true alignment threshold to $c = 0.2$. We generate 2000 samples and assign 80% to the test set, which is a largely unlabeled scenario. To construct the alignment predictor $g(X)$, we fit three standard models on (X, A) : Gradient Boosting Regression (GBR), Random Forest Regression (RF), and Support Vector Regression (SVR).

Figures 3 and 4 compare FDR and power under the homoscedastic setting (i) and the heterogeneous setting (ii). Across all base learners, CMS consistently achieves more stable FDR control and higher power in both settings. This reflects a key practical benefit of mirror statistics: unlike conformal p -values, CMS does not rely on exchangeability assumptions, and therefore remains reliable even when noise levels and signal structures vary across the covariate space. In contrast, conformal p -values exhibit clear FDR inflation and reduced power, especially in Setting (ii), where input-dependent variance breaks the calibration–test comparability required by the ranking-based construction. Overall, CMS is substantially more robust in nonlinear and heterogeneous environments, precisely where the conformal p -value method tends to degrade.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398

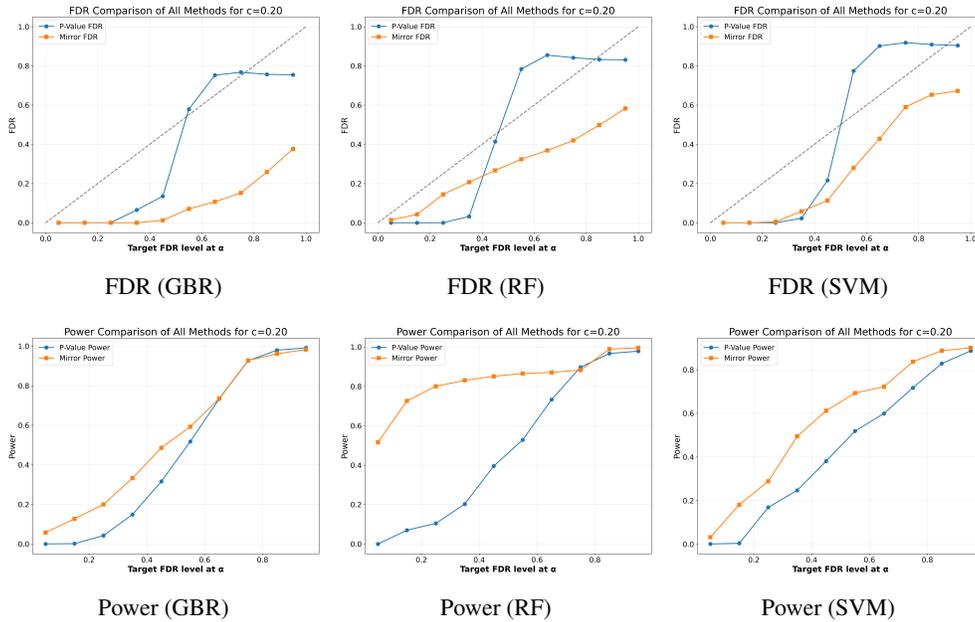


Figure 3: FDR (top) and power (bottom) using GBR, RF, and SVM under Setting (i) $\sigma(X) = 1.5$ at $c = 0.20$.

401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421

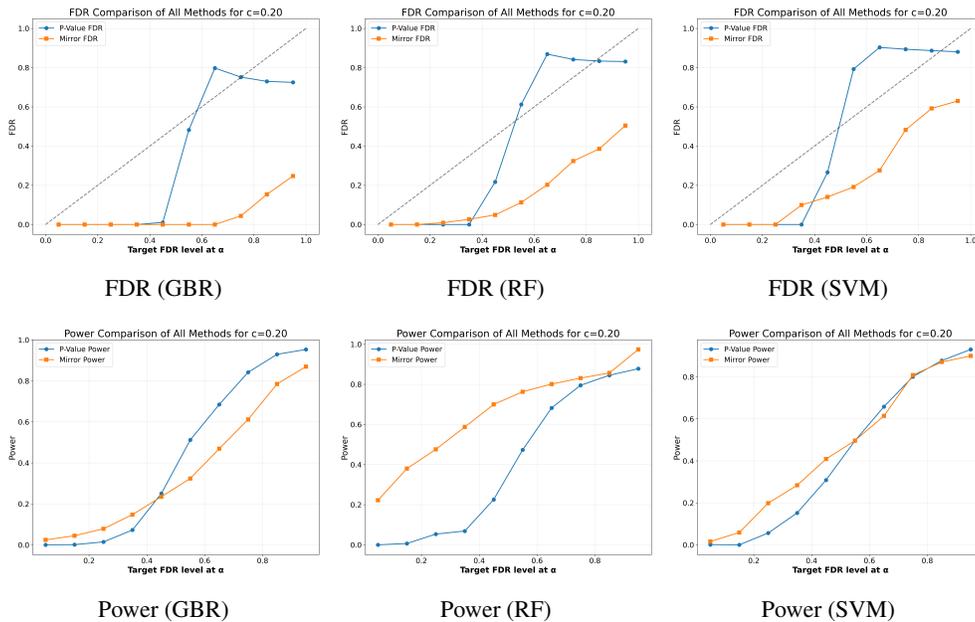


Figure 4: FDR (top) and power (bottom) using GBR, RF, and SVM under Setting (ii) $\sigma(X) = \frac{5.5 - |\mu(X)|}{2}$ at $c = 0.20$.

