

Constructing Prediction Intervals with Neural Networks: An Empirical Evaluation of Bootstrapping and Conformal Inference Methods

Alex Contarino

*Department of Mathematical Sciences
United States Air Force Academy*

ALEXANDER.CONTARINO@AFACADEMY.AF.EDU

**Christine Schubert Kabban
Chancellor Johnstone**

*Department of Mathematics and Statistics
Air Force Institute of Technology
Wright-Patterson AFB, OH 45433, USA*

CHRISTINE.SCHUBERTKABBAN@AFIT.EDU
CHANCELLOR.JOHNSTONE@AFIT.EDU

Fairul Mohd-Zaid

*711th Human Performance Wing
Machine Learning Section, Mission Analytics
Wright-Patterson AFB, OH 45433, USA*

FAIRUL.MOHD-ZAID@US.AF.MIL

Editor:

Abstract

Artificial neural networks (ANNs) are popular tools for accomplishing many machine learning tasks, including predicting continuous outcomes. However, the general lack of confidence measures provided with ANN predictions limit their applicability. Supplementing point predictions with prediction intervals (PIs) is common for other learning algorithms, but the complex structure and training of ANNs renders constructing PIs difficult. This work provides the network design choices and inferential methods for creating better performing PIs with ANNs. A two-step experiment is executed across 11 data sets, including an imaged-based data set. Two distribution-free methods for constructing PIs, bootstrapping and conformal inference, are considered. The results of the first experimental step reveal that the choices inherent to building an ANN affect PI performance. Guidance is provided for optimizing PI performance with respect to each network feature and PI method. In the second step, 20 algorithms for constructing PIs—each using the principles of bootstrapping or conformal inference—are implemented to determine which provides the best performance while maintaining reasonable computational burden. In general, this trade-off is optimized when implementing the cross-conformal method, which maintained interval coverage and efficiency with decreased computational burden.

Keywords: neural networks, prediction, inference, bootstrapping, conformal inference

The views expressed in this document are those of the authors and should not be interpreted as representing the official policy or position of the U.S. Air Force, the U.S. Department of Defense, or the U.S. Government.

1. Introduction

Neural networks, specifically multi-layer perceptrons (MLPs) and convolutional neural networks (CNNs), are popular tools used for many machine learning tasks. Comprised of connected nodes, neural networks aim to create a flexible means to mathematically learn representations of data. Within each node, vector-to-scalar and activation operations are performed. Each activation function $\phi(\cdot)$, e.g., the sigmoid (“logistic”) function, serves to map a scalar value to a desired range. The collection of nodes acting in parallel represent an automated method for mapping a set of features \mathbf{X} to a target outcome \mathbf{y} (Chollet, 2017). Neural networks have proved useful for a variety of tasks; one specific setting is natural language processing where manually engineering a mapping between \mathbf{X} and \mathbf{y} is difficult (Goodfellow et al., 2017). Other regression-based tasks in which neural networks have provided state-of-the-art performance include age prediction (Rothe et al., 2018), orientation estimation (Berg et al., 2020) and microscopy image processing (Xie et al., 2018).

MLPs and CNNs are “feed-forward” neural networks, with nodes arranged hierarchically in successive hidden layers (Goodfellow et al., 2017). The key difference between MLPs and CNNs is the choice of the vector-to-scalar operator $\Sigma(\cdot)$, with the former using vector multiplication. In contrast, CNNs are so-called because a convolution operation is used in the nodes of one or more hidden layers. Convolution effectively induces sparse connectivity among the nodes (generally called “filters” or “feature maps”), with favorable results when \mathbf{X} has high dimensionality, e.g., a set of images.

In regression, \mathbf{y} is continuous-valued variable, which can pose challenges when using neural networks for inference. For instance, suppose that an MLP is trained on \mathbf{X} and \mathbf{y} under mean squared error (MSE) loss, and a new instance of \mathbf{X} , \mathbf{x}_{test} , is observed. Suppose further that we wish to predict the value of the associated, unobserved target value, y_{test} . The trained network will provide a prediction for y_{test} given \mathbf{x}_{test} ; denote this value as $\hat{f}(\mathbf{x}_{test})$. However, with \mathbf{y} being real-valued, the probability that $\hat{f}(\mathbf{x}_{test}) = y_{test}$ is zero. Thus, associating $\hat{f}(\mathbf{x}_{test})$ with some measure of uncertainty is useful for enhancing the usability of a trained network. A prediction interval (PI) is one such measure, placing probabilistic bounds on the value of y_{test} (Casella and Berger, 2002). A $1 - \alpha$ PI is a set of target values $\Gamma \equiv \Gamma(\mathbf{x}_{test})$ such that

$$P(y_{test} \in \Gamma) = 1 - \alpha. \tag{1}$$

Prediction intervals must account for two sources of error inherent to training a regression model, including model-fitting error (or “prediction variance”) and irreducible error (or “data noise”) (Gareth et al., 2013). The performance of a prediction interval is generally evaluated according to two metrics: validity and efficiency (Fraser and Guttman, 1956). Validity is a function of the coverage of the PI, or the rate at which the PI correctly captures the estimated target value, and whether or not this rate matches the nominal confidence level $1 - \alpha$. The second metric of performance is efficiency which measures the width of the PI. Narrow prediction intervals give more information about y_{test} and are preferred, though not at the expense of coverage. Thus, the optimal PI for an estimated target y_{test} is one having the smallest average width while still satisfying (1).

Due to their complex structures and iterative training process, constructing PIs for neural networks is generally a difficult task. Factors such as performance and computational

time must be considered. The design choices inherent to constructing a neural network, such as the number of layers and nodes, have sizeable impacts on the network’s predictive ability, yet little is known regarding how such factors impact PI performance. Moreover, each method for constructing a PI has its own advantages and disadvantages. For instance, parametric approaches, such as maximum likelihood and approximate Bayes, estimate uncertainty by computing covariances among the network parameters (Papadopoulos et al., 2001). To do so, gradient matrices are computed at each training iteration, both complicating the optimization of the network and introducing a large computational burden (Khosravi et al., 2015).

In contrast to the uncertainty quantification methods mentioned above, methods leveraging bootstrap resampling are more easily implemented and reduce dependency on assumptions related to the distributions of parameter estimates. In a bootstrap method, an ensemble of B neural networks is trained, each to a different bootstrap resample of the data tuple (\mathbf{X}, \mathbf{y}) . For a given instance \mathbf{x}_i denote $\hat{f}_b(\mathbf{x}_i)$ as the predicted value of the target y_i from the network trained on the b^{th} bootstrap resample. The empirical distribution of values $\hat{f}_b(\mathbf{x}_i)$, $b = 1, \dots, B$, is then used to calculate the statistics needed to construct a PI (Efron and Tibshirani, 1993). As a flexible approach, both in assumptions and practical implementation, bootstrapping has been used extensively with neural networks (Papadopoulos et al. 2001; Khosravi et al. 2015; Anirudh and Thiagarajan 2017). However, the time and resources required to train and store modern MLPs and CNNs make bootstrapping methods increasingly difficult to implement. Moreover, to our knowledge, no experiment has examined how PI performance changes as neural network hyperparameters are tuned.

Conformal inference methods offer an intriguing alternative to bootstrapping. These approaches have many of the same advantages of bootstrapping but, for the most part, are less computationally expensive. Conformal inference determines whether a candidate value q should belong in Γ by evaluating how well the data pair (\mathbf{x}_{test}, q) conforms to the prior knowledge derived from \mathbf{X} and \mathbf{y} (Gammerman et al., 1998). While there has been research evaluating the efficacy of conformal inference approaches in deep learning settings (Papadopoulos 2008; Angelopoulos et al. 2020), these experiments have not, for example, compared the various conformal inference methods on differing data sets.

1.1 Contribution and Document Overview

The contribution of our research is two-fold. First, we show how the choice of hyperparameter settings in a MLP or CNN can affect the performance of PIs in terms of validity and efficiency using bootstrap and conformal inference methods. We specifically design MLPs and CNNs with varying hyperparameter settings, with each network trained to a common training set. PIs are then constructed for responses associated with pre-sequestered test observations. Differences in PI performance among the different hyperparameter settings are then identified using an analysis of variance (ANOVA) approach. These results are used to quantify the sensitivities of the bootstrap and conformal inference PI methods to various hyperparameter settings in a neural network.

The second contribution is a thorough evaluation of bootstrap and conformal inference methods to determine which provide the best performing PIs, while maintaining a reasonable computational burden. Within this analysis, we also conduct a cursory analysis of

conditional coverage. We use the findings of these analyses to provide guidance regarding the optimal PI method given varying time and resource constraints. Experiments are executed across eleven commonly-used data sets for regression tasks, including the image-based Rotated MNIST (“RotNIST”) data set.

The remainder of this document is organized as follows. Section 2 provides comprehensive overview of the methods examined in this research. Section 3 describes the two-step experiment conducted in order to evaluate PI method and performance. In Section 4, the empirical results of the two-step experiment are presented and analyzed. Section 5 concludes the paper with key findings.

2. Description of Methods

The following section describes the prediction interval methods examined in this research. These approaches fall into two classes: bootstrapping and conformal inference. These methods are discussed in terms of their theoretical underpinnings, practical implementation, as well as their advantages and shortcomings in various settings.

2.1 Bootstrapping

Bootstrapping is a resampling-based method for inference, used in computing measures of accuracy for statistics, such as bias and variance (Gentle, 2009). Specifically, suppose a random sample of size n , $\mathbf{z} = \{z_1, z_2, \dots, z_n\}$, is drawn from an unknown distribution and it is desired to use \mathbf{z} to calculate some statistic of interest, $T = g(\mathbf{z})$. Define \hat{F} as the empirical distribution of the observed data, where each z_i in \mathbf{z} occurs with probability $1/n$. Then, a “bootstrap resample” is a random sample of size n drawn (with replacement) from \hat{F} (Efron and Tibshirani, 1993). Each of the B resamples collected are generally denoted as \mathbf{z}_b^* , $b = 1, \dots, B$. The sampling distribution of T is then estimated by the empirical distribution of $T_b = g(\mathbf{z}_b)$, $b = 1, \dots, B$, calculated on each bootstrapped resample \mathbf{z}_b^* .

Bootstrapping is regularly applied in the construction of prediction intervals, with several different strategies available for constructing PIs. The two bootstrapping variations examined in this research are the pivot and percentile methods.

2.1.1 PIVOT BOOTSTRAP METHOD

A popular bootstrapping approach for constructing PIs is a pivotal quantity method where the distribution of the unknown target y_{test} , conditioned on the observed instance \mathbf{x}_{test} , is assumed to be approximately normal. In this approach, the half-width of the PI is a function of the estimates for the model-fitting and irreducible errors, with the interval centered at the bootstrap aggregated prediction for y_{test} . This prediction, denoted as $\hat{f}(\mathbf{x}_{test})$, is found through the simple averaging of each resampled model:

$$\hat{f}(\mathbf{x}_{test}) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(\mathbf{x}_{test}). \quad (2)$$

The pivot bootstrap PI, Γ_{pivot} , is the interval

$$\Gamma_{pivot} = \left[\hat{f}(\mathbf{x}_{test}) - z_{1-\alpha/2} \sqrt{\hat{\sigma}^2(\mathbf{x}_{test}) + \hat{\sigma}_\epsilon^2}, \hat{f}(\mathbf{x}_{test}) + z_{1-\alpha/2} \sqrt{\hat{\sigma}^2(\mathbf{x}_{test}) + \hat{\sigma}_\epsilon^2} \right], \quad (3)$$

where $\hat{\sigma}^2(\mathbf{x}_{test})$ is the estimated prediction variance, $\hat{\sigma}_\epsilon^2$ is the estimated irreducible error, and $z_{1-\alpha/2}$ is the $(1 - \alpha/2)$ -th percentile of the standard normal distribution. In general, using the corresponding percentile from the Student’s t distribution is a more mathematically rigorous method to capture the variability of the estimated target value. However, determining the degrees of freedom when modeling with a neural network is not a straightforward process (Gao and Jojic, 2016). The standard normal distribution is used here for simplicity.

