

---

# Calibrated Out-of-Distribution Detection with Conformal P-values

---

Stephen Bates<sup>1</sup> Emmanuel Candès<sup>2</sup> Lihua Lei<sup>3</sup> Yaniv Romano<sup>4</sup> Matteo Sesia<sup>5</sup>

## Abstract

This paper studies the construction of p-values for nonparametric out-of-distribution detection, taking a multiple-testing perspective. The goal is to test whether new independent samples belong to the same distribution as a reference data set or are outliers. We propose a solution based on conformal inference, a broadly applicable framework which yields p-values that are marginally valid but mutually dependent for different test points. We prove these p-values are positively dependent and enable exact false discovery rate control, although in a relatively weak marginal sense. We then introduce a new method to compute p-values that are both valid conditionally on the training data and independent of each other for different test points; this paves the way to stronger type-I error guarantees. Our results depart from classical conformal inference as we leverage concentration inequalities rather than combinatorial arguments to establish our finite-sample guarantees. Furthermore, our techniques also yield a uniform confidence bound for the false positive rate of any out-of-distribution detection algorithm, as a function of the threshold applied to its raw statistics. Finally, the relevance of our results is demonstrated by numerical experiments on real and simulated data.

## 1. Introduction

We consider the out-of-distribution (OOD) detection problem in which one observes an in-distribution data set  $\mathcal{D} = \{X_i\}_{i=1}^{2n}$  containing  $2n$  independent and identically distributed points  $X_i \in \mathbb{R}^d$  drawn from an unknown distribution

<sup>\*</sup>Equal contribution <sup>1</sup>Departments of Statistics and of EECS, UC Berkeley <sup>2</sup>Departments of Statistics and of Mathematics, Stanford University <sup>3</sup>Department of Statistics, Stanford University <sup>4</sup>Departments of Electrical Engineering and of Computer Science, Technion—Israel Institute of Technology <sup>5</sup>Department of Data Sciences and Operations, University of Southern California. Correspondence to: Lihua Lei <lihualei@stanford.edu>.

Presented at the ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning., Copyright 2021 by the author(s).

buton  $P_X$  (which may be continuous, discrete, or mixed). The goal is to test which among a new set of  $n_{\text{test}} \geq 1$  independent observations  $\mathcal{D}^{\text{test}} = \{X_{2n+i}\}_{i=1}^{n_{\text{test}}}$  are *outliers*, in the sense that they were not drawn from the same distribution  $P_X$ . By contrast, we refer to points drawn from  $P_X$  as *inliers*. A variety of machine-learning tools have been developed to address this classification task (e.g., Hendrycks & Gimpel, 2016; Liang et al., 2017; Lee et al., 2018b;a), which is sometimes referred to as *one-class classification* (Moya et al., 1993; Pimentel et al., 2014). They often produce a score for each test point measuring the abnormality. However, such algorithms are often complex and their scores are not directly covered by any precise statistical guarantees.

Fortunately, conformal inference (Vovk et al., 1999; 2005) allows one to practically convert the output of any one-class classifier (if it is invariant to the ordering of the training observations) into a provably valid p-value for the null hypothesis  $\mathcal{H}_{0,i} : X_i \sim P_X$ , for any  $X_i \in \mathcal{D}^{\text{test}}$ . The conformal p-value provides a calibrated measure of abnormality in finite samples without any distributional assumptions. Therefore, thresholding the p-value at level  $\alpha$  controls the Type-I error at the same level for each single test point. Conformal inference has been applied before in the context of outlier detection (Laxhammar & Falkman, 2015; Smith et al., 2015; Ishimtsev et al., 2017; Guan & Tibshirani, 2019; Cai & Koutsoukos, 2020; Haroush et al., 2021).

Despite the desirable property of conformal p-values, there is another challenge that has yet to be addressed. For modern machine learning applications, the number of outlier tests,  $n_{\text{test}}$ , is large and, therefore, it may be necessary to account for multiple comparisons to avoid making an excessive number of false discoveries. A meaningful error rate in this setting is the false discovery rate (FDR) (Benjamini & Hochberg, 1995): the expected proportion of true inliers among the test points reported as outliers. From a statistics perspective, multiple testing in this setting requires some care because classical conformal p-values corresponding to different values of  $i > 2n$  are independent of each other only conditional on  $\mathcal{D}$ , although they are valid only marginally over  $\mathcal{D}$ . This situation is delicate because FDR control typically requires p-values that either are mutually independent or follow certain patterns of dependence (Benjamini & Yekutieli, 2001; Clarke et al., 2009). Similarly, global testing (i.e., aggregating evidence from multiple ob-

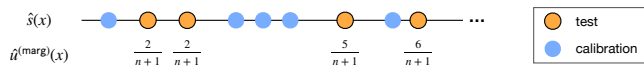


Figure 1: Visualization of the joint distribution of the conformal p-values. The distribution of  $\hat{s}(x)$  is the same for calibration and inlier test points. The conformal p-value for each test point is the number of calibration points to its left, divided by the total number of calibration points plus one, as in (1).

servations to test weaker batch-level hypotheses) may also require independent p-values. This paper addresses the above issues by carefully studying the theoretical properties of some standard multiple testing procedures applied to conformal p-values, and by developing new methods to compute p-values with stronger validity properties.

## 2. Marginal Conformal P-values

The conformal inference methods studied in this paper are statistical wrappers for any one-class classifiers. The latter are algorithms trained on data clean of any outliers to compute a score function  $\hat{s} : \mathbb{R}^d \rightarrow \mathbb{R}$  assigning a scalar value to any future data point, so that smaller (for example) values of  $\hat{s}(X)$  provide evidence that  $X$  may be an outlier.

After training  $\hat{s}$  on a subset of the observations in  $\mathcal{D}$ , namely those in  $\mathcal{D}^{\text{train}} = \{X_1, \dots, X_n\}$ , the scores are evaluated on the remaining  $n$  hold-out samples in  $\mathcal{D}^{\text{cal}} = \{X_{n+1}, \dots, X_{2n}\}$ . (Note that  $\mathcal{D}^{\text{train}}$  and  $\mathcal{D}^{\text{cal}}$  do not need to contain the same number of observations, although the current choice simplifies the notation without loss of generality). For a test point  $x$ , the marginal conformal p-value is then defined as

$$\hat{u}^{(\text{marg})}(x) = \frac{1 + |\{i \in \mathcal{D}^{\text{cal}} : \hat{s}(X_i) \leq \hat{s}(x)\}|}{n + 1}. \quad (1)$$

It is well-known (Vovk et al., 2005) that  $\hat{u}^{(\text{marg})}(X_{2n+1})$  is *marginally super-uniform* (conservative) in the sense that

$$\mathbb{P} [\hat{u}^{(\text{marg})}(X_{2n+1}) \leq t] \leq t, \quad (2)$$

for any  $t \in (0, 1)$ , whenever  $X_{2n+1}$  is an inlier. We say these p-values are marginally valid because they depend on the calibration data in  $\mathcal{D}^{\text{cal}}$ , and both  $\mathcal{D}^{\text{cal}}$  and  $X_{2n+1}$  are random in (2).

Notably, marginal p-values corresponding to different test points,  $\{\hat{u}^{(\text{marg})}(X)\}_{X \in \mathcal{D}^{\text{test}}}$ , are not mutually independent because they are all affected by  $\mathcal{D}^{\text{cal}}$ , see Figure 1 for a visualization of this dependence. This should be taken into account when adjusting for multiplicity in OOD detection applications because some common testing procedures are not generally valid for dependent p-values. We show that the dependence among marginal p-values invalidates Fisher’s combination test (Fisher, 1925) for the global null that there are no outliers in  $\mathcal{D}^{\text{test}}$ . By contrast, we can

prove the dependence between conformal p-values does not break the Benjamini-Hochberg (BH) procedure (Benjamini & Hochberg, 1995) for FDR control.

### 2.1. A negative result: global testing with conformal p-values can fail

Fisher’s combination test (Fisher, 1925) is a widely-used method to test the global null, in our case

$$H_0 : X_{2n+1}, \dots, X_{2n+m} \stackrel{\text{i.i.d.}}{\sim} P_X.$$

The idea is to aggregate the evidence from the individual tests, as follows. Given a p-value  $p_i$  for each null hypothesis  $i$ , Fisher’s test rejects the global null at level  $\alpha$  if

$$-2 \sum_{i=1}^m \log p_i \geq \chi^2(2m; 1 - \alpha),$$

where  $\chi^2(2m; 1 - \alpha)$  is the  $(1 - \alpha)$ -th quantile of the chi-square distribution with  $2m$  degrees of freedom. This test is valid if the p-values stochastically dominate  $\text{Unif}([0, 1])$  and are independent of each other. However, we prove in the following theorem that Fisher’s combination test becomes invalid when applied to marginal conformal p-values in the asymptotic regime where  $|\mathcal{D}^{\text{test}}|$  is proportional to  $|\mathcal{D}^{\text{cal}}|$ .

**Theorem 1.** *Assume that  $\hat{s}(X)$  is continuous and let  $p_i = \hat{u}^{(\text{marg})}(X_{2n+i})$ . Then, under the global null, if  $m = \lfloor \gamma n \rfloor$  for some  $\gamma \in (0, \infty)$ , as  $n$  tends to infinity,*

$$\mathbb{P} \left[ -2 \sum_{i=1}^m \log p_i \geq \chi^2(2m; 1 - \alpha) \right] \rightarrow \bar{\Phi} \left( \frac{z_{1-\alpha}}{\sqrt{1+\gamma}} \right),$$

where  $z_{1-\alpha}$  and  $\bar{\Phi}$  denote the  $(1 - \alpha)$ -th quantile and tail function of the standard normal distribution, respectively.

Since  $\gamma > 0$ , the marginal type-I error is always larger than  $\alpha$  whenever  $\alpha < 0.5$ . For illustration, consider  $\alpha = 5\%$ . When  $\gamma = 3$ , the type-I error is as large as 20.5%; when  $\gamma \rightarrow \infty$ , the type-I error is approaching 50%. This demonstrates the substantial adverse effect of dependence among marginal conformal p-values for Fisher’s combination test. In Appendix B.2, we show that corrections of Fisher’s combination test are possible.

### 2.2. A positive result: conformal p-values do not invalidate the BH procedure

Certain multiple testing methods, such as the BH procedure, are known to be robust to a particular type of mutual p-value dependence called *positive regression dependent on a subset* (PRDS) (Benjamini & Yekutieli, 2001).

**Definition 1** (PRDS). *A random vector  $X = (X_1, \dots, X_m)$  is PRDS on  $I_0 \subset \{1, \dots, m\}$  if for any  $i \in I_0$  and any increasing set  $D$ , the probability  $\mathbb{P}[X \in D \mid X_i = x]$  is increasing in  $x$ .*

Above, for vectors  $a$  and  $b$  of equal dimension, we say  $a \succeq b$  if every coordinate of  $a$  is no smaller than the corresponding coordinate of  $b$ , and a set  $D \subset \mathbb{R}^m$  is *increasing* if  $a \in D$  and  $b \succeq a$  implies  $b \in D$ . The PRDS property is a demanding form of positive dependence which can be interpreted, loosely speaking, as saying all pairwise correlations are positive. We show in the following theorem that marginal conformal p-values are PRDS on the set of inliers.

**Theorem 2.** *Assume that  $\hat{s}(X)$  is continuous. Consider  $m$  test points  $X_{2n+1}, \dots, X_{2n+m}$  such that the inliers are jointly independent of each other and of the data in  $\mathcal{D}$ . Then, the marginal conformal p-values  $(\hat{u}^{(\text{marg})}(X_{2n+1}), \dots, \hat{u}^{(\text{marg})}(X_{2n+m}))$  are PRDS on the set of inliers.*

It follows from Theorem 2 that marginal conformal p-values can be used to control the FDR with the BH procedure in finite samples for the null hypotheses

$$H_{0,i} : X_i \sim P_X, \quad i \in \{2n+1, \dots, 2n+m\}.$$

**Corollary 1** (Benjamini and Yekutieli (Benjamini & Yekutieli, 2001)). *In the setting of Theorem 2, the BH procedure applied to  $(\hat{u}^{(\text{marg})}(X_{2n+1}), \dots, \hat{u}^{(\text{marg})}(X_{2n+m}))$  at level  $\alpha \in (0, 1)$  controls the FDR at level  $\alpha$ . That is,*

$$\mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] \leq \alpha, \quad (3)$$

where  $\mathcal{H}_0 = \{i : H_{0,i} \text{ holds}\} \subseteq \{2n+1, \dots, 2n+m\}$  is the subset of true inliers in the test set, and  $\mathcal{R} \subseteq \{2n+1, \dots, 2n+m\}$  is the subset of test points reported as likely outliers.

### 3. Calibration-conditional conformal p-values

#### 3.1. Warm up: analyzing the false positive rate

Having noted that conformal inferences hold in theory only marginally over the calibration data, the first question one may ask is: how bad can these inferences be conditional on a particular calibration set? We will address this question by developing high-probability bounds for the conditional deviation from uniformity of marginal p-values, starting here from the simplest case of pointwise bounds. The purpose of a pointwise bound is to control the probability that a null p-value (corresponding to a true inlier) is smaller than  $\alpha$ , conditional on  $\mathcal{D}$ , for some *fixed* threshold  $\alpha \in (0, 1)$ . In other words, we wish to understand the conditional false positives rate (FPR) corresponding to the threshold  $\alpha$ ,

$$\text{FPR}(\alpha; \mathcal{D}) := \mathbb{P} \left[ \hat{u}^{(\text{marg})}(X_{2n+1}) \leq \alpha \mid \mathcal{D} \right], \quad (4)$$

beyond what we know from the marginal guarantee in (2), which is  $\mathbb{E}[\text{FPR}(\alpha; \mathcal{D})] \leq \alpha$ . The quantity in (4) can be studied precisely with existing results due to Vovk (2012),

which is formalized in Proposition 1 below. We revisit this topic here because it serves as an intuitive introduction to the more involved high-probability bounds that we will propose later.

**Proposition 1** (Pointwise FPR of marginal conformal p-values, adapted from Vovk (2012)). *Let  $\ell = \lfloor (n+1)\alpha \rfloor$ . If  $\hat{s}(X)$  is continuous,  $\text{FPR}(\alpha; \mathcal{D})$  follows a  $\text{BETA}(\ell, n+1-\ell)$  distribution.*

This shows precisely how a smaller  $\mathcal{D}^{\text{cal}}$  makes marginal p-values more conservative on average, but also more likely to be overly liberal on occasion. While this result is informative, it is limited for our purposes because it provides only a pointwise bound—it takes  $\alpha$  as fixed—whereas uniform bounds are needed to construct conditionally valid p-values that can be safely used with any multiple-testing procedure, as discussed in the next section.

#### 3.2. Conditionally valid conformal p-values

Proposition 1 implies marginal conformal p-values may be anti-conservative conditional on  $\mathcal{D}$ . We discuss a generic strategy to adjust the marginal conformal p-values into valid p-values conditional on  $\mathcal{D}$  in Appendix B.5. The following theorem describes a specific strategy based on the generalized Simes inequality (Sarkar et al., 2008).

**Theorem 3** (Simes adjustment). *Fix any integer  $k < n$ . Let  $h : [0, 1] \mapsto [0, 1]$  be a piece-wise constant function with*

$$h(t) = b_{\lceil (n+1)t \rceil}, \quad t \in [0, 1],$$

where, for any  $i = 1, \dots, n$ ,

$$b_{n+1-i} = 1 - \delta^{1/k} \left( \frac{i \cdots (i-k+1)}{n \cdots (n-k+1)} \right)^{1/k}.$$

Let also  $\hat{u}^{(\text{ccv})} = h \circ \hat{u}^{(\text{marg})}$ . Then

$$\mathbb{P} \left[ \mathbb{P} \left[ \hat{u}^{(\text{ccv})}(X_{2n+1}) \leq t \mid \mathcal{D} \right] \leq t, \forall t \in (0, 1) \right] \geq 1 - \delta.$$

Figure 2 illustrates the adjustment function  $h$  with  $n = 500$  and  $k = n/2$ . For example, a marginal p-value of  $\hat{u}^{(\text{marg})}(X) = 25/(n+1) \approx 0.05$  results in a CCV p-value of  $h(25/(n+1)) = b_{25} \approx 0.075$  in this case. While the adjustment is very conservative when the marginal conformal p-value is large, it typically does not matter in multiple testing problems, as it is the small ones that determine which hypotheses are rejected.

When  $X_{2n+1}, \dots, X_{2n+m}$  are independent, the conditional conformal p-values  $\hat{u}^{(\text{ccv})}(X_{2n+1}), \dots, \hat{u}^{(\text{ccv})}(X_{2n+m})$  are independent conditional on  $\mathcal{D}$ . Therefore, with probability at least  $1 - \delta$  in  $\mathcal{D}$ , the Fisher’s combination test would control the Type-I error for the global null and the BH procedure would control the FDR, conditional on  $\mathcal{D}$ .

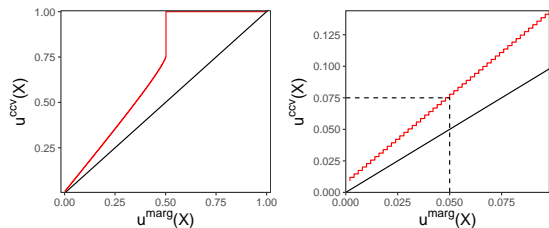


Figure 2: Illustration of Theorem 7. The red curve gives the sequence derived by Simes adjustment with  $k = n/2 = 250$  and the black curve gives the 45-degree line. The right panel zooms in on small indices.

## 4. Numerical Experiments

Suppose we can access  $J$  copies of inlier datasets  $\mathcal{D}_1, \dots, \mathcal{D}_J$  and  $L$  copies of test sets  $\mathcal{D}_{j,1}^{\text{test}}, \dots, \mathcal{D}_{j,L}^{\text{test}}$  for each  $j$ . The conditional FDR is defined as

$$\text{cFDR}(\mathcal{D}_j) := \mathbb{E} [\text{FDP}(\mathcal{D}^{\text{test}}; \mathcal{D}_j) \mid \mathcal{D}_j],$$

where  $\text{FDP}(\mathcal{D}^{\text{test}}; \mathcal{D}_j)$  is the proportion of inliers among the test points reported as outliers, based on the procedure calibrated on  $\mathcal{D}_j$ . This motivates the definition of the following performance measures. For any  $j \in [J]$ , we compute

$$\widehat{\text{cFDR}}(\mathcal{D}_j) := \frac{1}{L} \sum_{l=1}^L \text{FDP}(\mathcal{D}_{j,l}^{\text{test}}; \mathcal{D}_j), \quad (5)$$

$$\widehat{\text{cPower}}(\mathcal{D}_j) := \frac{1}{L} \sum_{l=1}^L \text{Power}(\mathcal{D}_{j,l}^{\text{test}}; \mathcal{D}_j), \quad (6)$$

where  $\text{Power}(\mathcal{D}_{j,l}^{\text{test}}; \mathcal{D}_j)$  is the proportion of outliers in  $\mathcal{D}_{j,l}^{\text{test}}$  correctly identified as such. Our experiments will demonstrate that the proposed simultaneous calibration method leads to sufficiently small  $\widehat{\text{cFDR}}(\mathcal{D}_j)$  for the desired fraction of  $\mathcal{D}_j$ 's, while the traditional point-wise calibration generally only leads to small values of the marginal FDR, namely  $\widehat{\text{mFDR}} := \frac{1}{J} \sum_{j=1}^J \widehat{\text{cFDR}}(\mathcal{D}_j)$ .

In Appendix D, we will investigate the performance of marginal and conditional conformal p-values on various simulated and real datasets. Here we present the results on the credit card data set<sup>1</sup> for illustration. The dataset involves 284,315 samples in  $\mathbb{R}^{30}$  with 492 outliers. To estimate the performance measures, we need to construct multiple training, calibration, and test sets by randomly splitting the  $n_{\text{inlier}}$  inlier examples into three disjoint subsets of size  $n_{\text{train}}$ ,  $n_{\text{cal}}$  and  $n_{\text{test}}$ , respectively. A total of  $n_{\text{inlier}}/2$  data points is used for training and calibration, i.e.,  $n_{\text{train}} + n_{\text{cal}} = n_{\text{inlier}}/2$  with  $n_{\text{cal}} = \min\{2000, n_{\text{train}}/2\}$ , while outlier examples are only included in the test sets. For

<sup>1</sup>The dataset is available at <http://www.ulb.ac.be/di/map/adalpozz/data/creditcard.Rdata>.

each training/calibration data subset, we sample 100 test sets of size  $n_{\text{test}} = \min\{2000, n_{\text{train}}/3\}$ . Each test set contains 90% of randomly chosen inliers, and 10% of outliers.

We utilize an isolation forest (Liu et al., 2008) machine-learning algorithms  $\hat{s}$  as the base method for detecting anomalies, available in the Python `sklearn` package. We rely on the default hyper-parameters, except for the ‘contamination’ parameter which we set equal to 0.1.

Figure 3 compares the performance of marginal and simultaneously calibrated p-values as a function of the nominal FDR level. Here, the BH procedure is applied. Note that the proposed Simes simultaneous calibration leads to FDR control for at least 90% of inlier datasets, as expected. This stands in contrast with the marginal calibration approach, which controls the FDR only marginally.

