

## PREDICTIVE INFERENCE WITH THE JACKKNIFE+

BY RINA FOYGEL BARBER<sup>1</sup>, EMMANUEL J. CANDÈS<sup>2</sup>, AADITYA RAMDAS<sup>3,\*</sup> AND RYAN J. TIBSHIRANI<sup>3,†</sup>

<sup>1</sup>*Department of Statistics, University of Chicago, [rina@uchicago.edu](mailto:rina@uchicago.edu)*

<sup>2</sup>*Departments of Statistics and Mathematics, Stanford University, [candes@stanford.edu](mailto:candes@stanford.edu)*

<sup>3</sup>*Department of Statistics and Data Science and Machine Learning Department, Carnegie Mellon University, [aramdas@cmu.edu](mailto:aramdas@cmu.edu); [†ryantibs@cmu.edu](mailto:ryantibs@cmu.edu)*

This paper introduces the *jackknife+*, which is a novel method for constructing predictive confidence intervals. Whereas the jackknife outputs an interval centered at the predicted response of a test point, with the width of the interval determined by the quantiles of leave-one-out residuals, the jackknife+ also uses the leave-one-out predictions at the test point to account for the variability in the fitted regression function. Assuming exchangeable training samples, we prove that this crucial modification permits rigorous coverage guarantees regardless of the distribution of the data points, for any algorithm that treats the training points symmetrically. Such guarantees are not possible for the original jackknife and we demonstrate examples where the coverage rate may actually vanish. Our theoretical and empirical analysis reveals that the jackknife and the jackknife+ intervals achieve nearly exact coverage and have similar lengths whenever the fitting algorithm obeys some form of stability. Further, we extend the jackknife+ to  $K$ -fold cross validation and similarly establish rigorous coverage properties. Our methods are related to *cross-conformal prediction* proposed by Vovk (*Ann. Math. Artif. Intell.* **74** (2015) 9–28) and we discuss connections.

**1. Introduction.** Suppose that we have i.i.d. training data  $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$ ,  $i = 1, \dots, n$ , and a new test point  $(X_{n+1}, Y_{n+1})$  drawn independently from the same distribution. We would like to fit a regression model to the training data, that is, a function  $\hat{\mu} : \mathbb{R}^d \rightarrow \mathbb{R}$  where  $\hat{\mu}(x)$  predicts  $Y_{n+1}$  given a new feature vector  $X_{n+1} = x$ , and then provide a prediction interval for the test point—an interval around  $\hat{\mu}(X_{n+1})$  that is likely to contain the true test response value  $Y_{n+1}$ . Specifically, given some target coverage level  $1 - \alpha$ , we would like to construct a prediction interval  $\hat{C}_{n,\alpha}$  as a function of the  $n$  training data points, such that

$$\mathbb{P}\{Y_{n+1} \in \hat{C}_{n,\alpha}(X_{n+1})\} \geq 1 - \alpha,$$

where the probability is taken with respect to a new test point  $(X_{n+1}, Y_{n+1})$  as well as with respect to the training data.

A naive solution might be to use the residuals on the training data,  $|Y_i - \hat{\mu}(X_i)|$ , to estimate the typical prediction error on the new test point—for instance, we might consider the prediction interval

$$(1.1) \quad \hat{\mu}(X_{n+1}) \pm (\text{the } (1 - \alpha) \text{ quantile of } |Y_1 - \hat{\mu}(X_1)|, \dots, |Y_n - \hat{\mu}(X_n)|).$$

However, in practice, this interval would typically undercover (meaning that the probability that  $Y_{n+1}$  lies in this interval would be lower than the target level  $1 - \alpha$ ), since due to overfitting, the residuals on the training data points  $i = 1, \dots, n$  are typically smaller than the residual on the previously unseen test point, that is,  $|Y_{n+1} - \hat{\mu}(X_{n+1})|$ .

Received May 2019; revised December 2019.

*MSC2020 subject classifications.* 62F40, 62G08, 62G09.

*Key words and phrases.* Distribution-free, jackknife, cross-validation, conformal inference, stability, leave-one-out.

In order to avoid the overfitting problem, the jackknife prediction method computes a margin of error with a leave-one-out construction:

- For each  $i = 1, \dots, n$ , fit the regression function  $\hat{\mu}_{-i}$  to the training data with the  $i$ th point removed, and compute the corresponding leave-one-out residual,  $|Y_i - \hat{\mu}_{-i}(X_i)|$ .
- Fit the regression function  $\hat{\mu}$  to the full training data, and output the prediction interval

$$(1.2) \quad \hat{\mu}(X_{n+1}) \pm (\text{the } (1 - \alpha) \text{ quantile of } |Y_1 - \hat{\mu}_{-1}(X_1)|, \dots, |Y_n - \hat{\mu}_{-n}(X_n)|).$$

Intuitively, this method should have the right coverage properties on average since it avoids overfitting—the leave-one-out residuals  $|Y_i - \hat{\mu}_{-i}(X_i)|$  reflect the typical magnitude of the error in predicting a new data point after fitting to a sample size  $n$  (or, almost equivalently,  $n - 1$ ), unlike the naive method where the residuals on the training data are likely to be too small due to overfitting.

However, the jackknife procedure does not have any universal theoretical guarantees. Although many results are known under asymptotic settings or under assumptions of stability of the regression algorithm  $\hat{\mu}$  (we will give an overview below), it is nonetheless the case that, for settings where  $\hat{\mu}$  is unstable, the jackknife method may lose predictive coverage. For example, we will see in our simulations in Section 7 that the jackknife can have extremely poor coverage using least squares regression when the sample size  $n$  is close to the dimension  $d$ .

In this paper, we introduce a new method, the *jackknife+*, that provides nonasymptotic coverage guarantees under no assumptions beyond the training and test data being exchangeable. We will see that the *jackknife+* offers, in the worst case, a  $1 - 2\alpha$  coverage rate (where  $1 - \alpha$  is the target), while the original jackknife may even have zero coverage in degenerate examples. On the other hand, empirically we often observe that the two methods yield nearly identical intervals and both achieve  $1 - \alpha$  coverage. Theoretically, we will see that under a suitable notion of stability, the *jackknife+* and jackknife both provably yield close to  $1 - \alpha$  coverage.

1.1. *Background.* The idea of resampling or subsampling from the available data, in order to assess the accuracy of our parameter estimates or predictions, has a rich history in the statistics literature. Early works developing the jackknife and bootstrap methods include Quenouille [18], Quenouille [19], Tukey [26], Miller [15], Efron [8], Stine [23]. Several papers from this period include leave-one-out methods for assessing or calibrating predictive accuracy, similar to the predictive interval constructed in (1.2) above, for example, Stone [24], Geisser [10], Butler and Rothman [4], generally using the term “cross-validation” to refer to this approach. (In this work, we will instead use the term “jackknife” to refer to the leave-one-out style of prediction methods, as is common in the modern literature.) Efron and Gong [9] provides an overview of the early literature on these types of methods.

While this rich literature demonstrated extensive evidence of the reliable performance of the jackknife in practice, relatively little has been known about the theoretical properties of this type of method until recently. Steinberger and Leeb [21], Steinberger and Leeb [22] have developed results proving valid predictive coverage of the jackknife under assumptions of algorithmic stability, meaning that the fitted model  $\hat{\mu}$  and its leave-one-out version  $\hat{\mu}_{-i}$  are required to give similar predictions at the test point. This work builds on earlier results by Bousquet and Elisseeff [2], which study generalization bounds for risk minimization through the framework of stability conditions; an earlier work in this line is that of Devroye and Wagner [7], which give analogous results for classification risk.

In contrast to cross-validation methods, which perform well but are difficult to analyze theoretically, we can instead consider a simple *validation* or *holdout* method. We first partition the training data as  $\{1, \dots, n\} = S_{\text{train}} \cup S_{\text{holdout}}$ , then fit  $\hat{\mu}_{\text{train}}$  on the subset  $S_{\text{train}}$  of the training data and construct a predictive interval

$$(1.3) \quad \hat{\mu}_{\text{train}}(X_{n+1}) \pm (\text{the } (1 - \alpha) \text{ quantile of } |Y_i - \hat{\mu}_{\text{train}}(X_i)|, i \in S_{\text{holdout}}).$$

Papadopoulos [16], Vovk [27], Lei et al. [14] study this type of method, under the name *split conformal prediction* or *inductive conformal prediction*, through the framework of exchangeability, and prove that  $1 - \alpha$  predictive coverage holds with no assumptions on the algorithm  $\mathcal{A}$  or on the distribution of the data (with a small correction to the definition of the quantile). This method is also computationally very cheap, as we only need to fit a single regression function  $\hat{\mu}_{\text{train}}$ —in contrast, jackknife and cross-validation methods require running the regression many times. However, these benefits come at a statistical cost. If the training size  $|S_{\text{train}}|$  is much smaller than  $n$ , then the fitted model  $\hat{\mu}_{\text{train}}$  may be a poor fit, leading to wide prediction intervals; if instead we decide to take  $|S_{\text{train}}| \approx n$  then instead  $|S_{\text{holdout}}|$  is very small, leading to high variability.

Finally, Vovk [28], Vovk et al. [30] proposed the *cross-conformal prediction* method, which is closely related to the jackknife+. We describe the cross-conformal method, and the previously known theoretical guarantees, in detail later on. Their work is based on the conformal prediction method (see Vovk, Gammerman and Shafer [29], Lei et al. [14] for background), which provably achieves distribution-free predictive coverage at the target level  $1 - \alpha$  but at an extremely high computational cost.

1.2. *Notation.* Before proceeding, we first define some notation. First, for any values  $v_i$  indexed by  $i = 1, \dots, n$ , define<sup>1</sup>

$$\hat{q}_{n,\alpha}^+ \{v_i\} = \text{the } \lceil (1 - \alpha)(n + 1) \rceil\text{-th smallest value of } v_1, \dots, v_n,$$

the  $1 - \alpha$  quantile of the empirical distribution of these values. Similarly, we will let  $\hat{q}_{n,\alpha}^- \{v_i\}$  denote the  $\alpha$  quantile of the distribution,

$$\hat{q}_{n,\alpha}^- \{v_i\} = \text{the } \lfloor \alpha(n + 1) \rfloor\text{-th smallest value of } v_1, \dots, v_n = -\hat{q}_{n,\alpha}^+ \{-v_i\}.$$

With this notation, the “naive” prediction interval in (1.1) can be defined as

$$(1.4) \quad \hat{C}_{n,\alpha}^{\text{naive}}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \hat{q}_{n,\alpha}^+ \{|Y_i - \hat{\mu}(X_i)|\}.$$

Second, we will write  $\mathcal{A}$  to denote the algorithm mapping a training data set of any size, to the fitted regression function. Formally,  $\mathcal{A}$  is a map from  $\bigcup_{m \geq 1} (\mathbb{R}^d \times \mathbb{R})^m$  (i.e., the collection of all training sets of any size  $m \geq 1$ ), to the space of functions  $\mathbb{R}^d \rightarrow \mathbb{R}$ . For example, when  $\hat{\mu}$  is the regression function fitted on the training data  $(X_1, Y_1), \dots, (X_n, Y_n)$ , we can write

$$(1.5) \quad \hat{\mu} = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n)).$$

Similarly, to compute the leave-one-out residuals for the jackknife, we let

$$(1.6) \quad \hat{\mu}_{-i} = \mathcal{A}((X_1, Y_1), \dots, (X_{i-1}, Y_{i-1}), (X_{i+1}, Y_{i+1}), \dots, (X_n, Y_n)),$$

and then the jackknife prediction interval (1.2) can be written as

$$(1.7) \quad \hat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1}) = \hat{\mu}(X_{n+1}) \pm \hat{q}_{n,\alpha}^+ \{R_i^{\text{LOO}}\},$$

where  $R_i^{\text{LOO}} = |Y_i - \hat{\mu}_{-i}(X_i)|$  denotes the  $i$ th leave-one-out residual.