5 REAL DATA EXPERIMENTS

5.1 EXPERIMENTAL SETUP

422
423
424
425
426
427
428
429
430
431

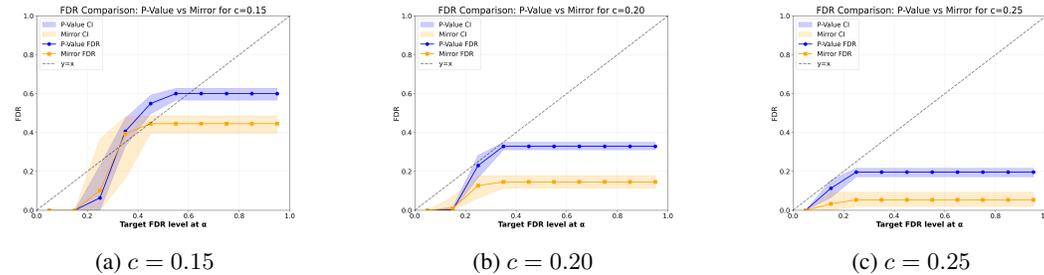
In this section, we will focus on the application on high-stakes clinical decision-making such as dynamic treatment regimes (DTRs). We evaluate the method using a cohort of sepsis patients derived from the Medical Information Mart for Intensive Care (MIMIC)-III dataset, focusing on early-stage

432 treatment planning within a 72-hour clinical window surrounding sepsis onset (Singer et al., 2016).
 433 For each patient we input the multimodel EHR data including the tabular data (e.g., vitals, labs, de-
 434 mographics) and clinical notes into the model. We use the multimodal recommendation model based
 435 on a Transformer architecture (Shen et al., 2025) to generate a treatment $f(X_i)$ among the treatment
 436 space, which consists of 25 discrete treatment combinations $\{T_i\}_{i=1}^{25}$ formed by intravenous fluid
 437 and vasopressor dosages.

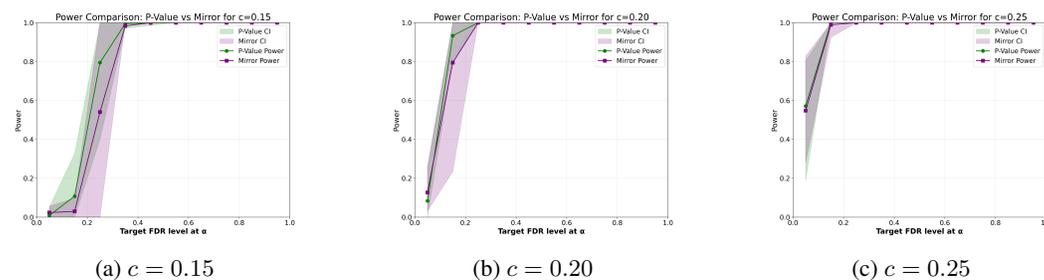
438 It is important to note that, unlike the idealized exchangeable assumption, real-world EHR data are
 439 not independent across patients. Patients treated in the same ICU may share correlated treatment
 440 policies, hospital practices, or clinical environments, leading to heteroskedasticity or weak depen-
 441 dence in their observed trajectories.

443 **Alignment Score $\varphi(X_i)$.** We follow Shen et al. (2025) and define the alignment score as the di-
 444 vergence between the base model’s predictive distribution and that of a student network fine-tuned
 445 only on survivor trajectories. Survivor trajectories are used not as unbiased ground truth but as a
 446 lower-noise reference: their treatment patterns tend to be more stable and less confounded by termi-
 447 nal deterioration than those from deceased patients. This design does not discard or ignore deceased
 448 patients. Instead, the divergence between the two models is explicitly interpreted as an uncertainty
 449 signal, reflecting label ambiguity in noisier mortality trajectories. In this way, the survivor-based
 450 refinement serves as a proxy for identifying unstable supervision rather than introducing selection
 451 bias, and the subsequent FDR-controlled selection steps ensure that such uncertainty is properly
 452 incorporated.

453 5.2 RESULTS ANALYSIS



464 Figure 5: FDR comparison across different c values.



476 Figure 6: Power comparison across different c values.

478 Across 100 independent experiments and a range of uncertainty thresholds c , we compare CMS
 479 with the standard conformal p -value method in terms of both FDR and power. Figure 5 shows the
 480 realized FDR across three representative thresholds $c \in \{0.15, 0.20, 0.25\}$. CMS exhibits stable
 481 and accurate FDR control across all target levels α , with empirical FDR curves closely tracking the
 482 nominal line. In contrast, conformal p -values tend to yield higher FDR, reflecting their sensitivity
 483 to heteroscedasticity and the violation of exchangeability.

484 Figure 6 presents the corresponding power results. The two methods achieve broadly comparable
 485 power across all examined c values, with only minor differences attributable to variance in the align-

486 [ment score or threshold structure. This confirms that CMS does not sacrifice detection ability while](#)
487 [providing more reliable FDR control.](#)

488
489 Together, these findings demonstrate that CMS provides accurate finite-sample FDR control and
490 competitive power across a wide range of uncertainty thresholds. This stability persists even in
491 heterogeneous settings where conformal p -values can become miscalibrated, highlighting the ro-
492 bustness of CMS in practical, non-exchangeable data environments.

493 6 CONCLUSION

494
495 In this work, we introduced Conformal Mirror Statistics (CMS), a flexible and distribution-free
496 framework for model alignment with rigorous FDR control. By leveraging mirror statistics, CMS
497 eliminates the need for large calibration sets and remains valid under far weaker conditions than
498 conventional conformal p -value methods. In particular, our theory accommodates non-exchangeable
499 and heteroskedastic data settings where exchangeability-based conformal p -values typically fail.
500 Empirical studies on simulations and sepsis treatment recommendation further demonstrate that
501 CMS achieves competitive efficiency while maintaining strict FDR guarantees.