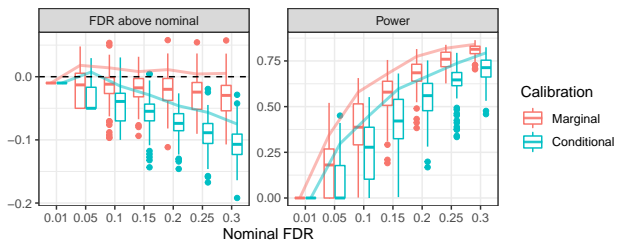


Figure 3: Performance of OOD detection on credit card fraud data. The BH procedure is applied on marginal and conditional conformal p-values to control the FDR over the set of test points. The box plots visualize the distribution of the (conditional) excess FDR and power, as defined in (41), conditional on 100 independent data sets. The solid curves indicate the 90-th quantile of the conditional FDR distribution. The nominal FDR 0.1, and the conditional method is applied with  $\delta = 0.1$ .

## 5. Discussion

We discussed statistical techniques to calibrate the Type-I errors for OOD detection. In some applications, false negatives are substantially more costly than false positives, in which case it is more appropriate to control Type-II error than Type-I error. Not surprisingly, it is generally impossible if only in-distribution data are available; in the extreme case the outlier distribution can be arbitrarily close to the inlier distribution to the extent that controlling Type-II error becomes impossible.

Nevertheless, Type-II error control is possible with labelled OOD samples. A naive idea is to flip the in-distribution and OOD samples and treat the latter as the null. We can further make the score more informative by training a classifier to distinguish the in-distribution and OOD samples. The caveat is that this approach guarantees Type-II error control only if the labelled OOD samples are sufficiently representative. We leave this investigation for future research.

## References

- Aggarwal, C. C. Outlier analysis. In *Data mining*, pp. 237–263. Springer, 2015.
- Agrawal, S. and Agrawal, J. Survey on anomaly detection using data mining techniques. *Procedia Computer Science*, 60:708–713, 2015.
- Angelopoulos, A. N., Bates, S., Malik, J., and Jordan, M. I. Uncertainty sets for image classifiers using conformal prediction. *preprint at arXiv:2009.14193*, 2020.
- Arnold, B. C., Balakrishnan, N., and Nagaraja, H. N. *A first course in order statistics*. SIAM, 2008.
- Barber, R. F., Candès, E. J., Ramdas, A., and Tibshirani, R. J. The limits of distribution-free conditional predictive inference. *preprint at arXiv:1903.04684*, 2019.
- Barber, R. F., Candès, E. J., Ramdas, A., Tibshirani, R. J., et al. Predictive inference with the jackknife+. *Annals of Statistics*, 49(1):486–507, 2021.
- Bates, S., Angelopoulos, A., Lei, L., Malik, J., and Jordan, M. I. Distribution-free, risk-controlling prediction sets. *arXiv preprint*, 2021. arXiv:2101.02703.
- Benjamini, Y. and Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*, 57(1):289–300, 1995.
- Benjamini, Y. and Yekutieli, D. The control of the false discovery rate in multiple testing under dependency. *Annals of Statistics*, pp. 1165–1188, 2001.
- Benjamini, Y., Krieger, A. M., and Yekutieli, D. Adaptive linear step-up procedures that control the false discovery rate. *Biometrika*, 93(3):491–507, 2006.
- Brown, M. B. 400: A method for combining non-independent, one-sided tests of significance. *Biometrics*, pp. 987–992, 1975.
- Cai, F. and Koutsoukos, X. Real-time out-of-distribution detection in learning-enabled cyber-physical systems. In *2020 ACM/IEEE 11th International Conference on Cyber-Physical Systems (ICCPs)*, pp. 174–183. IEEE, 2020.
- Campos, G. O., Zimek, A., Sander, J., Campello, R. J., Mícenková, B., Schubert, E., Assent, I., and Houle, M. E. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data mining and knowledge discovery*, 30(4):891–927, 2016.
- Cauchois, M., Gupta, S., and Duchi, J. Knowing what you know: valid confidence sets in multiclass and multilabel prediction. *preprint at arXiv:2004.10181*, 2020.
- Ceroli, A. Multivariate outlier detection with high-breakdown estimators. *Journal of the American Statistical Association*, 105(489):147–156, 2010.
- Chalapathy, R. and Chawla, S. Deep learning for anomaly detection: A survey. *preprint at arXiv:1901.03407*, 2019.
- Chernozhukov, V., Wüthrich, K., and Zhu, Y. Distributional conformal prediction. *preprint at arXiv:1909.07889*, 2019.
- Clarke, S., Hall, P., et al. Robustness of multiple testing procedures against dependence. *Annals of Statistics*, 37(1):332–358, 2009.
- Dempster, A. Generalized  $d_n^+$  statistics. *Ann. Math. Stat.*, 30(2):593–597, 1959.
- Durbin, J. *Distribution Theory for Tests Based on Sample Distribution Function*, volume 9. SIAM, 1973.
- Dvoretzky, A., Kiefer, J., and Wolfowitz, J. Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. *Ann. Math. Stat.*, pp. 642–669, 1956.
- Fedorova, V., Gammerman, A., Nouretdinov, I., and Vovk, V. Plug-in martingales for testing exchangeability on-line. *arXiv preprint arXiv:1204.3251*, 2012.
- Fisher, R. *Statistical methods for research workers*. Oliver & Boyd (Edinburgh), 1925.
- Friedman, J. On multivariate goodness-of-fit and two-sample testing. Technical report, No. SLAC-PUB-10325. Stanford Linear Accelerator Center, Menlo Park, CA (US), 2004.
- Guan, L. and Tibshirani, R. Prediction and outlier detection in classification problems. *arXiv preprint arXiv:1905.04396*, 2019.
- Gupta, C., Kuchibhotla, A. K., and Ramdas, A. K. Nested conformal prediction and quantile out-of-bag ensemble methods. *arXiv preprint arXiv:1910.10562*, 2019.
- Haroush, M., Frostig, T., Heller, R., and Soudry, D. Statistical testing for efficient out of distribution detection in deep neural networks. *arXiv preprint arXiv:2102.12967*, 2021.
- Hawkins, D. M. *Identification of outliers*, volume 11. Springer, 1980.
- Hechtlinger, Y., Póczos, B., and Wasserman, L. Cautious deep learning, 2018. arXiv:1805.09460.
- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.

- Hu, X. and Lei, J. A distribution-free test of covariate shift using conformal prediction. *arXiv preprint arXiv:2010.07147*, 2020.
- Ishimtsev, V., Bernstein, A., Burnaev, E., and Nazarov, I. Conformal  $k$ -nn anomaly detector for univariate data streams. In *Conformal and Probabilistic Prediction and Applications*, pp. 213–227. PMLR, 2017.
- Izbicki, R., Shimizu, G., and Stern, R. Flexible distribution-free conditional predictive bands using density estimators. In *International Conference on Artificial Intelligence and Statistics*, pp. 3068–3077. PMLR, 2020.
- Khan, S. S. and Madden, M. G. One-class classification: taxonomy of study and review of techniques. *The Knowledge Engineering Review*, 29(3):345–374, 2014.
- Kim, B., Xu, C., and Foygel Barber, R. Predictive inference is free with the jackknife+-after-bootstrap. *Advances in Neural Information Processing Systems*, 33, 2020.
- Kim, I., Ramdas, A., Singh, A., Wasserman, L., et al. Classification accuracy as a proxy for two-sample testing. *Annals of Statistics*, 49(1):411–434, 2021.
- Kivaranovic, D., Johnson, K. D., and Leeb, H. Adaptive, distribution-free prediction intervals for deep networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 4346–4356. PMLR, 2020.
- Kost, J. T. and McDermott, M. P. Combining dependent  $p$ -values. *Statistics & Probability Letters*, 60(2):183–190, 2002.
- Kotel’Nikova, V. and Chmaladze, E. On computing the probability of an empirical process not crossing a curvilinear boundary. *Theory of probability & its applications*, 27(3):640–648, 1983.
- Krishnamoorthy, K. and Mathew, T. *Statistical Tolerance Regions: Theory, Applications, and Computation*. Wiley Series in Probability and Statistics. Wiley, 2009. ISBN 9780470473894. URL <https://books.google.com/books?id=ljQh0miU6PQC>.
- Kuchibhotla, A. K. Exchangeability, conformal prediction, and rank tests. *arXiv preprint arXiv:2005.06095*, 2020.
- Laxhammar, R. and Falkman, G. Inductive conformal anomaly detection for sequential detection of anomalous sub-trajectories. *Annals of Mathematics and Artificial Intelligence*, 74(1-2):67–94, 2015.
- Lee, K., Lee, H., Lee, K., and Shin, J. Training confidence-calibrated classifiers for detecting out-of-distribution samples. In *International Conference on Learning Representations*, 2018a.
- Lee, K., Lee, K., Lee, H., and Shin, J. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *NeurIPS*, 2018b.
- Lei, J., Rinaldo, A., and Wasserman, L. A conformal prediction approach to explore functional data. *Annals of Mathematics and Artificial Intelligence*, 74, 02 2013. doi: 10.1007/s10472-013-9366-6.
- Liang, S., Li, Y., and Srikant, R. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.
- Lipták, T. On the combination of independent tests. *Magyar Tud Akad Mat Kutato Int Kozl*, 3:171–197, 1958.
- Liu, F. T., Ting, K. M., and Zhou, Z.-H. Isolation forest. In *2008 eighth ieee international conference on data mining*, pp. 413–422. IEEE, 2008.
- Lopez-Paz, D. and Oquab, M. Revisiting classifier two-sample tests. In *International Conference on Learning Representations*, 2017.
- Massart, P. The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *Annals of Probability*, pp. 1269–1283, 1990.
- Moran, P. The random division of an interval. *Supplement to the Journal of the Royal Statistical Society*, 9(1):92–98, 1947.
- Moya, M. M., Koch, M. W., and Hostetler, L. D. One-class classifier networks for target recognition applications. *NASA STI/Recon Technical Report N*, 93:24043, 1993.
- Papadopoulos, H., Proedrou, K., Vovk, V., and Gammerman, A. Inductive confidence machines for regression. In *Machine Learning: European Conference on Machine Learning ECML 2002*, pp. 345–356, 2002.
- Park, S., Bastani, O., Matni, N., and Lee, I. PAC confidence sets for deep neural networks via calibrated prediction. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=BJxVI04YvB>.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Petrov, V. Sums of independent random variables. *Yu. V. Prokhorov. V. Statulevičius (Eds.)*, 1975.

- Pimentel, M. A., Clifton, D. A., Clifton, L., and Tarassenko, L. A review of novelty detection. *Signal Processing*, 99: 215–249, 2014.
- Riani, M., Atkinson, A. C., and Cerioli, A. Finding an unknown number of multivariate outliers. *Journal of the Royal Statistical Society: series B (statistical methodology)*, 71(2):447–466, 2009.
- Romano, Y., Patterson, E., and Candès, E. Conformalized quantile regression. In *Advances in Neural Information Processing Systems 32*, pp. 3543–3553. 2019.
- Romano, Y., Sesia, M., and Candès, E. J. Classification with valid and adaptive coverage. *Advances in Neural Information Processing Systems*, 33, 2020.
- Sabokrou, M., Khalooei, M., Fathy, M., and Adeli, E. Adversarially learned one-class classifier for novelty detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3379–3388, 2018.
- Sarkar, S. K. et al. Generalizing simes’ test and hochberg’s stepup procedure. *Annals of Statistics*, 36(1):337–363, 2008.
- Siegmund, D. Boundary crossing probabilities and statistical applications. *Annals of Statistics*, pp. 361–404, 1986.
- Smith, J., Nouretdinov, I., Craddock, R., Offer, C., and Gammerman, A. Conformal anomaly detection of trajectories with a multi-class hierarchy. In *International symposium on statistical learning and data sciences*, pp. 281–290. Springer, 2015.
- Storey, J. D. A direct approach to false discovery rates. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3):479–498, 2002.
- Storey, J. D., Taylor, J. E., and Siegmund, D. Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(1):187–205, 2004.
- Tukey, J. W. Non-parametric estimation ii. statistically equivalent blocks and tolerance regions—the continuous case. *Ann. Math. Statist.*, 18(4):529–539, 12 1947. doi: 10.1214/aoms/1177730343. URL <https://doi.org/10.1214/aoms/1177730343>.
- Van Zwet, W. and Oosterhoff, J. On the combination of independent test statistics. *Ann. Math. Stat.*, 38(3):659–680, 1967.
- Vovk, V. Conditional validity of inductive conformal predictors. In *Proceedings of the Asian Conference on Machine Learning*, volume 25, pp. 475–490, 2012.
- Vovk, V. Cross-conformal predictors. *Annals of Mathematics and Artificial Intelligence*, 74(1-2):9–28, 2015.
- Vovk, V. Testing randomness. *arXiv preprint arXiv:1906.09256*, 2019.
- Vovk, V. Testing for concept shift online. *arXiv preprint arXiv:2012.14246*, 2020.
- Vovk, V. and Wang, R. Combining p-values via averaging. *Biometrika*, 107(4):791–808, 2020.
- Vovk, V., Gammerman, A., and Saunders, C. Machine-learning applications of algorithmic randomness. In *International Conference on Machine Learning*, pp. 444–453, 1999.
- Vovk, V., Nouretdinov, I., and Gammerman, A. Testing exchangeability on-line. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 768–775, 2003.
- Vovk, V., Gammerman, A., and Shafer, G. *Algorithmic learning in a random world*. Springer, 2005. doi: 10.1007/b106715.
- Vovk, V., Petej, I., Nouretdinov, I., Ahlberg, E., Carlsson, L., and Gammerman, A. Retrain or not retrain: Conformal test martingales for change-point detection. *arXiv preprint arXiv:2102.10439*, 2021.
- Wald, A. An extension of wilks’ method for setting tolerance limits. *Ann. Math. Statist.*, 14(1):45–55, 03 1943. doi: 10.1214/aoms/1177731491. URL <https://doi.org/10.1214/aoms/1177731491>.
- Wilcoxon, F. Individual comparisons by ranking methods. In *Breakthroughs in statistics*, pp. 196–202. Springer, 1992.
- Wilks, S. S. Determination of sample sizes for setting tolerance limits. *Ann. Math. Statist.*, 12(1):91–96, 03 1941. doi: 10.1214/aoms/1177731788. URL <https://doi.org/10.1214/aoms/1177731788>.
- Wilks, S. S. Statistical prediction with special reference to the problem of tolerance limits. *Ann. Math. Statist.*, 13(4):400–409, 12 1942. doi: 10.1214/aoms/1177731537. URL <https://doi.org/10.1214/aoms/1177731537>.
- Wilks, S. S. Multivariate statistical outliers. *Sankhyā: The Indian Journal of Statistics, Series A*, pp. 407–426, 1963.
- Zhang, Y. and Politis, D. N. Bootstrap prediction intervals with asymptotic conditional validity and unconditional guarantees. *arXiv preprint arXiv:2005.09145*, 2020.

## Appendix

### A. Related Work

The out-of-distribution detection problem considered in this paper is fully non-parametric, in the sense that we leverage the information contained in an external clean data set, and nothing else, to infer whether a future test point may be an outlier. This is in contrast with the more classical problem of multivariate outlier detection within a single data set, leveraging modeling assumptions rather than clean external samples (Wilks, 1963; Hawkins, 1980; Riani et al., 2009; Cerioli, 2010). A wealth of data mining and machine-learning methods have been developed to address our non-parametric task (Khan & Madden, 2014; Agrawal & Agrawal, 2015; Aggarwal, 2015; Sabokrou et al., 2018; Chalapathy & Chawla, 2019); these do not provide precise finite-sample guarantees on their own, but we can leverage them to compute scoring functions that powerfully separate outliers from inliers.

Our paper is based on conformal inference (Vovk et al., 1999; 2005), which has been applied before in the context of outlier detection (Laxhammar & Falkman, 2015; Smith et al., 2015; Ishimtsev et al., 2017; Guan & Tibshirani, 2019; Cai & Koutsoukos, 2020; Haroush et al., 2021). However, previous works did not study the implications of marginal p-values on the validity of multiple outlier testing procedures, nor did they seek the conditional guarantees obtained here. Another line of work applied conformal inference to test the global null for streaming data (Vovk et al., 2003; Fedorova et al., 2012; Vovk, 2019; 2020; Vovk et al., 2021). However, the guarantee no longer holds in the offline setting or beyond the global null. The most closely related work is that of (Vovk, 2012), which extends conformal inference to provide a form of calibration-conditional coverage. That paper focused explicitly on the prediction setting rather than on outlier detection, but is also directly relevant in our context, as discussed in Section 3.1. The main difference is that our novel high-probability bounds in Section 3 hold simultaneously for all possible coverage levels (in the language of (Vovk, 2012)) not just for a pre-specified one—this feature being necessary to obtain conditionally valid p-values for multiple outlier testing.

Other works on conformal inference focused on different types of conditional coverage. For example, (Barber et al., 2019) studied the difficulty of computing valid conformal predictions (in a supervised setting) conditional on the features of a new test point, while we are interested in conditioning on the calibration data (in an outlier detection setting). Other works have focused on seeking approximate feature-conditional coverage in multi-class classification (Hechtlinger et al., 2018; Romano et al., 2020; Cauchois et al., 2020; Angelopoulos et al., 2020) or in regression problems (Romano et al., 2019; Izbicki et al., 2020; Chernozhukov et al., 2019; Kivaranovic et al., 2020; Gupta et al., 2019). This paper is orthogonal, in the sense that our results could be applied to strengthen their coverage guarantees by conditioning on the calibration data. It should be noted that, although conformal inference can be based on different data hold-out strategies (Vovk, 2015; Barber et al., 2021; Kim et al., 2020), our paper focuses on sample splitting (Papadopoulos et al., 2002; Lei et al., 2013). The latter has the advantage of being the most computationally efficient option, and is necessary for us in theory because our high-probability bounds require the independence of the data points in addition to their exchangeability.

Further, the problem we consider is related to classical two-sample testing (Wilcoxon, 1992), although we take a different perspective. Two-sample testing compares two data sets to determine whether they were sampled from the same distribution, while our goal is to contrast many independent test points (or batches thereof) to the same reference set accounting for multiplicity. In any case, several recent works have explored the use of machine-learning and data hold-out methods for two-sample testing (Friedman, 2004; Lopez-Paz & Oquab, 2017; Kuchibhotla, 2020; Hu & Lei, 2020; Kim et al., 2021), which reinforces the connection with our work.

Finally, the duality between hypothesis testing and confidence intervals connects our conditionally calibrated p-values to the classical statistical topic of *tolerance regions*, which goes back to Wilks (Wilks, 1941; 1942), Wald (Wald, 1943), and Tukey (Tukey, 1947). See (Krishnamoorthy & Mathew, 2009) for an overview of the subject, (Vovk, 2012) for a discussion of their connection with conformal inference, and (Park et al., 2020; Bates et al., 2021) for modern examples using tolerance regions for predictive inference with neural networks. (Tolerance regions are predictive sets with a high-probability guarantee to contain the desired fraction of the population. For example, one can generate a tolerance region guaranteed to contain at least 80% of the population with probability 99%.) The construction of predictive intervals with (asymptotic) conditional validity in the aforementioned sense was also recently studied in (Zhang & Politis, 2020) with bootstrap rather than conformal inference methods.

## B. Theoretical Results and Technical Proofs

### B.1. Correlation structure of null marginal conformal p-values

For notational convenience, we write  $p_i$  instead of  $\hat{u}^{(\text{marg})}(X_{2n+i})$ . When  $X_{2n+1}, \dots, X_{2n+m}$  are all inliers which are drawn from  $P_X$ , the conformal p-values  $p_1, \dots, p_m$  are exchangeable. We prove in the following lemma that the standard (marginal) conformal p-values are positively correlated under arbitrary transformations, suggesting an inflation of the variance of the combination statistics.