From this point on, we will assume without comment that  $\mathcal{A}$  satisfies a symmetry condition, namely,  $\mathcal{A}$  must be invariant to reordering the data, that is,

$$(1.8) \quad \mathcal{A}((X_{\pi(1)}, Y_{\pi(1)}), \dots, (X_{\pi(m)}, Y_{\pi(m)})) = \mathcal{A}((X_1, Y_1), \dots, (X_m, Y_m))$$

for any sample size  $m \geq 1$ , any points  $(X_1, Y_1), \dots, (X_m, Y_m)$ , and any permutation  $\pi$  of the indices  $\{1, \dots, m\}$ .

<sup>1</sup>In defining the quantiles  $\hat{q}_{n,\alpha}^+$  of the residuals, we use  $(1 - \alpha)(n + 1)$  rather than  $(1 - \alpha)n$  to correct for the finite sample size—we will see later on why this correction is natural. For the jackknife, it is perhaps more common to see  $n$  in place of  $n + 1$ , that is, the residual quantile is defined slightly differently, but for large  $n$  the difference is negligible. Formally, if  $\alpha < \frac{1}{n+1}$  and so  $(1 - \alpha)(n + 1) > n$ , then we set  $\hat{q}_{n,\alpha}^+ \{v_i\} = \infty$ .

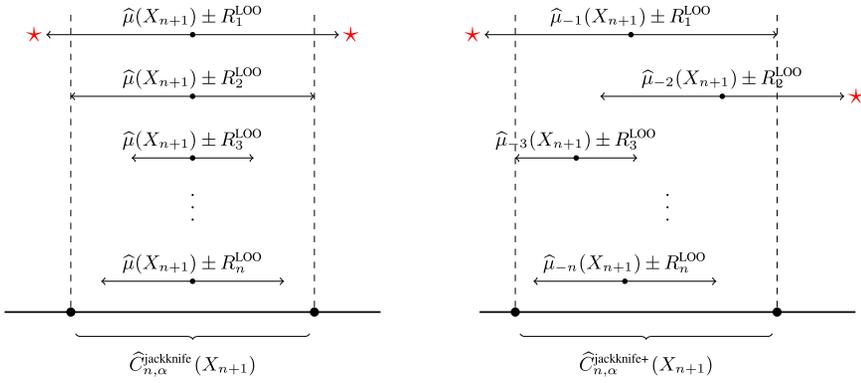


FIG. 1. Illustration of the usual jackknife and the new jackknife+. The resulting prediction intervals are chosen so that, on either side, the boundary is exceeded by a sufficiently small proportion of the two sided arrows—above, these are marked with a star.

**2. The jackknife+.** Our jackknife+ method is a modification of the jackknife (1.7). Defining  $\hat{\mu}_{-i}$  as in (1.6), the jackknife+ prediction interval is given by

$$(2.1) \quad \hat{C}_{n,\alpha}^{\text{jackknife+}}(X_{n+1}) = [\hat{q}_{n,\alpha}^- \{ \hat{\mu}_{-i}(X_{n+1}) - R_i^{\text{LOO}} \}, \hat{q}_{n,\alpha}^+ \{ \hat{\mu}_{-i}(X_{n+1}) + R_i^{\text{LOO}} \}].$$

To compare this to the usual jackknife, we observe that  $\hat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1})$  can equivalently be defined as

$$\hat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1}) = [\hat{q}_{n,\alpha}^- \{ \hat{\mu}(X_{n+1}) - R_i^{\text{LOO}} \}, \hat{q}_{n,\alpha}^+ \{ \hat{\mu}(X_{n+1}) + R_i^{\text{LOO}} \}].$$

The constructions of the usual jackknife and the new jackknife+ are compared in Figure 1. While both versions of jackknife use the leave-one-out residuals, the difference is that for jackknife, we center our interval on the predicted value  $\hat{\mu}(X_{n+1})$  fitted on the full training data, while for jackknife+ we use the leave-one-out predictions  $\hat{\mu}_{-i}(X_{n+1})$  for the test point.

Figure 1 illustrates that, if the leave-one-out fitted functions  $\hat{\mu}_{-i}$  are all quite similar to  $\hat{\mu}$ , which was fitted on the full training data, then the two methods should return nearly identical prediction intervals. On the other hand, in settings where the regression algorithm is extremely sensitive to the training data, such that removing one data point can substantially change the predicted value at  $X_{n+1}$ , the output may be quite different. In Section 5, we will examine the role of this type of instability in  $\hat{\mu}$  more closely.

To give one further interpretation of the difference between the two methods, while the jackknife interval  $\hat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1})$  is defined as a symmetric interval around the prediction  $\hat{\mu}(X_{n+1})$  for the test point (1.7), the jackknife+ interval can be interpreted as an interval around the median prediction,

$$\text{Median}(\hat{\mu}_{-1}(X_{n+1}), \dots, \hat{\mu}_{-n}(X_{n+1})),$$

which is guaranteed to lie inside  $\hat{C}_{n,\alpha}^{\text{jackknife+}}(X_{n+1})$  for any  $\alpha \leq \frac{1}{2}$  (in general, however, the jackknife+ interval will not be symmetric around this median prediction).

As detailed in Section 7, the jackknife and jackknife+ often perform nearly identically in practice (and generally achieve an empirical coverage level very close to the target  $1 - \alpha$ ), but in some more challenging examples where the regression algorithm is less stable, the original jackknife may lose coverage while jackknife+ still achieves the target coverage level.

Finally, we remark that in settings where the distribution of  $Y|X$  is highly skewed, it may be more natural to consider an asymmetric version of this method; we consider this extension in Appendix A.

2.1. *Assumption-free guarantees.* Remarkably, although the jackknife+ method appears to only be a slight modification of jackknife, our main result proves that the jackknife+ is guaranteed to achieve predictive coverage at the level  $1 - 2\alpha$ , without making any assumptions on the distribution of the data  $(X, Y)$  or the nature of the regression method  $\mathcal{A}$ .

For this theorem, and all results that follow, all probabilities are stated with respect to the distribution of the training data points  $(X_1, Y_1), \dots, (X_n, Y_n)$  and the test data point  $(X_{n+1}, Y_{n+1})$  drawn i.i.d. from an arbitrary distribution  $P$ , and we assume implicitly that the regression method  $\mathcal{A}$  is invariant to the ordering of the data (1.8). We will treat the sample size  $n \geq 2$  and the target coverage level  $\alpha \in [0, 1]$  as fixed throughout.

**THEOREM 1.** *The jackknife+ prediction interval satisfies*

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife+}}(X_{n+1})\} \geq 1 - 2\alpha.$$

This result is proved in Section 6 using the exchangeability of the  $n + 1$  data points  $(X_1, Y_1), \dots, (X_{n+1}, Y_{n+1})$ —we remark that this theorem actually holds more generally under the assumption that the  $n + 1$  data points are exchangeable, with the i.i.d. assumption as a special case.

In practice, we generally expect to achieve the target level  $1 - \alpha$  with either version of the jackknife. A natural question is whether the factor of 2 appearing in the coverage guarantee for jackknife+ is real, or is merely an artifact of the proof. We would also want to know whether analogous results may be possible for the original jackknife.

In fact, our next result constructs explicit pathological examples to see that, without making assumptions, we cannot improve our theoretical guarantee for the jackknife+ (i.e., we cannot remove the factor of 2 appearing in Theorem 1), and no guarantee at all is possible for the jackknife. For completeness, we also construct an example to show zero coverage for the naive method, although for that method we expect to see undercoverage in practice, not just in pathological examples.

**THEOREM 2.** *For any sample size  $n \geq 2$ , any  $\alpha \in [\frac{1}{n+1}, 1]$ , and any dimension  $d \geq 1$ , there exists a distribution on  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$  and a regression algorithm  $\mathcal{A}$ , such that the predictive coverage of the naive prediction interval (1.4) and the jackknife prediction interval (1.7) satisfy*

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{naive}}(X_{n+1})\} = \mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1})\} = 0.$$

*Furthermore, if  $\alpha \leq \frac{1}{2}$ , there exists a distribution on  $(X, Y) \in \mathbb{R}^d \times \mathbb{R}$  and a regression algorithm  $\mathcal{A}$ , such that the predictive coverage of jackknife+ satisfies*

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife+}}(X_{n+1})\} \leq 1 - 2\alpha + 6\sqrt{\frac{\log n}{n}}.$$

The proof of this theorem, and the proofs for all our theoretical results presented below, are deferred to the Supplementary Material [1]. The example for the original jackknife is simple—we choose the regression algorithm  $\mathcal{A}$  so that models fitted at sample size  $n$  are always less accurate than models fitted at sample size  $n - 1$ . The construction for jackknife+ is substantially more technical, and is similar in spirit to the example sketched in Vovk [28], Appendix A, for cross-conformal predictors in the setting of exchangeable data. (The constant 6 on the vanishing term in the bound for jackknife+ is simply an artifact of the proof, and can certainly be improved with a more careful construction.)

2.2. *The jackknife-minmax method.* We have seen that the best possible coverage guarantee for jackknife+, in the assumption-free setting, is  $1 - 2\alpha$  rather than the target level  $1 - \alpha$ . To address this gap, we can consider a more conservative alternative to the jackknife+, which will remove the factor of 2 from the theoretical bound. We define the jackknife-minmax method as follows:

$$(2.2) \quad \widehat{C}_{n,\alpha}^{\text{jack-mm}}(X_{n+1}) = \left[ \min_{i=1,\dots,n} \widehat{\mu}_{-i}(X_{n+1}) - \widehat{q}_{n,\alpha}^+ \{R_i^{\text{LOO}}\}, \max_{i=1,\dots,n} \widehat{\mu}_{-i}(X_{n+1}) + \widehat{q}_{n,\alpha}^+ \{R_i^{\text{LOO}}\} \right].$$

It is simple to verify that this interval is strictly more conservative than jackknife+, meaning that for any data set, we have

$$\widehat{C}_{n,\alpha}^{\text{jackknife+}}(X_{n+1}) \subseteq \widehat{C}_{n,\alpha}^{\text{jack-mm}}(X_{n+1}).$$

The advantage that jackknife-minmax provides is that, without any assumptions on the algorithm or distribution of the data, it always achieves the target coverage rate.