502
503 Looking forward, a promising direction is to integrate CMS with modern large-scale machine learn-
504 ing pipelines, for example in reinforcement learning or large language models. Also we could extend
505 CMS to online or streaming environments, where both calibration and test distributions evolve over
506 time. Such extensions would broaden CMS’s applicability in dynamically changing, real-world
507 decision-making systems.

508 ETHICS STATEMENT

509
510 This study uses the de-identified MIMIC-III database, which is accessible to credentialed researchers
511 who have completed the required data use agreement and human-subjects training. No personally
512 identifiable information is included. All experiments were conducted in compliance with institu-
513 tional and ethical guidelines. Our methods are intended for uncertainty quantification and model
514 alignment in predictive modeling. Deployment in clinical practice should be approached cautiously
515 to avoid misuse, such as replacing clinical judgment or reinforcing existing biases.

516 REPRODUCIBILITY STATEMENT

517
518 We provide detailed descriptions of our methodology, assumptions, and hyperparameter settings in
519 the main text and appendix. Pseudocode and theoretical proofs are included to ensure transparency.
520 The source code is currently under preparation for public release, and we aim to make it available in
521 the future. The MIMIC-III data cannot be shared directly due to access restrictions, but researchers
522 with appropriate credentials and data access should be able to reproduce our experiments by follow-
523 ing the provided procedures.
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- 540
541
542 Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and
543 distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- 544
545 Rina Foygel Barber and Emmanuel J. Candès. Controlling the false discovery rate via knockoffs. *The*
546 *Annals of Statistics*, 43(5):2055–2085, 2015. ISSN 00905364. URL [http://www.jstor.](http://www.jstor.org/stable/43818570)
547 [org/stable/43818570](http://www.jstor.org/stable/43818570).
- 548
549 Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: A practical and powerful
550 approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodolog-*
551 *ical)*, 57(1):289–300, 1995. ISSN 00359246. URL [http://www.jstor.org/stable/](http://www.jstor.org/stable/2346101)
552 [2346101](http://www.jstor.org/stable/2346101).
- 553
554 George Casella and Roger L Berger. *Statistical Inference*, volume 2. Duxbury Pacific Grove, CA,
555 2nd ed. edition, 2002.
- 556
557 Zichen Chen, Jiaao Chen, Jianda Chen, and Misha Sra. Standard benchmarks fail – auditing
558 llm agents in finance must prioritize risk, 2025. URL [https://arxiv.org/abs/2502.](https://arxiv.org/abs/2502.15865)
559 [15865](https://arxiv.org/abs/2502.15865).
- 560
561 Chanyeol Choi, Jihoon Kwon, Jaeseon Ha, Hojun Choi, Chaewoon Kim, Yongjae Lee, Jy-yong
562 Sohn, and Alejandro Lopez-Lira. Finder: Financial dataset for question answering and evaluating
563 retrieval-augmented generation, 2025. URL <https://arxiv.org/abs/2504.15800>.
- 564
565 Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E Ho. Large legal fictions: Profiling legal
566 hallucinations in large language models. *Journal of Legal Analysis*, 16(1):64–93, January 2024.
567 ISSN 1946-5319. doi: 10.1093/jla/laae003. URL [http://dx.doi.org/10.1093/jla/](http://dx.doi.org/10.1093/jla/laae003)
568 [laae003](http://dx.doi.org/10.1093/jla/laae003).
- 569
570 Chenguang Dai, Buyu Lin, Xin Xing, and Jun S. Liu. A scale-free approach for false discovery rate
571 control in generalized linear models. *Journal of the American Statistical Association*, 118(543):
572 1551–1565, 2023a.
- 573
574 Chenguang Dai, Buyu Lin, Xin Xing, and S. Jun Liu. False discovery rate control via data splitting.
575 *Journal of the American Statistical Association*, 118(544):2503–2520, 2023b.
- 576
577 Bradley Efron. Nonparametric estimates of standard error: The jackknife, the bootstrap and other
578 methods. *Biometrika*, 68(3):589–599, 1981. ISSN 00063444, 14643510. URL [http://www.](http://www.jstor.org/stable/2335441)
579 [jstor.org/stable/2335441](http://www.jstor.org/stable/2335441).
- 580
581 Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Der-
582 noncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. Bias and fairness in large language
583 models: A survey, 2024. URL <https://arxiv.org/abs/2309.00770>.
- 584
585 Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin.
586 Bayesian data analysis third edition. *Chapman and Hall/CRC*, 2013.
- 587
588 Yu Gui, Ying Jin, and Zhimei Ren. Conformal alignment: Knowing when to
589 trust foundation models with guarantees. In A. Globerson, L. Mackey, D. Bel-
590 grave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural In-*
591 *formation Processing Systems*, volume 37, pp. 73884–73919. Curran Associates, Inc.,
592 2024. URL [https://proceedings.neurips.cc/paper_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/870ccde24673d3970a680bb48496ed63-Paper-Conference.pdf)
593 [file/870ccde24673d3970a680bb48496ed63-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/870ccde24673d3970a680bb48496ed63-Paper-Conference.pdf).
- 594
595 Paul Hager, Friederike Jungmann, Robbie Holland, Kunal Bhagat, Inga Hubrecht, Manuel Knauer,
596 Jakob Vielhauer, Marcus Makowski, Rickmer Braren, Georgios Kaissis, and Daniel Rueck-
597 ert. Evaluation and mitigation of the limitations of large language models in clinical decision-
598 making. *Nature Medicine*, 30(9):2613–2622, September 2024. ISSN 1546-170X. doi:
599 [10.1038/s41591-024-03097-1](https://doi.org/10.1038/s41591-024-03097-1).
- 600
601 Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. Cuad: An expert-annotated nlp dataset
602 for legal contract review, 2021. URL <https://arxiv.org/abs/2103.06268>.