Suppose that the feature and target data, \mathbf{X} and \mathbf{y} , are observed. B bootstrap resamples are collected, and a neural network is trained on each to produce an estimated regression function \hat{f}_b , $b = 1, \dots, B$. Now, suppose a new feature vector, \mathbf{x}_{test} , is observed and it is desired to provide a point estimate and PI for its unknown target value, y_{test} . After calculating $\hat{f}(\mathbf{x}_{test})$, as in (2), constructing the PI is a matter of estimating the component model-fitting and irreducible errors. The latter prediction variance estimate for \mathbf{x}_{test} , $\hat{\sigma}^2(\mathbf{x}_{test})$, can be calculated from the distribution of predicted target values calculated from the bootstrap resamples,

$$\hat{\sigma}^2(\mathbf{x}_{test}) = \frac{1}{B-1} \sum_{b=1}^B \left[\hat{f}_b(\mathbf{x}_{test}) - \hat{f}(\mathbf{x}_{test}) \right]^2. \quad (4)$$

The irreducible error $\hat{\sigma}_\epsilon^2$ can be similarly estimated from the residuals of the out-of-sample predictions. In bootstrapping, out-of-sample observations are conveniently provided in the form of the “out-of-bag” sets. For a bootstrap resample \mathbf{z}_b^* , the corresponding set of out-of-bag observations is $\{z_i \in \mathbf{z} | z_i \notin \mathbf{z}_b^*\}$ and is denoted here as $\tilde{\mathbf{z}}_b^*$. Then, for each bootstrap resample, $\hat{\sigma}_{\epsilon(b)}^2$ is:

$$\hat{\sigma}_{\epsilon(b)}^2 = \frac{\sum_{\mathbf{x}_i^{(b)} \in \tilde{\mathbf{X}}_b^*} \left[y_i^{(b)} - \hat{f}_b(\mathbf{x}_i^{(b)}) \right]^2}{|\tilde{n}_b - 1|}, \quad (5)$$

where:

- \tilde{n}_b is the number of observations in the out-of-bag set $\tilde{\mathbf{X}}_b^*$
- $\mathbf{x}_i^{(b)}$ is the i^{th} observation in $\tilde{\mathbf{X}}_b^*$
- $y_i^{(b)}$ is the i^{th} observation in the out-of-bag set $\tilde{\mathbf{y}}_b^*$ and the target value of $\mathbf{x}_i^{(b)}$

The irreducible error estimate for the entire data set then is:

$$\hat{\sigma}_\epsilon^2 = \frac{1}{B} \sum_{b=1}^B \hat{\sigma}_{\epsilon(b)}^2. \quad (6)$$

Algorithm 1 provides pseudocode for calculating each component of the pivot bootstrap PI. While the approach described in Algorithm 1 is readily applied, the performance of the PIs constructed with the pivot bootstrap method have their own limitations. In particular,

confidence intervals constructed in this manner tend to be conservative (Papadopoulos et al., 2001), indicating the desired coverage could be achieved with a smaller prediction region. Further note that the performance of the PIs is dependent upon the sampling distribution of $\hat{f}(\mathbf{x}_{test})$; if it is skewed or biased, the validity and efficiency of the resulting PIs will deteriorate.

Algorithm 1: Pivot Bootstrap Method

Input: training data \mathbf{X} and \mathbf{y} , test observation \mathbf{x}_{test} , learning algorithm L , desired number of bootstrap resamples B , and desired coverage probability $1 - \alpha$

- 1 **for** $b = 1$ to B **do**
- 2 Generate bootstrap resamples of \mathbf{X} and \mathbf{y} ; denote them as \mathbf{X}_b^* and \mathbf{y}_b^*
- 3 Find the out-of-bag sets and denote them as $\tilde{\mathbf{X}}_b^*$ and $\tilde{\mathbf{y}}_b^*$
- 4 Train learning algorithm L on \mathbf{X}_b^* and \mathbf{y}_b^* ; denote the trained regressor as \hat{f}_b
- 5 Calculate $\hat{f}_b(\mathbf{x}_{test})$ and $\hat{f}_b(\mathbf{X}_b^*)$
- 6 Calculate $\hat{\sigma}_{\epsilon(b)}^2$ as in (5)
- 7 **end**
- 8 Calculate $\hat{f}(\mathbf{x}_{test})$ as in (2)
- 9 Calculate $\hat{\sigma}^2(\mathbf{x}_{test})$ as in (4)
- 10 Calculate $\hat{\sigma}_{\epsilon}^2$ as in (6)

Output: a $1 - \alpha$ prediction interval, as constructed in (3)

2.1.2 PERCENTILE BOOTSTRAP METHOD

A $1 - \alpha$ PI constructed using the percentile bootstrap method leverages the empirical distribution of the predicted values \hat{f}_b , $b = 1, \dots, B$ and a random sampling of the error terms. This procedure weakens assumptions regarding the sampling distribution of \hat{f} and does not require the explicit computation of component error terms.

The collection of values $\hat{f}_b(\mathbf{x}_{test})$ for $b = 1, \dots, B$ yield an empirical cumulative distribution function (ECDF), \hat{F} . Then, the percentile bootstrap algorithm computes the $\alpha/2$ and $1 - \alpha/2$ percentiles from \hat{F} , $\hat{F}_{(\alpha/2)}$ and $\hat{F}_{(1-\alpha/2)}$, respectively. Note that

$$\mathbb{P} \left(\hat{F}_{(\alpha/2)} < \hat{f}(\mathbf{x}_{test}) < \hat{F}_{(1-\alpha/2)} \right) \approx 1 - \alpha. \quad (7)$$

To accurately estimate the tail percentiles from sampling distributions requires specifying B to be relatively large (i.e., $B > 1,000$) (Efron and Tibshirani, 1993).

A valid confidence interval for the value $f(\mathbf{x}_{test})$, or the expected value of y_{test} given \mathbf{x}_{test} , can be constructed by inverting (7). In contrast, a valid PI for y_{test} must also account for the irreducible error. To do so, the values of \hat{F} are adjusted with a random error, sampled from the prediction errors of an out-of-sample data set (Davidson and Hinkley, 1997). When employing bootstrap resampling, such sets come in the form of the out-of-bag sets, $\tilde{\mathbf{X}}_b^*$ and $\tilde{\mathbf{y}}_b^*$, of each b^{th} resample of training data, \mathbf{X}_b^* and \mathbf{y}_b^* , respectively. Therefore, to construct the PI, prediction errors for each fitted regressor \hat{f}_b need to be calculated:

$$\mathbf{r}_b^* = \tilde{\mathbf{y}}_b^* - \hat{f}_b(\tilde{\mathbf{X}}_b^*). \quad (8)$$

From each collection \mathbf{r}_b^* , a random sample of size one is collected—these values are denoted here as e_b^* , $b = 1, \dots, B$. Then, for each $\hat{f}_b(\mathbf{x}_{test})$, the random-error-adjusted value, $\hat{g}_b(\mathbf{x}_{test})$, is calculated:

$$\hat{g}_b(\mathbf{x}_{test}) = \hat{f}_b(\mathbf{x}_{test}) + e_b^*. \quad (9)$$

Using the set of values $\hat{g}_b(\mathbf{x}_{test})$, $b = 1, \dots, B$, an ECDF, \hat{G} , can be constructed. Then, for any desired significance level α ,

$$P\left(\hat{G}_{(\alpha/2)} < y_{test} < \hat{G}_{(1-\alpha/2)}\right) \approx 1 - \alpha, \quad (10)$$

resulting in a $1 - \alpha$ PI for y_{test} of the form

$$\Gamma_{percentile} = \left[\hat{G}_{(\alpha/2)}, \hat{G}_{(1-\alpha/2)}\right]. \quad (11)$$

Algorithm 2 provides a pseudocode implementation of the percentile bootstrap PI method (Davidson and Hinkley, 1997).

Algorithm 2: Percentile Bootstrap Method

Input: training data \mathbf{X} and \mathbf{y} , test observation \mathbf{x}_{test} , learning algorithm L , desired number of bootstrap resamples B , and desired coverage probability $1 - \alpha$

- 1 **for** $b = 1$ to B **do**
- 2 Generate bootstrap resamples of \mathbf{X} and \mathbf{y} ; denote them as \mathbf{X}_b^* and \mathbf{y}_b^*
- 3 Find the out-of-bag sets and denote them as $\tilde{\mathbf{X}}_b^*$ and $\tilde{\mathbf{y}}_b^*$
- 4 Train learning algorithm L on \mathbf{X}_b^* and \mathbf{y}_b^* ; denote the trained regressor as \hat{f}_b
- 5 Calculate $\hat{f}_b(\mathbf{x}_{test})$ and $\hat{f}_b(\tilde{\mathbf{X}}_b^*)$
- 6 Calculate the prediction errors \mathbf{r}_b^* as in (8)
- 7 From \mathbf{r}_b^* , randomly sample a value e_b^*
- 8 Calculate $\hat{g}_b(\mathbf{x}_{test})$ as in (9)
- 9 **end**
- 10 Construct the ECDF \hat{G} from the values $\hat{g}_b(\mathbf{x}_{test})$, $b = 1, \dots, B$

Output: a $1 - \alpha$ prediction interval, as constructed in (11)

2.2 Conformal Inference

Conformal inference (CI) methods provide a potentially attractive alternative to the bootstrap method for constructing PIs from neural network outputs, with many of the same advantages as bootstrapping. With CI we can avoid some of the computational cost associated with training hundreds or thousands of networks. In particular, CI methods can be easily implemented with no changes to the underlying prediction algorithm (Papadopoulos, 2008). They are also agnostic to the quantitative and distributional properties of the target variable, applicable to both continuous and discrete outcomes, and thus, can be applied to any machine learning task (Schafer and Vovk, 2008). Like bootstrapping, CI does not require weighty assumptions regarding the distribution of the target variable. The only assumption guaranteeing the validity of the PIs constructed is that the data be exchangeable (Schafer and Vovk, 2008). A sequence of n random variables is said to be exchangeable if any of $n!$ possible permutations of observing them are equally likely.

The central idea of CI is to use a measure of “nonconformity” to evaluate how “strange” a potential test observation is compared to previously observed data (Gammerman et al., 1998). Typically, a context-appropriate distance function is chosen as the nonconformity measure; the absolute residual is generally used for regression tasks.

To understand how conformal PIs are constructed, suppose a value q is being considered for potential inclusion into the set Γ , a $1 - \alpha$ PI for y_{test} . The choice of whether to include q in Γ is based upon a statistical test of the null hypothesis $y_{test} = q$ (Lei et al., 2018).

The first step of this test is to calculate “conformity score” $R(\mathbf{x}_{test}, q)$, such that

$$R(\mathbf{x}_{test}, q) = d\left(\hat{f}(\mathbf{x}_{test}), q\right), \quad (12)$$

where $\hat{f}(\mathbf{x}_{test})$ is the prediction for y_{test} and d is a measurable function mapping to \mathbb{R} . The strangeness of q is determined by comparing the conformity score associated with (\mathbf{x}_{test}, q) to a set of other conformity scores, denoted here as C . With what data these other scores are calculated varies by the different conformal inference methods and is discussed further in the succeeding subsections. However, for all methods, a test statistic for $R(\mathbf{x}_{test}, q)$ can be calculated:

$$\pi(q) = \frac{\sum_{i=1}^l \mathbb{I}[R(\mathbf{x}_i, y_i) > R(\mathbf{x}_{test}, q)]}{l + 1}, \quad (13)$$

where \mathbb{I} is the indicator function and $l = |C|$. By exchangeability, $\pi(q)$ is sub-uniform over the set $\left\{\frac{0}{l+1}, \frac{1}{l+1}, \dots, \frac{l}{l+1}\right\}$, implying:

$$\mathbb{P}(\pi(q) \leq \alpha) \leq \alpha.$$

Since $\pi(q)$ is sub-uniform, it is a valid p -value. Thus, if $\pi(q) < \alpha$ then there is sufficient evidence that y_{test} will not equal q , and thus, q should not be included in Γ .

When \mathbf{y} is continuous, the number of candidates that can be considered for this process is infinite. To practically implement CI for regression tasks, a very fine grid of target values $Q = \{q_1, q_2, \dots, q_M\}$ is examined, where Q represents a size M partition of the candidate range. A conservative $1 - \alpha$ PI is then constructed as:

$$\Gamma = \{q_j \in Q \mid \pi(q_j) \geq \alpha\} \quad (14)$$

This research evaluates the performance and computational feasibility of four conformal inference algorithms: full conformal inference, split conformal inference, cross-conformal inference and bootstrap conformal inference. Operating in the trade-space of performance and computational efficiency, each has their own advantages and drawbacks. These methods, in addition to a modification in which a fitted kernel density function is used as the nonconformity measure, are discussed in the following sections.

2.2.1 FULL CONFORMAL INFERENCE METHOD

Given training data \mathbf{X} and \mathbf{y} of size n , the full CI method provides confidence estimation for “particular to particular” new observations (Saunders et al., 1999). That is, if a new

instance \mathbf{x}_{test} is observed, the algorithm learns a PI particular to that new instance, as opposed to a general rule to use for all future instances.