**THEOREM 3.** *The jackknife-minmax prediction interval satisfies*

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jack-mm}}(X_{n+1})\} \geq 1 - \alpha.$$

In practice, however, we will see that the jackknife-minmax prediction interval is generally too conservative.

**3. CV+ for K-fold cross-validation.** Suppose that we split the training sample into  $K$  disjoint subsets  $S_1, \dots, S_K$  each of size  $m = n/K$  (assumed to be an integer). Let

$$\widehat{\mu}_{-S_k} = \mathcal{A}((X_i, Y_i) : i \in \{1, \dots, n\} \setminus S_k)$$

be the regression function fitted onto the training data with the  $k$ th subset removed. To assess the quality of our regression algorithm using cross-validation (CV), we would consider the residuals from this  $K$ -fold process, namely,

$$R_i^{\text{CV}} = |Y_i - \widehat{\mu}_{-S_{k(i)}}(X_i)|, \quad i = 1, \dots, n,$$

where  $k(i) \in \{1, \dots, K\}$  identifies the subset that contains  $i$ , that is,  $i \in S_{k(i)}$ . Using these residuals, we can define the CV+ prediction interval as

$$(3.1) \quad \widehat{C}_{n,K,\alpha}^{\text{CV+}}(X_{n+1}) = \left[ \widehat{q}_{n,\alpha}^- \{ \widehat{\mu}_{-S_{k(i)}}(X_{n+1}) - R_i^{\text{CV}} \}, \widehat{q}_{n,\alpha}^+ \{ \widehat{\mu}_{-S_{k(i)}}(X_{n+1}) + R_i^{\text{CV}} \} \right].$$

Of course, jackknife+ can be viewed as a special case of CV+, by setting  $K = n$ . The advantage of the CV+ method, when we choose a smaller  $K$ , is that we only need to compute  $K$  rather than  $n$  models; however, this will likely come at the cost of slightly wider intervals, because the models  $\widehat{\mu}_{-S_k}$  are fitted using a lower sample size (i.e.,  $n(1 - 1/K)$ ) and will lead to slightly larger residuals.

3.1. *Assumption-free guarantee for CV+.* Our next result verifies that the CV+ prediction interval enjoys essentially the same worst-case coverage guarantee as jackknife+.

**THEOREM 4.** *The K-fold CV+ prediction interval satisfies the following coverage guarantees:*

(a) (Adapted from Vovk and Wang [31], Vovk et al. [30].)

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,K,\alpha}^{\text{CV}+}(X_{n+1})\} \geq 1 - 2\alpha - \frac{2(1 - 1/K)}{n/K + 1}.$$

(b)

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,K,\alpha}^{\text{CV}+}(X_{n+1})\} \geq 1 - 2\alpha - \frac{1 - K/n}{K + 1}.$$

Combining the two bounds, it follows that for all  $K$ ,

$$\begin{aligned} \mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,K,\alpha}^{\text{CV}+}(X_{n+1})\} &\geq 1 - 2\alpha - \min\left\{\frac{2(1 - 1/K)}{n/K + 1}, \frac{1 - K/n}{K + 1}\right\} \\ &\geq 1 - 2\alpha - \sqrt{2/n}. \end{aligned}$$

The first part of this result, part (a), is derived from the work of Vovk and Wang [31] and Vovk et al. [30]; we give more details on this in Section 3.2 below. This known result proves that the worst-case coverage is essentially  $1 - 2\alpha$  when  $K$  is sufficiently small, that is,  $K \ll n$ . Our new work, proving part (b), completes the picture by giving a meaningful bound for the case where  $K$  is large (at the extreme,  $K = n$  for leave-one-out methods). By combining the two bounds, we see that coverage is essentially  $1 - 2\alpha$  at any  $K$ , since the excess noncoverage is at most  $\sqrt{2/n}$  uniformly over any choice of  $K$ .

We can also compare our result to the holdout or split conformal method (1.3), which is equivalent to fitting a model  $\widehat{\mu}_{-S_1}$  and constructing the prediction interval using the quantile of the residuals  $R_i^{\text{CV}}$  for  $i \in S_1$ , but using only a single subset  $S_1$  (without repeating  $K$  times for each fold in the partition  $S_1, \dots, S_K$ ). As discussed earlier, this method offers an assumption-free guarantee of  $1 - \alpha$  coverage, but this comes at the cost of higher variance due to the single split—in contrast, CV+ reduces variance by averaging over all  $K$  splits, but at the cost of a weaker theoretical guarantee.

3.2. *Related method: Cross-conformal predictors.* Our proposed CV+ prediction interval is related to the *cross-conformal prediction* method of Vovk [28], Vovk et al. [30], which (in its symmetric version) returns the predictive set

$$\begin{aligned} &\widehat{C}_{n,K,\alpha}^{\text{cross-conf}}(X_{n+1}) \\ (3.2) \quad &= \left\{ y \in \mathbb{R} : \frac{\tau + \sum_{i=1}^n \mathbb{1}\{|y - \widehat{\mu}_{-S_{k(i)}}(X_{n+1})| < R_i^{\text{CV}}\} + \tau \mathbb{1}\{|y - \widehat{\mu}_{-S_{k(i)}}(X_{n+1})| = R_i^{\text{CV}}\}}{n + 1} > \alpha \right\}. \end{aligned}$$

Here,  $\tau \sim \text{Unif}[0, 1]$  introduces randomization into the method. By comparing to CV+, we can verify that

$$(3.3) \quad \widehat{C}_{n,K,\alpha}^{\text{cross-conf}}(X_{n+1}) \subseteq \widehat{C}_{n,K,\alpha}^{\text{CV}+}(X_{n+1})$$

deterministically (we demonstrate this in the Supplementary Material [1]). The two methods will sometimes produce the same output, but not always—in particular,  $\widehat{C}_{n,K,\alpha}^{\text{cross-conf}}(X_{n+1})$  may in principle return a predictive set that is a disjoint union of multiple intervals, while CV+ always returns an interval.

We next compare our theoretical findings with the known results for cross-conformal. Vovk et al. [30] show that the  $K$ -fold cross-conformal method has coverage at least<sup>2</sup>

$$(3.4) \quad 1 - 2\alpha - 2(1 - \alpha) \frac{1 - 1/K}{n/K + 1}.$$

When  $K$  is small, this additional term is negligible, and so we essentially have  $1 - 2\alpha$  coverage for cross-conformal. However, for large  $K$ , such as  $K = n$  for the leave-one-out method, their earlier result does not yield a meaningful guarantee—the guaranteed coverage level is zero. In contrast, our new assumption-free result in Theorem 1 proves  $1 - 2\alpha$  coverage for the jackknife+ method (i.e., with  $K = n$  folds), and Theorem 4 ensures  $1 - 2\alpha - \sqrt{2/n}$  coverage for  $K$ -fold CV+ at any choice of  $K$ .

REMARK 1. By examining the proofs of Theorems 1 and 4, we can see that the arguments apply directly to the  $K$ -fold (or  $n$ -fold) cross-conformal method; that is, our proofs for these theorems also establish that

$$\begin{aligned} \mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,K,\alpha}^{\text{CV}+}(X_{n+1})\} &\geq 1 - 2\alpha - \min\left\{\frac{2(1 - 1/K)}{n/K + 1}, \frac{1 - K/n}{K + 1}\right\} \\ &\geq 1 - 2\alpha - \sqrt{2/n} \end{aligned}$$

for  $K$ -fold cross-conformal with any  $K$ . In the special case that  $K = n$  we are guaranteed coverage  $\geq 1 - 2\alpha$ . The first term in the minimum was established by Vovk et al. [30] as presented in (3.4) above, but the second term (which allows for meaningful coverage for large values of  $K$ , for example,  $K = n$ ) is a new result.

3.3. *An alternative method: Conformal prediction.* The final related method we present is *conformal prediction* [29]. (We will sometimes refer to this method as “full” conformal prediction in order to distinguish it from the split conformal or cross-conformal methods described above.) Given the base algorithm  $\mathcal{A}$ , the full conformal prediction method outputs a prediction set (which consists of a union of one or more intervals) constructed as follows:

$$(3.5) \quad \widehat{C}_{n,\alpha}^{\text{conf}}(X_{n+1}) = \{y \in \mathbb{R} : |y - \widehat{\mu}^y(X_{n+1})| \leq \widehat{q}_{n,\alpha}^+ \{|Y_i - \widehat{\mu}^y(X_i)|\}\},$$

where

$$\widehat{\mu}^y = \mathcal{A}((X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, y))$$

denotes the output of the algorithm run on the training data augmented with the hypothesized test point  $(X_{n+1}, y)$ . In other words, to determine whether to include a value  $y$  in the prediction set at a new point  $X_{n+1}$ , we need to train the algorithm on the training+test data (as though  $Y_{n+1} = y$  were the true response value), and then see whether the residual of the test point “conforms” with the residuals on the remaining  $n$  points. The exchangeability of the test and training data ensures that

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{conf}}(X_{n+1})\} \geq 1 - \alpha,$$

that is, coverage at the target level. However, this desirable theoretical property comes at a high cost—we can see by construction of the prediction interval  $\widehat{C}_{n,\alpha}^{\text{conf}}(x)$  that the training algorithm  $\mathcal{A}$  needs to be rerun for every test point feature vector  $x$  we might consider, and for every possible response value  $y \in \mathbb{R}$  (or in practice, for each  $y$  in a fine grid over  $\mathbb{R}$ ).

<sup>2</sup>Vovk et al. [30] do not state this coverage result directly, but instead prove  $1 - 2\alpha$  coverage for a modification of the cross-conformal method; however, these two formulations can be shown to be equivalent. We give details in the Supplementary Material [1].

In certain special cases, there are computational tricks allowing for efficient calculation of the prediction set, for example, linear regression or ridge regression [3], and the Lasso [13]. Outside of these special cases, full conformal is prohibitively expensive in practice, even on moderately sized data sets; while it provides an extremely elegant and theoretically rigorous framework for distribution-free inference, it is not practical in many applied settings.

**4. Summary of coverage guarantees and computational costs.** In light of these theoretical results, which method should a statistician choose in practice? Table 1 summarizes the theoretical results behind each of the methods under consideration, and the typical empirical performance that we have observed.