- 594 Ying Jin and Emmanuel J Candès. Model-free selective inference under covariate shift via weighted
595 conformal p-values. *arXiv preprint arXiv:2307.09291*, 2023a.
- 596
- 597 Ying Jin and Emmanuel J Candès. Selection by prediction with conformal p-values. *Journal of*
598 *Machine Learning Research*, 24(244):1–41, 2023b.
- 599
- 600 Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-
601 free predictive inference for regression. *Journal of the American Statistical Association*, 113
602 (523):1094–1111, 2018.
- 603 Ziyi Liang, Matteo Sesia, and Wenguang Sun. Integrative conformal p-values for out-of-distribution
604 testing with labelled outliers. *Journal of the Royal Statistical Society Series B: Statistical Method-*
605 *ology*, 86(3):671–693, 01 2024.
- 606 Changdae Oh, Zhen Fang, Shawn Im, Xuefeng Du, and Yixuan Li. Understanding multimodal
607 LLMs under distribution shifts: An information-theoretic approach. In *Forty-second International*
608 *Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=iA3kafgHG1>.
- 609
- 610 Tobin Olatunji, Li Yao, Ben Covington, Alexander Rhodes, and Anthony Upton. Caveats in generat-
611 ing medical imaging labels from radiology reports, 2019. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1905.02283)
612 [1905.02283](https://arxiv.org/abs/1905.02283).
- 613
- 614 Rajesh Ranjan, Shailja Gupta, and Surya Narayan Singh. A comprehensive survey of bias in llms:
615 Current landscape and future directions, 2024. URL [https://arxiv.org/abs/2409.](https://arxiv.org/abs/2409.16430)
616 [16430](https://arxiv.org/abs/2409.16430).
- 617 Zhimei Ren and Rina Foygel Barber. Derandomised knockoffs: leveraging e-values for false dis-
618 covery rate control. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 86
619 (1):122–154, 2024.
- 620
- 621 Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning*
622 *Research*, 9(3), 2008.
- 623
- 624 Yishan Shen, Yuyang Ye, Hui Xiong, and Yong Chen. Safer: A calibrated risk-aware multimodal
625 recommendation model for dynamic treatment regimes, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2506.06649)
626 [abs/2506.06649](https://arxiv.org/abs/2506.06649).
- 627 Mervyn Singer, Clifford S. Deutschman, Christopher Warren Seymour, Manu Shankar-Hari, Djillali
628 Annane, Michael Bauer, Rinaldo Bellomo, Gordon R. Bernard, Jean-Daniel Chiche, Craig M.
629 Coopersmith, Richard S. Hotchkiss, Mitchell M. Levy, John C. Marshall, Greg S. Martin,
630 Steven M. Opal, Gordon D. Rubenfeld, Tom van der Poll, Jean-Louis Vincent, and Derek C.
631 Angus. The third international consensus definitions for sepsis and septic shock (sepsis-3).
632 *JAMA*, 315(8):801–810, 02 2016. ISSN 0098-7484. doi: 10.1001/jama.2016.0287. URL
633 <https://doi.org/10.1001/jama.2016.0287>.
- 634 Zhaoxue Tong, Zhanrui Cai, Songshan Yang, and Runze Li. Model-free conditional feature screen-
635 ing with fdr control. *Journal of the American Statistical Association*, 118(544):2575–2587, 2023.
- 636
- 637 Vladimir Vovk and Ruodu Wang. E-values: Calibration, combination and applications. *The Annals*
638 *of Statistics*, 49(3):1736–1754, 2021.
- 639
- 640 Ruodu Wang and Aaditya Ramdas. False discovery rate control with e-values. *Journal of the Royal*
641 *Statistical Society Series B: Statistical Methodology*, 84(3):822–852, 2022.
- 642
- 643
- 644
- 645
- 646
- 647

A USE OF LARGE LANGUAGE MODELS (LLMs)

During the preparation of this manuscript, we used an LLM-based writing assistant (OpenAI’s ChatGPT) for language polishing. Specifically, the tool was employed to improve grammar, rephrase sentences for clarity, and refine the overall readability of the text. All ideas, technical content, theoretical results, and experimental analyses are solely the work of the authors, who take full responsibility for the correctness and integrity of the paper.

B PROOF OF PROPOSITION 3.5

Proof. Write $\hat{\mu}_c := \frac{\sum_{i \in \mathcal{D}_l^{\text{cal}}} I\{A_i \leq c\} \hat{A}_i}{\sum_{i \in \mathcal{D}_l^{\text{cal}}} I\{A_i \leq c\}}$, we have

$$\begin{aligned} M_j &= \hat{A}_{n+j} - \frac{\sum_{i \in \mathcal{D}_l^{\text{cal}}} I\{A_i \leq c\} \hat{A}_i}{\sum_{i \in \mathcal{D}_l^{\text{cal}}} I\{A_i \leq c\}} \\ &= (\hat{A}_{n+j} - \mu_c) - (\hat{\mu}_c - \mu_c) \\ &= \hat{A}_{n+j} - \mu_c + O_p(n_{\text{cal}}^{-1/2}), \end{aligned}$$

where the last equality comes from Assumption 3.2 and Chebyshev inequality.

Assumption 3.2 implies that the unconditional density Q'_j is uniformly bounded. Since $H_{0,j} : A_{n+j} \leq c$ satisfies $\mathbb{P}(A_{n+j} \leq c) > 0$, conditioning on $H_{0,j}$ only rescales the unconditional density by the factor $1/\mathbb{P}(A_{n+j} \leq c)$, the conditional density under $H_{0,j}$ is also uniformly bounded. Hence each $G_{0,j}$ is Lipschitz with a universal constant L .