To construct a $1 - \alpha$ PI for the target y_{test} , each q_j in a set of candidate values Q is evaluated against the originally observed data \mathbf{X} and \mathbf{y} . In particular, \mathbf{X} and \mathbf{y} are augmented with the data pair (\mathbf{x}_{test}, q_j) . A learning algorithm, L , is trained on this augmented data set; denote the trained regressor as \hat{f}_{aug} . For every data pair in the augmented data set, conformity scores are computed as in (12) using \hat{f}_{aug} . The set of conformity scores against which every candidate value q is compared is then:

$$C_{full} = \{R(\mathbf{x}_i, y_i) \mid (\mathbf{x}_i, y_i) \in (\mathbf{X}_{aug}, \mathbf{y}_{aug})\}, \quad (15)$$

where $(\mathbf{X}_{aug}, \mathbf{y}_{aug})$ is the augmented training set. Next, $\pi(q_j)$ is calculated using (13); q_j 's inclusion into Γ is regulated by whether or not $\pi(q_j) \geq \alpha$. Repeating over each value in Q , the $1 - \alpha$ PI is constructed as in (14). Algorithm 3 provides a pseudocode implementation of the full CI method.

Algorithm 3: Full Conformal Inference Method

Input: training data \mathbf{X} and \mathbf{y} , test observation \mathbf{x}_{test} , learning algorithm L , a set of candidate values Q , nonconformity measure d , and desired coverage probability $1 - \alpha$

- 1 **for** $q_j \in Q$ **do**
- 2 Augment \mathbf{X} and \mathbf{y} with the data pair (\mathbf{x}_{test}, q_j) and denote them as \mathbf{X}_{aug} and \mathbf{y}_{aug}
- 3 Fit learning algorithm L to the augmented training set and denote the trained regressor as \hat{f}_{aug}
- 4 Compute $R(\mathbf{x}_i, y_i)$ for $(\mathbf{x}_i, y_i) \in (\mathbf{X}_{aug}, \mathbf{y}_{aug})$ as in (12) and \hat{f}_{aug}
- 5 Build C_{full} as in (15)
- 6 Compute the p -value of $R(\mathbf{x}_{test}, q_j)$ as in (13)
- 7 **end**

Output: a $1 - \alpha$ prediction interval, as defined in (14)

While experimental results suggest this approach produces efficient and valid PIs (Linnusson et al. 2014; Lei et al. 2018), its computational burden becomes undesirable when the training process for L is time-consuming. Since, in this approach, an augmented regression function is learned for each q_j in Q , L must be trained $|Q|$ times. While under certain regularity conditions we can avoid such repeated computations (Nouretdinov et al. 2001; Ndiaye 2022), we do not explore these approaches in our work.

2.2.2 SPLIT CONFORMAL INFERENCE METHOD

The split conformal inference method alters the CI approach to be inductive. That is, given training data \mathbf{X} and \mathbf{y} , split CI creates a general rule to construct a prediction interval for any future observation.

To do so, \mathbf{X} and \mathbf{y} are randomly partitioned into two subsets: $(\mathbf{X}_{proper}, \mathbf{y}_{proper})$ and $(\mathbf{X}_{cal}, \mathbf{y}_{cal})$, the ‘‘proper training’’ and the ‘‘calibration’’ sets, respectively. A learning

algorithm L is trained on $(\mathbf{X}_{proper}, \mathbf{y}_{proper})$; denote the estimated regression function as \hat{f}_{split} . The set of conformity scores C_{split} , used to evaluate candidate values, is built as:

$$C_{split} = \{R(\mathbf{x}_i, y_i) \mid (\mathbf{x}_i, y_i) \in (\mathbf{X}_{cal}, \mathbf{y}_{cal})\}, \quad (16)$$

where each conformity score $R(\mathbf{x}_i, y_i)$ is calculated as in (12) using \hat{f}_{split} . Now, for any future instance \mathbf{x}_{test} the scores in C_{split} are used to evaluate candidate values for potential inclusion in the interval, specifically by calculating a p -value as in (13). A pseudocode implementation of the split conformal inference method is provided in Algorithm 4 (Lei et al., 2018).

Algorithm 4: Split Conformal Inference Method

Input: training data \mathbf{X} and \mathbf{y} , test observation \mathbf{x}_{test} , learning algorithm L , a set of candidate values Q , nonconformity measure d , and desired coverage probability $1 - \alpha$

- 1 Evenly split \mathbf{X} and \mathbf{y} into two sub-sets: the proper training set $(\mathbf{X}_{proper}, \mathbf{y}_{proper})$, and the calibration set $(\mathbf{X}_{cal}, \mathbf{y}_{cal})$
- 2 Fit learning algorithm L to the proper training set and denote the trained regressor as \hat{f}
- 3 Build C_{split} as in (16)
- 4 **for** $q_j \in Q$ **do**
- 5 Compute $R(\mathbf{x}_{test}, q_j)$ as in (12) using \hat{f}
- 6 Compute the p -value of $R(\mathbf{x}_{test}, q_j)$ as in (13) using \hat{f}
- 7 **end**

Output: a $1 - \alpha$ prediction interval, as defined in (14)

Split CI presents a clear computational advantage over the full conformal inference and bootstrap methods. For learning algorithms relying on time-consuming gradient descent optimization (e.g., neural networks) the reduction in computation time is non-trivial. On the other hand, the PIs constructed with the split CI method are “less informationally efficient” by virtue of using only a subset of the original sample for training the learning algorithm (Vovk, 2015). Moreover, the randomness introduced by the single splitting of the data—along with, in this research, the use of a single, trained neural network for predictions—increases the variability of PI performance across test sets.

2.2.3 AGGREGATED CONFORMAL PREDICTORS

A compromise between the computationally-intensive full CI method and the cheap but high-variance split CI method is to leverage an “aggregated conformal predictor” (ACP) (Carlsson et al., 2014). An ACP uses a resampling procedure, such as k -fold partitioning or bootstrapping, to generate multiple resamples of the training data set. Given a test observation \mathbf{x}_{test} , candidate value q , and conformity score $R(\mathbf{x}_{test}, q)$, these resamples are leveraged to produce several estimates of the p -value for q . The idea of using an ACP is to calculate more accurate and stable p -values for constructing PIs. Two common aggregated

CI approaches, which are discussed here, are the cross-conformal inference (cross-CI) and bootstrap CI methods.

The implementation of the cross-CI is similar to split CI, as both methods are inductive. However the process of fitting a learning algorithm L and calculating $\pi(q_j)$ (for a candidate value $q_j \in Q$), is repeated across K splits of the data. The separate K estimates for $\pi(q_j)$, $\pi_k(q_j)$ for $k = 1, \dots, K$, are then aggregated to produce a single value, $\bar{\pi}(q_j)$ (Vovk, 2015). Assuming that the splits are approximately of equal size, and that K is significantly smaller than the number of observations in the original sample, the aggregation function can be the average of each $\pi_k(q_j)$:

$$\bar{\pi}(q_j) = \frac{1}{K} \sum_{k=1}^K \pi_k(q_j). \quad (17)$$

The $1 - \alpha$ PI can then be constructed for a test observation \mathbf{x}_{test} , with unknown target value y_{test} , as

$$\Gamma_{aggregated} = \{q_j \in Q \mid \bar{\pi}(q_j) \geq \alpha\}, \quad (18)$$

where Q is the set of candidates considered for y_{test} (Vovk, 2015). Algorithm 5 provides a pseudocode implementation of the cross-CI algorithm.

Algorithm 5: Cross-Conformal Inference Method

Input: training data \mathbf{X} and \mathbf{y} , test observation \mathbf{x}_{test} , desired number of folds K , learning algorithm L , a set of candidate values Q , nonconformity measure d , and desired coverage probability $1 - \alpha$

- 1 Evenly split \mathbf{X} and \mathbf{y} into K sub-sets, denoting each as $(\mathbf{X}_k, \mathbf{y}_k)$
- 2 **for** $k = 1$ to K **do**
- 3 The calibration set is $(\mathbf{X}_k, \mathbf{y}_k)$ and the proper training set is $(\mathbf{X}_{-k}, \mathbf{y}_{-k}) \equiv \{(\mathbf{x}_i, y_i) \in (\mathbf{X}, \mathbf{y}) \mid (\mathbf{x}_i, y_i) \notin (\mathbf{X}_k, \mathbf{y}_k)\}$
- 4 Fit learning algorithm L to the proper training set and denote the trained regressor as \hat{f}_k
- 5 **for** $q_j \in Q$ **do**
- 6 $C_k = \{R(\mathbf{x}_i, y_i); (\mathbf{x}_i, y_i) \in (\mathbf{X}, \mathbf{y})\}$ where each $R(\mathbf{x}_i, y_i)$ is calculated as in (12) using \hat{f}_k
- 7 Compute $R(\mathbf{x}_{test}, q_j)$ as in 12 using \hat{f}_k
- 8 Compute the p-value of $R(\mathbf{x}_{test}, q_j)$ as in (13)
- 9 **end**
- 10 **end**
- 11 Compute $\bar{\pi}(q_j)$ for $q_j \in Q$ as in (17)

Output: a $1 - \alpha$ prediction interval, as defined in (18)

An alternate approach, bootstrap CI, leverages the concepts of bootstrap resampling to generate several calculations of p -values which are then aggregated. The algorithmic approach of the bootstrap CI method is effectively the same as cross-CI, with the only change being the method of resampling.

Suppose the training data \mathbf{X} and \mathbf{y} are bootstrap resampled B times; denote each resample as \mathbf{X}_b^* and \mathbf{y}_b^* . Then, each \mathbf{X}_b^* and \mathbf{y}_b^* serves as the proper training set, with its corresponding set of out-of-bag samples, $\tilde{\mathbf{X}}_b^*$ and $\tilde{\mathbf{y}}_b^*$, serving as the calibration set. Following the same process as in cross-CI, conformity scores are calculated from each $\tilde{\mathbf{X}}_b^*$ and $\tilde{\mathbf{y}}_b^*$. Then, for a test observation \mathbf{x}_{test} , the p -value of a candidate value q_j can be calculated on each of the b sets of conformity scores: $\pi_b(q_j)$, $b = 1, \dots, B$ (Vovk, 2015). These values are then aggregated to a single value, $\bar{\pi}(q_j)$, calculated as their arithmetic mean. A $1 - \alpha$ PI is then constructed as in (18). Algorithm 6 provides a pseudocode implementation of the bootstrap CI method.

Algorithm 6: Bootstrap Conformal Inference Method

Input: training data \mathbf{X} and \mathbf{y} , test observation \mathbf{x}_{test} , desired number of resamples B , learning algorithm L , a set of candidate values Q , nonconformity measure d , and desired coverage probability $1 - \alpha$

- 1 Generate B resamples of \mathbf{X} and \mathbf{y} into K sub-sets, denoting each as $(\mathbf{X}_b^*, \mathbf{y}_b^*)$
- 2 **for** $b = 1$ to B **do**
- 3 The proper training set is $(\mathbf{X}_b^*, \mathbf{y}_b^*)$ and the calibration set is $(\tilde{\mathbf{X}}_b^*, \tilde{\mathbf{y}}_b^*)$
- 4 Fit learning algorithm L to the proper training set and denote the trained regressor as \hat{f}_b
- 5 **for** $q_j \in Q$ **do**
- 6 $C_b = \{R(\mathbf{x}_i, y_i); (\mathbf{x}_i, y_i) \in (\mathbf{X}, \mathbf{y})\}$ where each $R(\mathbf{x}_i, y_i)$ is calculated as in (12) using \hat{f}_b
- 7 Compute $R(\mathbf{x}_{test}, q_j)$ as in (12) using \hat{f}_b
- 8 Compute the p -value of $R(\mathbf{x}_{test}, q_j)$ as in (13)
- 9 **end**
- 10 **end**
- 11 Compute $\bar{\pi}(q_j)$ for $q_j \in Q$ as in (17)

Output: a $1 - \alpha$ prediction interval, as defined in (18)

2.2.4 CONFORMAL INFERENCE WITH KERNEL DENSITY ESTIMATION

While the absolute error is the typical choice for the nonconformity measure for regression tasks, the performance of these CI algorithms may suffer if the distribution of prediction errors does not behave as expected (e.g., if they are biased, skewed, or multi-modal). Neural networks compound this problem, as well. The parameter values of a neural network will generally converge to different values across different training sessions due to randomness in the optimization procedure, specifically random initialization and the existence of local minima. This effect adds an additional layer of noise to the network’s fit to the data and, by extension, the prediction errors.

One solution to this issue is kernel density estimation (KDE), a non-parametric method for estimating a probability density function (PDF) over a set of observed data (Rosenblatt, 1956). Estimating the PDF through KDE preserves the general shape of the data, while smoothing over spurious deviations arising from random sampling.

The two hyperparameters to be tuned for fitting the KDE to observed data are the kernel function (K) and bandwidth (h). The bandwidth parameter is effectively the strength of smoothing applied to the prediction errors. As $h \rightarrow 0$, the estimated PDF more closely resembles the observed data; as $h \rightarrow \infty$, it becomes smoother and flatter over the range of observed values (Rosenblatt, 1956). There exists several rules of thumb and analytic methods for tuning h (Lei et al., 2011). However, for most practitioners, constructing and evaluating several density estimates from a finite grid of potential bandwidths $H = \{h_1, h_2, \dots, h_l\}$ is the most readily applied approach.