Given the theoretical and empirical properties of the various options, we therefore recommend the jackknife+ as a practical alternative to the usual jackknife predictive intervals. On the one hand, the empirical performance of the jackknife+ is nearly identical to that of the original jackknife (assuming we avoid pathological examples), with both methods giving intervals of nearly the same width and achieving close to the target  $1 - \alpha$  coverage level. However, while the jackknife offers no theoretical guarantees in the absence of stability assumptions, the jackknife+ achieves at least  $1 - 2\alpha$  coverage in the worst possible case. On the other hand, the methods achieving  $1 - \alpha$  (rather than  $1 - 2\alpha$ ) coverage guarantees are either less statistically efficient in the sense of producing wider intervals (split conformal uses models fitted on a smaller portion of the data while jackknife-minmax is generally too conservative), or suffer from computational infeasibility (full conformal is computationally prohibitive aside from perhaps a few special cases).

We now turn to a direct comparison of the computational costs of these eight methods. The split conformal and naive methods require only one run of the regression algorithm  $\mathcal{A}$  (to fit  $\hat{\mu}$  on the full training data), while each of the jackknife methods requires  $n$  runs (to fit  $\hat{\mu}_{-i}$  for each  $i = 1, \dots, n$ —and one additional run to fit  $\hat{\mu}$ , in the case of the original jackknife). If the training sample size  $n$  is so large that fitting  $n$  regression functions is not feasible, we may instead prefer to use  $K$ -fold cross-validation. In contrast, the full conformal method must train  $\mathcal{A}$  many more times—once for each possible combination of a test point feature vector  $x$  and a possible response value  $y$ —except for special cases such as linear regression or ridge regression. These observations are summarized below:

Table 2 compares the computational cost (ignoring constants) of each method when run on a training sample of size  $n$ , for producing prediction sets on  $n_{\text{test}}$  many test points. The middle column (“Model training cost”) counts the number of times that the model fitting algorithm  $\mathcal{A}$  is run on a training data set<sup>3</sup> of size (up to)  $n$ . The value  $n_{\text{grid}}$  denotes the number of grid

TABLE 1

*Summary of distribution-free theoretical guarantees and typical empirical performance for all methods*

Method	Assumption-free theory	Typical empirical coverage
Naive (1.4)	No guarantee	$< 1 - \alpha$
Split conf. (holdout) (1.3)	$\geq 1 - \alpha$ coverage	$\approx 1 - \alpha$
Jackknife (1.7)	No guarantee	$\approx 1 - \alpha$ , or $< 1 - \alpha$ if $\hat{\mu}$ unstable
Jackknife+ (2.1)	$\geq 1 - 2\alpha$ coverage	$\approx 1 - \alpha$
Jackknife-minmax (2.2)	$\geq 1 - \alpha$ coverage	$> 1 - \alpha$
Full conformal (3.5)	$\geq 1 - \alpha$ coverage	$\approx 1 - \alpha$ , or $> 1 - \alpha$ if $\hat{\mu}$ overfits
K-fold CV+ (3.1)	$\geq 1 - 2\alpha$ coverage	$\gtrsim 1 - \alpha$
K-fold cross-conf. (3.2)	$\geq 1 - 2\alpha$ coverage	$\gtrsim 1 - \alpha$

<sup>3</sup>It is worth mentioning that for several common regression algorithms, the  $n$  leave-one-out residuals can be obtained without refitting  $n$  times, but by simply reweighting the in-sample training residuals. Examples include

TABLE 2  
*Summary of computational costs for all methods*

Method	Model training cost	Model evaluation cost
Naive (1.4)	1	$n + n_{\text{test}}$
Split conf. (holdout) (1.3)	1	$n + n_{\text{test}}$
Jackknife (1.7)	$n$	$n + n_{\text{test}}$
Jackknife+ (2.1)	$n$	$n_{\text{test}} \cdot n$
Jackknife-minmax (2.2)	$n$	$n_{\text{test}} \cdot n$
K-fold CV+ (3.1)	$K$	$n + n_{\text{test}} \cdot K$
K-fold cross-conf. (3.2)	$K$	$n + n_{\text{test}} \cdot K$
Full conformal (3.5)	$n_{\text{test}} \cdot n_{\text{grid}}$	$n_{\text{test}} \cdot n_{\text{grid}} \cdot n$

points of possible  $y$  values (a fine grid over  $\mathbb{R}$ ), used in the construction of the full conformal prediction method (3.5). The last column (“Model evaluation cost”) counts the number of times we evaluate a fitted model  $\hat{\mu}$  on a new data point. In most settings, the model training cost is dominant, for example, training a neural network is far more costly than evaluating the prediction of a trained network on a new example. There are important exceptions, however, such as  $K$ -nearest neighbors, where computing a prediction incurs the cost of identifying the  $K$  neighbors of the test point.

**5. Guarantees under stability assumptions.** Next, we consider how adding stability assumptions—conditions that ensure that the fitted regression function  $\hat{\mu}$  is not too sensitive to perturbations of the training data set—can improve the theoretical guarantees of the jackknife and its variants. (For simplicity, we only consider leave-one-out methods, and do not examine  $K$ -fold cross-validation here.)

5.1. *In-sample and out-of-sample stability.* Fix any  $\epsilon \geq 0$ ,  $\nu \in [0, 1]$ , any sample size  $n \geq 2$ , and any distribution  $P$  on  $(X, Y)$ . We say that a regression algorithm  $\mathcal{A}$  satisfies  $(\epsilon, \nu)$  out-of-sample stability with respect to the distribution  $P$  and sample size  $n$  if, for all  $i \in \{1, \dots, n\}$ ,

$$(5.1) \quad \mathbb{P}\{|\hat{\mu}(X_{n+1}) - \hat{\mu}_{-i}(X_{n+1})| \leq \epsilon\} \geq 1 - \nu,$$

for  $\hat{\mu}$  and  $\hat{\mu}_{-i}$  defined as before in (1.5) and (1.6). The probability above is taken with respect to the distribution of the data points  $(X_1, Y_1), \dots, (X_n, Y_n), (X_{n+1}, Y_{n+1})$  drawn i.i.d. from  $P$ . Similarly,  $\mathcal{A}$  satisfies  $(\epsilon, \nu)$  in-sample stability with respect to the distribution  $P$  and sample size  $n$  if, for all  $i \in \{1, \dots, n\}$ ,

$$(5.2) \quad \mathbb{P}\{|\hat{\mu}(X_i) - \hat{\mu}_{-i}(X_i)| \leq \epsilon\} \geq 1 - \nu.$$

Naturally, since the data points are exchangeable, if (5.1) or (5.2) holds for any single  $i \in \{1, \dots, n\}$  then it holds for all  $i \in \{1, \dots, n\}$ . These types of conditions appear elsewhere in the literature, for example, Bousquet and Elisseeff [2] define similar conditions, termed “hypothesis stability” and “pointwise hypothesis stability.”

While the out-of-sample and in-sample stability properties may at first appear similar, they are extremely different in practice. Out-of-sample stability requires that, for a test point that is *independent of the training data*, the predicted value does not change much if we remove

---

linear smoothing methods like ordinary least squares, kernel ridge regression, kernel smoothing, thin plate splines and smoothing splines. Another interesting example is random forests, where the  $i$ th leave-one-out fit can be obtained by simply ignoring all trees containing the  $i$ th point.

one point from the training data. In contrast, in-sample stability requires that, *for a point in the training data set*, the predicted value does not change much if we remove this point itself from the training data set. In a scenario where the model fitting algorithm suffers from strong overfitting, we would expect in-sample stability to be very poor, while out-of-sample stability may still hold, for example, we will see in Section 5.5 that this is the case for  $K$ -nearest neighbor methods. On the other hand, strongly convex regularization, such as ridge regression, induces both in- and out-of-sample stability [2], Example 3. This is not the case, however, for sparse regression methods (e.g.,  $\ell_1$  regularization), which are proved by Xu, Caramanis and Mannor [32] to be incompatible with in-sample stability.

5.2. *Summary of stability results.* Before giving the details of our theoretical results, we summarize our findings on the various methods’ predictive coverage guarantees, with and without stability assumptions (Table 3).

The assumption free results are the same as those discussed in Section 4, while the results under stability assumptions are presented next in Theorems 5 and 6.

5.3. *Out-of-sample stability and the jackknife.* We will next prove that out-of-sample stability is sufficient for the jackknife and jackknife+ methods to achieve the target coverage rate, with a slight modification. Define

$$\widehat{C}_{n,\alpha}^{\text{jackknife},\epsilon}(X_{n+1}) = \widehat{\mu}(X_{n+1}) \pm (\widehat{q}_{n,\alpha}^+ \{R_i^{\text{LOO}}\} + \epsilon),$$

and similarly, define

$$\widehat{C}_{n,\alpha}^{\text{jackknife+},\epsilon}(X_{n+1}) = [\widehat{q}_{n,\alpha}^- \{\widehat{\mu}_{-i}(X_{n+1}) - R_i^{\text{LOO}}\} - \epsilon, \widehat{q}_{n,\alpha}^+ \{\widehat{\mu}_{-i}(X_{n+1}) + R_i^{\text{LOO}}\} + \epsilon],$$

which we refer to as the  $\epsilon$ -inflated versions of the jackknife and jackknife+ predictive intervals.

**THEOREM 5.** *Suppose that the regression algorithm  $A$  satisfies the  $(\epsilon, \nu)$  out-of-sample stability property (5.1) with respect to the data distribution  $P$  and the sample size  $n$ . Then the  $\epsilon$ -inflated jackknife prediction interval satisfies*

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife},\epsilon}(X_{n+1})\} \geq 1 - \alpha - 2\sqrt{\nu}.$$

Similarly, the  $2\epsilon$ -inflated jackknife+ prediction interval satisfies

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife+},2\epsilon}(X_{n+1})\} \geq 1 - \alpha - 4\sqrt{\nu}.$$

(The different amounts of inflation,  $\epsilon$  for jackknife versus  $2\epsilon$  for jackknife+, are simply an artifact of the particular definition of out-of-sample stability that we use, and should not be interpreted as a meaningful difference between these two methods.)