Thus for any t ,

$$P(M_j \leq t) = G_{0,j}(t) + O_p(n_{\text{cal}}^{-1/2}), \quad P(M_j \leq -t) = 1 - G_{0,j}^-(t) + O_p(n_{\text{cal}}^{-1/2}).$$

From Assumption 3.1, we have

$$\sup_t |P(M_j \leq t) - (1 - P(M_j \leq -t))| \eta_s + O_p(n_{\text{cal}}^{-1/2}),$$

where $\eta_s \rightarrow 0$ when $n_{\text{cal}} \rightarrow \infty$. This proves the approximate symmetry of M_j . \square

C PROOF OF THEOREM 3.6

We begin by proving an essential proposition, which serves as a key step toward Theorem 3.6.

Proposition C.1. (Weak dependence) *Under Assumptions 3.1 and 3.2, there exist constants $\omega > 0$ and $\gamma \in (0, 2)$ such that*

$$\text{Var} \left(\sum_{j \in \mathcal{H}_0} I(M_j > \tau) \right) \leq \omega m_0^\gamma, \quad \forall \tau \in \mathbb{R}.$$

Proof of Proposition C.1. Let $S(\tau) := \sum_{j \in \mathcal{H}_0} I\{M_j > \tau\}$. By the law of total variance,

$$\text{Var}(S(\tau)) = \mathbb{E}[\text{Var}(S(\tau) \mid \mathcal{D}_l^{\text{cal}})] + \text{Var}(\mathbb{E}[S(\tau) \mid \mathcal{D}_l^{\text{cal}}]).$$

Given $\mathcal{D}_l^{\text{cal}}$, $\Delta := \hat{\mu}_c - \mu_c$ is a constant and $M_j = (\hat{A}_{n+j} - \mu_c) - \Delta$ depends only on the j -th test sample. By Assumption 3.2, the test samples are independent, hence $\{I(M_j > \tau)\}_{j \in \mathcal{H}_0}$ are conditionally independent Bernoulli with success probability $p_{j,\Delta}(\tau) := \mathbb{P}(\hat{A}_{n+j} - \mu_c > \tau + \Delta \mid \mathcal{D}_l^{\text{cal}})$. Therefore,

$$\text{Var}(S(\tau) \mid \mathcal{D}_l^{\text{cal}}) = \sum_{j \in \mathcal{H}_0} p_{j,\Delta}(\tau)(1 - p_{j,\Delta}(\tau)) \leq \frac{m_0}{4},$$

702 and thus $\mathbb{E}[\text{Var}(S(\tau) \mid \mathcal{D}_l^{\text{cal}})] \leq \frac{m_0}{4}$.

703
704 Remember $Q_j(\delta) = \mathbb{P}(\widehat{A}_j - \mu_c > \delta)$, so that $p_{j,\Delta}(\tau) = Q_j(\tau + \Delta)$. By Assumption 3.2, Q_j is
705 differentiable with bounded derivative Q'_j . A first-order expansion yields

$$706 \quad p_{j,\Delta}(\tau) = Q_j(\tau) + Q'_j(\tau)\Delta + r_j(\Delta), \quad |r_j(\Delta)| \leq C \Delta^2,$$

707
708 for some constant C . Hence

$$709 \quad \mathbb{E}[S(\tau) \mid \mathcal{D}_l^{\text{cal}}] = \sum_{j \in \mathcal{H}_0} p_{j,\Delta}(\tau),$$

710
711 and therefore

$$712 \quad \text{Var}(\mathbb{E}[S(\tau) \mid \mathcal{D}_l^{\text{cal}}]) = \text{Var}\left(\sum_{j \in \mathcal{H}_0} p_{j,\Delta}(\tau)\right).$$

713
714 Since each $p_{j,\Delta}(\tau) = Q_j(\tau + \Delta)$ is uniformly Lipschitz in Δ ,

$$715 \quad \text{Var}\left(\sum_{j \in \mathcal{H}_0} p_{j,\Delta}(\tau)\right) \lesssim m_0^2 \text{Var}(\Delta).$$

716
717 Since $\Delta = \widehat{\mu}_c - \mu_c$ is the sample mean of $Y_i := \widehat{A}_i I(A_i \leq c)$ over the $n_{\text{cal},c} := \sum_{i \in \mathcal{D}_{\text{cal}}} I(A_i \leq c)$
718 calibration points with $A_i \leq c$, we have under Assumption 3.2, the calibration estimated scores $\{\widehat{A}_i\}$
719 are independent with a uniform second-moment bound $\sup_i \mathbb{E}[\widehat{A}_i^2] < \infty$. Therefore,

$$720 \quad \text{Var}(\widehat{\mu}_c) = \frac{1}{n_{\text{cal},c}^2} \sum_{i: A_i \leq c} \text{Var}(Y_i) = O\left(\frac{1}{n_{\text{cal},c}}\right),$$

721
722 which gives $\text{Var}(\Delta) = O(1/n_{\text{cal},c})$.