A natural choice for K in the context of evaluating regression estimates is the Gaussian kernel. As an example, suppose a KDE with Gaussian kernel and bandwidth h is fitted to a set of observed values $\mathbf{z} = \{z_1, z_2, \dots, z_n\}$. Then for an arbitrary value u , the estimated density, $\hat{p}(u)$, under the fitted KDE is

$$\hat{p}(u) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{u - z_i}{h}\right),$$

where the Gaussian kernel K is calculated as

$$K\left(\frac{u - z}{h}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{u-z}{h}\right)^2}$$

for any z (Casella and Berger, 2002). Thus, for every potential value of u , the fitted KDE outputs the estimated density at u . The KDE function fitted to \mathbf{z} does not necessarily resemble the original kernel function; so long as the smoothing parameter h is not overly strong, the modalities and asymmetries in the empirical distribution of \mathbf{z} will still be present.

KDE gives a methods for estimating the likelihood (i.e., density) associated with a particular observation. Values in regions of high density are more likely, while values in regions of low density are less likely. Given the goal of CI to construct prediction intervals based on the "strangeness" of some candidate value, we can use KDE to construct a conformity measure where strangeness is inherently related to the estimated distribution of residuals, as opposed to the absolute value of residuals. Specifically, for a candidate value q_j , we can use $-\ln(R(\mathbf{x}_i, y_i))$ as our conformity measure, or any monotone function thereof.

3. Experimental Setup

This research executes a two-step experiment on eleven data sets to evaluate how the PI methods discussed in Section 2 perform when used in conjunction with MLPs and CNNs. The first experimental step seeks to understand how changes in the settings of key neural network hyperparameters affect PI performance. In the second step, bootstrapping and CI PI methods are compared to determine which better optimizes the trade-off between performance and computational burden. In this section, we first describe the breadth of the data sets, and therefore data structures, used. Then, we provide further details about the experiments to be performed.

3.1 Data Sets

The data sets chosen for this research are intended to capture a variety of data structures, such that the expected performance of different PI methods can be well understood regard-

less of the specific task at hand. Additionally, the data sets used are familiar in the neural networks literature, used to support a variety of research goals (Hernandez-Lobato and Adams 2015; Gal and Ghahramani 2016; Foong et al. 2019). Ten of these data sets (summarized in Table 1) have one-dimensional feature vectors, for which MLPs are constructed during experimentation.

Data Set	Samples (n)	Features (d)	n/d
Boston Housing	506	13	38.9
Wine Quality*	1,599	11	145.9
Concrete Strength*	1,030	8	128.8
Energy Efficiency*	768	8	96.0
Kinematics	8,192	8	1024.0
Naval Propulsion*	11,934	16	745.9
Power Plant*	9,568	4	2392.0
Protein Structure	45,730	9	5081.1
Yacht Hydrodynamics*	308	6	51.3
Year Prediction MSD*	515,345	90	5726.1

Table 1: Summary of benchmark data sets. An asterisk (*) denotes data sets accessed from the UCI Machine Learning Repository (Dua and Graff, 2021).

The eleventh data set is the Rotated Modified National Institute of Standards (RotNIST) data set of handwritten digits. The MNIST data set has been used regularly as a benchmark data set to evaluate the performance of novel network architectures (Jarrett et al. 2009; Glorot et al. 2011; Lin et al. 2014). While the original MNIST data set is generally used for the classification of the handwritten digits, RotNIST extends its scope to regression tasks by applying random rotations to the digits. The task is then to train a neural network to predict these angles of rotation, which are applied uniformly from -45 and 45 degrees.

The particular iteration of this data set used in this paper is downloaded from the MATLAB code platform (link: [instructions to download data set](#)). It contains 10,000 observations of $28 \times 28 \times 1$ images. Pixel values are scaled such that every entry is between zero and one. For this analysis, images are resized to $14 \times 14 \times 1$ using bilinear interpolation. Figure 1 displays four sample images from the processed data set.

3.2 Design of Experiment

The first experimental step seeks to determine the effects of neural network hyperparameter selection on PI performance using an ANOVA approach. To do so, experimental designs of various MLP and CNN hyperparameters are constructed. These hyperparameter designs are then used to fit corresponding neural networks. The networks are trained on each data set in conjunction with the pivot bootstrap (1,000 resamples) and split conformal inference methods to construct 95% PIs (i.e., $\alpha = 0.05$) for observations in a pre-sequestered test set (20% of original data set). The two PI methods are chosen as exemplars of the bootstrap and

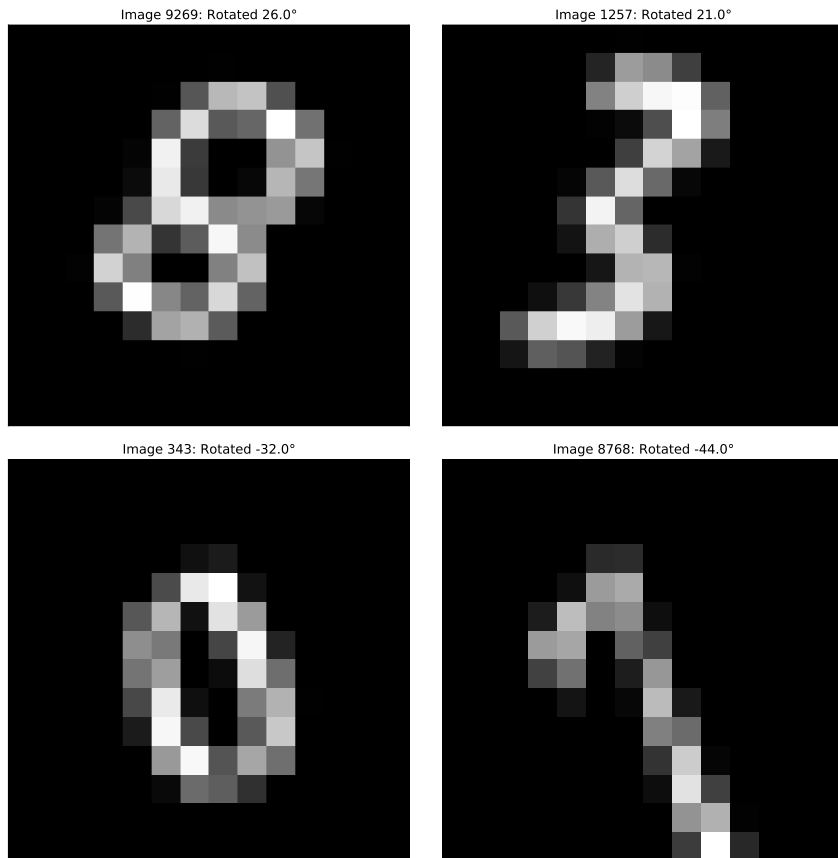


Figure 1: Example images from the RotNIST dataset.

conformal inference approaches, respectively. The validity and efficiency of the constructed PIs are evaluated across the modeled network architectures. For robustness, a five-fold cross validation approach is implemented, with each instance in a data set used as a test observation exactly once.

The experimental design for evaluating MLPs is summarized in Table 2. Three factors are considered: activation function, number of hidden layers, and number of nodes in each hidden layers. Every combination of levels is explored, resulting in a total of $6 \times 3 \times 3 = 54$ network architectures. MLPs typically used to model the benchmark data sets utilize one or two hidden layers, each with 50 to 100 nodes with ReLU or Tanh activation functions (Hernandez-Lobato and Adams 2015; Gal and Ghahramani 2016; Foong et al. 2019). By examining other popular activation functions (i.e., sigmoid) and constructing both small and large networks, the experimental design covers the decision space typical of these key hyperparameters. Other network hyperparameters not included in this design are either assigned to constant values for each data set, such as the settings of the Adam optimization algorithm for training (Kingma and Ba, 2014), or are outside of the scope of the analysis. In particular, regularization schemes are not considered for these MLPs. Algorithm 8 provides a pseudocode execution for implementing this first experimental design.

Factor	Levels
Activation	ReLU, Sigmoid, Tanh
Layers	1, 2, 3
Nodes	5, 10, 25, 50, 75, 100

Table 2: Design of experiment factor and levels for the benchmark data sets.

Algorithm 7: First Experimental Step (Benchmark Data Sets)

```

1 for each benchmark data set do
2   for each design point in Table 2 do
3     Build network using combination of hyperparameters; designate the
       untrained network as  $L$ 
4     Use 5-fold cross validation to create 5 pairs of training and test sets
5     for each training/test set pair do
6       Run Algorithm 1 to implement Pivot Bootstrap method ( $B = 1,000$ );
       calculate network ensemble’s RMSE using out-of-bag sets
7       Run Algorithm 4 to implement Split Conformal Inference method;
       calculate network’s RMSE using validation set
8     end
9   end
10  Perform ANOVA modeling to assess PI performance across network design
      choices
11 end

```

A similar experimental methodology to that portrayed in Algorithm 8 is pursued for the RotNIST dataset. However, this experimental design uses the relevant hyperparameters for a CNN; namely, the number of convolutional layers, the number of filters for each layer, and the kernel size. These experimental factors and their corresponding levels are shown in Table 3. All the CNNs fit in the experiment also employ regularization to improve test set performance. In particular, batch normalization is executed after each convolutional layer, as well as dropout before the output layer (probability= 0.05) (Srivastava et al., 2014).

Factor	Levels
Convolutional Layers	2, 3, 4
Kernel Size	1×1 , 3×3 , 5×5
Pooling	Average, Maximum

Table 3: Design of experiment factors and levels for the RotNIST data set.

At the conclusion of this first experimental step, further analysis is performed to determine the sensitivity of network hyperparameters to PI performance and optimal network structure for each data set. Architectures are assessed according to the validity and efficiency of the PIs constructed with them. The optimal network architecture results in the

most efficient PIs (smallest average widths) that maintain validity. The validity of a set of PIs is tested using the Agresti and Coull method for binomial proportions (Agresti and Coull, 1998). Valid PIs are those whose average coverage is not statistically different from 0.95. The average width of the PIs for the optimal network architecture and PI method is recorded as W_j for the j^{th} data set.

After determining the optimal neural network architecture for each data set, the second experimental step begins. In this step, the optimal network structure is used to implement each of the PI methods described in Section 2, such that the validity, efficiency, and computational burden of each can be evaluated and compared. For a given data set, 100 observations are randomly selected and sequestered as the test set. 95% PIs are constructed for the test set using each PI method, with coverage and width recorded for each test observation. This portion of the experiment is repeated 10 times, after which summary statistics for PI coverage and average width are calculated using each repetition. The examined methods are summarized in Table 4. All CI methods are implemented using both the absolute residual and a fitted KDE as the conformity measure, with the results of both presented in Section 4. The pseudocode to execute this second experimental step is provided in Algorithm 8. Here, computational burden is assessed as the number of times the architecture has to be trained in order to implement the PI method, as the training of neural networks is generally the most time and resource intensive task for inference within deep learning tools.

Method	Execution	Computational Burden (number of networks trained)
Pivot Bootstrap	100 resamples	100
Pivot Bootstrap	500 resamples	500
Pivot Bootstrap	1,000 resamples	1,000
Percentile Bootstrap	1,000 resamples	1,000
Bootstrap Conformal Inference	5 resamples	5
Bootstrap Conformal Inference	10 resamples	10
Bootstrap Conformal Inference	20 resamples	20
Cross-Conformal Inference	5 folds	5
Cross-Conformal Inference	10 folds	10
Cross-Conformal Inference	20 folds	20
Full Conformal Inference		10,001
Split Conformal Inference		1

Table 4: Execution and computation burden of methods examined during experimentation.

Recall from Section 2 that the conformal inference methods require the examination of a fine grid of candidate target values. With the scale and variability of observations varying by data set, and considering the computational burden of full conformal inference, a systematic approach is needed to ensure that the location, size, and fineness of Q considered for each observation is appropriate and experimentally feasible. To that end, the fineness of Q is determined by the conformal inference method used. PIs constructed with full

Algorithm 8: Second Experimental Step

```

1 for each data set do
2   for trial 1 to 10 do
3     Randomly sample and remove 100 test observations from the original data
       set; designate the remaining observations as the training set
4     for each method  $A$  in Table 4 do
5       Implement  $A$  using the data set’s optimal neural network structure
6       Record the performance metrics of the constructed PIs
7     end
8   end
9   Average replicates of PI performance metrics
10  Analyze PIs
11 end

```

conformal inference examine a search grid of 100 values, while other conformal inference methods examined grids of 1,000 values. The width of Q is set as twice the value of W_j , the average width of the optimal set of PIs for data set j found in the first experimental step. Lastly, for each test observation \mathbf{x}_{test} , Q is constructed such that it is centered around the point prediction for its target value. This prediction, $\hat{f}(\mathbf{x}_{test})$, is calculated from the underlying learning algorithm fit to the training data. For the full conformal inference method, point predictions for each test observation are found by training a neural network on the original training set, before the PI algorithm is implemented. Thus, the search grid for \mathbf{x}_{test} from dataset j is the set of 100 (or 1,000) evenly-spaced values in the interval $[\hat{f}(\mathbf{x}_{test}) - W_j, \hat{f}(\mathbf{x}_{test}) + W_j]$.