**Lemma 1.** *Assume that  $\hat{s}(X)$  is continuous. Then, for any function  $G : [0, 1] \mapsto \mathbb{R}$ , and for any pair of nulls  $(i, j)$ ,*

$$\text{Cor} [G(\hat{u}^{(\text{marg})}(X_{2n+i})), G(\hat{u}^{(\text{marg})}(X_{2n+j}))] = \frac{1}{n+2}.$$

Lemma 1 suggests that the variance of the combination statistic with any transformation  $G(\cdot)$  is  $(1 + \gamma)$  times as large as that when the p-values are i.i.d.. In fact, when  $G(\cdot)$  is square-integrable, under the global null,

$$\begin{aligned} \text{Var} \left[ \sum_{i=1}^m G(p_i) \right] &= m \text{Var} [G(p_1)] + m(m-1) \text{Cov} [G(p_1), G(p_2)] \\ &= \left( m + \frac{m(m-1)}{n+2} \right) \text{Var} [G(p_1)] \\ &\approx (1 + \gamma)m \text{Var} [G(p_1)]. \end{aligned}$$

*Proof of Lemma 1.* Without loss of generality, assume  $i = 1$  and  $j = 2$ . Let  $(R_1, \dots, R_n, R_{n+1}, R_{n+2})$  be the rank of  $(S_1, \dots, S_{n+2}) \stackrel{d}{=} (\hat{s}(X_{n+1}), \dots, \hat{s}(X_{2n}), \hat{s}(X_{2n+1}), \hat{s}(X_{2n+2}))$  in the ascending order. Then  $S_1, \dots, S_{n+2}$  are i.i.d. draws from a non-atomic distribution,  $(R_1, \dots, R_{n+2})$  are mutually distinct almost surely and for any permutation  $\pi : \{1, \dots, n+2\} \mapsto \{1, \dots, n+2\}$ ,

$$(S_{\pi(1)}, \dots, S_{\pi(n+2)}) \stackrel{d}{=} (S_1, \dots, S_{n+2}).$$

Therefore,

$$(R_1, \dots, R_{n+2}) \sim \text{Unif}(\{1, \dots, n+2\}).$$

By definition,

$$p_1 = \frac{1}{n+1} \sum_{i=1}^{n+1} I(S_i \leq S_{n+1}), \quad R_{n+1} = \sum_{i=1}^{n+2} I(S_i \leq S_{n+1}).$$

Thus,

$$p_1 = \frac{R_{n+1} - I(S_{n+2} \leq S_{n+1})}{n+1}.$$

Similarly,

$$p_2 = \frac{R_{n+2} - I(S_{n+1} \leq S_{n+2})}{n+1}.$$

For any  $j \in \{1, \dots, n+1\}$ ,

$$\begin{aligned} \mathbb{P} \left[ p_1 = p_2 = \frac{j}{n+1} \right] &= \mathbb{P} [R_{n+1} = j+1, R_{n+2} = j] + \mathbb{P} [R_{n+1} = j, R_{n+2} = j+1] \\ &= 2\mathbb{P} [R_{n+1} = j+1, R_{n+2} = j] = \frac{2}{(n+2)(n+1)}. \end{aligned}$$

For any  $1 \leq j < k \leq n+1$ ,

$$\mathbb{P} \left[ p_1 = \frac{j}{n+1}, p_2 = \frac{k}{n+1} \right] = \mathbb{P} [R_{n+1} = j, R_{n+2} = k+1] = \frac{1}{(n+2)(n+1)}.$$

By symmetry,

$$\mathbb{P}\left[p_1 = \frac{k}{n+1}, p_2 = \frac{j}{n+1}\right] = \frac{1}{(n+2)(n+1)}.$$

As a result,

$$\begin{aligned} \mathbb{E}[G(p_1)G(p_2)] &= \frac{2}{(n+2)(n+1)} \sum_{j=1}^{n+1} G^2\left(\frac{j}{n+1}\right) + \frac{1}{(n+2)(n+1)} \sum_{j \neq k} G\left(\frac{j}{n+1}\right) G\left(\frac{k}{n+1}\right) \\ &= \frac{1}{(n+2)(n+1)} \sum_{j=1}^{n+1} G^2\left(\frac{j}{n+1}\right) + \frac{1}{(n+2)(n+1)} \left\{ \sum_{j=1}^{n+1} G\left(\frac{j}{n+1}\right) \right\}^2. \end{aligned}$$

On the other hand, since  $p_1$  is uniformly distributed on  $\{1/(n+1), 2/(n+1), \dots, 1\}$ ,

$$\mathbb{E}[G(p_1)] = \frac{1}{n+1} \sum_{j=1}^{n+1} G\left(\frac{j}{n+1}\right), \quad \mathbb{E}[G^2(p_1)] = \frac{1}{n+1} \sum_{j=1}^{n+1} G^2\left(\frac{j}{n+1}\right).$$

Note that  $\mathbb{E}[G^2(p_1)] < \infty$  since  $G(i/(n+1)) \in \mathbb{R}$ . As a result,

$$\begin{aligned} \text{Cov}[G(p_1), G(p_2)] &= \frac{1}{(n+2)(n+1)} \sum_{j=1}^{n+1} G^2\left(\frac{j}{n+1}\right) - \frac{1}{(n+2)(n+1)^2} \left\{ \sum_{j=1}^{n+1} G\left(\frac{j}{n+1}\right) \right\}^2 \\ &= \frac{1}{n+2} \{ \mathbb{E}[G^2(p_1)] - (\mathbb{E}[G(p_1)])^2 \} = \frac{1}{n+2} \text{Var}[G(p_1)]. \end{aligned}$$

Therefore,

$$\text{Cor}[G(p_1), G(p_2)] = \frac{\text{Cov}[G(p_1), G(p_2)]}{\sqrt{\text{Var}[G(p_1)]\text{Var}[G(p_2)]}} = \frac{1}{n+2}.$$

□

## B.2. Failure of type-I error control with combination tests

We state a theorem for general (adjusted) combination tests which reject the global null if

$$\sum_{i=1}^m G(\hat{u}^{(\text{marg})}(Z_{2n+i})) \geq \xi c_{1-\alpha}(G) - m(\xi - 1) \int_0^1 G(u) du, \quad (7)$$

where  $\xi > 0$  is a pre-specified constant and

$$c_{1-\alpha}(G) \triangleq \text{Quantile}_{1-\alpha} \left( \sum_{i=1}^m G(U_i) \right), \quad U_i \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0, 1]).$$

**Theorem 4.** Assume  $\hat{s}(X)$  is continuous and  $G(\cdot) : [0, 1] \mapsto \mathbb{R}$  is a non-constant function satisfying

- (i)  $\int_0^1 G^{2+\eta}(u) du < \infty$ ;
- (ii)  $\left| \frac{1}{n+1} \sum_{j=1}^{n+1} G^k(j/(n+1)) - \int_0^1 G^k(u) du \right| = o(1/\sqrt{n})$ , for  $k \in \{1, 2\}$ ;
- (iii)  $\max_{j \in \{1, \dots, n+1\}} G(j/(n+1)) = o(\sqrt{n})$ .

Then, under the global null, if  $m = \lfloor \gamma n \rfloor$  for some  $\gamma \in (0, \infty)$ , as  $n \rightarrow \infty$ ,

$$\mathbb{P} \left[ \sum_{i=1}^m G(\hat{u}^{(\text{marg})}(X_{2n+i})) \geq \xi c_{1-\alpha}(G) - m(\xi - 1) \int_0^1 G(u) du \right] \rightarrow \bar{\Phi} \left( \frac{\xi z_{1-\alpha}}{\sqrt{1+\gamma}} \right), \quad (8)$$

where  $z_{1-\alpha}$  and  $\bar{\Phi}$  denote the  $(1 - \alpha)$ -th quantile and the tail function of the standard normal distribution, respectively. Furthermore, under the same asymptotic regime, for  $W \sim N(0, 1)$ ,

$$\mathbb{P} \left[ \sum_{i=1}^m G(\hat{u}^{(\text{marg})}(X_{2n+i})) \geq \xi c_{1-\alpha}(G) - m(\xi - 1) \int_0^1 G(u) du \mid \mathcal{D} \right] \xrightarrow{d} \bar{\Phi}(\xi z_{1-\alpha} + \sqrt{\gamma}W). \quad (9)$$

**Remark 1.** For Fisher's combination test,  $G(u) = -2 \log u$ . Since  $G(U) \sim \chi^2(2)$ , condition (i) is clearly satisfied. To verify (ii), we note that  $G(u)$  is decreasing and  $|G'(u)| = 2/u$  is decreasing. Thus, for  $u \in [(j-1)/(n+1), j/(n+1)]$ , for  $k \in \{1, 2\}$ ,

$$0 \leq G^k(u) - G^k\left(\frac{j}{n+1}\right) \leq \frac{k}{n+1} G^{k-1}\left(\frac{j}{n+1}\right) G'\left(\frac{j}{n+1}\right) \leq \frac{8 \log(n+1)}{j}.$$

As a result,

$$\begin{aligned} \left| \frac{1}{n+1} \sum_{j=1}^{n+1} G^k(j/(n+1)) - \int_0^1 G^k(u) du \right| &\leq \sum_{j=1}^{n+1} \left| \frac{1}{n+1} G^k\left(\frac{j}{n+1}\right) - \int_{(j-1)/(n+1)}^{j/(n+1)} G^k(u) du \right| \\ &\leq \sum_{j=1}^{n+1} \int_{(j-1)/(n+1)}^{j/(n+1)} |G^k(j/(n+1)) - G^k(u)| du \leq \frac{1}{n+1} \sum_{j=1}^{n+1} \frac{8 \log(n+1)}{j} = O\left(\frac{\log^2 n}{n}\right). \end{aligned}$$

Thus, (ii) is proved. Finally, (iii) is satisfied because  $G(j/(n+1)) \leq G(1/(n+1)) = O(\log n)$ . Therefore, Theorem 1 is a special case of Theorem 4 with  $\xi = 1$ . In general, it is easy to verify (i)–(iii) for various other combination functions (Lipták, 1958; Van Zwet & Oosterhoff, 1967; Vovk & Wang, 2020).

**Remark 2.** By (8), the limiting marginal type-I error is  $\alpha$  when  $\xi = \sqrt{1+\gamma}$ . This implies (7) by noting that  $\int_0^1 (-2 \log u) du = 2$ . By (9), since the random variable  $\bar{\Phi}(\xi z_{1-\alpha} + \sqrt{\gamma}W)$  has a positive density everywhere, the  $(1 - \delta)$ -th quantile of the conditional type-I error converges to the  $(1 - \delta)$ -th quantile of  $\bar{\Phi}(\xi z_{1-\alpha} + \sqrt{\gamma}W)$ , which is  $\bar{\Phi}(\xi z_{1-\alpha} - \sqrt{\gamma}z_{1-\delta})$ . Thus, the conditional type-I error is controlled at level  $\alpha$  with probability at least  $1 - \delta$  asymptotically if  $\xi = 1 + \sqrt{\gamma}z_{1-\delta}/z_{1-\alpha}$ .

**Remark 3.** To confirm our theory, we run Monte-Carlo simulations with  $n = 10^5$  and  $\gamma \in \{2^{-3}, 2^{-2}, \dots, 2^3\}$ , estimating the average type-I error across  $10^4$  samples. Since  $\hat{s}(X)$  is continuous, we can assume that  $\hat{s}(X) \sim \text{Unif}([0, 1])$  without loss of generality, as we will show in the proof. Figure A1 presents the simulated and asymptotic type-I errors for both the unadjusted ( $\xi = 1$ ) and adjusted ( $\xi = \sqrt{1+\gamma}$ ) Fisher's combination test given by (7) with  $G(u) = -2 \log u$ .

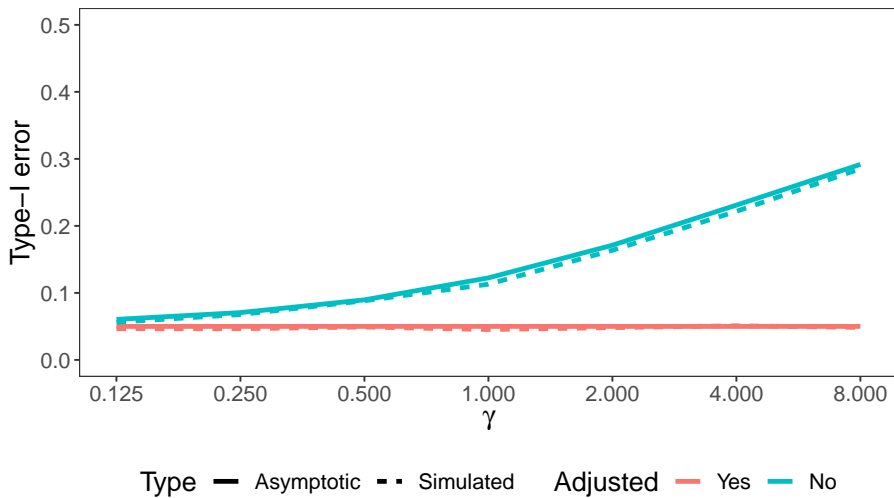


Figure A1: Type-I errors of unadjusted and adjusted Fisher's combination test.

**Remark 4.** If the  $p_i$ 's are dependent, [Brown \(1975\)](#) and [Kost & McDermott \(2002\)](#) approximate the null distribution by a rescaled chi-square distribution  $c\chi^2(f)$ , where  $c$  and  $f$  are chosen to match the mean and variance of the Fisher's combination statistic  $S_{\text{Fisher}}$ . Specifically,

$$c = \frac{\text{Var}[S_{\text{Fisher}}]}{2\mathbb{E}[S_{\text{Fisher}}]}, \quad f = \frac{2\mathbb{E}[S_{\text{Fisher}}]^2}{\text{Var}[S_{\text{Fisher}}]}.$$

In our case, it is easy to see that

$$\mathbb{E}[S_{\text{Fisher}}] \approx 2m, \quad \text{Var}[S_{\text{Fisher}}] \approx 4m(1 + \gamma).$$

As a result, the null distribution is approximated by  $(1 + \gamma)\chi^2(2m/(1 + \gamma))$ . The central limit theorem implies that  $\chi^2(f) \approx N(f, 2f)$  when  $f$  is large. Thus, the critical value for this approximation is

$$(1 + \gamma)\chi^2(2m/(1 + \gamma); 1 - \alpha) \approx (1 + \gamma) \left( \frac{2m}{1 + \gamma} + \sqrt{\frac{2m}{1 + \gamma}} z_{1-\alpha} \right) = 2m + \sqrt{2m(1 + \gamma)} z_{1-\alpha}.$$

Similarly, the critical value for our correction (7) is

$$\sqrt{1 + \gamma}\chi^2(2m; 1 - \alpha) - 2(\sqrt{1 + \gamma} - 1)m \approx \sqrt{1 + \gamma}(2m + \sqrt{2m} z_{1-\alpha}) - 2(\sqrt{1 + \gamma} - 1)m \approx 2m + \sqrt{2m(1 + \gamma)} z_{1-\alpha}.$$

Therefore, both corrections are asymptotically equivalent.

To prove Theorem 4, we start by stating two lemmas. The first lemma is a general Berry-Esseen bound for sums of independent (but not necessarily identically distributed) random variables with potentially infinite third moments.

**Lemma 2.** [[Petrov \(1975\)](#), p. 112, Theorem 5] Let  $X_1, X_2, \dots, X_n$  be independent random variables such that  $\mathbb{E}[X_j] = 0$ , for all  $j$ . Assume also  $\mathbb{E}[X_j^2 g(X_j)] < \infty$  for some function  $g$  that is non-negative, even, and non-decreasing in the interval  $x > 0$ , with  $x/g(x)$  being non-decreasing for  $x > 0$ . Write  $B_n = \sum_j \text{Var}[X_j]$ . Then,

$$d_K \left( \mathcal{L} \left( \frac{1}{\sqrt{B_n}} \sum_{j=1}^n X_j \right), N(0, 1) \right) \leq \frac{A}{B_n g(\sqrt{B_n})} \sum_{j=1}^n \mathbb{E}[X_j^2 g(X_j)],$$

where  $A$  is a universal constant,  $\mathcal{L}(\cdot)$  denotes the probability law,  $d_K$  denotes the Kolmogorov-Smirnov distance (i.e., the  $\ell_\infty$ -norm of the difference of CDFs)

The second lemma is a well-known representation of the spacing between consecutive order statistics.

**Lemma 3** (From [Moran \(1947\)](#); see also Section 4 of [Arnold et al. \(2008\)](#)). Let  $U_1, \dots, U_n \stackrel{i.i.d.}{\sim} \text{Unif}([0, 1])$  and  $U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(n)}$  be their order statistics. Then

$$(U_{(1)} - U_{(0)}, \dots, U_{(n+1)} - U_{(n)}) \stackrel{d}{=} \left( \frac{V_1}{\sum_{k=1}^{n+1} V_k}, \dots, \frac{V_{n+1}}{\sum_{k=1}^{n+1} V_k} \right),$$

where  $U_{(0)} = 0, U_{(n+1)} = 1$ , and  $V_1, \dots, V_{n+1} \stackrel{i.i.d.}{\sim} \text{Exp}(1)$ .

**Proof of Theorem 4.** We first prove the limiting conditional type-I error (9). For convenience, we write  $p_i$  instead of  $\hat{u}^{(\text{marg})}(X_{2n+i})$  and  $S_j$  instead of  $\hat{s}(X_{n+j})$ . Since  $\hat{s}(X)$  is continuous,

$$p_i = \frac{1 + |\{j \in \mathcal{D}^{\text{cal}} : S_j \leq \hat{s}(X_{2n+i})\}|}{n + 1} = \frac{1 + |\{j \in \mathcal{D}^{\text{cal}} : F_S(S_j) \leq F_S(\hat{s}(X_{2n+i}))\}|}{n + 1}$$

where  $F_S$  denotes the CDF of  $\hat{s}(X)$  conditional on  $\mathcal{D}$ . As a result, we can assume  $\hat{s}(X) \sim \text{Unif}([0, 1])$  without loss of generality. Conditional on  $\mathcal{D}$ ,  $p_1, \dots, p_m$  are i.i.d. random variables with

$$\mathbb{P} \left[ p_i = \frac{j}{n + 1} \mid \mathcal{D} \right] = S_{(j)} - S_{(j-1)}, \quad j = 1, \dots, n + 1,$$

where  $S_{(1)} < S_{(2)} < \dots < S_{(n)}$  denote the order statistics of  $(S_1, \dots, S_n)$ , and  $S_{(0)} = 0, S_{(n+1)} = 1$ . By Lemma 3, we can reformulate the distribution of  $p_i$  conditional on  $\mathcal{D}$  as

$$\mathbb{P} \left[ p_i = \frac{j}{n+1} \mid \mathcal{D} \right] = \frac{V_j}{\sum_{k=1}^{n+1} V_k}, \quad j = 1, \dots, n+1. \quad (10)$$

As a result, for  $k \in \{1, 2\}$ ,

$$\mathbb{E} [G^k(p_i) \mid \mathcal{D}] = \frac{\sum_{j=1}^{n+1} G^k \left( \frac{j}{n+1} \right) V_j}{\sum_{j=1}^{n+1} V_j} = \frac{(n+1)^{-1} \sum_{j=1}^{n+1} G^k \left( \frac{j}{n+1} \right) V_j}{(n+1)^{-1} \sum_{j=1}^{n+1} V_j}. \quad (11)$$

By the strong law of large number,

$$\frac{1}{n+1} \sum_{j=1}^{n+1} V_j \xrightarrow{\text{a.s.}} \mathbb{E}[V_1] = 1. \quad (12)$$

Let  $g_n = \max_{j \in \{1, \dots, n+1\}} G(j/(n+1))$ . Since  $V_1, \dots, V_{n+1}$  are independent,

$$\begin{aligned} \text{Var} \left[ \frac{1}{n+1} \sum_{j=1}^{n+1} G^k \left( \frac{j}{n+1} \right) V_j \right] &= \sum_{j=1}^{n+1} \frac{1}{(n+1)^2} \mathbb{E} \left[ G^{2k} \left( \frac{j}{n+1} \right) (V_j - 1)^2 \right] \\ &= \frac{1}{(n+1)^2} \sum_{j=1}^{n+1} G^{2k} \left( \frac{j}{n+1} \right) \leq \frac{g_n^{2k-2}}{(n+1)} \frac{1}{n+1} \sum_{j=1}^{n+1} G^2 \left( \frac{j}{n+1} \right). \end{aligned} \quad (13)$$

By condition (ii),

$$\left| \frac{1}{n+1} \sum_{j=1}^{n+1} G^2 \left( \frac{j}{n+1} \right) - \int_0^1 G^2(u) du \right| = o(1),$$

and thus

$$\frac{1}{n+1} \sum_{j=1}^{n+1} G^2 \left( \frac{j}{n+1} \right) = O(1).$$

By condition (iii),  $g_n = o(\sqrt{n})$ . Together with (13), we obtain that for  $k \in \{1, 2\}$ ,

$$\text{Var} \left[ \frac{1}{n+1} \sum_{j=1}^{n+1} G^k \left( \frac{j}{n+1} \right) V_j \right] = o(1).$$

By Chebyshev's inequality,

$$\frac{1}{n+1} \sum_{i=1}^n G^k \left( \frac{j}{n+1} \right) (V_j - 1) = o_P(1).$$

Applying the condition (ii) again, we arrive at

$$\frac{1}{n+1} \sum_{i=1}^n G^k \left( \frac{j}{n+1} \right) V_j - \int_0^1 G^k(u) du = o_P(1).$$

By (11),

$$\mathbb{E} [G^k(p_i) \mid \mathcal{D}] - \int_0^1 G^k(u) du = o_P(1), \quad k \in \{1, 2\}. \quad (14)$$

Let  $a_n$  be a deterministic sequence such that  $a_n < 1/2$ , and  $U \sim \text{Unif}([0, 1])$ . Let also  $\mathcal{E}_n$  be the event that  $\mathcal{D}$  is such that

$$\mathcal{E}_n = \left\{ \mathcal{D} : \frac{\text{Var}[G(p_i) \mid \mathcal{D}]}{\text{Var}[G(U)]} \in [1 - a_n, 1 + a_n] \right\}. \quad (15)$$

Since  $G$  is a non-constant function,  $\text{Var}[G(U)] > 0$ . By (14), we can choose  $a_n = o(1)$  such that

$$\mathbb{P}[\mathcal{E}_n^c] = o(1).$$

Let

$$W_m = \frac{\sum_{i=1}^m \{G(p_i) - \mathbb{E}[G(p_i) \mid \mathcal{D}]\}}{\sqrt{m \text{Var}[G(p_i) \mid \mathcal{D}]}}.$$

By Lemma 2 with  $g(x) = x$ ,

$$d_K(\mathcal{L}(W_m \mid \mathcal{D}), N(0, 1)) \leq \frac{A}{\sqrt{m}} \frac{\mathbb{E}[|G(p_i) - \mathbb{E}[G(p_i) \mid \mathcal{D}]|^3]}{\text{Var}[G(p_i) \mid \mathcal{D}]^{3/2}},$$

where  $A$  is a universal constant. Since  $G(p_i) \leq g_n$  almost surely, by condition (iii),

$$\mathbb{E}[|G(p_i) - \mathbb{E}[G(p_i) \mid \mathcal{D}]|^3] \leq 2g_n \text{Var}[G(p_i) \mid \mathcal{D}].$$

Thus,

$$d_K(\mathcal{L}(W_m \mid \mathcal{D}), N(0, 1)) \leq \frac{2A}{\sqrt{m}} \frac{g_n}{\text{Var}[G(p_i) \mid \mathcal{D}]^{1/2}}.$$

On the event  $\mathcal{E}_n$ , the condition (iii) and that  $n = O(m)$  imply that

$$d_K(\mathcal{L}(W_m \mid \mathcal{D}), N(0, 1)) \leq \frac{4Ag_n}{\sqrt{m \text{Var}[G(U)]}} = o(1).$$

Since the Kolmogorov distance is invariant under rescalings, we have

$$d_K\left(\mathcal{L}\left(\sqrt{\frac{\text{Var}[G(p_i) \mid \mathcal{D}]}{\text{Var}[G(U)]}} W_m \mid \mathcal{D}\right), N\left(0, \frac{\text{Var}[G(p_i) \mid \mathcal{D}]}{\text{Var}[G(U)]}\right)\right) = o(1).$$

Since  $\text{Var}[G(p_i) \mid \mathcal{D}]/\text{Var}[G(U)] \in [1 - a_n, 1 + a_n] \rightarrow 1$ ,

$$d_K\left(N\left(0, \frac{\text{Var}[G(p_i) \mid \mathcal{D}]}{\text{Var}[G(U)]}\right), N(0, 1)\right) = o(1).$$

Let

$$K_n \triangleq d_K\left(\mathcal{L}\left(\sqrt{\frac{\text{Var}[G(p_i) \mid \mathcal{D}]}{\text{Var}[G(U)]}} W_m \mid \mathcal{D}\right), N(0, 1)\right). \quad (16)$$

The above arguments show that  $K_n = o(1)$  on the event  $\mathcal{E}_n$ .