723
724 By the law of large numbers, $n_{\text{cal},c} \asymp n_{\text{cal}} \cdot \mathbb{P}(A \leq c)$, so that

$$725 \quad \text{Var}(\mathbb{E}[S(\tau) \mid \mathcal{D}_l^{\text{cal}}]) = O\left(\frac{m_0^2}{n_{\text{cal}}}\right).$$

726
727 Combining with the conditional variance bound, we obtain

$$728 \quad \text{Var}(S(\tau)) \leq C_1 m_0 + C_2 \frac{m_0^2}{n_{\text{cal}}}.$$

729
730 Since we have the calibration size grows polynomially with m_0 :

$$731 \quad n_{\text{cal}} \asymp m_0^\beta, \quad \beta > 0.$$

732
733 Then the second term satisfies

$$734 \quad \frac{m_0^2}{n_{\text{cal}}} = O(m_0^{2-\beta}),$$

735
736 so that

$$737 \quad \text{Var}(S(\tau)) \lesssim m_0 + m_0^{2-\beta}.$$

738
739 In particular, when $\beta = 1$ (e.g., a fixed split ratio), we obtain $\text{Var}(S(\tau)) \lesssim m_0$, corresponding to
740 $\gamma = 1$. More generally, for any $\beta > 0$, the bound holds with $\gamma = \max\{1, 2 - \beta\} \in (0, 2)$, i.e.

$$741 \quad \text{Var}\left(\sum_{j \in \mathcal{H}_0} I\{M_j > \tau\}\right) \leq \omega m_0^\gamma, \quad \forall \tau \in \mathbb{R}.$$

742
743
744
745
746
747
748
749
750
751
752
753 \square

754
755 With Propositions 3.5 and C.1 in place, Proposition 2.2 of Dai et al. (2023b) directly yields Theorem 3.6.

D PROOF OF PROPOSITION 3.8

Proof. By Theorem 3.6, under Assumption 3.1 and i.i.d. the CMS estimator $\widehat{\text{FDP}}(\tau)$ is a consistent estimator of $\text{FDR}(\tau)$ at every fixed $\tau > 0$. Hence, for any fixed threshold τ ,

$$\widehat{\text{FDP}}(\tau) \xrightarrow{P} \text{FDR}(\tau) \quad (n_{\text{cal}}, m \rightarrow \infty).$$

Assumption 3.2 ensures that the mirror statistic M has a continuous distribution, so its empirical distribution satisfies Glivenko–Cantelli uniform convergence. Consequently,

$$\sup_{\tau \in [0, T]} \left| \widehat{\text{FDP}}(\tau) - \text{FDR}(\tau) \right| \xrightarrow{P} 0, \quad \text{for every fixed } T > 0.$$

Define

$$\psi(\tau) = \text{FDR}(\tau) - \alpha, \quad \widehat{\psi}(\tau) = \widehat{\text{FDP}}(\tau) - \alpha.$$

By the uniform convergence above, $\widehat{\psi} \rightarrow \psi$ uniformly on compact intervals. The mild technical condition that $\text{FDR}(\tau)$ is strictly decreasing at

$$\tau^*(\alpha) = \inf\{\tau : \text{FDR}(\tau) \leq \alpha\}$$

ensures that $\tau^*(\alpha)$ is the unique first zero-crossing of ψ . Classical stability results for zero-crossings under uniform convergence now yield

$$\tau_\alpha = \inf\{\tau : \widehat{\text{FDP}}(\tau) \leq \alpha\} \xrightarrow{P} \tau^*(\alpha).$$

Next consider the empirical proportion of true aligned selections under the random threshold τ_α . Since $\tau_\alpha \rightarrow \tau^*(\alpha)$ and M has a continuous distribution, the probability that M lies in any shrinking neighborhood of $\tau^*(\alpha)$ vanishes. Thus the difference

$$I\{M_j > \tau_\alpha, A_{n+j} > c\} - I\{M_j > \tau^*(\alpha), A_{n+j} > c\}$$

contributes negligibly in the limit. Because $M_j = g(X_{n+j}) - \widehat{\mu}_c$ and the empirical baseline $\widehat{\mu}_c$ satisfies $\widehat{\mu}_c \xrightarrow{P} \mu_c := \mathbb{E}[g(X) \mid A \leq c]$, we obtain the convergence

$$M_j \xrightarrow{P} M := g(X) - \mu_c.$$

Consequently, the indicators $I\{M_j > \tau_\alpha, A_{n+j} > c\}$ converge in distribution to $I\{M > \tau^*(\alpha), A > c\}$. By the law of large numbers,

$$\frac{1}{m} \sum_{j=1}^m I\{M_j > \tau_\alpha, A_{n+j} > c\} \xrightarrow{P} \mathbb{P}(M > \tau^*(\alpha), A > c).$$

Finally, the CMS power satisfies

$$\text{Power}(\tau_\alpha) = \mathbb{E} \left[\frac{m^{-1} \sum_{j=1}^m I\{M_j > \tau_\alpha, A_{n+j} > c\}}{m^{-1} \sum_{j=1}^m I\{A_{n+j} > c\}} \right].$$

The denominator converges to $\mathbb{P}(A > c) > 0$, and the numerator converges as above. Slutsky’s theorem yields

$$\lim_{n_{\text{cal}}, m \rightarrow \infty} \text{Power}(\tau_\alpha) = \mathbb{P}(M > \tau^*(\alpha) \mid H_1),$$

together with

$$\lim_{n_{\text{cal}}, m \rightarrow \infty} \frac{1}{m} \sum_{j=1}^m I\{M_j > \tau_\alpha, A_{n+j} > c\} = \mathbb{P}(M > \tau^*(\alpha), A > c).$$

Since no valid level- α procedure can exceed the oracle power $\mathbb{P}(M > \tau^*(\alpha) \mid H_1)$, CMS achieves the maximal asymptotic power. \square

E STRUCTURAL EQUIVALENCE BETWEEN CONFORMAL ALIGNMENT AND CMS

First, we will rewrite the power of conformal p-value.