3.3 Code Implementation and Other Details

Neural networks are trained in the Python coding language using the Keras library, with TensorFlow as the backend platform. Notebooks are executed on Google Colaboratory (“Google Colab”). The code is publicly available at <https://github.com/alexcontarino/Constructing-Prediction-Intervals-for-Neural-Networks>.

For the 10 benchmark data sets shown in Table 1, the Adam optimizer is used for training neural networks. Optimization settings, such as the learning rate, batch size, and number of training epochs, are tuned for each data set to ensure networks fit well to the training data and in a timely fashion. Other settings for the Adam optimizer function are left to their default value assigned in Keras. The optimization settings used for each data set are summarized in Table 5.

Parameters of the training algorithm for the CNNs tasked for learning the RotNIST data set are similarly tuned. For that data set, the stochastic gradient descent (SGD) algorithm with a learning rate of 0.001 is used to train networks for 20 epochs. The batch size for training is 64 samples. Other settings for the SGD optimizer are left to their Keras-assigned default values.

Data Set	Epochs	Batch Size	Learning Rate
Boston Housing	200	32	0.001
Wine Quality	75	32	0.001
Concrete Strength	200	32	0.001
Energy Efficiency	250	32	0.001
Kinematics	150	32	0.001
Naval Propulsion	80	256	0.001
Power Plant	350	600	0.01
Protein Structure	75	1024	0.025
Yacht Hydrodynamics	500	64	0.001
Year Predictions MSD	15	4096	0.1

Table 5: Optimization settings for benchmark data sets.

4. Results

This section provides results and discussion related to our two-stage experiment. We first present the effects of neural network hyperparameters on PI performance. The statistical significance of effects are measured and evaluated using an Analysis of Variance (ANOVA) approach. Next, we compare the performance of the PI methods discussed in Section 2, examining the differences in the coverage and average widths on similar data sets. Lastly, we perform a case study assessing the conditional coverage of select PI methods on various data sets.

4.1 Effects of Neural Network Hyperparameters

The results of the first experimental step reveals how changes in neural network hyperparameter settings affect PI performance. Performance in terms of coverage and efficiency is discussed in the succeeding subsections.

4.1.1 COVERAGE

We analyze PI coverage of both the pivot bootstrap and split conformal inference methods is analyzed using an ANOVA approach. With the factors identified in Table 2 and Table 3 for MLPs and CNNs, respectfully, full factorial models are built to estimate coverage for each data set. These models include main effects, as well as all second- and third-order effects, in order to determine which hyperparameters are significant and, in particular, understand how neural network architecture affects PI coverage.

Table 6 summarizes the significance of effects for the benchmark regression data sets. A significance level of 0.05 was chosen for this analysis. The most striking result is in the significance in effects when comparing the bootstrapping and split CI methods. For the former, in which it is expected that coverage remain around the nominal level, none of the effects are significant in any of the experimental data sets. However, a majority of effects are significant for modeling pivot bootstrap PI coverage in a given data set. For the latter PIs, coverage generally improves with the fit of the network. This result is intuitive given

the pivot bootstrap method relies on the distribution of the target outcome, conditioned on the observed feature vector, being approximately Normal.

	Boston Housing	Wine Quality	Concrete Strength	Energy Efficiency	Kinematics	Naval Propulsion	Power Plant	Protein Structure	Yacht Hydrodynamics	Year Prediction
Effect	Pivot Bootstrap (1,000 Resamples)									
Activation	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Not Significant	Significant
Layers	Significant	Not Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant
Nodes	Significant	Significant	Not Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant
Activation × Layers	Significant	Not Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant
Activation × Nodes	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant
Layers × Nodes	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Not Significant
Activation × Layers × Nodes	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Significant	Not Significant
Effect	Split Conformal Inference									
Activation	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant
Layers	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant
Nodes	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant
Activation × Layers	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant
Activation × Nodes	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant
Layers × Nodes	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant
Activation × Layers × Nodes	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant	Not Significant

	Significant at $\alpha=0.05$
	Not significant at $\alpha=0.05$

Table 6: Significance of effects modeling PI coverage.

Figure 2 provides a comparison between the two PI methods for the Boston Housing data set. The significant differences in coverage among the activation functions can be observed for pivot bootstrap PIs, indicated by the non-overlapping error bars in the trend lines. Note too that these significant differences result in relatively conservative PIs; the empirical coverages are above the desired confidence level of 0.95. The analogous plot of split conformal inference PI coverage in Figure 2 reveals no design choice in the network architecture which yields a statistically significant change in coverage. This pattern extends to the image-based RotNIST data set (Figure 3), with the modeled PI coverage plotted by the relevant CNN hyperparameters. Graphical results of PI coverage for the remaining data sets are provided in Appendix A.

The results presented here reveal that the coverage of the pivot bootstrap method is sensitive to the underlying neural network structure. The apparent mechanisms for this relationship are the fit of the network (measured by root mean squared error on the test set) and the response variable being modeled. Figure 4 provides an example plot of coverages of the bootstrap PIs against the test RMSEs from the networks used to construct the PIs. In this example from the Energy Efficiency data set, PI coverage is highest for networks at the extremes of observed RMSE. However, the relationship between coverage of the bootstrapped PIs and RMSE is not consistent across the response variables of the modeled

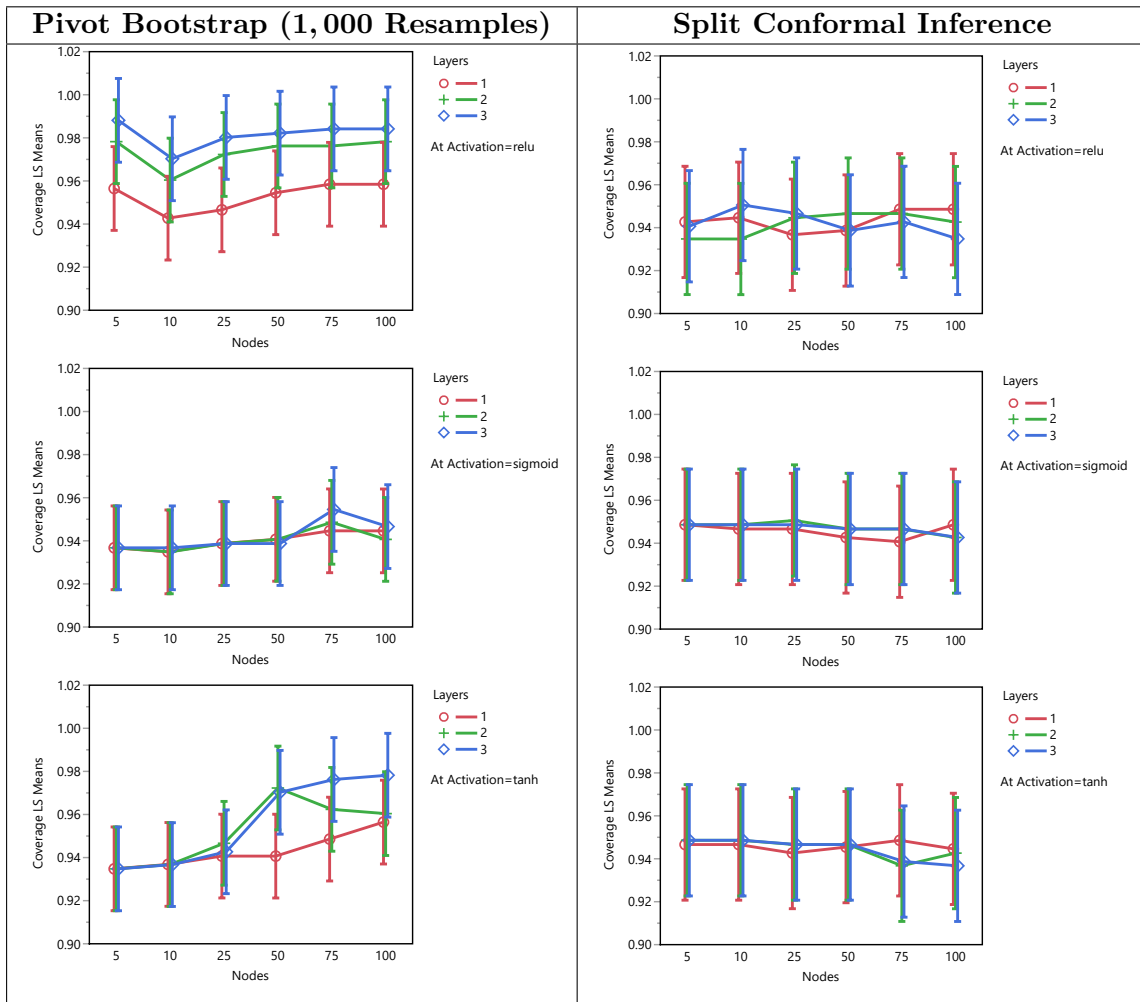


Figure 2: PI Coverage in the Boston Housing Data Set.

data sets. Plots of bootstrapped PI coverage versus RMSE for other data sets are provided in Appendix C.

Meanwhile, the split conformal inference method is generally insensitive to both structure and fit of the underlying networks. The PIs constructed with this method have an average coverage near the nominal level of 0.95, independent of any factors relating to the learning algorithm used.

4.1.2 AVERAGE WIDTH

A similar analysis approach is taken to examine PI average width as a function of neural network hyperparameter effects. In particular, full factorial models are built from the factors shown in Table 2 (for benchmark data sets) and Table 3 (for RotNIST data set) to estimate pivot bootstrap and split conformal inference PI average widths. As before, these models include all main effects, as well as second- and third-order interactions. In the models for

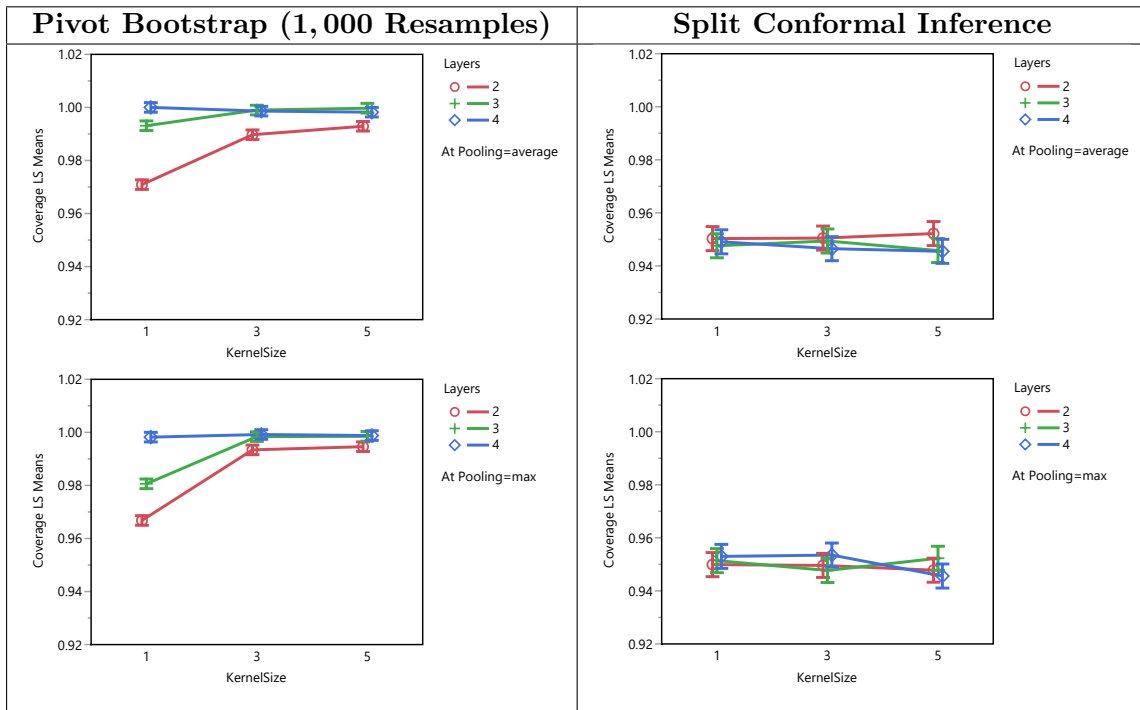


Figure 3: PI Coverage in the RotNIST Data Set.

all eleven data sets, as well as for both pivot bootstrap and split conformal inference PIs, all modeled effects are statistically significant at $\alpha = 0.05$. Moreover, for every data set, the constructed models had a R^2 of at least 0.98, indicating they accounted for approximately 98% of the observed variance in average width.