TABLE 3  
Summary of theoretical guarantees for all methods under various stability assumptions

Method	Assumption-free theory	Out-of-sample stability	In-sample and out-of-sample stability
Naive (1.4)	No guarantee	No guarantee	$\approx 1 - \alpha$
Jackknife (1.7)	No guarantee	$\approx 1 - \alpha$	$\approx 1 - \alpha$
Jackknife+ (2.1)	$1 - 2\alpha$	$\approx 1 - \alpha$	$\approx 1 - \alpha$
Jackknife-minmax (2.2)	$1 - \alpha$	$1 - \alpha$	$1 - \alpha$

We remark that if we additionally assume that, in the data distribution,  $Y|X$  has a bounded conditional density (e.g.,  $Y = \mu(X) + \mathcal{N}(0, \sigma^2)$  for some unknown true mean function  $\mu(\cdot)$ ), then the result of Theorem 5 is sufficient to ensure that the (noninflated) jackknife and jackknife+ intervals achieve close to target coverage. The reason is this: if the conditional density of  $Y|X$  is bounded by some constant  $c < \infty$ , then very little probability can be captured by inflating the interval. Specifically,

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife},\epsilon}(X_{n+1}) \setminus \widehat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1})\} \leq 2\epsilon c.$$

Combined with the result of Theorem 5, this proves that

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1})\} \geq 1 - \alpha - 2\sqrt{v} - 2\epsilon c.$$

Similarly, for jackknife+, we have

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife}+}(X_{n+1})\} \geq 1 - \alpha - 4\sqrt{v} - 4\epsilon c.$$

5.4. *In-sample stability and overfitting.* To contrast the scenarios of in-sample and out-of-sample stability, we will next demonstrate that adding the in-sample stability assumption would in fact be sufficient for the “naive” prediction interval, defined earlier in (1.4), to have coverage at roughly the target level. Its  $\epsilon$ -inflated version is defined as

$$(5.3) \quad \widehat{C}_{n,\alpha}^{\text{naive},\epsilon}(X_{n+1}) = \widehat{\mu}(X_{n+1}) \pm (\widehat{q}_{n,\alpha}^+ \{|Y_i - \widehat{\mu}(X_i)|\} + \epsilon).$$

Recall from Section 1 that we would typically expect  $\widehat{C}_{n,\alpha}^{\text{naive}}(X_{n+1})$  to undercover severely due to the overfitting problem (thus inspiring the use of the jackknife to avoid this issue), and similarly  $\widehat{C}_{n,\alpha}^{\text{naive},\epsilon}(X_{n+1})$  will also undercover whenever  $\epsilon$  is too small to correct for overfitting. This is often the case even when out-of-sample stability is satisfied. With in-sample stability, however, this is no longer the case. In other words, the in-sample stability property is essentially assuming that inflation by  $\epsilon$  is sufficient to correct for overfitting.

**THEOREM 6.** *Suppose that the regression algorithm  $\mathcal{A}$  satisfies both the  $(\epsilon, v)$  in-sample stability property (5.2) and the  $(\epsilon, v)$  out-of-sample stability property (5.1) with respect to the data distribution  $P$  and the sample size  $n$ . Then the naive prediction interval satisfies*

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{naive},2\epsilon}(X_{n+1})\} \geq 1 - \alpha - 4\sqrt{v}.$$

5.5. *Example:  $K$ -nearest neighbors.* To give an illustrative example, consider a  $K$ -nearest neighbor ( $K$ -NN) method. This style of example is also considered in Steinberger and Leeb [22], Example 4.1, and was studied earlier by Devroye and Wagner [7] in the context of estimating the error of a classifier, and by Bousquet and Elisseeff [2] in the context of error in regression. Given a training data set  $(X_1, Y_1), \dots, (X_n, Y_n)$  and a new test point  $x$ , our prediction is

$$\widehat{\mu}(x) = \frac{1}{K} \sum_{i \in N(x)} Y_i,$$

where  $N(x) \subset \{1, \dots, n\}$  is the set of the  $K$  nearest neighbors to the test point  $x$ , that is, the  $K$  indices  $i$  giving the smallest values of  $\|X_i - x\|_2$  (of course, we can replace the  $\ell_2$  norm with any other metric). We will assume for simplicity that there are no ties between these distances (for instance, the  $X_i$  points might be continuously distributed on  $\mathbb{R}^d$ , or we apply a random tie-breaking rule). Now consider out-of-sample stability. Let  $N(X_{n+1})$  and  $N_{-i}(X_{n+1})$  be the  $K$ -nearest neighbor sets for the test point  $X_{n+1}$  given the full training data, or the training data with data point  $i$  removed, respectively. Then we can easily verify that

$$i \notin N(X_{n+1}) \iff N(X_{n+1}) = N_{-i}(X_{n+1}) \implies \widehat{\mu}(X_{n+1}) = \widehat{\mu}_{-i}(X_{n+1}).$$

Therefore,

$$\mathbb{P}\{|\widehat{\mu}(X_{n+1}) - \widehat{\mu}_{-i}(X_{n+1})| = 0\} \geq \mathbb{P}\{i \notin N(X_{n+1})\} = 1 - \frac{K}{n},$$

where the last step holds by exchangeability of the  $n$  training points. This proves that the  $K$ -NN method satisfies  $(\epsilon, \nu)$ -out-of-sample stability with  $\epsilon = 0$  and  $\nu = K/n$ . (In contrast, we cannot hope for a similar argument to guarantee in-sample stability, since we will always have  $i \in N(X_i)$ ; that is,  $X_i$  is one of its own nearest neighbors—and so in general we will have  $\widehat{\mu}(X_i) \neq \widehat{\mu}_{-i}(X_i)$ .)

Applying the conclusion of Theorem 5 to this setting, then we see that  $K$ -NN leads to a coverage rate at least

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife}}(X_{n+1})\} \geq 1 - \alpha - 2\sqrt{K/n}$$

for the jackknife, and

$$\mathbb{P}\{Y_{n+1} \in \widehat{C}_{n,\alpha}^{\text{jackknife+}}(X_{n+1})\} \geq 1 - \alpha - 4\sqrt{K/n}$$

for the jackknife+. These results hold with no assumptions whatsoever on the distribution of the data, and in particular, we do not need to assume that the  $K$ -NN prediction is accurate or consistent on the given data.

*5.6. Comparison to existing results.* As mentioned above, Bousquet and Elisseeff [2] study stability in the context of generalization bounds for regression, with the aim of bounding risk rather than predictive inference. The predictive accuracy of the jackknife under assumptions of algorithm stability was explored by Steinberger and Leeb [21] for the linear regression setting, and in a more general setting by Steinberger and Leeb [22]. Their stability assumption (see, e.g., Steinberger and Leeb [22], Definition 1) is essentially equivalent to our out-of-sample stability condition (5.1). However, the theory obtained in their work is asymptotic, and relies also on distributional assumptions (see Steinberger and Leeb [22], (C1)), namely that  $Y_i = \mathbb{E}[Y_i | X_i] + v_i$  where the noise  $v_i$  is continuously distributed and is independent of  $X_i$  (e.g., this does not allow for heteroskedasticity). In contrast, our guarantee, in Theorem 5, offers a simple finite-sample coverage guarantee with no distributional assumptions, requiring only algorithm stability.

**6. Proof of Theorem 1.** Suppose for a moment that we have access to the test point  $(X_{n+1}, Y_{n+1})$  as well as the training data. For any indices  $i, j \in \{1, \dots, n + 1\}$  with  $i \neq j$ , let  $\widetilde{\mu}_{-(i,j)}$  define the regression function fitted on the training plus test data, with points  $i$  and  $j$  removed. (We use  $\widetilde{\mu}$  rather than  $\widehat{\mu}$  to remind ourselves that the model is fitted on a subset of the combined training and test data  $i = 1, \dots, n + 1$ , rather than a subset of only the training data.) Note that  $\widetilde{\mu}_{-(i,j)} = \widetilde{\mu}_{-(j,i)}$  for any  $i \neq j$ , and that  $\widetilde{\mu}_{-(i,n+1)} = \widehat{\mu}_{-i}$  for any  $i = 1, \dots, n$ .

Next, we define a matrix of residuals,  $R \in \mathbb{R}^{(n+1) \times (n+1)}$ , with entries

$$R_{ij} = \begin{cases} +\infty & i = j, \\ |Y_i - \widetilde{\mu}_{-(i,j)}(X_i)| & i \neq j, \end{cases}$$

that is, the off-diagonal entries represent the residual for the  $i$ th point when both  $i$  and  $j$  are left out of the regression. We also define a comparison matrix,  $A \in \{0, 1\}^{(n+1) \times (n+1)}$ , with entries

$$A_{ij} = \mathbb{1}\{R_{ij} > R_{ji}\}.$$

In other words,  $A_{ij}$  is the indicator for the event that, when excluding data points  $i$  and  $j$  from the regression, data point  $i$  has higher residual than data point  $j$ . Naturally, we see that

$A_{ij} = 1$  implies  $A_{ji} = 0$ , for any  $i, j$ . We note that this comparison matrix construction is also examined by Vovk [28], Appendix A, where it is used to establish that leave-one-out conformal methods fail to achieve  $1 - \alpha$  coverage.

Next, we are interested in finding data points  $i$  with unusually large residuals—the ones that are hardest to predict. We will define a set  $\mathcal{S}(A) \subseteq \{1, \dots, n + 1\}$  of “strange” points,<sup>4</sup>

$$\mathcal{S}(A) = \{i \in \{1, \dots, n + 1\} : A_{i\bullet} \geq (1 - \alpha)(n + 1)\},$$

where  $A_{i\bullet} = \sum_{j=1}^{n+1} A_{ij}$  is the  $i$ th row sum of  $A$ . In other words, the  $i$ th point is “strange” (i.e.,  $i \in \mathcal{S}(A)$ ) if it holds that, when we compare the residual  $R_{ij}$  of the  $i$ th point against residual  $R_{ji}$  for the  $j$ th point (for each  $j \neq i$ ), the residual  $R_{ij}$  for the  $i$ th point is the larger one, for a sufficiently high fraction of these comparisons.

From this point on, the proof will proceed as follows:

- Step 1: we will establish deterministically that  $|\mathcal{S}(A)| \leq 2\alpha(n + 1)$ , that is, for any comparison matrix  $A$  it is impossible to have more than  $2\alpha(n + 1)$  many strange points.
- Step 2: using the fact that the data points are i.i.d. (or more generally exchangeable), we will show that the probability that the test point  $n + 1$  is strange (i.e.,  $n + 1 \in \mathcal{S}(A)$ ) is therefore bounded by  $2\alpha$ .
- Step 3: finally, we will verify that the jackknife+ interval can only fail to cover the test response value  $Y_{n+1}$  if  $n + 1$  is a strange point.

*Step 1: Bounding the number of strange points.* This bound is essentially a consequence of Landau’s theorem for tournaments [12]. For data points  $i$  and  $j$ , we say that data point  $i$  “wins” its game against data point  $j$ , if  $A_{ij} = 1$ ; that is, point  $i$  has a higher residual than point  $j$ , under the corresponding regression  $\tilde{\mu}_{-(i,j)}$ . Note that each strange point  $i \in \mathcal{S}(A)$  can lose against at most  $\alpha(n + 1) - 1$  other strange points—this is because point  $i$  must win against at least  $(1 - \alpha)(n + 1)$  points in total since it is strange, and as we have defined it, point  $i$  cannot win against itself.

Let  $s = |\mathcal{S}(A)|$  denote the number of strange points. The key realization is now that, if we think about grouping each *pair* of strange points by the losing point, then we see that there are at most

$$s \cdot (\alpha(n + 1) - 1)$$

pairs of strange points. This is because there are at most  $s$  unique possibilities for the loser, and for each such loser, it can only lose against at most  $\alpha(n + 1) - 1$  other strange points, as argued above. In other words, we have established

$$\frac{s(s - 1)}{2} \leq s \cdot (\alpha(n + 1) - 1),$$

and rearranging gives  $s \leq 2\alpha(n + 1) - 1 < 2\alpha(n + 1)$ , as desired.