Proposition E.1. *Consider the conformal p-value*

$$p_j = \frac{1 + \sum_{i \in \mathcal{D}_{\text{cal}}} I\{A_i \leq c, g(X_i) \geq g(X_{n+j})\}}{|\mathcal{D}_{\text{cal}}| + 1}.$$

The asymptotic power of conformal p-value admits the expression

$$\lim_{n, m \rightarrow \infty} \text{Power} = \frac{\mathbb{P}\{g(X) \geq T^*, H_1\}}{\mathbb{P}(A > c)} \quad \text{for some constant } T^*.$$

Proof. Let $V(x, a) = 2\bar{M} \cdot I\{a > c\} - g(x)$ for some constant $\bar{M} > \sup_x |g(x)|$. We begin by rewriting the p-value as

$$p_j = \frac{1 + \sum_{i \in \mathcal{D}_{\text{cal}}} I\{V_i \leq \widehat{V}_{n+j}\}}{|\mathcal{D}_{\text{cal}}| + 1},$$

so that p_j is a monotone nondecreasing function of \widehat{V}_{n+j} .

apply the Proposition 2.10 in Jin & Candès (2023b), we can get the asymptotic power taking the form

$$\frac{\mathbb{P}\{g(X) \geq T^*, A > c\}}{\mathbb{P}(A > c)},$$

for some constant T^* , which completes the proof. \square

Because the CMS statistic satisfies

$$M(X) = g(X) - \frac{\sum_{i \in \mathcal{D}_i^{\text{cal}}} I\{A_i \leq c\} \widehat{A}_i}{\sum_{i \in \mathcal{D}_i^{\text{cal}}} I\{A_i \leq c\}},$$

the rejection event $M(X) > \tau^*(\alpha)$ is equivalent to

$$g(X) > \tau^*(\alpha) + \frac{\sum_{i \in \mathcal{D}_i^{\text{cal}}} I\{A_i \leq c\} \widehat{A}_i}{\sum_{i \in \mathcal{D}_i^{\text{cal}}} I\{A_i \leq c\}}.$$

Hence the asymptotic CMS power can be written as

$$\lim_{n_{\text{cal}}, m \rightarrow \infty} \text{Power}_{\text{CMS}} = \mathbb{P}\left(g(X) > \tau^*(\alpha) + \frac{\sum_{i \in \mathcal{D}_i^{\text{cal}}} I\{A_i \leq c\} \widehat{A}_i}{\sum_{i \in \mathcal{D}_i^{\text{cal}}} I\{A_i \leq c\}} \mid H_1\right).$$

In summary, both conformal p-values and CMS admit the same structural form for their asymptotic power. Conformal p-values yield a selection rule of the form

$$g(X) \geq \tau_{\text{cp}}^*, \quad \implies \quad \lim_{n, m \rightarrow \infty} \text{Power}_{\text{cp}} = \mathbb{P}(g(X) \geq \tau_{\text{cp}}^* \mid H_1),$$

while CMS induces a shifted threshold,

$$g(X) > \tau_{\text{cms}}^* + C_{\text{cal}}, \quad \implies \quad \lim_{n_{\text{cal}}, m \rightarrow \infty} \text{Power}_{\text{cms}} = \mathbb{P}(g(X) > \tau_{\text{cms}}^* + C_{\text{cal}} \mid H_1).$$

Thus both procedures share the same functional form

$$\mathbb{P}\{g(X) > \text{threshold} \mid H_1\},$$

differing only in the data-dependent threshold. Since the CMS threshold explicitly incorporates the alignment boundary c through the mirror construction, it is generally more adaptive to the underlying distributional structure than the conformal p-value threshold.

F ROBUST CONFORMAL ALIGNMENT

Although Algorithm 1 controls the FDR at the desired level α , two important concerns remain. First, splitting the data inflates the variance of the predicted alignment scores compared to competing methods that fully utilize the data. Second, the conformal alignment selections produced by Algorithm 1 may be unstable and vary substantially across different random splits.

To address these issues, we propose a multiple data-splitting procedure that aggregates selection results from independent replications of Algorithm 1. Specifically, we run the algorithm K times with random sample splits, obtaining conformal mirror statistics $\{M_j^{(k)}\}_{j \in [m]}$, thresholds $\tau_\alpha^{(k)}$ and the selected set $\widehat{S}^{(k)}$. We then aggregate these results using the framework of e -values (Vovk & Wang, 2021), a recently proposed alternative to p -values that allows valid inference under arbitrary dependence. A nonnegative random variable e is called an e -value if $\mathbb{E}[e] \leq 1$ under the null, with larger values indicating stronger evidence against it. For instance, we should reject when $e \geq 1/\alpha$ controls type-I error at level α .

For each test unit X_{n+j} , we construct its aggregated e -value as the average of the K replications:

$$e_j = \frac{1}{K} \sum_{k=1}^K e_j^{(k)}, \quad e_j^{(k)} = \text{weight}_j^{(k)} \cdot I\{j \in \widehat{S}^{(k)}\}. \quad (6)$$

Inspired by Ren & Barber (2024) and Dai et al. (2023b), we construct two e -value forms tailored to our setting: the derandomized version, $e_j^{(k)} = \frac{m \cdot I(M_j^{(k)} > \tau_\alpha^{(k)})}{1 + \sum_{j \in [m]} I(M_j^{(k)} \leq -\tau_\alpha^{(k)})}$, and the inclusion-rate, $e_j^{(k)} = \frac{m \cdot I(M_j^{(k)} > \tau_\alpha^{(k)})}{\alpha \cdot \sum_{j \in [m]} I\{M_j^{(k)} \geq \tau_\alpha^{(k)}\}}$, which differ in how they normalize by the number of selected units.

Given the aggregated e_j values, we apply the e-BH procedure (Algorithm 2) to control FDR. By Wang & Ramdas (2022), e-BH ensures $\text{FDR} \leq \alpha \cdot m_0/m \leq \alpha$, as summarized in the following theorem.