In general, design choices which increase network capacity (i.e., ability to encode relationships between the features and target variable) result in better fitting networks (as measured by test RMSE) and smaller prediction intervals. Thus, networks having more layers or nodes will typically provide narrower PIs than networks having fewer. Similarly, networks utilizing the ReLU activation also provide an advantage in average width over those using Tanh or Logistic. However, these advantages often diminish as the size of the network grows. For example, Figure 5 for the Power Plant data set, shows the difference in average width between networks using five and ten nodes is generally larger than the difference between networks using 75 and 100.

Similarly, the advantage of using three layers instead of two is insignificant regardless of the chosen activation for many data sets. In general, for the benchmark regression data sets, increasing network capacity beyond two layers and fifty nodes provides statistically insignificant changes in PI average width. Exceptions to this rule include the Energy Efficiency and Protein Structure data sets, which have target variables with more complex behaviors.

Focusing solely on the CNNs constructed for RotNIST data set (Figure 6), a similar pattern in the effects of their hyperparameters emerges. The network’s prediction performance

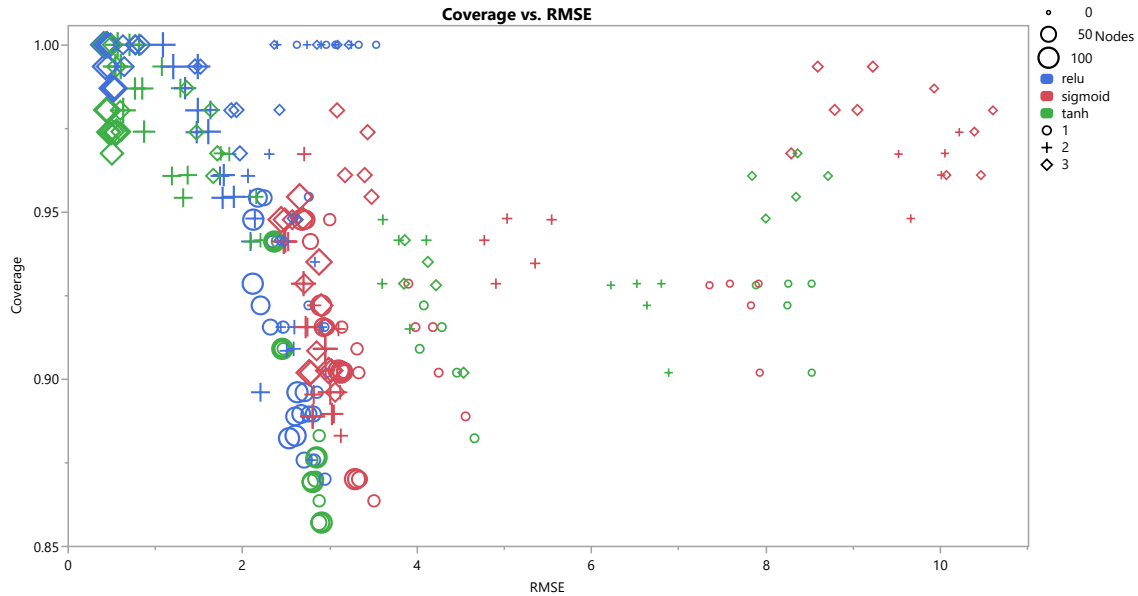


Figure 4: PI Coverage versus RMSE for the Energy Efficiency Data Set.

measured by RMSE and the average widths of the corresponding PIs improve significantly as layers increases from two to three, and as kernel size increases from 1×1 to 3×3 . Other changes in hyperparameters yield less noticeable and generally insignificant differences in width. Graphical results of PI average width for the remaining data sets are provided in Appendix B.

However, the results presented thus far should not be interpreted to say that arbitrarily increasing network size will always result in better PI performance. Rather, in select data sets, increasing network capacity yields poor performing PIs. The experimental results reveal that in none of the eleven data sets modeled did the network (split conformal inference method) or network ensemble (pivot bootstrap method) minimizing RMSE also produce the most efficient set of PIs (that also maintained coverage). Thus, optimal neural network performance does not necessarily translate to optimal average widths for its corresponding PIs. These results suggests that further heuristics in the optimization of neural networks—beyond, say, monitoring RMSE on a validation set—are needed if the end goal of the network is to provide efficient PIs.

4.1.3 OTHER RESULTS AND DISCUSSION

The empirical results presented thus far illustrate how the choice of neural network hyperparameters affect PI performance. Most notably, the coverage of the pivot bootstrap PIs are sensitive to changes in network architecture. Larger networks, or more generally, those that are able to better learn the training data, lead to bootstrapped PIs with higher coverage than smaller, poorer fitting networks. Conversely, the coverage of split conformal

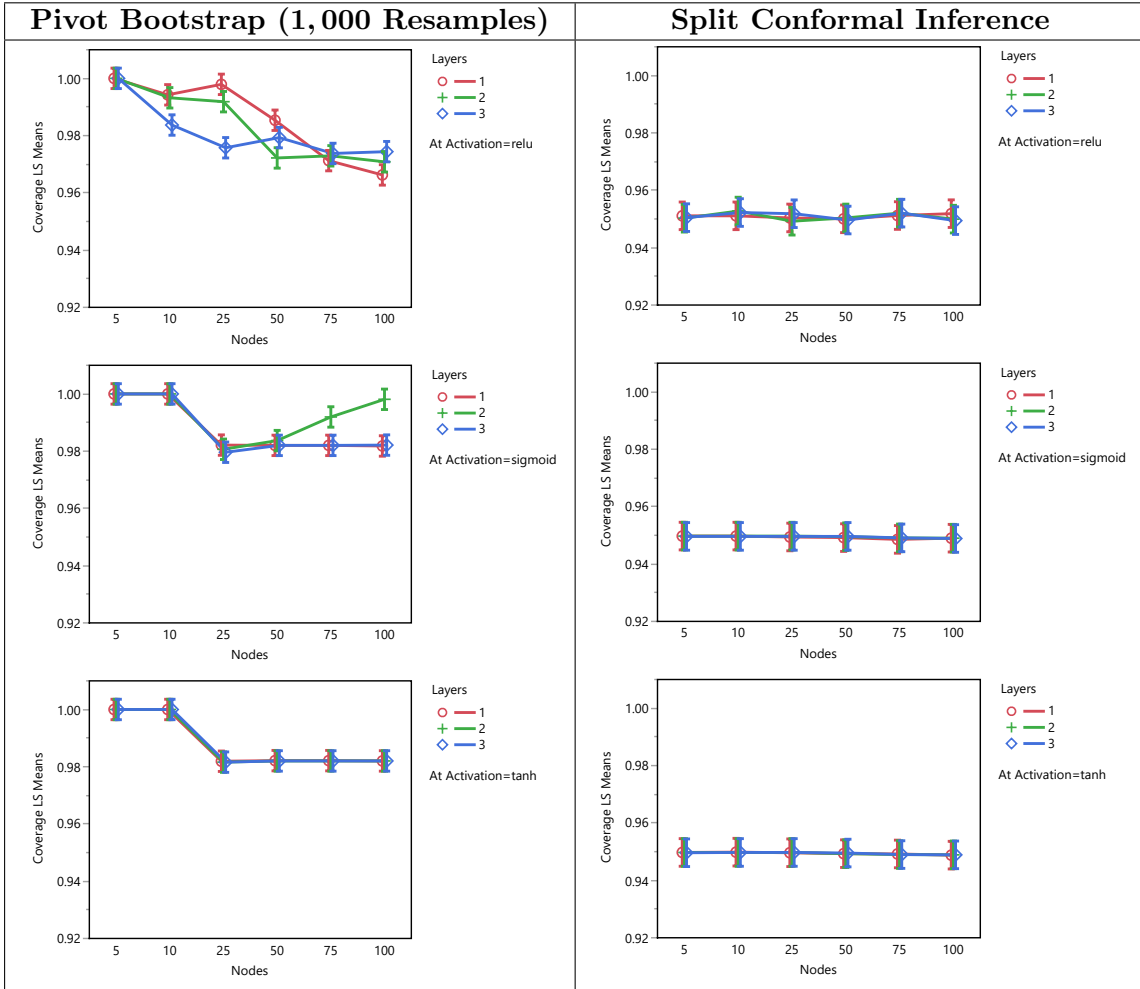


Figure 5: PI Average Width in the Power Plant Data Set.

inference PI sets remain close to the nominal level regardless of how the size or fit of the neural network used for regression learning, which is not surprising.

Moreover, PI average width is closely related to the fit of the neural network. This is evidenced in summary plots for each data set: PI average width and test set RMSE move in tandem with one another. Despite this persistent relationship between PI average width and RMSE, it is not necessarily the case that the neural network (or ensemble of networks) which provide the minimum test set RMSE also provide the PIs with the smallest average widths, particularly for the networks constructed in conjunction with the bootstrap PI method.

The exact mechanism for this somewhat contradictory result is unknown. Further exploration is needed to understand the relationship between the fit of the network and the distribution of its resulting residuals, especially as the former improves from highly biased to minimally bias.

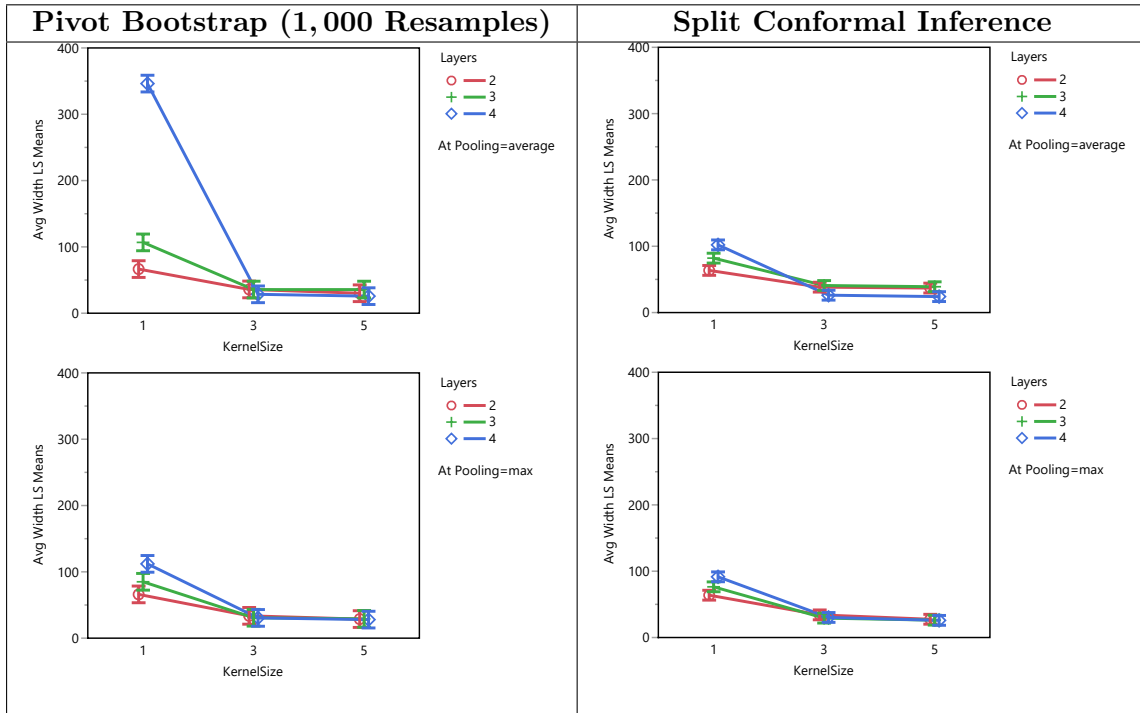


Figure 6: PI Average Width in the RotNIST Data Set.

Table 7 provides network architecture for each data set producing the optimal PIs, along with their average width and the PI algorithm used. Split conformal inference produces the optimal set of PIs for all but one of the data set. This result is driven largely by the conservative nature of the bootstrap method; the result is PIs whose coverage is much higher than nominal. The difference in overall performance between the split conformal inference method and pivot bootstrap method is large, with the former providing roughly nominal coverage while being, on average, 3.4% narrower. A more detailed comparison of PI performance is provided in Section 4.2. The lone exception to this rule is seen in the Year Prediction data set. The PIs of this data set generally large (usually more than half of the original target range from the training set); the lack of precision in these estimates may mean that the relative advantage of using one PI method over another is negligible.