*Step 2: Exchangeability of the data points.* We next leverage the exchangeability of the data points to show that, since there are at most  $2\alpha(n + 1)$  strange points among a total of  $n + 1$  points, it follows that the test point has at most  $2\alpha$  probability of being strange—this reasoning uses the exchangeability of the data in exactly the same way as the conformal prediction literature [29].

To establish this formally, since the data points  $(X_1, Y_1), \dots, (X_{n+1}, Y_{n+1})$  are exchangeable and the regression fitting algorithm  $\mathcal{A}$  is invariant to the ordering of the data points (the symmetry condition (1.8)), it follows that  $A \stackrel{d}{=} \Pi A \Pi^\top$  for any  $(n + 1) \times (n + 1)$  permutation matrix  $\Pi$ , where  $\stackrel{d}{=}$  denotes equality in distribution. In particular, for any index

<sup>4</sup>The authors thank an anonymous reviewer for suggesting this presentation of the proof.

$j \in \{1, \dots, n+1\}$ , suppose we take  $\Pi$  to be any permutation matrix with  $\Pi_{j,n+1} = 1$  (i.e., corresponding to a permutation mapping  $n+1$  to  $j$ ). Then, deterministically, we have

$$n+1 \in \mathcal{S}(A) \Leftrightarrow j \in \mathcal{S}(\Pi A \Pi^\top)$$

and, therefore,

$$\mathbb{P}\{n+1 \in \mathcal{S}(A)\} = \mathbb{P}\{j \in \mathcal{S}(\Pi A \Pi^\top)\} = \mathbb{P}\{j \in \mathcal{S}(A)\}$$

for all  $j = 1, \dots, n+1$ . In other words, if we compare an arbitrary training point  $j$  versus the test point  $n+1$ , these two points are equally likely to be strange, by exchangeability of the data. We can then calculate

$$\mathbb{P}\{n+1 \in \mathcal{S}(A)\} = \frac{1}{n+1} \sum_{j=1}^{n+1} \mathbb{P}\{j \in \mathcal{S}(A)\} = \frac{\mathbb{E}[|\mathcal{S}(A)|]}{n+1} \leq 2\alpha,$$

where the last step applies the result of Step 1.

*Step 3: Connecting to jackknife+.* Finally, we need to relate the question of coverage of the jackknife+ interval, to our notion of strange points. Suppose that  $Y_{n+1} \notin \widehat{C}_{n,\alpha}^{\text{jackknife}^+}(X_{n+1})$ . This means that either

$$Y_{n+1} > \widehat{q}_{n,\alpha}^+ \{\widehat{\mu}_{-i}(X_{n+1}) + R_i^{\text{LOO}}\},$$

which implies that  $Y_{n+1} > \widehat{\mu}_{-j}(X_{n+1}) + R_j^{\text{LOO}}$  for at least  $(1-\alpha)(n+1)$  many indices  $j \in \{1, \dots, n\}$ , or otherwise

$$Y_{n+1} < \widehat{q}_{n,\alpha}^- \{\widehat{\mu}_{-i}(X_{n+1}) - R_i^{\text{LOO}}\},$$

which implies that  $Y_{n+1} < \widehat{\mu}_{-j}(X_{n+1}) - R_j^{\text{LOO}}$  for at least  $(1-\alpha)(n+1)$  many indices  $j \in \{1, \dots, n\}$ . In either case, then we have

$$\begin{aligned} (1-\alpha)(n+1) &\leq \sum_{j=1}^n \mathbb{1}\{Y_{n+1} \notin \widehat{\mu}_{-j}(X_{n+1}) \pm R_j^{\text{LOO}}\} \\ &= \sum_{j=1}^n \mathbb{1}\{|Y_j - \widehat{\mu}_{-j}(X_j)| < |Y_{n+1} - \widehat{\mu}_{-j}(X_{n+1})|\} \\ &= \sum_{j=1}^{n+1} \mathbb{1}\{R_{j,n+1} < R_{n+1,j}\} = \sum_{j=1}^{n+1} A_{n+1,j}, \end{aligned}$$

and, therefore,  $n+1 \in \mathcal{S}(A)$ , that is, point  $n+1$  is strange. Combining this with the result of Step 2, we have

$$\mathbb{P}\{Y_{n+1} \notin \widehat{C}_{n,\alpha}^{\text{jackknife}^+}(X_{n+1})\} \leq \mathbb{P}\{n+1 \in \mathcal{S}(A)\} \leq 2\alpha.$$

**7. Empirical results.** In this section, we compare seven methods—naive (1.4), jackknife (1.7), jackknife+ (2.1), jackknife-minmax (2.2), CV+ (3.1), split conformal (1.3), and full conformal (3.5)—on simulated and real data. Code for reproducing all results and figures is available online.<sup>5</sup>

<sup>5</sup><http://www.stat.uchicago.edu/~rina/jackknife.html>

7.1. *Simulations.* We first examine the performance of the various prediction intervals on a simulated example, using least squares as our regression method. We will see that when the training sample size  $n$  is equal or approximately equal to the dimension  $d$ , the instability of the least squares method (due to poor conditioning of the  $n \times d$  design matrix) leads to a wide disparity in performance between the various methods. This simulation is thus designed to demonstrate the role of stability in the performance of these various methods.

7.1.1. *Data and methods.* Our target coverage level is  $1 - \alpha = 0.9$ . We use training sample size  $n = 100$ , and repeat the experiment at each dimension  $d = 5, 10, \dots, 200$ , with i.i.d. data points  $(X_i, Y_i)$  generated as

$$X_i \sim \mathcal{N}(0, I_d) \quad \text{and} \quad Y_i | X_i \sim \mathcal{N}(X_i^\top \beta, 1).$$

The true coefficient vector  $\beta$  is drawn as  $\beta = \sqrt{10} \cdot u$  for a uniform random unit vector  $u \in \mathbb{R}^d$ . The regression method  $\mathcal{A}$  is simply least squares, with the convention that if the linear system is underdetermined then we take the solution with the lowest  $\ell_2$  norm (the limit of ridge regression as the regularization tends to zero). Specifically, for training data  $(X_1, Y_1), \dots, (X_n, Y_n)$ , we return the regression function  $\hat{\mu}(x) = x^\top \hat{\beta}$ ,

$$\hat{\beta} = X_{\text{mat}}^\dagger Y_{\text{vec}},$$

where  $X_{\text{mat}}$  denotes the  $n \times d$  matrix of covariates,  $Y_{\text{vec}}$  the vector of responses, and  $\dagger$  the Moore–Penrose pseudoinverse.

We then generate 100 test data points from the same distribution, and calculate the empirical probability of coverage (i.e., the proportion of test points for which the prediction interval computed at the  $X$  value contains the  $Y$  value) and the average width of the prediction interval.

7.1.2. *Results.* Figure 2 displays the results of the simulation, averaged over 50 trials (where each trial has an independent draw of the training sample of  $n = 100$  and the test sample of size 100).

When  $d < n$ , the jackknife and jackknife+ show very similar performance, with approximately the right coverage level  $1 - \alpha = 0.9$  and with nearly identical interval width. For  $d \approx n$  (the regime where least squares is quite unstable), the jackknife has substantial undercoverage—at  $d = n$  the jackknife shows coverage rate around 0.5, and continues to show substantial undercoverage when  $d$  is slightly larger than  $n$ . In this regime, the jackknife+ continues to show the right coverage level, at the cost of a prediction interval that is only slightly wider than the jackknife. For large  $d$ , the jackknife and jackknife+ again show very similar performance. In fact, this connects to recent work on interpolation methods (methods that achieve zero training error). Specifically, Hastie et al. [11] study “ridgeless” regression (i.e., the least squares solution with the lowest  $\ell_2$  norm, as in our simulation), and demonstrate that this provides a stable solution with good test error as long as  $d$  is either sufficiently small or sufficiently large relative to  $n$ . We see a similar phenomenon in the predictive coverage performance of the jackknife.

As expected, the jackknife-minmax is over-conservative, with typical coverage higher than  $1 - \alpha = 0.9$  across all dimensions  $d$ , while the naive method drastically undercovers due to increasing overfitting as  $d$  grows (and in fact, at  $d \geq n$ , the training error is exactly zero, so the prediction intervals have width zero and coverage zero.)

When  $d > n$ , we note that full conformal prediction will always have infinite length intervals since for every potential  $y$  in the  $(X_{n+1}, y)$  pair, all  $n + 1$  residuals will equal zero. Naturally, in such a situation, full conformal will have coverage equal to one deterministically. In practice, it is common to modify the conformal prediction method by truncating to

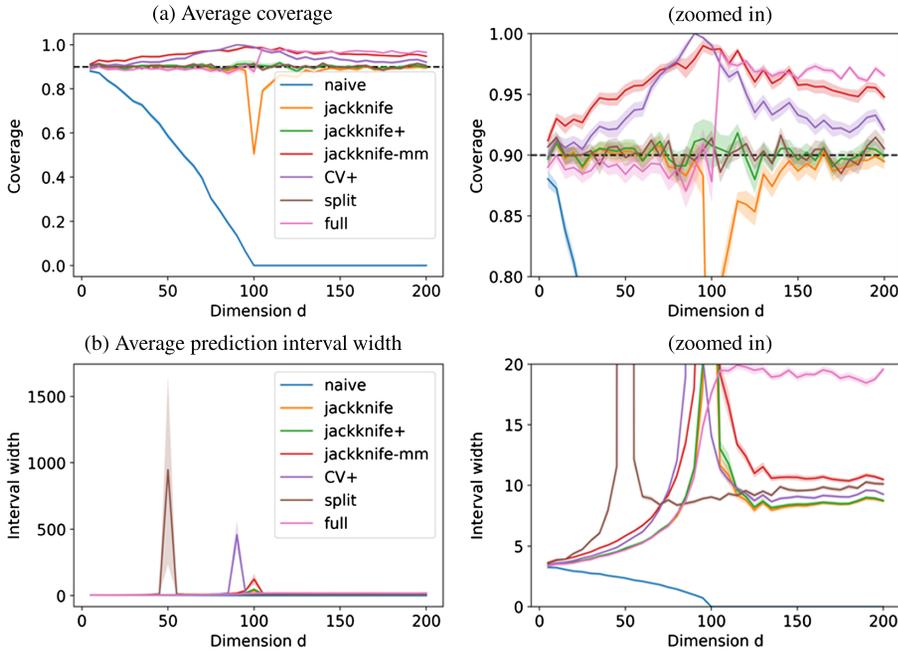


FIG. 2. Simulation results, showing the coverage and width of the predictive interval for all methods. The solid lines show the mean over 50 independent trials, with shading to show  $\pm$  one standard error. We observe that the jackknife undercovers around  $d = 100$  due to instability (since  $n = 100$ ). Jackknife+ and split conformal are the only two methods that maintain the correct coverage level throughout without under- or over-covering, but we can observe that jackknife+ often produces shorter intervals than split conformal. (See text for more details.)

a finite range, for example, to the observed range of  $Y$  values in the training data (which has minimal effect on the coverage guarantee [6]); this is why we see finite length intervals for full conformal in our simulation results.