On the other hand, let

$$c_m = \frac{c_{1-\alpha}(G) - m\mathbb{E}[G(U)]}{\sqrt{m \text{Var}[G(U)]}},$$

and

$$\tilde{W}_n = \frac{\sqrt{n+1}(\mathbb{E}[G(p_i) \mid \mathcal{D}] - \mathbb{E}[G(U)])}{\sqrt{\text{Var}[G(U)]}}.$$

Then

$$\mathbb{P}\left[\sum_{i=1}^m G(p_i) \geq \xi c_{1-\alpha}(G) - m(\xi - 1)\mathbb{E}[G(U)] \mid \mathcal{D}\right] = \mathbb{P}\left[\sqrt{\frac{\text{Var}[G(p_i) \mid \mathcal{D}]}{\text{Var}[G(U)]}} W_m + \sqrt{\frac{m}{n+1}} \tilde{W}_n \geq \xi c_m \mid \mathcal{D}\right].$$

By (16),

$$\left| \mathbb{P} \left[ \sqrt{\frac{\text{Var}[G(p_i) | \mathcal{D}]}{\text{Var}[G(U)]}} W_m + \sqrt{\frac{m}{n+1}} \tilde{W}_n \geq \xi c_m \mid \mathcal{D} \right] - \bar{\Phi} \left( \xi c_m - \sqrt{\frac{m}{n+1}} \tilde{W}_n \right) \right| \leq K_n.$$

Since  $K_n = o(1)$  on  $\mathcal{E}_n$  and  $\mathbb{P}[\mathcal{E}_n^c] = o(1)$ , we obtain that

$$\left| \mathbb{P} \left[ \sum_{i=1}^m G(p_i) \geq c_{1-\alpha}(G) \mid \mathcal{D} \right] - \bar{\Phi} \left( \xi c_m - \sqrt{\frac{m}{n+1}} \tilde{W}_n \right) \right| = o_P(1). \quad (17)$$

Since  $\bar{\Phi}$  is a continuous function and  $m/n \rightarrow \gamma$ , to prove (9), it remains to prove

$$c_m \xrightarrow{P} z_{1-\alpha}, \quad \tilde{W}_n \xrightarrow{d} N(0, 1). \quad (18)$$

Without loss of generality, we assume that  $\eta \leq 1$  in the condition (i). By Lemma 2 with  $g(x) = x^\eta$ , which clearly fulfills the criteria, we have that

$$d_K \left( \frac{\sum_{j=1}^m G(U_j) - \mathbb{E}[G(U)]}{\sqrt{m \text{Var}[G(U)]}}, N(0, 1) \right) \leq \frac{A}{m^{\eta/2}} \frac{\mathbb{E}[|G(U) - \mathbb{E}[G(U)]|^{2+\eta}]}{\text{Var}[G(U)]^{1+\eta/2}} = o(1). \quad (19)$$

By definition,  $c_m$  is the  $(1 - \alpha)$ -th quantile of  $\left( \sum_{j=1}^m G(U_j) - \mathbb{E}[G(U)] \right) / \sqrt{m \text{Var}[G(U)]}$ . By (19),

$$|\bar{\Phi}(c_m) - \alpha| = |\bar{\Phi}(c_m) - \bar{\Phi}(z_{1-\alpha})| = o(1).$$

Since  $\bar{\Phi}'(z_{1-\alpha}) > 0$ , it implies the first part of (18).

To prove the second part of (18), we recall (11) with  $k = 1$  that

$$\tilde{W}_n = \frac{(n+1)^{-1/2} \sum_{j=1}^{n+1} \left\{ G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right\} V_j}{\sqrt{\text{Var}[G(U)]} \left( \sum_{j=1}^{n+1} V_j \right) / (n+1)}.$$

Set  $X_j = (n+1)^{-1/2} \left\{ G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right\} (V_j - 1)$  and  $g(x) = x$  in Lemma 2. Then

$$B_n = \frac{1}{n+1} \sum_{j=1}^{n+1} \left\{ G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right\}^2.$$

By the condition (ii), we have that

$$B_n = \frac{1}{n+1} \sum_{j=1}^n G^2\left(\frac{j}{n+1}\right) - \frac{2\mathbb{E}[G(U)]}{n+1} \sum_{j=1}^n G\left(\frac{j}{n+1}\right) + (\mathbb{E}[G(U)])^2 \rightarrow \text{Var}[G(U)]. \quad (20)$$

By the condition (i), (iii) and (20),

$$\begin{aligned} \sum_{j=1}^{n+1} \mathbb{E}|X_j|^3 &\leq \frac{1}{(n+1)^{3/2}} \sum_{j=1}^{n+1} \left| G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right|^3 \\ &\leq \frac{g_n + \mathbb{E}[G(U)]}{\sqrt{n+1}} \frac{1}{n+1} \sum_{j=1}^{n+1} \left( G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right)^2 \\ &= \frac{g_n + \mathbb{E}[G(U)]}{\sqrt{n+1}} B_n = o(1). \end{aligned}$$

Let

$$\tilde{W}'_n = \frac{1}{\sqrt{B_n}} \sum_{j=1}^{n+1} X_j = \frac{1}{\sqrt{(n+1)B_n}} \sum_{j=1}^{n+1} \left\{ G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right\} (V_j - 1).$$

Then Lemma 2 implies that

$$d_K \left( \tilde{W}'_n, N(0, 1) \right) \leq \frac{A \sum_{j=1}^{n+1} \mathbb{E}|X_j|^3}{B_n^{3/2}} = o(1). \quad (21)$$

By definition,

$$\tilde{W}_n = \left( \tilde{W}'_n + \frac{1}{\sqrt{(n+1)B_n}} \sum_{j=1}^{n+1} \left\{ G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right\} \right) \sqrt{\frac{B_n}{\text{Var}[G(U)]}} \frac{1}{\left( \sum_{j=1}^n V_j \right) / (n+1)}.$$

The condition (ii) with  $k = 1$  implies that

$$\frac{1}{\sqrt{(n+1)}} \sum_{j=1}^{n+1} \left\{ G\left(\frac{j}{n+1}\right) - \mathbb{E}[G(U)] \right\} = o(1). \quad (22)$$

By (12), (20), (21), (22) and Slutsky's Lemma, we prove the second part of (18). Therefore, the limiting conditional type-I error (9) is proved.

Next, we prove the limiting marginal type-I error (8). Since

$$\mathbb{P} \left[ \sum_{i=1}^m G(p_i) \geq \xi_{c_{1-\alpha}}(G) - m(\xi - 1)\mathbb{E}[G(U)] \mid \mathcal{D} \right]$$

is bounded almost surely, the convergence in distribution implies the convergence in expectation. Therefore,

$$\mathbb{P} \left[ \sum_{i=1}^m G(p_i) \geq \xi_{c_{1-\alpha}}(G) - m(\xi - 1)\mathbb{E}[G(U)] \right] \rightarrow \mathbb{E}[\bar{\Phi}(\xi z_{1-\alpha} + \sqrt{\gamma}W)].$$

Let  $W'$  be an independent copy of  $W$ . Then

$$\bar{\Phi}(\xi z_{1-\alpha} + \sqrt{\gamma}W) = \mathbb{P}[W' \geq \xi z_{1-\alpha} + \sqrt{\gamma}W \mid W].$$

As a result,

$$\mathbb{E}[\bar{\Phi}(\xi z_{1-\alpha} + \sqrt{\gamma}W)] = \mathbb{P}[W' \geq \xi z_{1-\alpha} - \sqrt{\gamma}W] = \mathbb{P}[W' - \sqrt{\gamma}W \geq \xi z_{1-\alpha}].$$

The proof is completed by the fact that  $W' - \sqrt{\gamma}W \sim N(0, 1 + \gamma)$ . □

### B.3. Conformal p-values are PRDS

*Proof of Theorem 2.* Let  $Z = (S_{(1)}, \dots, S_{(n)})$  be the order statistics of  $(\hat{s}(X_i))_{i \in \{n+1, \dots, 2n\}}$ , the conformal scores evaluated on the calibration set. Let  $Y = (p_1, \dots, p_m)$  be the conformal p-values evaluated on the test set (i.e.,  $p_j = \hat{u}^{(\text{marg})}(X_{2n+j})$ ). Then,

$$\begin{aligned} \mathbb{P}[Y \in D \mid Y_i = y] &= \int \mathbb{P}[Y \in D \mid Z = z] \mathbb{P}[Z = z \mid Y_i = y] dz \\ &= \mathbb{E}_{Z \mid Y_i = y} [\mathbb{P}[Y \in D \mid Z]]. \end{aligned}$$

With this representation, the conclusion will be implied by the following two lemmas.

**Lemma 4.** For a non-decreasing set  $D$  and vectors  $z, z'$  such that  $z \preceq z'$ , then

$$\mathbb{P}[Y \in D \mid Z = z] \geq \mathbb{P}[Y \in D \mid Z = z'].$$

**Lemma 5.** For  $y \geq y'$ , there exists  $Z_1 \sim Z \mid Y_i = y$  and  $Z_2 \sim Z \mid Y_i = y'$  such that  $\mathbb{P}[Z_1 \preceq Z_2] = 1$ .

In words, Lemma 4 states that the conformal p-values increase as the conformal scores on the calibration set decrease, while Lemma 5 states that a larger conformal p-value indicates the calibration conformal scores are smaller. The proof follows easily from these. Take any  $y \geq y'$  and let  $Z_1$  and  $Z_2$  be as in the statement of Lemma 5. Then,

$$\begin{aligned} \mathbb{P}[Y \in D \mid Y_i = y] &= \mathbb{E}_{Z_1} [\mathbb{P}[Y \in D \mid Z = Z_1]] \\ &\geq \mathbb{E}_{Z_2} [\mathbb{P}[Y \in D \mid Z = Z_2]] \\ &= \mathbb{P}[Y \in D \mid Y_i = y']. \end{aligned}$$

The inequality follows from Lemma 4 and the fact that  $\mathbb{P}[Z_1 \preceq Z_2] = 1$ , which comes from Lemma 5. □

Lemma 4 follows immediately from the definition of marginal conformal p-values in (1). Lemma 5 is proved below.

*Proof of Lemma 5, continuous case.* As in the proof of Theorem 4, since  $\hat{s}(X)$  is continuous, we can assume without loss of generality that the scores  $S_i$  follow the uniform distribution on  $[0, 1]$ . Let  $S'_{(1)} \leq S'_{(2)} \leq \dots \leq S'_{(n+1)}$  be the order statistics of  $(\hat{s}(X_{n+1}), \dots, \hat{s}(X_{2n+1}))$  and  $R_{2n+1}$  be the rank of  $\hat{s}(X_{2n+1})$  among these. By definition,

$$\left\{ (S_{(1)}, \dots, S_{(n)}) \mid R_{2n+1} = k, S'_{(1)}, \dots, S'_{(n+1)} \right\} = (S'_{(1)}, \dots, S'_{(k-1)}, S'_{(k+1)}, \dots, S'_{(n+1)}).$$

Since  $\hat{s}(X)$  is continuous,  $R_{2n+1}$  is independent of  $(S'_{(1)}, S'_{(2)}, \dots, S'_{(n+1)})$ . As a result, for any positive integer  $k \leq n+1$ ,

$$\left\{ (S_{(1)}, \dots, S_{(n)}) \mid R_{2n+1} = k \right\} \stackrel{d}{=} (S'_{(1)}, \dots, S'_{(k-1)}, S'_{(k+1)}, \dots, S'_{(n+1)}).$$

The right-hand-side is clearly entry-wise non-increasing in  $k$ . Since  $p_1 = R_{2n+1}/(n+1)$ , Lemma 5 is proved for  $i = 1$ . The same proof carries over to other indices  $i$ . □

**Extension to non-continuous scores.** When  $\hat{s}(X)$  has atoms, the set of conformity scores  $\{\hat{s}(X_i) : i \in \mathcal{D}^{\text{cal}}\}$  have ties with non-zero probability. In this case, we replace the marginal conformal p-value (2) by a randomized version, i.e.,

$$p_j = \frac{|\{i \in \mathcal{D}^{\text{cal}} : \hat{s}(X_i) < \hat{s}(X_{2n+j})\}| + \lceil (1 + |\{i \in \mathcal{D}^{\text{cal}} : \hat{s}(X_i) = \hat{s}(X_{2n+j})\}|)U_j \rceil}{n+1}, \quad (23)$$

where  $U_1, U_2, \dots$  are i.i.d. random variables drawn from  $\text{Unif}([0, 1])$  which are independent of the data. Note that (23) is identical to (2) almost surely when  $\hat{s}(X)$  is continuous. Now we prove that the marginal conformal p-values defined in (23) satisfy the PRDS property.

**Proposition 2** (Theorem 2 for the non-continuous case). *Consider the setting of Theorem 2, but where  $\hat{s}(\cdot)$  is not assumed to be continuous. Define the randomized marginal p-values as in (23). Then, the marginal conformal p-values  $(p_1, \dots, p_m)$  are PRDS.*

The proof follows as above, once we verify Lemma 4 and Lemma 5 in the more general setting.

*Proof of Lemma 4, general case.* Let  $U = (U_1, \dots, U_m)$ . By definition,  $U$  is independent of  $(Y, Z)$ , and thus

$$\mathbb{P}[Y \in D \mid Z = z] = \mathbb{P}[Y \in D \mid Z = z, U], \quad \text{a.s.}$$

Let  $p_j(x; z, u)$  denote the mapping from  $(X_{2n+j}, Z, U)$  to  $p_j$ . Then

$$p_j(x; z, u) = \frac{m_{<}(x; z) + \lceil \{1 + m_{=}(x; z)\}u \rceil}{n+1},$$

where

$$m_{<}(x, z) = \sum_{i=1}^n I(z_i < x), \quad m_{=}(x, z) = \sum_{i=1}^n I(z_i = x).$$

If  $z \preceq z'$ ,

$$m_{<}(x, z) \geq m_{<}(x, z'), \quad m_{<}(x, z) + m_{=}(x, z) \geq m_{<}(x, z') + m_{=}(x, z'). \quad (24)$$

We claim that the mapping  $p_j(x; z, u)$  is non-increasing in  $z$  for every  $x$  and  $u$ . Equivalently, we will show that for any  $x$  and  $u \in [0, 1]$ ,

$$m_{<}(x, z) + \lceil \{1 + m_{=}(x, z)\}u \rceil \geq m_{<}(x, z') + \lceil \{1 + m_{=}(x, z')\}u \rceil. \quad (25)$$

We consider three cases.

Case 1: if  $m_{<}(x, z) = m_{<}(x, z')$ , (24) implies that  $m_{=}(x, z) \geq m_{=}(x, z')$ . Thus, (25) is obviously true.

Case 2: if  $m_{<}(x, z) + m_{=}(x, z) = m_{<}(x, z') + m_{=}(x, z')$ , let  $a = 1 + m_{=}(x, z)$  and  $b = m_{<}(x, z) - m_{<}(x, z')$ . Then (25) is equivalent to

$$b \geq \lceil (a + b)u \rceil - \lceil au \rceil.$$

This can be proved using the fact that  $\lceil (a + b)u \rceil \leq \lceil au \rceil + \lceil bu \rceil$ .

Case 3: if  $m_{<}(x, z) > m_{<}(x, z')$  and  $m_{<}(x, z) + m_{=}(x, z) > m_{<}(x, z') + m_{=}(x, z')$ , then  $m_{<}(x, z) \geq m_{<}(x, z') + 1$  and  $m_{<}(x, z) + m_{=}(x, z) \geq m_{<}(x, z') + m_{=}(x, z') + 1$  since  $m_{<}(x, z)$ ,  $m_{<}(x, z')$ ,  $m_{=}(x, z)$ , and  $m_{=}(x, z')$  are all integers. Then

$$\begin{aligned} m_{<}(x, z) + \lceil \{1 + m_{=}(x, z)\}u \rceil &\geq m_{<}(x, z) + \{1 + m_{=}(x, z)\}u \\ &= m_{<}(x, z)(1 - u) + \{1 + m_{=}(x, z) + m_{<}(x, z)\}u \\ &\geq \{1 + m_{<}(x, z')\}(1 - u) + \{2 + m_{=}(x, z') + m_{<}(x, z')\}u \\ &= m_{<}(x, z') + \{1 + m_{=}(x, z')\}u + 1 \\ &\geq m_{<}(x, z') + \lceil \{1 + m_{=}(x, z')\}u \rceil. \end{aligned}$$

Therefore, (25) is proved. As a result, the mapping from  $(X_{2n+1}, \dots, X_{2n+m}, Z, U)$  to  $Y$  is entry-wise non-increasing in  $Z$  given  $(X_{2n+j}, \dots, X_{2n+m}, U)$ . Since  $\{X_{2n+j} : j = 1, \dots, m\}$ ,  $Z$ , and  $U$  are mutually independent, we arrive at

$$\mathbb{P}[Y \in D \mid Z = z, U] \geq \mathbb{P}[Y \in D \mid Z = z', U], \quad \text{a.s.}$$

The independence between  $U$  and  $Z$  implies that  $(U \mid Z = z) \stackrel{d}{=} (U \mid Z = z')$ . Lemma 4 then follows from the above inequality.  $\square$

*Proof of Lemma 5, general case.* Let  $R_{2n+j} = (n+1)p_j$ . Note that  $R_{2n+j}$  can be interpreted as the rank with ties broken randomly. As in the proof for the continuous case, we first prove that

$$\left\{ (S_{(1)}, \dots, S_{(n)}) \mid R_{2n+1} = k, S'_{(1)}, \dots, S'_{(n+1)} \right\} = (S'_{(1)}, \dots, S'_{(k-1)}, S'_{(k+1)}, \dots, S'_{(n+1)}). \quad (26)$$

Let  $k_- = \max\{\ell : S'_{(\ell)} < S_{2n+1}\}$  and  $k_+ = \min\{\ell : S'_{(\ell)} > S_{2n+1}\}$ . Then  $S'_\ell = S_{2n+1}$  for any  $k_- < \ell < k_+$ . Since there exists at least one  $\ell$  with  $S'_{(\ell)} = S_{2n+1}$ , i.e., the index corresponding to  $S_{2n+1}$ , we have  $k_+ - k_- \geq 2$ . By definition,

$$1 + |\{i \in \mathcal{D}^{\text{cal}} : \hat{s}(X_i) = \hat{s}(X_{2n+j})\}| = |\{i \in \mathcal{D}^{\text{cal}} \cup \{2n+1\} : \hat{s}(X_i) = \hat{s}(X_{2n+j})\}| = k_+ - k_- - 1.$$

As a result,

$$k = k_- + \lceil (k_+ - k_- - 1)U_1 \rceil \in (k_-, k_+).$$

Therefore,  $\hat{s}(X_{2n+1}) = S'_{(k)}$  and (26) is proved.