Considering the observed PI performance in coverage, width, as well as the potential effects of over-fitting, the recommended benchmark neural network structure to begin analysis of data set is an architecture consisting of two hidden layers, with 50 nodes in each, and employing the ReLU activation. Additionally, single-layer networks generally do not provide enough opportunity for the network to map the feature space to the response effectively. Specifically, that lack of depth must be overcome by prolonging the training time or by adding nodes to the structure. A similar modification must also be implemented if using Tanh or Sigmoid activation function, as these functions require more parameters and longer training times to effectively learn the patterns in the data when compared to the ReLU activation. For data sets containing more complicated patterns (e.g., the Energy Effi-

Benchmark Data Sets	MLP Settings			Method	Width
	Layers	Nodes	Activation		
Boston Housing	3	100	ReLU	Split Conformal	13.135
Wine Quality	2	50	ReLU	Split Conformal	2.560
Concrete Strength	3	100	ReLU	Split Conformal	24.836
Energy Efficiency	3	100	Tanh	Split Conformal	4.243
Kinematics	3	75	Tanh	Split Conformal	0.303
Naval Propulsion	3	100	ReLU	Split Conformal	0.126
Power Plant	3	50	ReLU	Split Conformal	16.005
Protein Structure	3	50	Tanh	Split Conformal	18.218
Yacht Hydrody- namics	3	100	ReLU	Split Conformal	4.560
Year Prediction MSD	1	50	ReLU	Pivot Bootstrap	59.917
Image-Based Data Set	CNN Settings			Method	Width
	Layers	Kernel Size	Pooling		
RotNIST	3	5×5	Maximum	Split Conformal	26.170

Table 7: Optimal neural network hyperparameter settings found for each data set. Also shown is the PI method producing the optimally performing PIs with their corresponding average width.

ciency data set) increasing network capacity—through the addition of layers or nodes—will improve both metrics.

An analogous approach can be taken with CNNs. Based on the results from the RotNIST data set, deeper networks (three or four hidden layers) provide a marked improvement in PI and network point prediction performance as compared to CNNs having just two hidden layers. In the relatively limited design of architectures examined here, use of 3×3 and 5×5 convolution kernels provided a similar advantage over 1×1 kernels. However, design space of CNN depth and kernel size remains an area of active research (Simonyan and Zisserman 2014, Howard et al. 2020).

4.2 Comparison of PI Methods

Table 8 provides the computed average widths for each method implemented across each data set (standard errors in parentheses) from the second experiment. Table 9 similarly displays coverage values, which are evaluated according to their closeness to the nominal coverage level of 0.95. To enable easier interpretation of the average width results, the values are normalized by data set. In particular, every average width is expressed as a ratio to the narrowest PIs produced for that data sets with the most efficient PIs for each data set having a value of 1.00. These values are shown in Table 10.

Method	Boston Housing		Wine Quality		Concrete Strength		Energy Efficiency		Kinematics		Naval Propulsion		Power Plant		Protein Structure		Yacht Hydrodynamics		Year Prediction		RoMIST
	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	
Pivot Bootstrap (100 resamples)	16.2 ± 0.8	2.64 ± 0.02	25.8 ± 0.4	2.31 ± 0.05	2.32 ± 0.04	0.332 ± 0.0001	0.111 ± 0.0004	19.8 ± 0.3	18.5 ± 0.3	5.71 ± 0.51	56.9 ± 0.2	28.9 ± 0.5									
Pivot Bootstrap (500 resamples)	16.2 ± 0.7	2.64 ± 0.01	25.8 ± 0.3	2.32 ± 0.04	2.37 ± 0.06	0.333 ± 0.0001	0.112 ± 0.0002	19.7 ± 0.2	18.7 ± 0.1	5.74 ± 0.46	56.9 ± 0.3	28.9 ± 0.2									
Pivot Bootstrap (1,000 resamples)	16.2 ± 0.2	2.64 ± 0.05	26.1 ± 0.6	2.37 ± 0.06	2.37 ± 0.06	0.351 ± 0.0001	0.112 ± 0.0001	18.7 ± 0.0	18.9 ± 0.1	5.03 ± 0.15	59.9 ± 0.1	29.1 ± 0.1									
Percentile Bootstrap (1,000 resamples)	14.9 ± 0.4	2.67 ± 0.03	25.6 ± 0.3	2.36 ± 0.02	2.36 ± 0.02	0.329 ± 0.0001	0.09 ± 0.0002	19.7 ± 0.1	19.7 ± 0.1	6.11 ± 0.49	58.0 ± 0.2	22.0 ± 0.1									
Bootstrap Conformal Inference (5 resamples)	13.8 ± 1.0	2.65 ± 0.09	24.1 ± 0.9	2.27 ± 0.07	2.27 ± 0.07	0.300 ± 0.0006	0.097 ± 0.0033	17.2 ± 0.8	18.6 ± 0.6	5.95 ± 1.50	59.1 ± 1.0	23.7 ± 2.6									
Bootstrap Conformal Inference (5 resamples) - KDE	13.5 ± 0.9	2.58 ± 0.06	22.0 ± 2.4	2.19 ± 0.05	2.19 ± 0.05	0.298 ± 0.0006	0.091 ± 0.026	16.7 ± 1.0	18.2 ± 0.3	5.26 ± 1.09	52.5 ± 1.2	22.2 ± 2.3									
Bootstrap Conformal Inference (10 resamples)	13.5 ± 0.6	2.67 ± 0.04	21.1 ± 0.6	2.22 ± 0.09	2.22 ± 0.09	0.297 ± 0.0002	0.084 ± 0.0004	17.3 ± 0.6	18.3 ± 0.4	5.38 ± 0.65	57.6 ± 1.3	22.5 ± 2.8									
Bootstrap Conformal Inference (10 resamples) - KDE	13.3 ± 0.6	2.61 ± 0.05	21.8 ± 2.8	2.17 ± 0.08	2.17 ± 0.08	0.296 ± 0.0002	0.082 ± 0.0006	17.0 ± 0.7	18.1 ± 0.3	4.76 ± 0.64	52.6 ± 1.1	22.6 ± 2.4									
Bootstrap Conformal Inference (20 resamples)	13.6 ± 0.6	2.66 ± 0.03	24.3 ± 0.5	2.27 ± 0.06	2.27 ± 0.06	0.297 ± 0.0003	0.089 ± 0.0006	17.5 ± 0.4	18.3 ± 0.1	5.45 ± 0.61	58.1 ± 0.7	23.4 ± 0.8									
Bootstrap Conformal Inference (20 resamples) - KDE	13.3 ± 0.5	2.61 ± 0.05	21.9 ± 2.7	2.18 ± 0.06	2.18 ± 0.06	0.296 ± 0.0003	0.088 ± 0.0006	17.2 ± 0.4	18.2 ± 0.2	4.76 ± 0.49	54.2 ± 0.7	22.8 ± 0.7									
Cross-Conformal Inference (5 folds)	12.7 ± 0.8	2.62 ± 0.05	23.9 ± 0.9	2.16 ± 0.06	2.16 ± 0.06	0.282 ± 0.0003	0.080 ± 0.0009	16.7 ± 0.5	18.1 ± 0.4	4.39 ± 0.26	60.9 ± 1.3	22.3 ± 1.6									
Cross-Conformal Inference (5 folds) - KDE	12.5 ± 0.7	2.58 ± 0.05	20.8 ± 3.3	2.00 ± 0.06	2.00 ± 0.06	0.282 ± 0.0004	0.077 ± 0.0007	16.2 ± 0.4	17.7 ± 0.4	3.94 ± 0.43	53.6 ± 2.1	20.9 ± 1.5									
Cross-Conformal Inference (10 folds)	12.3 ± 0.6	2.62 ± 0.04	23.2 ± 0.9	2.12 ± 0.07	2.12 ± 0.07	0.279 ± 0.0004	0.087 ± 0.0006	17.5 ± 0.7	18.4 ± 0.3	3.95 ± 0.41	59.5 ± 0.6	30.5 ± 5.5									
Cross-Conformal Inference (10 folds) - KDE	12.3 ± 0.6	2.58 ± 0.04	19.9 ± 4.0	2.00 ± 0.08	2.00 ± 0.08	0.278 ± 0.0003	0.085 ± 0.0006	16.7 ± 0.7	18.2 ± 0.4	3.49 ± 0.49	54.0 ± 1.3	28.5 ± 5.1									
Cross-Conformal Inference (20 folds)	12.7 ± 0.6	2.61 ± 0.04	22.4 ± 0.7	2.16 ± 0.06	2.16 ± 0.06	0.281 ± 0.0003	0.084 ± 0.0006	17.3 ± 0.3	18.2 ± 0.2	3.99 ± 0.33	58.4 ± 0.8	20.5 ± 1.2									
Cross-Conformal Inference (20 folds) - KDE	12.6 ± 0.6	2.60 ± 0.03	24.8 ± 0.3	2.42 ± 0.06	2.42 ± 0.06	0.326 ± 0.0006	0.097 ± 0.0002	16.1 ± 0.2	18.5 ± 0.3	2.72 ± 0.16	63.5 ± 3.4	25.4 ± 1.6									
Full Conformal Inference	12.4 ± 0.4	2.58 ± 0.02	24.8 ± 0.3	2.26 ± 0.35	2.26 ± 0.35	0.317 ± 0.0005	0.090 ± 0.0001	16.5 ± 0.0	17.5 ± 0.2	2.72 ± 0.16	50.5 ± 0.3	19.5 ± 0.2									
Split Conformal Inference	13.4 ± 1.6	2.63 ± 0.15	24.9 ± 1.7	3.48 ± 0.89	3.48 ± 0.89	0.305 ± 0.0009	0.113 ± 0.0023	16.1 ± 0.2	18.6 ± 1.3	6.38 ± 1.99	65.6 ± 0.9	30.5 ± 7.6									
Split Conformal Inference - KDE	12.6 ± 1.6	2.58 ± 0.14	22.8 ± 2.3	3.26 ± 0.80	3.26 ± 0.80	0.295 ± 0.0012	0.104 ± 0.0018	15.8 ± 0.3	18.0 ± 0.2	4.81 ± 1.34	50.1 ± 3.1	21.8 ± 1.4									

Table 8: Average widths of PI methods with calculated standard errors.