Split conformal is the only method other than jackknife+ to maintain coverage at 0.9 throughout. (Note that since split conformal trains on half of the data, its length spikes near  $d = 50$ , rather than  $d = 100$  as for the other methods; this is simply due to the change in sample size  $n/2 = 50$  used in training. This is a result of instability of OLS when  $n \approx d$  and is not reflective of comparisons between holdout and jackknife+.)

**7.2. Real data.** We next compare the various methods on three real data sets. We will try three regression algorithms: ridge regression, random forests and neural networks (details given below). Our aim in these experiments is to demonstrate the typical performance of the various prediction interval methods in a real data setting; we do not seek to optimize the base methods used as our regression algorithms, but are only interested in how the various prediction interval methods behave in comparison to each other. Due to the high computational cost of the full conformal method, we do not include it in the comparison.

**7.2.1. Data.** The Communities and Crime data set<sup>6</sup> [20] contains information on 1994 communities, with covariates such as median income, distribution of ages, family size, etc., and the goal of predicting a response variable defined as the per capita violent crime rate. After removing categorical variables and variables with missing data,  $d = 99$  covariates remain.

The BlogFeedback data set<sup>7</sup> [5] contains 52,397 data points, each corresponding to a single blog post. The goal is to predict the response variable of the number of comments left on the

<sup>6</sup><http://archive.ics.uci.edu/ml/datasets/communities+and+crime>

<sup>7</sup><https://archive.ics.uci.edu/ml/datasets/BlogFeedback>

blog post in the following 24 hours, using  $d = 280$  covariates such as the length of the post, the number of comments on previous posts, etc. Since the distribution of the response is extremely skewed, we transform it as  $Y = \log(1 + \# \text{ comments})$ .

The Medical Expenditure Panel Survey 2016 data set,<sup>8</sup> provided by the Agency for Healthcare Research and Quality, contains data on individuals' utilization of medical services such as visits to the doctor, hospital stays, etc. Details on the data collection for older versions of this data set are described in Trena et al. [25]. We select a subset of relevant features, such as age, race/ethnicity, family income, occupation type, etc. After splitting categorical features into dummy variables to encode each category separately, the resulting dimension is  $d = 107$ . The goal is to predict the health care system utilization of each individual, which is a composite score reflecting the number of visits to a doctor's office, hospital visits, days in nursing home care, etc. With missing data removed, this data set contains 33,005 data points. Since the distribution of the response is highly skewed, we transform it as  $Y = \log(1 + (\text{utilization score}))$ .

**7.2.2. Methods.** Our procedure is the same for each of the three data sets. We randomly sample  $n = 200$  data points from the full data set, to use as the training data. The remaining points form the test set.

We run our experiment using three different regression algorithms  $\mathcal{A}$ —namely, ridge regression, random forests and neural networks. The details of these algorithms are as follows:

- For ridge regression, we define  $\hat{\mu}(x) = \hat{\beta}_0 + x^\top \hat{\beta}$  for

$$\hat{\beta}_0, \hat{\beta} = \arg \min_{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^d} \left\{ \frac{1}{2} \sum_{i=1}^n (Y_i - \beta_0 - X_i^\top \beta)^2 + \lambda \|\beta\|_2^2 \right\},$$

where the penalty parameter is chosen as  $\lambda = 0.001 \|X_{\text{mat}}\|^2$ , where  $X_{\text{mat}} \in \mathbb{R}^{n \times d}$  is the covariate matrix of the training data, and  $\|X_{\text{mat}}\|$  is its spectral norm. This choice is to accommodate situations in which the matrix  $X_{\text{mat}}$  does not have full column rank as in the case where  $d > n$ . In such cases, the solution above is nearly the least-squares solution with minimum  $\ell_2$  norm.

- For random forests, we use the `RandomForestRegressor` method from the `scikit-learn` package [17] in Python, with 20 trees grown for each random forest using the mean absolute error criterion, and with default settings otherwise.
- For neural networks, we use the `MLPRegressor` method also from `scikit-learn`, run with the L-BFGS solver and the logistic activation function, and with default settings otherwise.

For each choice of  $\mathcal{A}$ , we construct six prediction intervals (naive, jackknife, jackknife+, jackknife-minmax, CV+, split conformal), and calculate their empirical coverage rate and their average width on the test set. We then repeat this procedure 20 times, with the train/test split formed randomly each time, and report the mean and standard error over these 20 trials.

**7.2.3. Results.** Figure 3 displays the results of the real data experiments. For each data set, each regression algorithm, and each one of the six prediction interval methods, the figure plots the average coverage and average width, together with their standard errors across the 20 independent trials.

We see that the jackknife and jackknife+ methods both yield empirical coverage extremely close to the target level of 90%, and have very similar predictive interval widths. However,

<sup>8</sup>[https://meps.ahrq.gov/mepsweb/data\\_stats/download\\_data\\_files\\_detail.jsp?cboPufNumber=HC-192](https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-192)

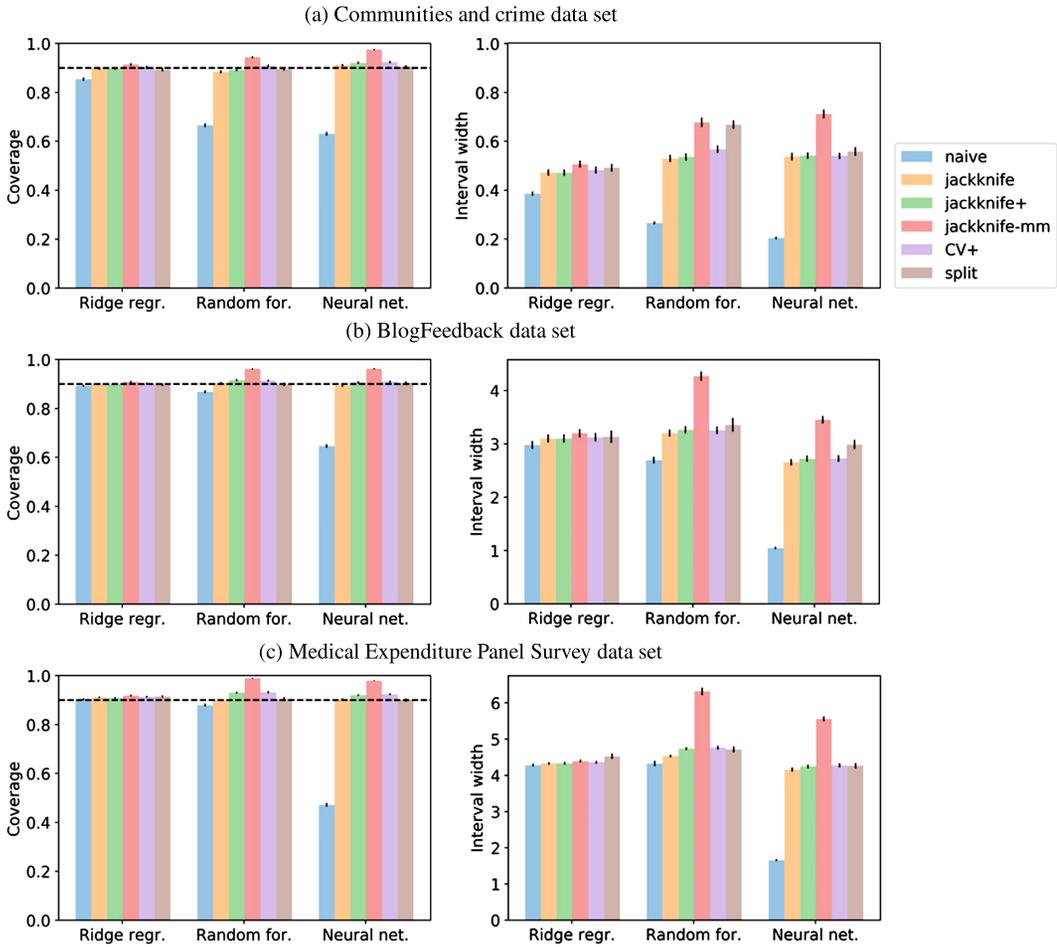


FIG. 3. Results on three real data sets, using either ridge regression, random forests, or neural networks as the regression algorithm. The bar plots show the coverage and the width of the predictive interval for all methods. The figures display the mean over 20 independent trials (i.e., splits into training and test data), with error bars to show  $\pm$  one standard error. In general, the naive method undercovers while jackknife-minmax overcovers, and the remaining methods have well calibrated coverage. In terms of their interval lengths, we typically (but not necessarily) get the expected order: jackknife < jackknife+ < 10-fold CV+ < split conformal.

in some settings, the jackknife+ shows slightly higher coverage than jackknife, and slightly wider prediction intervals. These settings correspond to regression methods with greater instability. As expected, the naive method undercovers in some settings and the jackknife-minmax is generally overly conservative. Split conformal performs reasonably well: its length and coverage is sometimes comparable to the jackknife+, but is also significantly wider in some instances. Intuitively, if the best regression function in the considered function class is simple and the dataset is large, split conformal should perform fine even though it uses  $n/2$  points for training and  $n/2$  for calibration; however, in settings where the dataset is small relative to the complexity of the best regressor, then we should observe significant gains in using  $n - 1$  points for training and  $n$  points for calibration. One phenomenon that is not visible in the empirical results is that split conformal is a randomized method, with output varying slightly depending on the random split, while jackknife+ is a deterministic method on any fixed training data set.

**8. Summary.** The jackknife+ differs from the jackknife in that it uses the quantiles of

$$\widehat{\mu}_{-i}(X_{n+1}) \pm R_i^{\text{LOO}} = \widehat{\mu}(X_{n+1}) + (\widehat{\mu}_{-i}(X_{n+1}) - \widehat{\mu}(X_{n+1})) \pm R_i^{\text{LOO}},$$

instead of those of  $\widehat{\mu}(X_{n+1}) \pm R_i^{\text{LOO}}$ , to build predictive intervals. By applying the shifts  $\widehat{\mu}_{-i}(X_{n+1}) - \widehat{\mu}(X_{n+1})$ , the jackknife+ effectively accounts for the (possible) algorithm instability, yielding rigorous coverage guarantees under no assumptions other than exchangeable samples. This, together with its empirical performance on real data, makes it a better choice than the jackknife in practice. In cases where the jackknife+ is computationally prohibitive,  $K$ -fold CV+ offers an attractive alternative. Here, it would be interesting to see if the coverage guarantees for the latter method can be somewhat sharpened.