It remains to prove that  $R_{2n+1}$  is independent of  $(S'_{(1)}, S'_{(2)}, \dots, S'_{(n+1)})$ . For any non-decreasing sequence  $a_1 \leq \dots \leq a_{n+1}$ , let  $1 = n_0 < n_1 < \dots < n_m = n + 1$  be integers such that

$$a_{n_{j-1}} = \dots = a_{n_j-1} < a_{n_j}, \quad j = 1, \dots, m-1, \quad a_{n_{m-1}-1} < a_{n_{m-1}} = \dots = a_{n_m}$$

Let  $\pi : \{1, \dots, n+1\} \mapsto \{1, \dots, n+1\}$  be a uniform random permutation. Since  $X_{n+1}, \dots, X_{2n+1}$  are i.i.d., Conditioning on the event that,

$$\left\{ (\hat{s}(X_{n+1}), \dots, \hat{s}(X_{2n+1})) \mid (S'_{(1)}, \dots, S'_{(n+1)}) = (a_1, \dots, a_{n+1}) \right\} \stackrel{d}{=} (a_{\pi(1)}, \dots, a_{\pi(n+1)}).$$

For any  $j = 1, \dots, m-1$ , if  $\pi(n+1) \in [n_{j-1}, n_j]$ ,

$$|\{i : a_{\pi(i)} = a_{\pi(n+1)}\}| = n_j - n_{j-1}, \quad |\{i : a_{\pi(i)} < a_{\pi(n+1)}\}| = n_{j-1} - 1,$$

and thus,

$$R_{2n+1} = n_{j-1} - 1 + \lceil (n_j - n_{j-1})U_j \rceil.$$

Similarly, if  $\pi(n+1) \in [n_{m-1}, n_m]$ ,

$$R_{2n+1} = n_{m-1} - 1 + \lceil (n_m - n_{m-1} + 1)U_j \rceil.$$

For any  $k$ , let  $j_k = \max\{j : n_j \leq k\}$ , and  $\mathcal{I}_k$  be the set  $\{n_{j_k-1}, \dots, n_{j_k} - 1\}$  if  $j_k < m$  and  $\{n_{j_k-1}, \dots, n_{j_k}\}$  otherwise. Then

$$\begin{aligned} & \mathbb{P}(R_{2n+1} = k \mid (S'_{(1)}, \dots, S'_{(n+1)}) = (a_1, \dots, a_{n+1})) \\ &= \mathbb{P}\left(\pi(n+1) \in \mathcal{I}_k, U_1 \in \left(\frac{k - n_{j_k-1}}{|\mathcal{I}_k|}, \frac{k + 1 - n_{j_k-1}}{|\mathcal{I}_k|}\right]\right) \\ &= \mathbb{P}(\pi(n+1) \in \mathcal{I}_k) \mathbb{P}\left(U_1 \in \left(\frac{k - n_{j_k-1}}{|\mathcal{I}_k|}, \frac{k + 1 - n_{j_k-1}}{|\mathcal{I}_k|}\right]\right) \\ &= \frac{|\mathcal{I}_k|}{n+1} \frac{1}{|\mathcal{I}_k|} = \frac{1}{n+1}. \end{aligned}$$

Therefore,  $R_{2n+1}$  is independent of  $(S'_{(1)}, \dots, S'_{(n+1)})$ . The proof of Lemma 5 is then completed.  $\square$

#### B.4. Storey's correction does not break FDR control

When the proportion of nulls is much smaller than 1, as it may be the case in many out-of-distribution detection problems, the BH procedure controls the FDR at level  $\pi_0\alpha$  (Benjamini & Yekutieli, 2001), where  $\pi_0$  is the proportion of nulls, and hence is conservative. If  $\pi_0$  is known, a simple remedy is to replace the target FDR level with  $\alpha/\pi_0$ . However,  $\pi_0$  is rarely known in practice and hence it needs to be estimated. Given p-values  $p_i$  for all null hypotheses, it was proposed by Storey et al. in Storey (2002); Storey et al. (2004) to estimate  $\pi_0$  as

$$\hat{\pi}_0 = \frac{1 + \sum_{i=1}^m I(p_i > \lambda)}{m(1 - \lambda)},$$

and then to apply the BH procedure at level  $\alpha/\hat{\pi}_0$ . Specifically, given a p-value  $p_i$  for the  $i$ -th null hypothesis, let  $p_{(1)} \leq \dots \leq p_{(m)}$  be the ordered statistics. Given a target FDR level  $\alpha$  and a scalar  $\lambda \in (0, 1)$ , the rejection set of the Storey-BH procedure is

$$\mathcal{R} = \left\{ i : p_i \leq \frac{\alpha R}{m\hat{\pi}_0}, p_i < \lambda \right\},$$

where

$$\hat{\pi}_0 = \frac{1 + \sum_{i=1}^m I(p_i \geq \lambda)}{m(1 - \lambda)} \triangleq \frac{1 + A}{m(1 - \lambda)}$$

and

$$R = \max \left\{ r : p_{(r)} \leq \frac{\alpha r}{m\hat{\pi}_0}, p_{(r)} < \lambda \right\}.$$

The parameter  $\lambda$  is often chosen as 0.5 or  $\alpha$  or  $1 - \alpha$ .

If the null p-values are super-uniform in the sense of (2), mutually independent, and independent of the non-null p-values, this provably controls the FDR in finite samples (Storey et al., 2004). However, unlike in its standard version, the BH procedure with Storey's correction may fail to control the FDR if the p-values are PRDS; see Section 6.3 of Benjamini et al. (2006). Surprisingly, we show below that the positive correlation (Lemma 1) among the marginal conformal p-values does not break the FDR control at all.

**Theorem 5.** Set  $\lambda = K/(n+1)$  for any integer  $K$ . Assume  $\hat{s}(X)$  is continuous. In the setting of Corollary 1, the BH procedure with Storey's correction applied at level  $\alpha \in (0, 1)$  to the marginal  $p$ -values  $(\hat{u}^{(\text{marg})}(X_{2n+1}), \dots, \hat{u}^{(\text{marg})}(X_{2n+m}))$  controls the FDR at level  $\alpha$ . That is,

$$\mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] \leq \alpha. \quad (27)$$

The proof of Theorem 5 rests on a novel FDR bound for the BH procedure with Storey's correction applied to PRDS  $p$ -values, formalized in the following theorem.

**Theorem 6.** Assume that  $(p_1, \dots, p_n)$  is PRDS and each null  $p$ -value is super-uniform with an almost sure lower bound  $p_{\min} \in [0, 1]$ . Then

$$\mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] \leq \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \mathbb{E} \left[ \frac{1}{1+A} \mid p_i \leq p_* \right],$$

where

$$p_* = \max \left\{ \frac{\alpha(1-\lambda)}{m}, p_{\min} \right\}.$$

*Proof.* Let

$$V_i = I(H_i \text{ is rejected}) \leq I \left( p_i \leq \alpha(1-\lambda) \frac{R}{1+A} \right).$$

Then

$$\begin{aligned} \mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] &= \sum_{i \in \mathcal{H}_0} \mathbb{E} \left[ \frac{V_i}{R \vee 1} \right] = \sum_{i \in \mathcal{H}_0} \sum_{r=1}^m \frac{1}{r} \mathbb{P} \left( p_i \leq \alpha(1-\lambda) \frac{r}{1+A}, R=r \right) \\ &= \sum_{i \in \mathcal{H}_0} \sum_{r=1}^m \sum_{a=1}^m \frac{1}{r} \mathbb{P} \left( p_i \leq \alpha(1-\lambda) \frac{r}{1+a}, R=r, A=a \right). \end{aligned}$$

Let  $r_0(a) = \max\{1, \lceil (1+a)p_{\min}/(1-\lambda)\alpha \rceil\}$ . By definition, the summand for a given  $a$  is non-zero only if  $r \geq r_0(a)$ . Thus,

$$\begin{aligned} \mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] &= \sum_{i \in \mathcal{H}_0} \sum_{a=1}^m \sum_{r=r_0(a)}^m \frac{1}{r} \mathbb{P} \left( p_i \leq \alpha(1-\lambda) \frac{r}{1+a} \right) \mathbb{P} \left( R=r, A=a \mid p_i \leq \alpha(1-\lambda) \frac{r}{1+a} \right) \\ &\stackrel{(i)}{\leq} \sum_{i \in \mathcal{H}_0} \sum_{a=1}^m \sum_{r=r_0(a)}^m \frac{1}{r} \cdot \alpha(1-\lambda) \frac{r}{1+a} \mathbb{P} \left( R=r, A=a \mid p_i \leq \alpha(1-\lambda) \frac{r}{1+a} \right) \\ &= \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \sum_{a=1}^m \sum_{r=r_0(a)}^m \frac{1}{1+a} \mathbb{P} \left( R=r, A=a \mid p_i \leq \alpha(1-\lambda) \frac{r}{1+a} \right), \end{aligned}$$

where (i) uses the super-uniformity of the null  $p$ -value. Let  $\mathcal{T}$  denote the set of all possible values that  $r/(1+a)$  can take such that  $\mathbb{P}(p_i \leq \alpha(1-\lambda)r/(1+a)) > 0$ , i.e.

$$\mathcal{T} = \left\{ \frac{r}{1+a} : a \in \{1, \dots, m\}, r \in \{r_0(a), \dots, m\}, a+r \leq m \right\}.$$

Clearly,  $\mathcal{T}$  is a finite set. Let  $t_1 \leq t_2 \leq \dots \leq t_M$  denote the values of  $\mathcal{T}$ . It is easy to see that

$$\alpha(1-\lambda)t_1 \geq \max \left\{ p_{\min}, \frac{\alpha(1-\lambda)}{m} \right\} = p_*. \quad (28)$$

Then

$$\begin{aligned}
 \mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] &\leq \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \sum_{j=1}^M \sum_{a=1}^m \frac{1}{1+a} \mathbb{P}(R = (1+a)t_j, A = a \mid p_i \leq \alpha(1-\lambda)t_j) \\
 &= \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \sum_{j=1}^M \mathbb{E} \left[ \frac{I\{R = (1+A)t_j\}}{1+A} \mid p_i \leq \alpha(1-\lambda)t_j \right] \\
 &= \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \sum_{j=1}^M \left\{ \mathbb{E}[H_j(p) \mid p_i \leq \alpha(1-\lambda)t_j] - \mathbb{E}[H_{j+1}(p) \mid p_i \leq \alpha(1-\lambda)t_j] \right\} \\
 &= \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \left\{ \mathbb{E}[H_1(p) \mid p_i \leq \alpha(1-\lambda)t_1] \right. \\
 &\quad \left. - \sum_{j=1}^{M-1} (\mathbb{E}[H_{j+1}(p) \mid p_i \leq \alpha(1-\lambda)t_j] - \mathbb{E}[H_{j+1}(p) \mid p_i \leq \alpha(1-\lambda)t_{j+1}]) \right\},
 \end{aligned}$$

where

$$H_j(p) = \frac{I\{R \geq (1+A)t_j\}}{1+A}, \quad H_{M+1}(p) = 0.$$

Since  $A$  is an increasing function of  $p$  and  $R$  is a decreasing function of  $p$ ,  $H_j(p)$  is decreasing in  $p$ . The PRDS property implies that for any  $j = 1, \dots, M-1$ ,

$$\mathbb{E}[H_{j+1}(p) \mid p_i \leq \alpha(1-\lambda)t_j] - \mathbb{E}[H_{j+1}(p) \mid p_i \leq \alpha(1-\lambda)t_{j+1}] \geq 0.$$

Therefore,

$$\begin{aligned}
 \mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] &\leq \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \mathbb{E}[H_1(p) \mid p_i \leq \alpha(1-\lambda)t_1] \\
 &\leq \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \mathbb{E} \left[ \frac{1}{1+A} \mid p_i \leq \alpha(1-\lambda)t_1 \right] \\
 &\leq \alpha(1-\lambda) \sum_{i \in \mathcal{H}_0} \mathbb{E} \left[ \frac{1}{1+A} \mid p_i \leq p_* \right],
 \end{aligned}$$

where the last step follows from (28), the PRDS property, and the fact that  $p \mapsto 1/(1+A)$  is decreasing.  $\square$

To prove Theorem 5, we present an additional lemma.

**Lemma 6.** [Lemma 1 from Benjamini et al. (2006)] If  $Y \sim \text{Binom}(k-1, p)$ , then  $\mathbb{E}[1/(1+Y)] \leq 1/kp$ .

*Proof of Theorem 5.* As in the proof of Theorem 4, since  $\hat{s}(X)$  is continuous, we can assume  $\hat{s}(X) \sim \text{Unif}([0, 1])$  without loss of generality. We write  $p_i$  instead of  $\hat{u}^{(\text{marg})}(X_{2n+i})$  and  $S_j$  instead of  $\hat{s}(X_{n+j})$ . Then

$$p_j = \frac{1 + \sum_{i=1}^n I(S_i \leq S_{n+j})}{n+1}.$$

Then  $p_j \geq 1/(n+1)$  almost surely. Let  $m_0 = |\mathcal{H}_0|$  and we assume that  $\mathcal{H}_0 = \{1, \dots, m_0\}$  without loss of generality. Since  $p = (p_1, \dots, p_m)$  are PRDS and exchangeable, Theorem 7 implies that

$$\mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] \leq \alpha(1-\lambda)m_0 \mathbb{E} \left[ \frac{1}{1+A} \mid p_1 \leq \max \left\{ \frac{1}{n+1}, \frac{\alpha(1-\lambda)}{m} \right\} \right].$$

Since  $1/(1+A)$  is decreasing in  $p$ , using the PRDS property again, we have

$$\mathbb{E} \left[ \frac{|\mathcal{R} \cap \mathcal{H}_0|}{\max\{1, |\mathcal{R}|\}} \right] \leq \alpha(1-\lambda)m_0 \mathbb{E} \left[ \frac{1}{1+A} \mid p_1 \leq \frac{1}{n+1} \right] = \alpha(1-\lambda)m_0 \mathbb{E} \left[ \frac{1}{1+A} \mid p_1 = \frac{1}{n+1} \right]. \quad (29)$$

Let  $A_0 = \sum_{j=2}^{m_0} I(p_j \geq \lambda)$ . Then

$$\mathbb{E} \left[ \frac{1}{1+A} \mid p_1 = \frac{1}{n+1} \right] \leq \mathbb{E} \left[ \frac{1}{1+A_0} \mid p_1 = \frac{1}{n+1} \right].$$

Let  $S_{(1)} \leq S_{(2)} \leq \dots \leq S_{(n+1)}$  denote the order statistics of  $S_1, \dots, S_{n+1}$  and  $R_{n+1}$  denote the rank of  $S_{n+1}$ . Since  $S_1 \sim \text{Unif}([0, 1])$ , there is no tie almost surely.

Now we compute

$$\mathbb{E} \left[ \frac{1}{1+A_0} \mid p_1 = \frac{1}{n+1}, S_{(1)}, \dots, S_{(n+1)} \right] = \mathbb{E} \left[ \frac{1}{1+A_0} \mid R_{n+1} = 1, S_{(1)}, \dots, S_{(n+1)} \right]. \quad (30)$$

By definition,

$$p_2, \dots, p_{m_0} \mid S_1, \dots, S_{n+1} \stackrel{i.i.d.}{\sim} \frac{1 + \sum_{j=1}^n I(S_j \leq U)}{n+1}$$

where  $U \sim \text{Unif}([0, 1])$ . Note that there is a bijection between  $(S_1, \dots, S_{n+1})$  and  $(S_{(1)}, \dots, S_{(n+1)}, R_1, \dots, R_{n+1})$  for vectors without ties. The above distributional equivalence can be rewritten as

$$p_2, \dots, p_{m_0} \mid R_1, \dots, R_{n+1}, S_{(1)}, \dots, S_{(n+1)} \stackrel{i.i.d.}{\sim} \frac{1 + \sum_{j=1}^{n+1} I(S_{(j)} \leq U) - I(S_{(R_{n+1})} \leq U)}{n+1}.$$

Since the RHS does not depend on  $(R_1, \dots, R_n)$ ,  $(p_2, \dots, p_{m_0})$  is independent of  $(R_1, \dots, R_n)$  conditional on  $(R_{n+1}, S_{(1)}, \dots, S_{(n+1)})$ . As a result,

$$p_2, \dots, p_{m_0} \mid R_{n+1} = 1, S_{(1)}, \dots, S_{(n+1)} \stackrel{i.i.d.}{\sim} \frac{1 + \sum_{j=2}^{n+1} I(S_{(j)} \leq U)}{n+1}.$$

Recall  $K = (n+1)\lambda \in \mathbb{Z}$ . Then

$$\begin{aligned} \mathbb{P}(p_2 \geq \lambda \mid R_{n+1} = 1, S_{(1)}, \dots, S_{(n+1)}) &= \mathbb{P} \left( \sum_{j=2}^{n+1} I(S_{(j)} \leq U) \geq K - 1 \mid S_{(2)}, \dots, S_{(n+1)} \right) \\ &= \mathbb{P}(U \geq S_{(K)} \mid S_{(2)}, \dots, S_{(n+1)}) \\ &= 1 - S_{(K)}. \end{aligned}$$

Therefore,

$$I(p_2 \geq \lambda), \dots, I(p_{m_0} \geq \lambda) \mid R_{n+1} = 1, S_{(1)}, \dots, S_{(n+1)} \stackrel{i.i.d.}{\sim} \text{Ber}(1 - S_{(K)}).$$

This implies that

$$A_0 \mid R_{n+1} = 1, S_{(1)}, \dots, S_{(n+1)} \sim \text{Binom}(m_0 - 1, 1 - S_{(K)}).$$

By Lemma 6,

$$\mathbb{E} \left[ \frac{1}{1+A_0} \mid R_{n+1} = 1, S_{(1)}, \dots, S_{(n+1)} \right] \leq \frac{1}{m_0 \{1 - S_{(K)}\}}.$$

Since  $R_{n+1}$  is independent of  $(S_{(1)}, \dots, S_{(n+1)})$ ,

$$\mathbb{E} \left[ \frac{1}{1+A_0} \mid R_{n+1} = 1 \right] \leq \mathbb{E} \left[ \frac{1}{m_0 \{1 - S_{(K)}\}} \right]. \quad (31)$$

By symmetry and the property of order statistics,

$$1 - S_{(K)} \stackrel{d}{=} S_{(n+2-K)} \sim \text{Beta}(n+2-K, K).$$

Thus,

$$\begin{aligned}
 \mathbb{E} \left[ \frac{1}{1 - S_{(K)}} \right] &= \int_0^1 \frac{1}{x} \frac{\Gamma(n+2)}{\Gamma(n+2-K)\Gamma(K)} x^{n+1-K} (1-x)^{K-1} dx \\
 &= \int_0^1 \frac{\Gamma(n+2)}{\Gamma(n+2-K)\Gamma(K)} x^{n-K} (1-x)^{K-1} dx \\
 &= \frac{\Gamma(n+2)\Gamma(n+1-K)}{\Gamma(n+2-K)\Gamma(n+1)} \\
 &= \frac{n+1}{n+1-K} = \frac{1}{1-\lambda}.
 \end{aligned} \tag{32}$$

Putting (29), (31) and (32) together, we prove the result.  $\square$

### B.5. Conditional p-value adjustment

We say that  $\hat{u}^{(\text{ccv})}(X)$  satisfies the *calibration-conditional validity* (CCV) property if

$$\mathbb{P} \left[ \mathbb{P} \left[ \hat{u}^{(\text{ccv})}(X_{2n+1}) \leq t \mid \mathcal{D} \right] \leq t \text{ for all } t \in (0, 1] \right] \geq 1 - \delta. \tag{33}$$

The following theorem suggests a generic strategy through a simultaneous upper confidence bound for order statistics.

**Theorem 7** (Conditional p-value adjustment). *Let  $U_1, \dots, U_n \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0, 1])$ , with order statistics  $U_{(1)} \leq U_{(2)} \leq \dots \leq U_{(n)}$ , and fix any  $\delta \in (0, 1)$ . Suppose  $0 \leq b_1 \leq b_2 \leq \dots \leq b_n \leq 1$  are  $n$  reals such that*

$$\mathbb{P} \left[ U_{(1)} \leq b_1, \dots, U_{(n)} \leq b_n \right] \geq 1 - \delta. \tag{34}$$

Let also  $b_0 = 0, b_{n+1} = 1$ , and  $h : [0, 1] \mapsto [0, 1]$  be a piece-wise constant function such that

$$h(t) = b_{\lceil (n+1)t \rceil}, \quad t \in [0, 1].$$

Then,  $\hat{u}^{(\text{ccv})} = h \circ \hat{u}^{(\text{marg})}$  satisfies (33), i.e.,  $\hat{u}^{(\text{ccv})}(X_{2n+1})$  is a calibration-conditional valid p-value.