Algorithm 2 e-BH Procedure for FDR Control

Require: e -values: e_1, e_2, \dots, e_m , FDR level: α .

Ensure: Final selected set \widehat{S} .

- 1: Sort the e -values in decreasing order: $e_{(1)} \geq e_{(2)} \geq \dots \geq e_{(m)}$,
 - 2: Set $d = \max \left\{ j : e_{(j)} \geq \frac{m}{\alpha j} \right\}$,
 - 3: Return $\widehat{S} = \{j : e_j \geq \frac{m}{\alpha d}\}$.
-

Theorem F.1 (e-BH FDR control). *Suppose the values e_1, e_2, \dots, e_m we construct for aggregation satisfy the e -value criterion $\mathbb{E}(e_j) \leq 1$ for all $j \in \mathcal{H}_0$ (or the relaxed e -value criterion $\sum_{j \in \mathcal{H}_0} \mathbb{E}(e_j) \leq m$). Then the e-BH procedure in Algorithm 2 guarantees $\text{FDR} \leq \alpha$ for any level $\alpha \in (0, 1)$.*

We demonstrate the simulation results for multiple-splitting mirror statistics here, with the same setting from Section 4.

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

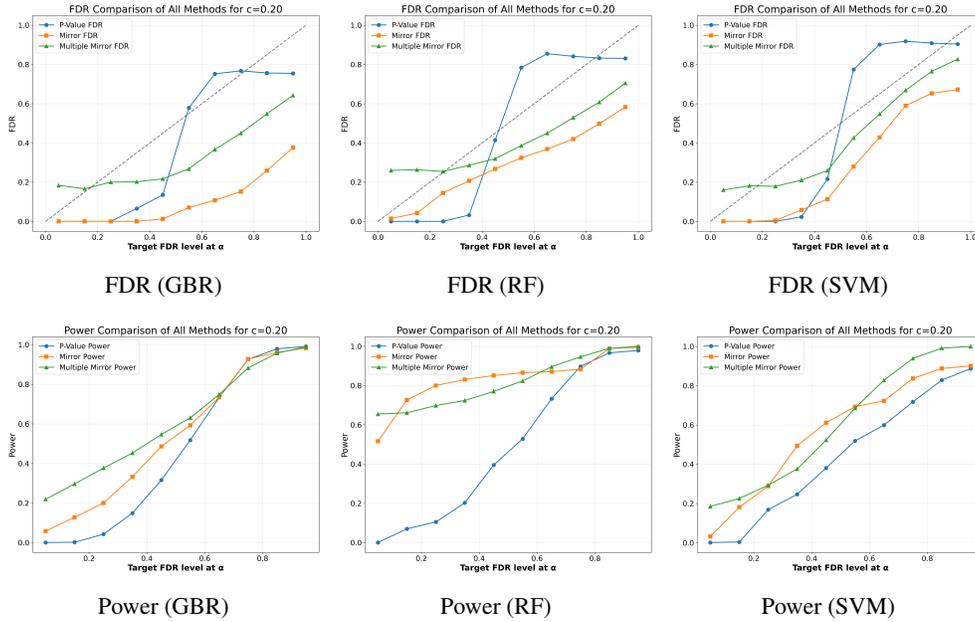


Figure 7: FDR (top) and power (bottom) using GBR, RF, and SVM under Setting (i) $\sigma(X) = 1.5$ at $c = 0.20$.

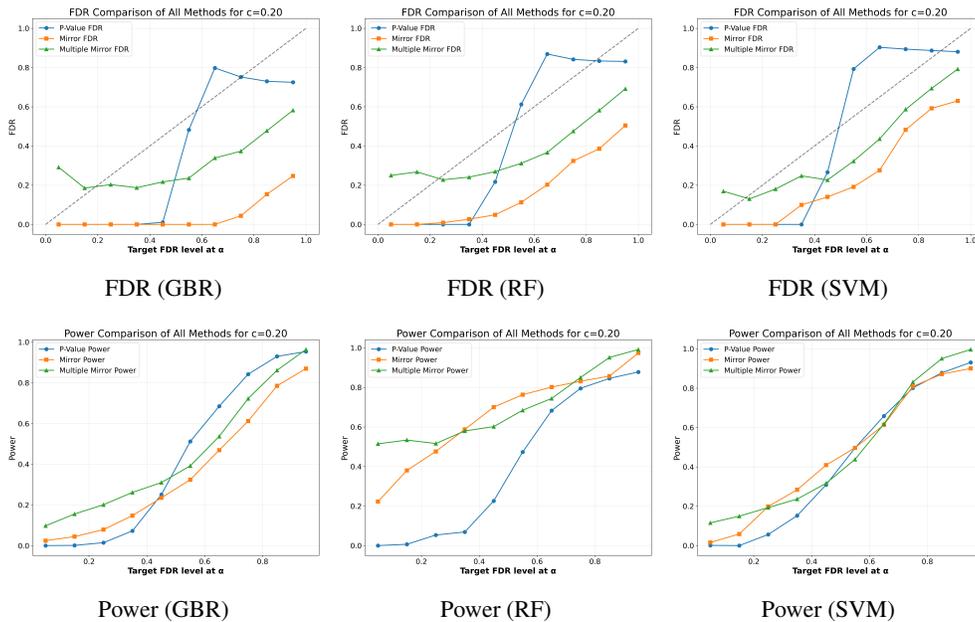


Figure 8: FDR (top) and power (bottom) using GBR, RF, and SVM under Setting (ii) $\sigma(X) = \frac{5.5 - |\mu(X)|}{2}$ at $c = 0.20$.