APPENDIX A: ASYMMETRIC JACKKNIFE+ AND CV+

In settings where the distribution of  $Y$  given  $X$  appears to have symmetric noise, it is natural to construct predictions intervals symmetrically, which is why we can consider absolute values of residuals. If the data is likely to be skewed, however, we may want to consider an asymmetric construction. Fix any  $\alpha_+, \alpha_- > 0$  with  $\alpha_+ + \alpha_- = \alpha$ , and let

$$(A.1) \quad \widehat{C}_{n,\alpha_{\pm}}^{\text{jackknife}+}(X_{n+1}) = [\widehat{q}_{n,\alpha_-}^- \{ \widehat{\mu}_{-i}(X_{n+1}) + R_i^{\text{sgn,LOO}} \}, \widehat{q}_{n,\alpha_+}^+ \{ \widehat{\mu}_{-i}(X_{n+1}) + R_i^{\text{sgn,LOO}} \}],$$

where the signed residuals are

$$R_i^{\text{sgn,LOO}} = Y_i - \widehat{\mu}_{-i}(X_i).$$

We can of course consider the analogous asymmetric version of the original jackknife,

$$(A.2) \quad \widehat{C}_{n,\alpha_{\pm}}^{\text{jackknife}}(X_{n+1}) = [\widehat{\mu}(X_{n+1}) + \widehat{q}_{n,\alpha_-}^- \{ R_i^{\text{sgn,LOO}} \}, \widehat{\mu}(X_{n+1}) + \widehat{q}_{n,\alpha_+}^+ \{ R_i^{\text{sgn,LOO}} \}].$$

This type of asymmetric jackknife was considered by Steinberger and Leeb [22]. Similarly, we can define an asymmetric version of jackknife-minmax or of CV+. We remark that, even if we were to choose  $\alpha_- = \alpha_+ = \alpha/2$ , these asymmetric constructions would not necessarily be equal to the original jackknife, jackknife+, jackknife-minmax and CV+ intervals, because the empirical distribution of the signed residuals will in general be asymmetric even if only due to random chance.

All of the coverage guarantees that we have proved for the various symmetric methods, hold also for their asymmetric counterparts. For example, to verify  $1 - 2\alpha$  coverage for the asymmetric jackknife+ in the assumption-free setting, the proof of Theorem 1 proceeds identically except that the matrix of residuals  $R \in \mathbb{R}^{(n+1) \times (n+1)}$  constructed in the proof is replaced with two matrices

$$(R_{\pm})_{ij} = \begin{cases} +\infty & i = j, \\ \pm(Y_i - \widetilde{\mu}_{-(i,j)}(X_i)) & i \neq j, \end{cases}$$

where  $R_+$  (resp.,  $R_-$ ) is used to bound the probability of noncoverage in the right (resp., left) tail by  $\alpha_+$  (resp.,  $\alpha_-$ ).

**Acknowledgments.** The authors are grateful to the American Institute of Mathematics for supporting and hosting our collaboration. Rina Foygel Barber was partially supported by the National Science Foundation via grant DMS-1654076 and by an Alfred P. Sloan fellowship. Emmanuel J. Candès was partially supported by the Office of Naval Research under grant N00014-16-1-2712, by the National Science Foundation via grant DMS-1712800, and by a generous gift from TwoSigma.

The authors are grateful to an anonymous reviewer for helpful suggestions on the presentation of the proof of Theorem 1. Emmanuel J. Candès thanks Yaniv Romano for help with some experiments. Aaditya Ramdas thanks Arun Kumar Kuchibhotla for discussions regarding cross-conformal prediction.

## SUPPLEMENTARY MATERIAL

**Supplement to “Predictive inference with the jackknife+”** (DOI: [10.1214/20-AOS1965SUPP](https://doi.org/10.1214/20-AOS1965SUPP); .pdf). In the Supplementary Material, we provide proofs for theoretical results.

## REFERENCES

- [1] BARBER, R. F., CANDÈS, E. J., RAMDAS, A. and TIBSHIRANI, R. J. (2021). Supplement to “Predictive inference with the jackknife+.” <https://doi.org/10.1214/20-AOS1965SUPP>
- [2] BOUSQUET, O. and ELISSEFF, A. (2002). Stability and generalization. *J. Mach. Learn. Res.* **2** 499–526. [MR1929416 https://doi.org/10.1162/153244302760200704](https://doi.org/10.1162/153244302760200704)
- [3] BURNAEV, E. and VOVK, V. (2014). Efficiency of conformalized ridge regression. In *Conference on Learning Theory* 605–622.
- [4] BUTLER, R. and ROTHMAN, E. D. (1980). Predictive intervals based on reuse of the sample. *J. Amer. Statist. Assoc.* **75** 881–889. [MR0600971](https://doi.org/10.1080/01621459.1980.1056032)
- [5] BUZA, K. (2014). Feedback prediction for blogs. In *Data Analysis, Machine Learning and Knowledge Discovery* 145–152. Springer, Berlin.
- [6] CHEN, W., CHUN, K.-J. and BARBER, R. F. (2018). Discretized conformal prediction for efficient distribution-free inference. *Stat* **7** e173. [MR3769053 https://doi.org/10.1002/sta.4.173](https://doi.org/10.1002/sta.4.173)
- [7] DEVROYE, L. P. and WAGNER, T. J. (1979). Distribution-free inequalities for the deleted and holdout error estimates. *IEEE Trans. Inf. Theory* **25** 202–207. [MR0521311 https://doi.org/10.1109/TIT.1979.1056032](https://doi.org/10.1109/TIT.1979.1056032)
- [8] EFRON, B. (1979). Bootstrap methods: Another look at the jackknife. *Ann. Statist.* **7** 1–26. [MR0515681](https://doi.org/10.2307/2685844)
- [9] EFRON, B. and GONG, G. (1983). A leisurely look at the bootstrap, the jackknife, and cross-validation. *Amer. Statist.* **37** 36–48. [MR0694281 https://doi.org/10.2307/2685844](https://doi.org/10.2307/2685844)
- [10] GEISSER, S. (1975). The predictive sample reuse method with applications. *J. Amer. Statist. Assoc.* **70** 320–328.
- [11] HASTIE, T., MONTANARI, A., ROSSET, S. and TIBSHIRANI, R. J. (2019). Surprises in high-dimensional ridgeless least squares interpolation. Preprint. Available at [arXiv:1903.08560](https://arxiv.org/abs/1903.08560).
- [12] LANDAU, H. G. (1953). On dominance relations and the structure of animal societies. III. The condition for a score structure. *Bull. Math. Biophys.* **15** 143–148. [MR0054933 https://doi.org/10.1007/bf02476378](https://doi.org/10.1007/bf02476378)
- [13] LEI, J. (2019). Fast exact conformalization of the lasso using piecewise linear homotopy. *Biometrika* **106** 749–764. [MR4031197 https://doi.org/10.1093/biomet/asz046](https://doi.org/10.1093/biomet/asz046)
- [14] LEI, J., G’SELL, M., RINALDO, A., TIBSHIRANI, R. J. and WASSERMAN, L. (2018). Distribution-free predictive inference for regression. *J. Amer. Statist. Assoc.* **113** 1094–1111. [MR3862342 https://doi.org/10.1080/01621459.2017.1307116](https://doi.org/10.1080/01621459.2017.1307116)
- [15] MILLER, R. G. (1974). The jackknife—a review. *Biometrika* **61** 1–15. [MR0391366 https://doi.org/10.1093/biomet/61.1.1](https://doi.org/10.1093/biomet/61.1.1)
- [16] PAPAPOPOULOS, H. (2008). Inductive conformal prediction: Theory and application to neural networks. In *Tools in Artificial Intelligence* InTech, Vienna.
- [17] PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A. et al. (2011). Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12** 2825–2830. [MR2854348](https://doi.org/10.26434/chemrxiv-2012-06)
- [18] QUENOUILLE, M. H. (1949). Approximate tests of correlation in time-series. *J. Roy. Statist. Soc. Ser. B* **11** 68–84. [MR0032176](https://doi.org/10.2307/2343176)
- [19] QUENOUILLE, M. H. (1956). Notes on bias in estimation. *Biometrika* **43** 353–360. [MR0081040 https://doi.org/10.1093/biomet/43.3-4.353](https://doi.org/10.1093/biomet/43.3-4.353)
- [20] REDMOND, M. and BAVEJA, A. (2002). A data-driven software tool for enabling cooperative information sharing among police departments. *European J. Oper. Res.* **141** 660–678.
- [21] STEINBERGER, L. and LEEB, H. (2016). Leave-one-out prediction intervals in linear regression models with many variables. Preprint. Available at [arXiv:1602.05801](https://arxiv.org/abs/1602.05801).
- [22] STEINBERGER, L. and LEEB, H. (2018). Conditional predictive inference for high-dimensional stable algorithms. Preprint. Available at [arXiv:1809.01412](https://arxiv.org/abs/1809.01412).
- [23] STINE, R. A. (1985). Bootstrap prediction intervals for regression. *J. Amer. Statist. Assoc.* **80** 1026–1031. [MR0819610](https://doi.org/10.2307/2343176)
- [24] STONE, M. (1974). Cross-validated choice and assessment of statistical predictions. *J. Roy. Statist. Soc. Ser. B* **36** 111–147. [MR0356377](https://doi.org/10.2307/2343176)
- [25] TRENA, M., EZZATI-RICE, ROHDE, F. and GREENBLATT, J. (2008). Sample design of the medical expenditure panel survey household component, 1998–2007. MEPS Methodology Report No. 22, Agency for Healthcare Research and Quality.

- [26] TUKEY, J. (1958). Bias and confidence in not quite large samples. *Ann. Math. Stat.* **29** 614.
- [27] VOVK, V. (2012). Conditional validity of inductive conformal predictors. In *Asian Conference on Machine Learning* 475–490.
- [28] VOVK, V. (2015). Cross-conformal predictors. *Ann. Math. Artif. Intell.* **74** 9–28. [MR3353894](#)  
<https://doi.org/10.1007/s10472-013-9368-4>
- [29] VOVK, V., GAMMERMAN, A. and SHAFER, G. (2005). *Algorithmic Learning in a Random World*. Springer, New York. [MR2161220](#)
- [30] VOVK, V., NOURETDINOV, I., MANOKHIN, V. and GAMMERMAN, A. (2018). Cross-conformal predictive distributions. In *Conformal and Probabilistic Prediction and Applications* 37–51.
- [31] VOVK, V. and WANG, R. (2012). Combining p-values via averaging. Preprint. Available at [arXiv:1212.4966](#).
- [32] XU, H., CARAMANIS, C. and MANNOR, S. (2012). Sparse algorithms are not stable: A no-free-lunch theorem. *IEEE Trans. Pattern Anal. Mach. Intell.* **34** 187–193.