*Proof of Theorem 7.* Let  $S_i = \hat{s}(X_{n+i})$  for  $i = 1, \dots, n$  with  $F^-(t) = \mathbb{P}[S_i < t \mid \mathcal{D}^{\text{train}}]$ , and  $S_{(1)} \leq S_{(2)} \leq \dots \leq S_{(n)}$  be the order statistics. Then it is easy to see that

$$(F^-(S_{(1)}), \dots, F^-(S_{(n)})) \preceq (U_{(1)}, \dots, U_{(n)}),$$

where  $\preceq$  denotes the entry-wise stochastic dominance in the sense that  $(A_1, \dots, A_n) \preceq (B_1, \dots, B_n)$  iff

$$\mathbb{P}[A_1 \leq z_1, \dots, A_n \leq z_n] \geq \mathbb{P}[B_1 \leq z_1, \dots, B_n \leq z_n], \quad \forall (z_1, \dots, z_n) \in \mathbb{R}^n.$$

When  $F$  is continuous, the equality in distribution holds. Let  $\mathcal{E}_n$  denote the event on which  $F^-(S_{(i)}) \leq b_i$  for all  $i = 1, \dots, n$ . Then

$$\mathbb{P}[\mathcal{E}_n] \geq 1 - \delta.$$

Now we prove the following claim, which directly yields the theorem:

$$\mathbb{P} \left[ \hat{u}^{(\text{ccv})}(X_{2n+1}) \leq t \mid \mathcal{D} \right] \leq t, \quad \forall t \in [0, 1], \quad \text{if } \mathcal{D} \in \mathcal{E}_n. \tag{35}$$

Note that the image of  $\hat{u}^{(\text{ccv})}$  is  $\{b_1, \dots, b_n, 1\}$ , it remains to prove (35) with  $t \in \{b_1, \dots, b_n, 1\}$ . When  $t = 1$ , it clearly holds. When  $t = b_i$ ,

$$\hat{u}^{(\text{ccv})}(X_{2n+1}) \leq b_i \iff \hat{u}^{(\text{marg})}(X_{2n+1}) \leq \frac{i}{n+1} \iff \hat{s}(X_{2n+1}) < S_{(i)}.$$

Thus,

$$\mathbb{P} \left[ \hat{u}^{(\text{ccv})}(X_{2n+1}) \leq b_i \mid \mathcal{D} \right] = \mathbb{P} \left[ \hat{s}(X_{2n+1}) < S_{(i)} \mid \mathcal{D} \right] = F^-(S_{(i)}).$$

By definition of  $\mathcal{E}_n$ , (35) holds for all  $t \in \{b_1, \dots, b_n\}$ .  $\square$

The larger p-values typically do not matter in multiple testing problems, as it is the small ones that determine which hypotheses are rejected. Therefore, to maximize power, we would like the  $b_i$  values in Theorem 7 to be as small as possible for low indices  $i$ , while we may be satisfied with letting  $b_i = 1$  for large  $i$ . The generalized Simes inequality yields a desirable class of  $(b_1, \dots, b_n)$  sequences with this property.

**Proposition 3** (Generalized Simes Inequality, from Equation (3.5) in Sarkar et al. (2008)). *For any positive integer  $k \leq n$ , the uniform bound (34) in Theorem 7 holds with*

$$b_{n+1-i} = 1 - \delta^{1/k} \left( \frac{i \cdots (i-k+1)}{n \cdots (n-k+1)} \right)^{1/k}, \quad i = 1, \dots, n.$$

The original motivation of Sarkar et al. (2008) was to compute thresholds for step-up procedure to achieve  $k$ -FWER control; there, the parameter  $k$  was set to be a small integer. Here, we exploit Proposition 3 differently, choosing  $k = n/2$  so that the  $b_i$  values with lower indices  $i$  are as small as possible while those with larger indices  $i$  may be uninformative (note that  $b_{n-k+2} = \dots = b_n = 1$ ). In particular, our choice corresponds to

$$b_1 = 1 - \delta^{2/n} = 1 - \exp \left\{ -\frac{2 \log(1/\delta)}{n} \right\} \approx \frac{2 \log(1/\delta)}{n}.$$

Therefore, the smallest possible marginal p-value equal to  $1/(n+1)$  would be mapped to  $h(1/(n+1)) \approx 2 \log(10)/n = 4.61/n$ , if  $\delta = 0.1$ , for example, since  $\hat{u}^{(\text{ccv})}(X) = h(\hat{u}^{(\text{marg})}(X))$ . If  $n = 10000$ , then  $h(1/(n+1)) < 0.0005$ , which is larger than the marginal p-value, but much smaller than what one would obtain from other standard uniform bounds. For example, the DKWM inequality (Dvoretzky et al., 1956; Massart, 1990) would imply a result similar to that of Proposition 3 but with  $b_i = \min\{i/n + \sqrt{\log(2/\delta)/2n}, 1\}$ ; this would map the smallest possible marginal p-value to  $1/(n+1) + \sqrt{\log(2/\delta)/2n} > 0.1$ , in the above example. The comparison between the generalized Simes inequality and the DKWM inequality is expanded in Appendix C, where we also consider an additional uniform bound based on the linear-boundary crossing probability for the empirical CDF (Dempster, 1959). This comparison confirms the generalized Simes inequality yields the most powerful adjustment for our multiple testing purposes. In practice, we find that  $k = n/2$  works well, as motivated empirically in Appendix D. (Note that larger values of  $k$  would lower further the smallest possible adjusted p-value, but at the cost of raising other small p-values).

## B.6. Simultaneous confidence bounds for the false positive rate

Some practitioners may be accustomed to thinking about outlier detection in terms of FPR—the probability of incorrectly reporting as outlier any true inlier—rather than p-values. In particular, they may wonder what the FPR can be if they report  $X_{2n+1}$  as likely to be an outlier whenever the classification score  $\hat{s}(X_{2n+1})$  (computed by some black-box outlier detection algorithm) is above a threshold  $t$ , as a function of  $t$ , so that they may choose a posteriori which value of  $t$  to adopt. This question is closely related to the problem of constructing CCV p-values, so our method provides an answer. In fact, the next result shows Theorem 7 also yields a simultaneous upper confidence bound for the CDF.

**Proposition 4.** *Let  $F$  denote the true CDF of some distribution from which  $n$  i.i.d. samples,  $Z_1, \dots, Z_n$ , are drawn, and denote by  $\hat{F}_n$  the corresponding empirical CDF. With the same notation as in Theorem 7,*

$$\mathbb{P} \left[ F(z) \leq h(\hat{F}_n(z)), \quad \forall z \in \mathbb{R} \right] \geq 1 - \delta. \quad (36)$$

*Proof of Proposition 4.* Note that  $h(i/n) = b_{\lceil i+1/n \rceil} = b_{i+1}$  where we let  $b_{n+1} = 1$  for convenience. Then, the event that  $F(Z_{(i)}) \leq h((i-1)/n) = b_i$  for all  $i \in \{1, \dots, n\}$  occurs with probability at least  $1 - \delta$ , where  $Z_{(1)} \leq \dots \leq Z_{(n)}$  are the order statistics. Under this event, for any  $z \in [Z_{(i-1)}, Z_{(i)})$ , where we let  $Z_{(0)} = \infty$  and  $Z_{(n+1)} = \infty$  for convenience,  $\hat{F}_n(z) = (i-1)/n$  and thus

$$F(z) \leq F(Z_{(i)}) \leq b_i = h(\hat{F}_n(z)).$$

On the other hand, if  $h : [0, 1] \rightarrow [0, 1]$  is a function such that  $h(\hat{F}_n(z))$  is a uniform upper confidence band of  $F$  for any CDF  $F$ , then (34) holds with  $b_i = h(i/n)$ .  $\square$

Applying Proposition 4 to the CDF of the scores  $\hat{s}$  computed by any one-class classification algorithm provides a uniform upper confidence bound for its FPR, namely  $\text{FPR}(t) := \mathbb{P}[\hat{s}(X_{2n+1}) \leq t]$ , as a function of the detection threshold  $t$ . In

words, this guarantees that reporting as outliers an observation with black-box score equal to  $z$  is likely (with probability at least  $1 - \delta$ ) to result in a FPR no greater than  $h(\hat{F}_n(z))$ , where  $\hat{F}_n(z)$  is the empirical CDF of the analogous scores computed on a calibration data set of size  $n$ . Figure A2 shows a practical example of this upper bound based on the empirical distribution of scores evaluated on 1000 calibration points, with  $\delta = 0.1$  and  $k = n/2$  (the exact details of this example are the same as those of the numerical experiments presented later in Section D.2). For instance, this plot informs us that reporting as outliers future samples with scores below  $-0.5$  is likely to result in an FPR below 0.025.

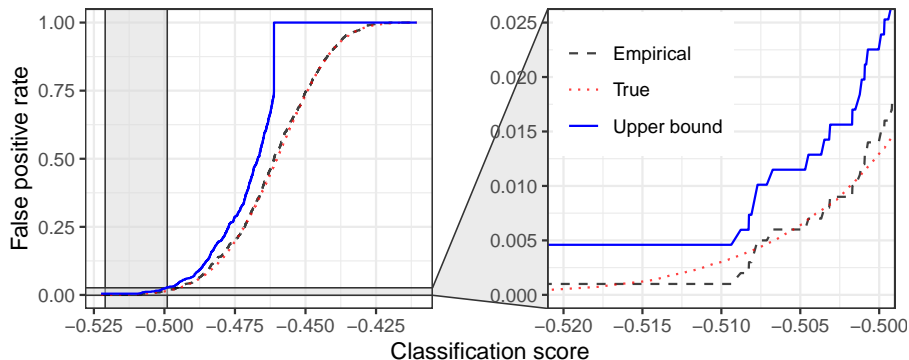


Figure A2: FPR calibration curve for an isolation forest one-class classifier on simulated data, as a function of the reporting threshold for the classification scores. The upper bound (solid blue) is guaranteed to lie above the true FPR curve (dotted red) with probability at least 90%. The dashed black curve corresponds to the empirical FPR. The panel on the right zooms in on small values (likely outliers).

Note that the construction of a uniform confidence band for an unknown CDF is a widely studied problem. For example, the DKWM inequality (Dvoretzky et al., 1956; Massart, 1990) implies the bound in (36) with  $h(z) = \min\{z + \sqrt{\log(2/\delta)/2n}, 1\}$ . However, the DKWM bound is tightest at  $z = 1/2$  and loose near 0, which would limit the power to detect outliers. Therefore, it is preferable for our purposes to have a function  $h(z)$  that is as close as possible to the identity for small values of  $z$ , as discussed earlier in Section B.5.

### B.7. Simultaneously-valid prediction sets

Lastly, CCV p-values can be easily re-purposed to strengthen the marginal guarantees generally obtainable for conformal predictions. In particular, for each  $\alpha \in (0, 1)$ , one can define a predictive set

$$\hat{C}^\alpha := \{x : \hat{u}^{(\text{ccv})}(x) > \alpha\}. \quad (37)$$

These sets are simultaneously valid for all  $\alpha$ , conditional on the calibration data. That is, they satisfy

$$\mathbb{P} \left[ \mathbb{P} [X_{2n+1} \in \hat{C}^\alpha \mid \mathcal{D}] \geq 1 - \alpha \text{ for all } \alpha \in (0, 1) \right] \geq 1 - \delta. \quad (38)$$

In words, if we use CCV p-values to construct prediction sets, the probability that a new observation falls within  $\hat{C}^\alpha$  is at least  $1 - \alpha$ , simultaneously for all  $\alpha \in (0, 1)$  with high probability. This is stronger than the usual conformal guarantee, as the latter holds marginally over  $\mathcal{D}$  and only for a single pre-specified  $\alpha$ .

## C. Numerical Comparisons of Different Adjustment Functions

In addition to the adjustment functions derived from the generalized Simes inequality and the DKWM inequality, we consider here another class of simultaneous bounds based on the so-called *boundary crossing probability* (Dempster, 1959; Durbin, 1973; Kotel’Nikova & Chmaladze, 1983; Siegmund, 1986)—the probability that  $F(z)$  ever crosses  $h(\hat{F}_n(z))$  for a fixed function  $h(\cdot)$ . This probability is generally difficult to compute analytically, but the special case of a linear  $h(\cdot)$  is an exception. Assuming that  $F$  is the CDF of  $\text{Unif}([0, 1])$ , let  $\hat{F}_n(z)$  is the empirical CDF of  $S_1, \dots, S_n \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0, 1])$ . Then, Dempster (1959) proved that

$$\mathbb{P} \left[ \hat{F}_n(z) \leq b + \frac{1-b}{1-a}z, \forall z \in (0, 1) \right] = 1 - \Delta_{\text{Dempster}}(a, b; n),$$

for any  $a, b \in (0, 1)$ , where

$$\Delta_{\text{Dempster}}(a, b; n) := a \sum_{j=0}^{\lfloor n(1-b) \rfloor} \frac{n!}{j!(n-j)!} \left( a + \frac{1-a}{1-b} \frac{j}{n} \right)^{j-1} \left( 1 - a - \frac{1-a}{1-b} \frac{j}{n} \right)^{n-j}. \quad (39)$$

If we replace  $S_i$  with  $1 - S_i$ , then  $\hat{F}_n(z)$  becomes  $1 - \hat{F}_n(1 - z)$ . Further, replacing  $z$  by  $1 - z$  leads to

$$\mathbb{P} \left[ z \leq \frac{1-a}{1-b} \hat{F}_n(z) + a, \forall z \in (0, 1) \right] = 1 - \Delta_{\text{Dempster}}(a, b; n). \quad (40)$$

For any pair  $(a, b)$  with  $\Delta_{\text{Dempster}}(a, b; n) = \delta$ , we obtain a function  $h(z) = a + (1-a)z/(1-b)$  satisfying (36), which yields the following sequence satisfying (34):

$$b_i = a + \frac{1-a}{1-b} \frac{i}{n}.$$

Given any  $a$ , it is easy to compute the corresponding  $b$  such that  $\Delta_{\text{Dempster}}(a, b; n) = \delta$  via a binary search.

Note that this leads to adjusted p-values that cannot be lower than  $b_1 = a + (1-a)/(1-b)n$ . To ensure a fair comparison with the method based on the generalized Simes inequality, we choose  $a$  via another binary search such that the resulting  $b_1$  matches that given by the Simes inequality for a particular value of  $k$ . If there exists no value of  $a$  yielding the same  $b_1$  as the Simes method, we set  $a$  as to minimize  $b_1$ . Figure A3 compares the adjustment functions yielded by the generalized Simes inequality, the DKWM inequality, and the Dempster exact linear-boundary crossing probability with  $k \in \{n/4, n/2\}$  and  $n \in \{300, 1000, 3000, 10000\}$  for small marginal p-values within  $[0, 0.05]$ . It is clear that the Simes adjustment function is the best option in most scenarios, except when  $n = 10000$  and  $\hat{u}^{(\text{marg})}(X) > 0.03$ , in which case the DKWM bound is tighter. Nonetheless, for the purpose of multiple testing, we would rarely expect p-values above 0.03 to be significant.

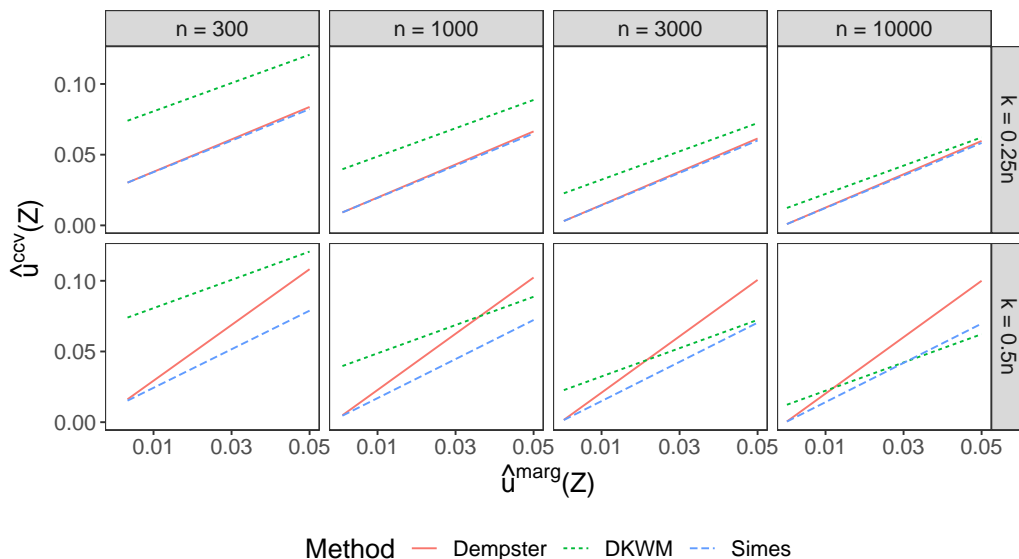


Figure A3: Comparison of different adjustment functions.

## D. Numerical Experiments

### D.1. Setup

The following experiments are designed to simulate a world in which our methods are independently applied by  $J$  practitioners. Each practitioner  $j \in [J]$  has an independent data set  $\mathcal{D}_j$  (to train and calibrate the method), and  $L$  test sets  $\mathcal{D}_{j,l}^{\text{test}}$  (to compute p-values and evaluate performance), each corresponding to different possible future scenarios  $l \in [L]$ . The data sets contain  $2n$  observations each ( $|\mathcal{D}_j| = 2n$ ), and the test sets contain  $n_{\text{test}}$  observations each ( $|\mathcal{D}_{j,l}^{\text{test}}| = n_{\text{test}}$ ).

Imagine that, from the practitioner’s present point of view, the data set  $\mathcal{D}_j$  is fixed but the test set is random, so that  $\mathcal{D}_{j,l}^{\text{test}}$  represents the test set for practitioner  $j$  under future scenario  $l$ . Then, as discussed in Section 2, practitioner  $j$  is most interested in the FDR (or other measures of type-I errors, alternatively) conditional on  $\mathcal{D}_j$ , i.e., in the random variable

$$\text{cFDR}(\mathcal{D}_j) := \mathbb{E} [\text{FDP}(\mathcal{D}^{\text{test}}; \mathcal{D}_j) \mid \mathcal{D}_j],$$

where  $\text{FDP}(\mathcal{D}^{\text{test}}; \mathcal{D}_j)$  is the proportion of inliers among the test points reported as outliers, based on the procedure calibrated on  $\mathcal{D}_j$ . This motivates the definition of the following performance measures. For any  $j \in [J]$ , we compute

$$\widehat{\text{cFDR}}(\mathcal{D}_j) := \frac{1}{L} \sum_{l=1}^L \text{FDP}(\mathcal{D}_{j,l}^{\text{test}}; \mathcal{D}_j), \quad \widehat{\text{cPower}}(\mathcal{D}_j) := \frac{1}{L} \sum_{l=1}^L \text{Power}(\mathcal{D}_{j,l}^{\text{test}}; \mathcal{D}_j), \quad (41)$$

where  $\text{Power}(\mathcal{D}_{j,l}^{\text{test}}; \mathcal{D}_j)$  is the proportion of outliers in  $\mathcal{D}_{j,l}^{\text{test}}$  correctly identified as such by practitioner  $j$ .

## D.2. Outlier detection on simulated data

### D.2.1. DATA DESCRIPTION

We begin to investigate the empirical performance of different methods for calibrating conformal p-values on synthetic data. The data are generated by sampling each data point  $X_i \in \mathbb{R}^{50}$  from a multivariate Gaussian mixture model  $P_X^a$ , such that  $X_i = \sqrt{a} V_i + W_i$ , for some constant  $a \geq 1$  and appropriate random vectors  $V_i, W_i \in \mathbb{R}^{50}$ . Here,  $V_i$  has independent standard Gaussian components, and each coordinate of  $W_i$  is independent and uniformly distributed on a discrete set  $\mathcal{W} \subseteq \mathbb{R}^{50}$  with cardinality  $|\mathcal{W}| = 50$ . The vectors in  $\mathcal{W}$  are sampled independently from the uniform distribution on  $[-3, 3]^{50}$ , before the beginning of our experiments, and then held constant thereafter. (Therefore, each coordinate of  $W_i$  is uniformly distributed on  $[-3, 3]$ , but it is not the case that the different  $W_i$ ’s are independent and identically distributed on  $[-3, 3]^{50}$ ; instead, the fixed set  $\mathcal{W}$  makes this a mixture model.)

The data sets  $\mathcal{D}_j$  are sampled from  $P_X^a$  with  $a = 1$  and  $n = 1000$ . The total  $2n$  observations in each  $\mathcal{D}_j$  are further divided into  $n_{\text{train}} = 1000$  observations used to fit a one-class SVM classifier scoring function  $\hat{s}$  (implemented in the Python package `scikit-learn` (Pedregosa et al., 2011)), and  $n_{\text{cal}} = 1000$  observations used to calibrate the conformal p-values, leading to a valid p-value  $\hat{u}(X_{n+1}) \in [0, 1]$  for any new data point  $X_{n+1}$ . The total number of data sets is  $J = 100$ , each of which is associated with  $L = 100$  test sets. A random subset of the observations in each test set  $\mathcal{D}_{j,l}^{\text{test}}$  is sampled from  $P_X^a$  with  $a = 1$ , while the others are outliers, in the sense that they are sampled from  $P_X^a$  with  $a > 1$ , as specified below.

### D.2.2. INDIVIDUAL OUTLIER DETECTION

First, we focus on a data generating model under which 90% of the  $n_{\text{test}} = 1000$  observations in each  $\mathcal{D}_{j,l}^{\text{test}}$  are sampled from  $P_X^a$  with  $a = 1$ , and we seek to identify the remaining 10% of outliers. For this purpose, we calibrate a conformal p-value for all observations in  $\mathcal{D}_{j,l}^{\text{test}}$ , and then we apply the BH procedure at some nominal FDR level  $\alpha$  to account for the multiple comparisons, with and without Storey’s correction based on the estimated null proportion. In the following, we apply our conditional calibration method with the parameters  $\delta = 0.1$  and  $k = n_{\text{cal}}/2$  (see below for comments about the choice of  $k$ ).

Figure A4 shows the distribution of  $\widehat{\text{cFDR}}(\mathcal{D}_s)$  and  $\widehat{\text{cPower}}(\mathcal{D}_s)$ , corresponding to  $\alpha = 0.1$ , for different values of the signal strength  $a$  (recall that here  $a = 1$  corresponds to no signal), when the BH procedure is utilized to account for the multiple comparisons. The results confirm the calibration-conditional p-values control the conditional FDR for at least 90% of practitioners, while the marginal p-values do not. In fact, marginal p-values only control the conditional FDR if the number of samples in the calibration data set is very large; see Figure A5. Furthermore, we note that both methods control the marginal FDR, as also predicted by our theoretical results. Figure A6 presents the results obtained by applying Storey’s correction to the BH procedure, while Figure A7 summarizes additional experiments in which our conditional calibration method is applied with  $\delta = 0.25$ . Finally, Figure A8 visualizes the effect of different values of the Simes parameter  $k$  on our calibration-conditional p-value, showing that  $k = n_{\text{cal}}/2$  works relatively well, although the performance does not appear to be extremely sensitive to this choice.

### D.2.3. BATCH OUTLIER DETECTION

We now consider the global testing problem of detecting whether a batch of new observations contains any outliers. For this purpose, we follow the same approach as before, with the only difference that the  $n_{\text{test}} = 1000$  observations in each test

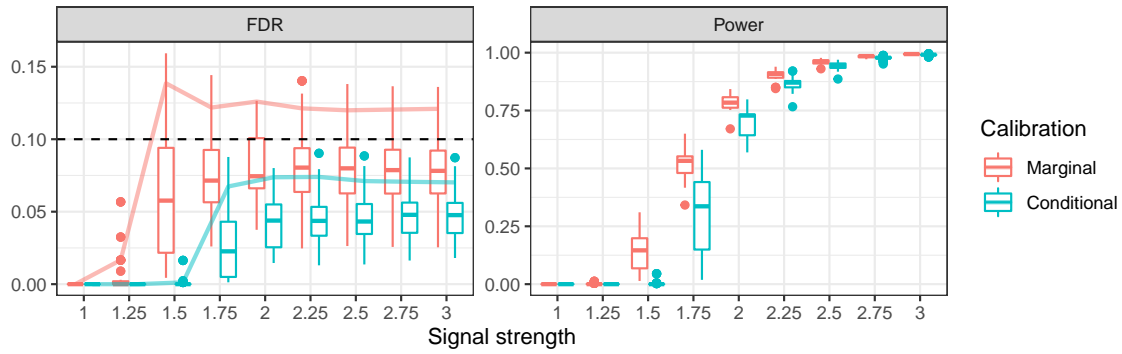


Figure A4: Performance of different methods for calibrating conformal p-values in a simulated outlier detection problem, as a function of the signal strength. The box plots visualize the distribution of FDR and power, as defined in (41), conditional on 100 independent data sets. The solid curves indicate the 90-th quantile of the conditional FDR distribution. The nominal FDR 0.1, and the conditional method is applied with  $\delta = 0.1$ .

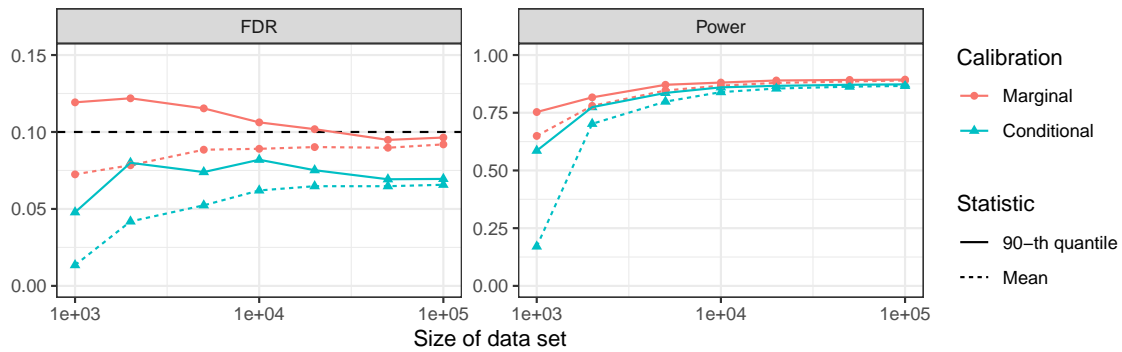


Figure A5: FDR and power in a simulated outlier detection problem as a function of the number of samples in the data set (half of which are utilized for calibration). Other details are as in Figure A4.

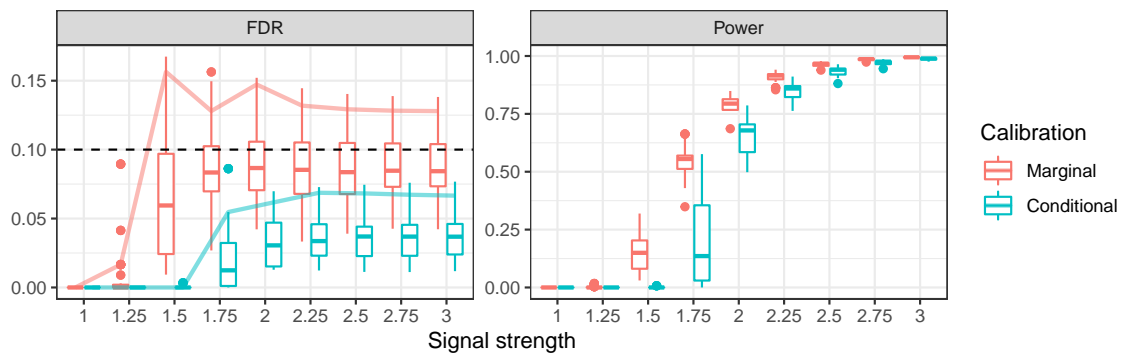


Figure A6: FDR and power in a simulated outlier detection problem, using the BH procedure with Storey's correction. Other details are as in Figure A4.

set are now sub-divided into 100 batches of size 10. The 10 calibrated p-values in each batch are combined with Fisher's method to test the batch-specific global null. Then, the BH procedure with Storey's correction is applied to control the FDR over all batches. This simulation is designed such that 90% of the batches contain no outliers (i.e., all samples are drawn from  $P_X^a$  with  $a = 1$ ), while 50% of the samples in the remaining batches are outliers (i.e., they are drawn from  $P_X^a$  with  $a = 2$ ). Of course, batched testing is less informative than the precise identification of outliers discussed in the previous section, but the advantage now is that we can achieve higher power. Figure A9 shows that, even though this problem is relatively easy (the power is almost equal to 1), the use of marginal p-values may still lead to a conditional FDR that is

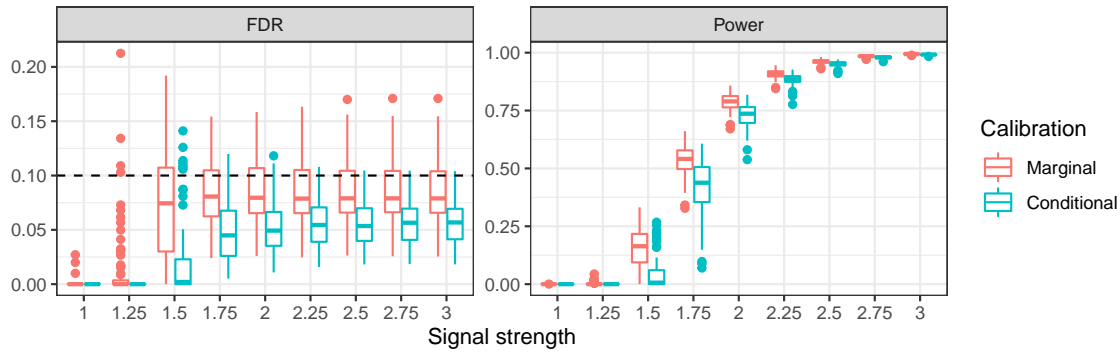


Figure A7: FDR and power in a simulated outlier detection problem, as a function of the signal strength. The conditional calibration method is applied with  $\delta = 0.25$  instead of  $\delta = 0.05$ . Other details are as in Figure A4.

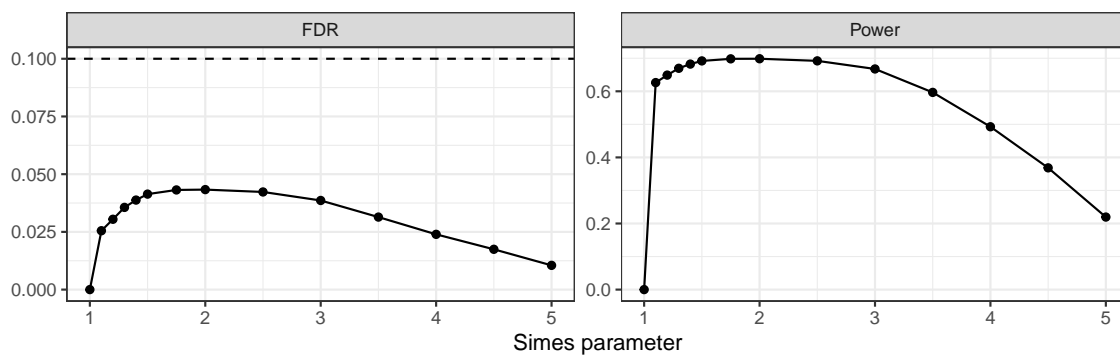


Figure A8: Performance of simultaneously calibrated conformal p-values as a function of the Simes parameter  $n/k$ . The signal strength is equal to 2. Other details are as in Figure A4.

noticeably higher than expected for many researchers. By contrast, simultaneous calibration appears to be conservative for all of them, without much sacrifice in power.

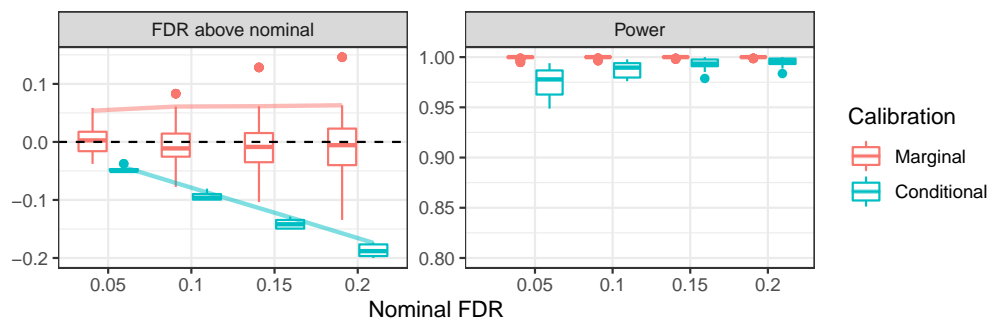


Figure A9: Performance of different methods for calibrating conformal p-values in a simulated outlier batch detection problem, as a function of the nominal FDR level. The excess FDR is defined as the difference between the empirical FDR and the nominal FDR. Note that both methods achieve power close to one in this example. Other details are as in Figure A4.

Finally, we study the effect of the batch size on the performance of different calibration methods under the global null hypothesis (i.e., when there are no outliers in the test set). As before, the p-values in each batch are combined with Fisher’s method and the global null is rejected if the resulting p-value is smaller than 0.1. As before, the experiment is repeated for 100 independent data sets and 1000 test sets. Figure A10 shows that marginal p-values do not lead to valid inferences, especially if the batch size is large. By contrast, the calibration-conditional method always remains valid.

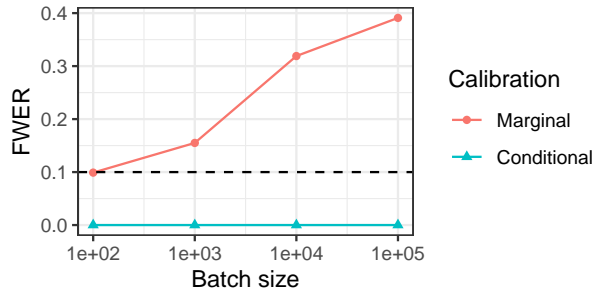


Figure A10: Family-wise error rate (FWER) in a simulated outlier batch detection problem under the global null hypothesis, using different calibration methods for the conformal p-values. The results are shown as a function of the batch size. The global null is rejected if the Fisher’s combined p-value is below 0.1, which means the nominal FWER is 10% (horizontal dashed line).

### D.3. Outlier detection on real data

#### D.3.1. DATA DESCRIPTION

Table A1: Summary of the benchmark data sets for outlier detection utilized in our applications.

	ALOI <sup>2</sup>	Cover <sup>3</sup>	Credit card <sup>4</sup>	KDDCup99 <sup>5</sup>	Mammography <sup>6</sup>	Digits <sup>7</sup>	Shuttle <sup>8</sup>
Features $d$	27	10	30	40	6	16	9
Inliers $n_{\text{inliers}}$	283301	286048	284315	47913	10923	6714	45586
Outliers $n_{\text{outliers}}$	1508	2747	492	200	260	156	3511

We turn to study the performance of the calibration schemes from Section D.2 on several benchmark data sets for outlier detection, summarized in Table A1. As before, the Simes simultaneous calibration is applied with  $\delta = 0.1$  and  $k = n_{\text{cal}}/2$ . We utilize an isolation forest (Liu et al., 2008) machine-learning algorithms  $\hat{s}$  as the base method for detecting anomalies, available in the Python `sklearn` package. We rely on the default hyper-parameters, except for the ‘contamination’ parameter which we set equal to 0.1. Additional experiments based on one-class SVM and Local Outlier Factor (LOF) algorithms are presented in Tables A5–A6.

#### D.3.2. INDIVIDUAL OUTLIER DETECTION

Here, we follow the experimental setup of Section 4. Figure A11 compares the performance of marginal and simultaneously calibrated p-values on the credit card data set with the BH procedure with Storey’s correction. The results are consistent with Figure 3.

Consistent conclusion can be drawn from Table A2, which compares the two calibration procedures on all benchmark data sets at the nominal FDR level of 0.2. Additional results corresponding to different outlier detection algorithms (one-class SVM and LOF) can be found in Table A4. In all cases, we adopt the `sklearn` default parameters. Finally, Table A5 summarizes the performance of different calibration and detection methods across all data sets when the BH procedure is applied without Storey’s correction.

<sup>2</sup>The dataset is available at <https://www.dbs.ifi.lmu.de/research/outlier-evaluation/DAMI/literature/ALOI>; see Campos et al. (2016)

<sup>3</sup>The dataset is available at <http://odds.cs.stonybrook.edu/forestcovercovertype-dataset>

<sup>4</sup>The dataset is available at <https://www.kaggle.com/mlg-ulb/creditcardfraud>

<sup>5</sup>The dataset is available at <https://www.kaggle.com/mlg-ulb/creditcardfraud>; see Campos et al. (2016)

<sup>6</sup>The dataset is available at <http://odds.cs.stonybrook.edu/mammography-dataset/>

<sup>7</sup>The dataset is available at <http://odds.cs.stonybrook.edu/pendigits-dataset>

<sup>8</sup>The dataset is available at <http://odds.cs.stonybrook.edu/shuttle-dataset>

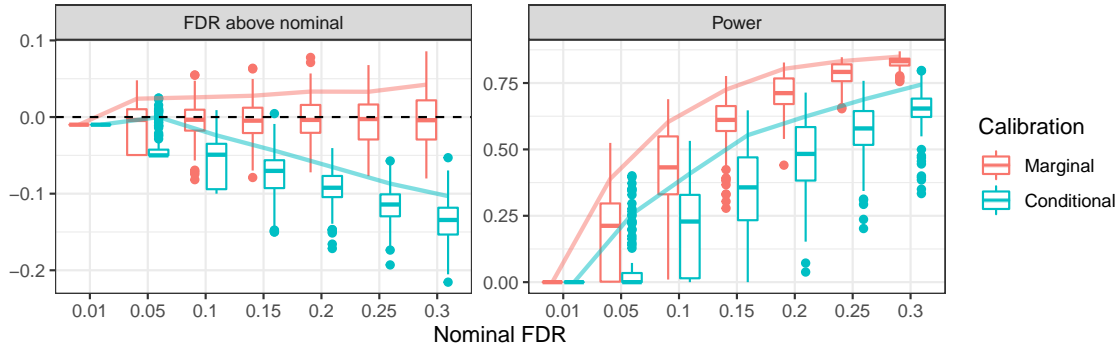


Figure A11: Outlier detection performance on credit card fraud data. Conformal p-values based on an isolation forest model are calibrated using different methods. The BH procedure with Storey’s correction is then applied to control the FDR over the set of test points. Other details are as in Figure A6.

Table A2: Outlier detection performance on different data sets, using alternative methods for calibrating conformal p-values. The FDR and power diagnostics are defined conditional on the training and calibration data, as defined in Section D.1. The nominal marginal FDR level is 0.2. Empirical FDR values larger than the nominal level are colored in orange; values at least one standard deviation above it are colored in red.

Dataset	FDR				Power			
	Mean		90th percentile		Mean		90-th quantile	
	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
ALOI	0.025	0	0.02	0	0	0	0	0
Cover	0.08	0.013	0.277	0.049	0.008	0.002	0.03	0.004
Credit card	0.197	0.106	0.233	0.135	0.712	0.469	0.803	0.624
KDDCup99	0.196	0.105	0.234	0.135	0.755	0.62	0.825	0.713
Mammography	0.18	0.031	0.282	0.112	0.167	0.036	0.342	0.155
Digits	0.177	0.029	0.27	0.116	0.347	0.056	0.603	0.213
Shuttle	0.196	0.107	0.234	0.138	0.981	0.976	0.984	0.981

### D.3.3. BATCH OUTLIER DETECTION

We now focus on global testing for outlier batch detection, similarly to Section D.2.3. The available data are divided into training, calibration, and test sets according to the same scheme as in Section D.3.2; the only difference is that the size of the test sets is now equal to 1000, so as to follow as closely as possible the same experimental protocol as in Section D.2.3.

Figure A12 compares the performance of the different calibration methods as a function of the nominal FDR level. The p-values in each batch are combined with Fisher’s method, and then the BH procedure is applied with Storey’s correction. Again, we observe that simultaneous calibration is required to ensure the conditional FDR is controlled in at least 90% of the applications, although it involves some power loss. Both calibration methods control the marginal FDR.

Table A3 summarizes the performance of the two alternative calibration methods on all data sets. Here, the nominal FDR level is 0.1 and the BH procedure is applied with the Storey correction. Again, the results show that the Simes method controls the conditional FDR 90% of the time, although at some cost in power, while the marginal calibration method does not. See Table A6 for additional results that, in addition to the isolation forest, include also the one-class SVM and LOF algorithms for outlier detection. Finally, Table A7 summarizes performance of the different methods on all data sets when the BH procedure is applied without the Storey correction.

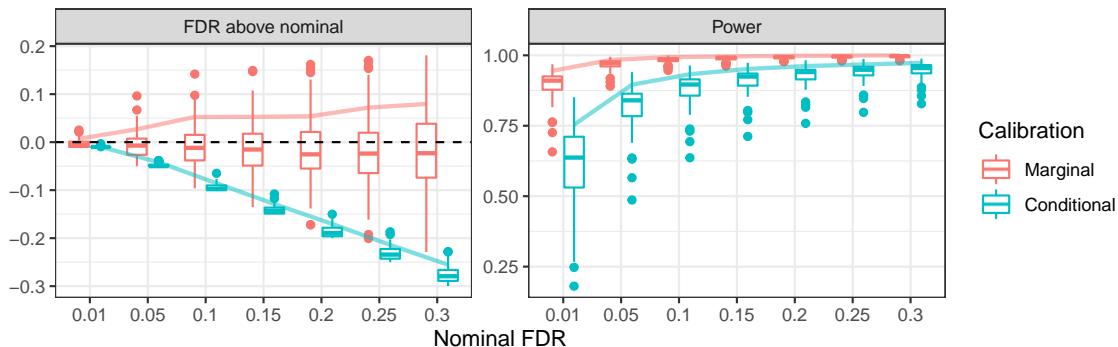


Figure A12: Outlier batch detection performance on credit card fraud data. Conformal p-values are computed based on an isolation forest model and calibrated using different methods. Other details are as in Figure A9.

Table A3: Outlier batch detection performance on different data sets, using alternative methods for calibrating conformal p-values. The nominal FDR level is 0.1. Other details are as in Table A2.

Data set	FDR				Power			
	Mean		90-th quantile		Mean		90-th quantile	
	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
ALOI	0.072	0.003	0.178	0	0.002	0	0.004	0
Cover	0.092	0.006	0.183	0.01	0.18	0.017	0.359	0.034
Credit card	0.092	0.005	0.153	0.014	0.983	0.885	0.993	0.933
KDDCup99	0.088	0.005	0.129	0.013	0.999	0.979	1	0.994
Mammography	0.072	0.004	0.126	0.016	0.61	0.21	0.765	0.361
Digits	0.09	0.005	0.148	0.014	0.97	0.626	0.999	0.836
Shuttle	0.087	0.006	0.137	0.013	1	1	1	1

Table A4: Outlier detection performance on real data, using different data sets, machine learning models, and nominal FDR levels. Other details are as in Table A2.

Model	Nominal	FDR				Power			
		Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
<b>ALOI</b>									
IForest	0.05	0	0	0	0	0	0	0	0
	0.10	0.001	0	0	0	0	0	0	0
	0.20	0.025	0	0.02	0	0	0	0	0
LOF	0.05	0	0	0	0	0	0	0	0
	0.10	0.005	0	0	0	0	0	0	0
	0.20	0.056	0.003	0.176	0	0.002	0	0.003	0
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0.005	0	0.012	0	0	0	0.001	0
	0.20	0.069	0.001	0.228	0	0.003	0	0.01	0
<b>Cover</b>									
IForest	0.05	0	0	0	0	0	0	0	0
	0.10	0.011	0	0.032	0	0.002	0	0.003	0
	0.20	0.08	0.013	0.277	0.049	0.008	0.002	0.03	0.004
LOF	0.05	0.05	0.026	0.069	0.041	0.949	0.91	0.968	0.943
	0.10	0.1	0.056	0.126	0.075	0.973	0.955	0.98	0.969
	0.20	0.198	0.111	0.23	0.138	0.987	0.976	0.991	0.982

(Continued on Next Page...)

**Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning**

Table A4: Outlier detection performance on real data, using different data sets, machine learning models, and nominal FDR levels. Other details are as in Table A2. (continued)

Model	FDR					Power			
	Nominal	Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Credit card</b>									
IForest	0.05	0.037	0.012	0.074	0.05	0.185	0.062	0.389	0.256
	0.10	0.095	0.042	0.126	0.076	0.426	0.207	0.603	0.409
	0.20	0.197	0.106	0.233	0.135	0.712	0.469	0.803	0.624
LOF	0.05	0.001	0	0	0	0	0	0	0
	0.10	0.018	0.001	0.022	0	0.001	0	0	0
	0.20	0.087	0.021	0.278	0.031	0.007	0.002	0.022	0.001
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>KDDCup99</b>									
IForest	0.05	0.04	0.013	0.074	0.036	0.378	0.208	0.515	0.44
	0.10	0.095	0.043	0.125	0.077	0.594	0.397	0.703	0.524
	0.20	0.196	0.105	0.234	0.135	0.755	0.62	0.825	0.713
LOF	0.05	0.001	0	0	0	0	0	0	0
	0.10	0.038	0	0.159	0	0.012	0	0.055	0
	0.20	0.141	0.037	0.27	0.158	0.039	0.011	0.07	0.055
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Mammography</b>									
IForest	0.05	0.011	0	0.051	0	0.012	0	0.035	0
	0.10	0.067	0	0.163	0	0.061	0	0.184	0
	0.20	0.18	0.031	0.282	0.112	0.167	0.036	0.342	0.155
LOF	0.05	0.003	0	0	0	0.001	0	0	0
	0.10	0.032	0	0.169	0	0.023	0	0.112	0
	0.20	0.167	0.018	0.272	0.061	0.195	0.017	0.316	0.036
SVM	0.05	0.011	0	0.031	0	0.004	0	0.015	0
	0.10	0.075	0	0.196	0	0.042	0	0.086	0
	0.20	0.186	0.003	0.256	0.002	0.169	0.001	0.267	0.001
<b>Digits</b>									
IForest	0.05	0.006	0	0.006	0	0.007	0	0.005	0
	0.10	0.049	0.002	0.159	0	0.073	0.003	0.245	0
	0.20	0.177	0.029	0.27	0.116	0.347	0.056	0.603	0.213
LOF	0.05	0.01	0	0.045	0	0.038	0.005	0.149	0
	0.10	0.06	0.003	0.142	0.001	0.282	0.017	0.775	0.005
	0.20	0.191	0.059	0.245	0.144	0.821	0.297	0.984	0.795
SVM	0.05	0.001	0	0	0	0	0	0	0
	0.10	0.045	0	0.152	0	0.018	0	0.048	0
	0.20	0.167	0.005	0.253	0.007	0.24	0.004	0.475	0.002
<b>Shuttle</b>									
IForest	0.05	0.048	0.023	0.068	0.035	0.946	0.872	0.977	0.97
	0.10	0.097	0.052	0.13	0.071	0.975	0.953	0.981	0.977
	0.20	0.196	0.107	0.234	0.138	0.981	0.976	0.984	0.981

(Continued on Next Page...)

**Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning**

---

Table A4: Outlier detection performance on real data, using different data sets, machine learning models, and nominal FDR levels. Other details are as in Table A2. (*continued*)

Model	FDR					Power			
	Nominal	Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
LOF	0.05	0.051	0.026	0.07	0.044	0.991	0.857	0.998	0.988
	0.10	0.099	0.055	0.125	0.075	0.999	0.992	1	0.999
	0.20	0.197	0.109	0.236	0.137	1	0.999	1	1
SVM	0.05	0.047	0.007	0.067	0.034	0.904	0.152	0.998	0.91
	0.10	0.101	0.052	0.12	0.072	0.999	0.953	1	0.998
	0.20	0.202	0.112	0.232	0.13	1	1	1	1

**Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning**

Table A5: Outlier detection performance on real data, using different data sets, machine learning models, and nominal FDR levels. Other details are as in Table A4.

Model	Nominal	FDR				Power			
		Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
<b>ALOI</b>									
IForest	0.05	0	0	0	0	0	0	0	0
	0.10	0.001	0	0	0	0	0	0	0
	0.20	0.027	0.003	0.03	0	0	0	0	0
LOF	0.05	0	0	0	0	0	0	0	0
	0.10	0.005	0	0	0	0	0	0	0
	0.20	0.054	0.007	0.175	0	0.002	0	0.003	0
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0.006	0	0.015	0	0	0	0.001	0
	0.20	0.083	0.013	0.263	0.031	0.004	0.001	0.012	0.002
<b>Cover</b>									
IForest	0.05	0	0	0	0	0	0	0	0
	0.10	0.008	0	0.031	0	0.001	0	0.002	0
	0.20	0.073	0.021	0.244	0.085	0.007	0.003	0.025	0.009
LOF	0.05	0.045	0.031	0.063	0.049	0.943	0.922	0.965	0.955
	0.10	0.09	0.065	0.115	0.085	0.971	0.961	0.978	0.971
	0.20	0.179	0.129	0.209	0.156	0.985	0.979	0.989	0.984
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Credit card</b>									
IForest	0.05	0.032	0.016	0.068	0.057	0.164	0.085	0.344	0.297
	0.10	0.084	0.051	0.114	0.085	0.384	0.248	0.58	0.452
	0.20	0.178	0.126	0.211	0.154	0.678	0.539	0.775	0.667
LOF	0.05	0.001	0	0	0	0	0	0	0
	0.10	0.017	0.003	0.013	0	0.001	0	0	0
	0.20	0.085	0.03	0.274	0.052	0.006	0.002	0.022	0.003
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>KDDCup99</b>									
IForest	0.05	0.034	0.018	0.067	0.047	0.355	0.238	0.501	0.455
	0.10	0.086	0.055	0.116	0.088	0.561	0.453	0.687	0.617
	0.20	0.176	0.123	0.212	0.156	0.738	0.666	0.813	0.735
LOF	0.05	0.001	0	0	0	0	0	0	0
	0.10	0.037	0.004	0.159	0.004	0.012	0.001	0.055	0
	0.20	0.138	0.051	0.264	0.177	0.038	0.015	0.069	0.055
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Mammography</b>									
IForest	0.05	0.008	0	0.024	0	0.009	0	0.014	0
	0.10	0.056	0.002	0.155	0.001	0.05	0.003	0.16	0.002
	0.20	0.161	0.053	0.252	0.165	0.147	0.059	0.315	0.217
LOF	0.05	0.001	0	0	0	0	0	0	0
	0.10	0.029	0	0.154	0	0.019	0	0.057	0
	0.20	0.141	0.035	0.259	0.178	0.161	0.036	0.307	0.164

*(Continued on Next Page...)*

**Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning**

Table A5: Outlier detection performance on real data, using different data sets, machine learning models, and nominal FDR levels. Other details are as in Table A4. (continued)

Model	FDR					Power			
	Nominal	Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
SVM	0.05	0.006	0	0.009	0	0.002	0	0.004	0
	0.10	0.065	0	0.185	0	0.036	0	0.068	0
	0.20	0.17	0.053	0.245	0.173	0.146	0.033	0.244	0.091
<b>Digits</b>									
IForest	0.05	0.005	0	0.002	0	0.006	0	0.001	0
	0.10	0.042	0.003	0.143	0	0.057	0.005	0.181	0
	0.20	0.159	0.05	0.257	0.172	0.296	0.093	0.541	0.368
LOF	0.05	0.009	0	0.029	0	0.029	0.005	0.111	0
	0.10	0.046	0.007	0.129	0.017	0.209	0.034	0.687	0.066
	0.20	0.173	0.086	0.226	0.162	0.77	0.429	0.975	0.861
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0.042	0	0.149	0	0.015	0	0.047	0
	0.20	0.145	0.02	0.242	0.052	0.168	0.018	0.417	0.048
<b>Shuttle</b>									
IForest	0.05	0.043	0.028	0.061	0.043	0.939	0.891	0.976	0.973
	0.10	0.088	0.061	0.117	0.087	0.973	0.963	0.98	0.978
	0.20	0.176	0.124	0.209	0.158	0.981	0.978	0.983	0.982
LOF	0.05	0.045	0.032	0.065	0.049	0.988	0.934	0.998	0.992
	0.10	0.089	0.065	0.111	0.09	0.998	0.995	1	0.999
	0.20	0.178	0.127	0.211	0.156	1	0.999	1	1
SVM	0.05	0.039	0.015	0.061	0.043	0.827	0.358	0.997	0.993
	0.10	0.091	0.063	0.111	0.081	0.999	0.997	1	0.999
	0.20	0.183	0.13	0.216	0.156	1	1	1	1

Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning

Table A6: Outlier batch detection performance on real data, using Storey’s correction to control the FDR. Other details are as in Table A3.

Model	Nominal	FDR				Power			
		Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
<b>ALOI</b>									
IForest	0.05	0.035	0.003	0.1	0	0.001	0	0.002	0
	0.10	0.071	0.003	0.164	0	0.001	0	0.004	0
	0.20	0.141	0.014	0.278	0.058	0.004	0	0.008	0.001
LOF	0.05	0.034	0	0.091	0	0.023	0.003	0.039	0.007
	0.10	0.082	0.003	0.161	0	0.047	0.005	0.071	0.012
	0.20	0.185	0.009	0.284	0.032	0.106	0.011	0.153	0.02
SVM	0.05	0.032	0.004	0.097	0	0.003	0	0.007	0.002
	0.10	0.064	0.006	0.173	0	0.006	0.001	0.01	0.003
	0.20	0.152	0.013	0.289	0.08	0.012	0.002	0.022	0.005
<b>Cover</b>									
IForest	0.05	0.036	0.006	0.088	0.003	0.086	0.013	0.183	0.029
	0.10	0.08	0.008	0.158	0.035	0.163	0.025	0.328	0.048
	0.20	0.173	0.016	0.25	0.066	0.301	0.049	0.536	0.106
LOF	0.05	0.035	0.004	0.059	0.012	1	1	1	1
	0.10	0.074	0.01	0.115	0.022	1	1	1	1
	0.20	0.163	0.021	0.225	0.039	1	1	1	1
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Credit card</b>									
IForest	0.05	0.039	0.004	0.066	0.011	0.967	0.863	0.983	0.917
	0.10	0.083	0.009	0.123	0.021	0.982	0.916	0.993	0.953
	0.20	0.168	0.023	0.238	0.045	0.992	0.951	0.997	0.973
LOF	0.05	0.034	0.005	0.109	0	0.03	0.005	0.047	0.01
	0.10	0.071	0.009	0.153	0.005	0.055	0.009	0.09	0.017
	0.20	0.158	0.017	0.248	0.088	0.109	0.016	0.155	0.029
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>KDDCup99</b>									
IForest	0.05	0.035	0.004	0.062	0.01	0.998	0.971	1	0.989
	0.10	0.077	0.009	0.11	0.021	0.999	0.988	1	0.997
	0.20	0.167	0.019	0.215	0.037	1	0.996	1	1
LOF	0.05	0.032	0.003	0.09	0	0.061	0.013	0.087	0.024
	0.10	0.072	0.005	0.14	0.011	0.103	0.022	0.144	0.036
	0.20	0.16	0.014	0.258	0.069	0.178	0.036	0.236	0.052
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Mammography</b>									
IForest	0.05	0.036	0.003	0.072	0.009	0.489	0.177	0.667	0.311
	0.10	0.073	0.007	0.116	0.022	0.615	0.266	0.767	0.433
	0.20	0.143	0.018	0.22	0.045	0.743	0.384	0.856	0.564
LOF	0.05	0.031	0.002	0.065	0.006	0.448	0.14	0.571	0.234
	0.10	0.066	0.005	0.118	0.017	0.58	0.228	0.699	0.35
	0.20	0.135	0.015	0.209	0.039	0.713	0.352	0.8	0.485

(Continued on Next Page...)

**Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning**

Table A6: Outlier batch detection performance on real data, using Storey’s correction to control the FDR. Other details are as in Table A3. (continued)

Model	FDR					Power			
	Nominal	Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
SVM	0.05	0.011	0.001	0.031	0	0.377	0.095	0.458	0.144
	0.10	0.024	0.002	0.054	0.003	0.492	0.157	0.568	0.221
	0.20	0.053	0.005	0.097	0.018	0.613	0.248	0.688	0.324
<b>Digits</b>									
IForest	0.05	0.04	0.003	0.074	0.01	0.924	0.56	0.988	0.783
	0.10	0.079	0.008	0.127	0.019	0.968	0.728	0.997	0.903
	0.20	0.161	0.02	0.234	0.042	0.99	0.86	1	0.962
LOF	0.05	0.042	0.004	0.08	0.012	0.999	0.945	1	0.999
	0.10	0.087	0.009	0.14	0.021	1	0.984	1	1
	0.20	0.178	0.022	0.258	0.045	1	0.997	1	1
SVM	0.05	0.041	0.004	0.079	0.017	0.803	0.367	0.88	0.511
	0.10	0.086	0.009	0.145	0.028	0.889	0.528	0.942	0.677
	0.20	0.179	0.019	0.265	0.039	0.95	0.691	0.978	0.808
<b>Shuttle</b>									
IForest	0.05	0.036	0.004	0.062	0.01	1	1	1	1
	0.10	0.077	0.009	0.116	0.019	1	1	1	1
	0.20	0.163	0.021	0.222	0.036	1	1	1	1
LOF	0.05	0.041	0.003	0.066	0.012	1	1	1	1
	0.10	0.085	0.009	0.128	0.023	1	1	1	1
	0.20	0.172	0.023	0.236	0.043	1	1	1	1
SVM	0.05	0.031	0.003	0.055	0.009	1	1	1	1
	0.10	0.067	0.008	0.107	0.018	1	1	1	1
	0.20	0.151	0.018	0.21	0.034	1	1	1	1

**Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning**

Table A7: Outlier batch detection performance on real data, using different data sets, machine learning models, and nominal FDR levels. Other details are as in Table A6.

Model	FDR					Power			
	Nominal	Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
<b>ALOI</b>									
IForest	0.05	0.037	0.004	0.1	0	0.001	0	0.002	0
	0.10	0.09	0.01	0.188	0.028	0.001	0	0.004	0
	0.20	0.173	0.017	0.316	0.09	0.003	0	0.007	0.001
LOF	0.05	0.034	0.002	0.098	0	0.021	0.003	0.038	0.007
	0.10	0.07	0.006	0.185	0.004	0.042	0.005	0.071	0.013
	0.20	0.153	0.014	0.294	0.082	0.095	0.01	0.143	0.019
SVM	0.05	0.036	0.002	0.092	0	0.003	0	0.006	0.001
	0.10	0.068	0.007	0.175	0.008	0.006	0.001	0.012	0.002
	0.20	0.156	0.012	0.283	0.074	0.013	0.001	0.023	0.004
<b>Cover</b>									
IForest	0.05	0.041	0.002	0.107	0	0.09	0.013	0.159	0.029
	0.10	0.074	0.01	0.15	0.04	0.162	0.027	0.273	0.06
	0.20	0.151	0.022	0.242	0.078	0.292	0.051	0.485	0.093
LOF	0.05	0.04	0.004	0.07	0.011	1	1	1	1
	0.10	0.083	0.009	0.121	0.023	1	1	1	1
	0.20	0.173	0.023	0.234	0.045	1	1	1	1
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Credit card</b>									
IForest	0.05	0.043	0.005	0.07	0.014	0.966	0.862	0.986	0.923
	0.10	0.087	0.012	0.133	0.027	0.983	0.914	0.995	0.957
	0.20	0.179	0.026	0.256	0.049	0.992	0.949	0.999	0.977
LOF	0.05	0.029	0.001	0.093	0	0.028	0.005	0.044	0.011
	0.10	0.06	0.004	0.134	0.005	0.051	0.009	0.08	0.016
	0.20	0.14	0.012	0.243	0.06	0.101	0.016	0.15	0.026
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>KDDCup99</b>									
IForest	0.05	0.037	0.004	0.061	0.012	0.998	0.972	1	0.991
	0.10	0.079	0.009	0.117	0.02	0.999	0.988	1	0.997
	0.20	0.166	0.021	0.232	0.038	1	0.996	1	1
LOF	0.05	0.036	0.001	0.102	0	0.062	0.011	0.094	0.021
	0.10	0.076	0.004	0.153	0.01	0.104	0.02	0.154	0.034
	0.20	0.161	0.012	0.264	0.053	0.18	0.035	0.261	0.059
SVM	0.05	0	0	0	0	0	0	0	0
	0.10	0	0	0	0	0	0	0	0
	0.20	0	0	0	0	0	0	0	0
<b>Mammography</b>									
IForest	0.05	0.031	0.004	0.057	0.01	0.472	0.146	0.611	0.256
	0.10	0.065	0.006	0.11	0.023	0.601	0.234	0.726	0.36
	0.20	0.134	0.014	0.197	0.038	0.732	0.352	0.825	0.49
LOF	0.05	0.033	0.004	0.061	0.009	0.434	0.127	0.552	0.216
	0.10	0.067	0.008	0.12	0.023	0.571	0.212	0.683	0.328
	0.20	0.138	0.016	0.204	0.037	0.707	0.331	0.802	0.447

*(Continued on Next Page...)*

**Submission and Formatting Instructions for ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning**

Table A7: Outlier batch detection performance on real data, using different data sets, machine learning models, and nominal FDR levels. Other details are as in Table A6. (continued)

Model	FDR					Power			
	Nominal	Mean		90th percentile		Mean		90-th quantile	
		Marg.	Cond.	Marg.	Cond.	Marg.	Cond.	Marg.	Cond.
SVM	0.05	0.012	0.001	0.03	0	0.389	0.095	0.484	0.137
	0.10	0.027	0.002	0.054	0.002	0.506	0.158	0.595	0.22
	0.20	0.06	0.004	0.098	0.016	0.627	0.248	0.709	0.323
<b>Digits</b>									
IForest	0.05	0.035	0.002	0.061	0.008	0.918	0.523	0.979	0.773
	0.10	0.075	0.007	0.119	0.017	0.966	0.7	0.997	0.872
	0.20	0.163	0.017	0.253	0.038	0.989	0.841	1	0.951
LOF	0.05	0.04	0.002	0.072	0.008	0.999	0.941	1	0.998
	0.10	0.083	0.006	0.127	0.018	1	0.983	1	1
	0.20	0.169	0.017	0.241	0.039	1	0.996	1	1
SVM	0.05	0.037	0.002	0.063	0.009	0.807	0.347	0.886	0.487
	0.10	0.082	0.007	0.121	0.019	0.894	0.517	0.94	0.659
	0.20	0.169	0.015	0.234	0.028	0.951	0.686	0.977	0.797
<b>Shuttle</b>									
IForest	0.05	0.042	0.004	0.069	0.012	1	1	1	1
	0.10	0.086	0.011	0.127	0.023	1	1	1	1
	0.20	0.176	0.025	0.244	0.045	1	1	1	1
LOF	0.05	0.039	0.004	0.063	0.013	1	1	1	1
	0.10	0.079	0.009	0.118	0.023	1	1	1	1
	0.20	0.164	0.022	0.226	0.041	1	1	1	1
SVM	0.05	0.032	0.003	0.061	0.01	1	1	1	1
	0.10	0.069	0.008	0.111	0.02	1	1	1	1
	0.20	0.153	0.018	0.22	0.038	1	1	1	1