Method	Boston Housing		Wine Quality		Concrete Strength		Energy Efficiency		Kanamataks		Naval Population		Power Plant		ProteinStructure		Yacht Hydrodynamics		Year Prediction		MeanST	
	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE	Mean	SE		
Bootstrap (100 resamples)	0.983 ± 0.009	0.969 ± 0.021	0.985 ± 0.013	0.987 ± 0.012	0.978 ± 0.009	0.981 ± 0.009	0.993 ± 0.008	0.992 ± 0.012	1.0 ± 0.0	1.0 ± 0.0	0.983 ± 0.016	0.971 ± 0.023	0.991 ± 0.012	0.941 ± 0.016	0.991 ± 0.012	0.991 ± 0.013	0.990 ± 0.013	0.991 ± 0.012	0.991 ± 0.012	0.991 ± 0.012	0.991 ± 0.012	1.0 ± 0.0
Bootstrap (500 resamples)	0.984 ± 0.009	0.969 ± 0.021	0.984 ± 0.011	0.985 ± 0.011	0.982 ± 0.011	0.979 ± 0.011	0.992 ± 0.012	0.992 ± 0.012	1.0 ± 0.0	1.0 ± 0.0	0.979 ± 0.006	0.972 ± 0.002	0.994 ± 0.014	0.939 ± 0.018	0.994 ± 0.014	0.994 ± 0.014	0.994 ± 0.014	0.994 ± 0.014	0.994 ± 0.014	0.994 ± 0.014	0.994 ± 0.014	1.0 ± 0.0
Bootstrap (1,000 resamples)	0.984 ± 0.011	0.955 ± 0.021	0.984 ± 0.011	0.985 ± 0.011	0.982 ± 0.011	0.979 ± 0.011	0.989 ± 0.008	0.989 ± 0.008	1.0 ± 0.0	1.0 ± 0.0	0.950 ± 0.021	0.979 ± 0.019	0.985 ± 0.019	0.948 ± 0.016	0.985 ± 0.019	0.985 ± 0.019	0.985 ± 0.019	0.985 ± 0.019	0.985 ± 0.019	0.985 ± 0.019	0.985 ± 0.019	1.0 ± 0.0
Percentile Bootstrap (1,000 resamples)	0.978 ± 0.011	0.968 ± 0.019	0.981 ± 0.012	0.981 ± 0.012	0.980 ± 0.012	0.978 ± 0.012	0.989 ± 0.008	0.989 ± 0.008	1.0 ± 0.0	1.0 ± 0.0	0.980 ± 0.005	0.977 ± 0.013	0.996 ± 0.005	0.939 ± 0.021	0.996 ± 0.005	0.996 ± 0.005	0.996 ± 0.005	0.996 ± 0.005	0.996 ± 0.005	0.996 ± 0.005	0.996 ± 0.005	1.0 ± 0.0
Bootstrap Conformal (5 resamples) - KDE	0.947 ± 0.021	0.964 ± 0.023	0.945 ± 0.025	0.966 ± 0.013	0.968 ± 0.016	0.968 ± 0.016	0.978 ± 0.011	0.978 ± 0.011	0.994 ± 0.007	0.994 ± 0.007	0.959 ± 0.027	0.972 ± 0.020	0.964 ± 0.023	0.921 ± 0.022	0.964 ± 0.023	0.964 ± 0.023	0.964 ± 0.023	0.964 ± 0.023	0.964 ± 0.023	0.964 ± 0.023	0.964 ± 0.023	0.990 ± 0.020
Bootstrap Conformal (10 resamples) - KDE	0.955 ± 0.020	0.966 ± 0.023	0.947 ± 0.035	0.973 ± 0.010	0.965 ± 0.012	0.965 ± 0.012	0.973 ± 0.018	0.973 ± 0.018	0.993 ± 0.010	0.993 ± 0.010	0.964 ± 0.022	0.972 ± 0.018	0.965 ± 0.025	0.925 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.988 ± 0.018
Bootstrap Conformal (20 resamples) - KDE	0.952 ± 0.023	0.965 ± 0.023	0.947 ± 0.038	0.967 ± 0.008	0.957 ± 0.014	0.957 ± 0.014	0.974 ± 0.013	0.974 ± 0.013	0.993 ± 0.010	0.993 ± 0.010	0.960 ± 0.021	0.974 ± 0.013	0.965 ± 0.025	0.925 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.965 ± 0.025	0.986 ± 0.021
Bootstrap Conformal (50 resamples) - KDE	0.959 ± 0.019	0.965 ± 0.022	0.967 ± 0.019	0.966 ± 0.019	0.966 ± 0.019	0.966 ± 0.019	0.976 ± 0.014	0.976 ± 0.014	0.997 ± 0.005	0.997 ± 0.005	0.968 ± 0.023	0.964 ± 0.023	0.968 ± 0.023	0.941 ± 0.018	0.968 ± 0.023	0.968 ± 0.023	0.968 ± 0.023	0.968 ± 0.023	0.968 ± 0.023	0.968 ± 0.023	0.968 ± 0.023	0.995 ± 0.009
Bootstrap Conformal (100 resamples) - KDE	0.959 ± 0.016	0.965 ± 0.024	0.948 ± 0.032	0.963 ± 0.032	0.963 ± 0.016	0.963 ± 0.016	0.977 ± 0.013	0.977 ± 0.013	0.998 ± 0.004	0.998 ± 0.004	0.966 ± 0.021	0.964 ± 0.022	0.966 ± 0.021	0.933 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.994 ± 0.007
Cross-conformal (5 folds) - KDE	0.953 ± 0.025	0.967 ± 0.024	0.968 ± 0.013	0.968 ± 0.013	0.960 ± 0.017	0.960 ± 0.017	0.965 ± 0.014	0.965 ± 0.014	0.986 ± 0.010	0.986 ± 0.010	0.956 ± 0.027	0.972 ± 0.018	0.962 ± 0.022	0.939 ± 0.020	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.987 ± 0.015
Cross-conformal (10 folds) - KDE	0.948 ± 0.022	0.966 ± 0.024	0.925 ± 0.050	0.943 ± 0.021	0.943 ± 0.021	0.943 ± 0.021	0.967 ± 0.012	0.967 ± 0.012	0.987 ± 0.010	0.987 ± 0.010	0.950 ± 0.020	0.972 ± 0.018	0.962 ± 0.022	0.918 ± 0.023	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.962 ± 0.022	0.985 ± 0.019
Cross-conformal (20 folds) - KDE	0.946 ± 0.019	0.967 ± 0.022	0.966 ± 0.014	0.921 ± 0.057	0.945 ± 0.016	0.945 ± 0.016	0.965 ± 0.022	0.965 ± 0.022	0.996 ± 0.005	0.996 ± 0.005	0.967 ± 0.021	0.962 ± 0.023	0.966 ± 0.021	0.939 ± 0.018	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.997 ± 0.009
Cross-conformal (50 folds) - KDE	0.946 ± 0.019	0.966 ± 0.024	0.921 ± 0.057	0.921 ± 0.057	0.940 ± 0.021	0.940 ± 0.021	0.966 ± 0.021	0.966 ± 0.021	0.996 ± 0.005	0.996 ± 0.005	0.965 ± 0.021	0.960 ± 0.026	0.965 ± 0.021	0.929 ± 0.016	0.965 ± 0.021	0.965 ± 0.021	0.965 ± 0.021	0.965 ± 0.021	0.965 ± 0.021	0.965 ± 0.021	0.965 ± 0.021	0.996 ± 0.012
Cross-conformal (100 folds) - KDE	0.950 ± 0.020	0.966 ± 0.024	0.928 ± 0.058	0.928 ± 0.058	0.961 ± 0.020	0.961 ± 0.020	0.970 ± 0.015	0.970 ± 0.015	0.996 ± 0.005	0.996 ± 0.005	0.968 ± 0.024	0.965 ± 0.022	0.968 ± 0.024	0.942 ± 0.018	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.996 ± 0.012
Cross-conformal (20 folds) - KDE	0.951 ± 0.020	0.966 ± 0.024	0.928 ± 0.058	0.928 ± 0.058	0.961 ± 0.020	0.961 ± 0.020	0.970 ± 0.015	0.970 ± 0.015	0.996 ± 0.005	0.996 ± 0.005	0.968 ± 0.024	0.965 ± 0.022	0.968 ± 0.024	0.942 ± 0.018	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.968 ± 0.024	0.996 ± 0.012
Full Conformal Inference	0.945 ± 0.021	0.963 ± 0.017	0.954 ± 0.018	0.947 ± 0.029	0.947 ± 0.029	0.947 ± 0.029	0.958 ± 0.022	0.958 ± 0.022	0.985 ± 0.019	0.985 ± 0.019	0.945 ± 0.024	0.945 ± 0.024	0.945 ± 0.024	0.939 ± 0.032	0.945 ± 0.024	0.945 ± 0.024	0.945 ± 0.024	0.945 ± 0.024	0.945 ± 0.024	0.945 ± 0.024	0.945 ± 0.024	0.984 ± 0.022
Split Conformal Inference - KDE	0.952 ± 0.023	0.963 ± 0.017	0.953 ± 0.019	0.948 ± 0.022	0.948 ± 0.022	0.948 ± 0.022	0.951 ± 0.018	0.951 ± 0.018	0.985 ± 0.019	0.985 ± 0.019	0.928 ± 0.028	0.931 ± 0.030	0.949 ± 0.016	0.873 ± 0.037	0.928 ± 0.028	0.928 ± 0.028	0.928 ± 0.028	0.928 ± 0.028	0.928 ± 0.028	0.928 ± 0.028	0.928 ± 0.028	0.987 ± 0.019
Split Conformal Inference - KDE	0.916 ± 0.031	0.961 ± 0.019	0.918 ± 0.047	0.917 ± 0.053	0.917 ± 0.053	0.917 ± 0.053	0.949 ± 0.023	0.949 ± 0.023	0.944 ± 0.013	0.944 ± 0.013	0.944 ± 0.019	0.944 ± 0.019	0.944 ± 0.019	0.876 ± 0.027	0.944 ± 0.019	0.944 ± 0.019	0.944 ± 0.019	0.944 ± 0.019	0.944 ± 0.019	0.944 ± 0.019	0.944 ± 0.019	0.985 ± 0.016

Table 9: Average coverage of PI methods with calculated standard errors.

Method	Boston Housing		Wine Quality		Concrete Strength		Energy Efficiency		Kinematics		Naval Propulsion		Power Plant		Yacht Hydrostructure		Year Prediction		RoMNIST		
	Best	Worst	Best	Worst	Best	Worst	Best	Worst	Best	Worst	Best	Worst	Best	Worst	Best	Worst	Best	Worst	Best	Worst	
Pivot Bootstrap (100 resamples)	1.320	1.035	1.404	1.252	1.196	1.436	1.273	1.059	2.096	1.136	1.497										
Pivot Bootstrap (500 resamples)	1.321	1.035	1.402	1.259	1.198	1.428	1.271	1.067	2.110	1.135	1.500										
Pivot Bootstrap (1,000 resamples)	1.317	1.035	1.420	1.287	1.262	1.452	1.202	1.081	1.846	1.196	1.510										
Percentile Bootstrap (1,000 resamples)	1.214	1.045	1.390	1.278	1.184	1.157	1.022	1.126	2.244	1.157	1.143										
Bootstrap Conformal Inference (5 resamples)	1.122	1.039	1.313	1.234	1.081	1.254	1.109	1.061	2.185	1.179	1.231										
Bootstrap Conformal Inference (5 resamples) - KDE	1.097	1.009	1.196	1.187	1.073	1.176	1.077	1.042	1.932	1.049	1.151										
Bootstrap Conformal Inference (10 resamples)	1.099	1.044	1.312	1.205	1.070	1.085	1.114	1.047	1.977	1.150	1.221										
Bootstrap Conformal Inference (10 resamples) - KDE	1.086	1.024	1.187	1.179	1.067	1.055	1.097	1.038	1.748	1.050	1.171										
Bootstrap Conformal Inference (20 resamples)	1.105	1.040	1.319	1.232	1.069	1.151	1.125	1.049	2.003	1.159	1.211										
Bootstrap Conformal Inference (20 resamples) - KDE	1.088	1.021	1.190	1.184	1.067	1.141	1.110	1.042	1.749	1.082	1.180										
Cross-Conformal Inference (5 folds)	1.038	1.028	1.297	1.171	1.016	1.032	1.074	1.033	1.613	1.216	1.155										
Cross-Conformal Inference (5 folds) - KDE	1.022	1.011	1.130	1.086	1.014	1.000	1.045	1.014	1.446	1.069	1.086										
Cross-Conformal Inference (10 folds)	1.001	1.027	1.263	1.149	1.003	1.119	1.124	1.051	1.451	1.188	1.584										
Cross-Conformal Inference (10 folds) - KDE	1.000	1.011	1.083	1.088	1.000	1.096	1.074	1.040	1.283	1.077	1.478										
Cross-Conformal Inference (20 folds)	1.034	1.023	1.218	1.171	1.010	1.086	1.115	1.042	1.466	1.165	1.062										
Cross-Conformal Inference (20 folds) - KDE	1.026	1.020	1.000	1.000	1.007	1.078	1.062	1.025	1.313	1.080	1.000										
Full Conformal Inference	1.008	1.009	1.351	1.316	1.175	1.257	1.036	1.060	1.000	1.268	1.315										
Full Conformal Inference - KDE	1.008	1.009	1.351	1.316	1.175	1.257	1.036	1.060	1.000	1.268	1.315										
Split Conformal Inference	1.095	1.030	1.354	1.888	1.098	1.457	1.034	1.066	2.343	1.309	1.582										
Split Conformal Inference - KDE	1.030	1.000	1.226	1.770	1.060	1.342	1.018	1.030	1.768	1.000	1.129										
% Difference Between Best and Worst Methods	32%	5%	42%	89%	26%	46%	27%	13%	134%	31%	58%										

Table 10: Relative efficiencies of PI methods.

For all eleven data sets examined, using a CI method with KDE yields the most efficient PIs. Across all CI methods and data sets, the use of KDE improves PI efficiency by an average of 5.8%. When examining Table 10, cross-CI with $K = 20$, appears to be the most efficient method overall. While the full CI method does produce the most efficient set of PIs for three data sets, its performance varies substantially in other data sets. Note that the bootstrap methods produce generally the least efficient PIs compared to the CI methods.

Supplementing CI method with KDE further improves their efficiency, however in certain cases this comes at the cost of unacceptable variance or coverage. In particular, the narrowest PIs for the Concrete Strength and Energy Efficiency data sets, each produced using cross-CI ($K = 20$) with KDE, also have coverages well below the nominal level of 95%. Moreover, these PIs also the highest variance in coverage and average width across test sets when compared to the other methods. The extra computation time required for KDE varies by data set, increasing as the number of training samples grows. For the larger data sets—in particular, the Year Prediction data set—this extra time is substantial. The average 5.8% improvement in efficiency suggests that the additional computation time remains worthwhile in most cases so long as coverage is acceptable. If time or computational resources are limited, work-arounds such as allowing non-zero error in the KDEs (VanderPlas, 2013), or reducing the number of folds implemented in cross-CI, can be pursued with minimal degradation in PI performance.

In terms of validity, full and split CI provide the PIs whose coverage most closely match the nominal level of 95%. The PIs of the aggregated CI methods are generally conservative, although the implementing with KDE reduces coverage (in some cases below 95%, as discussed). Bootstrapped PIs, regardless of the number of resamples, produce coverages exceeding 98% for many data sets, as seen previously in the first experimental step.

Across test sets, the bootstrap method provides the most stable PI performance, both in average width and coverage, while the set of CI methods exhibit higher levels of variation. Among the CI methods, the full conformal inference method has the lowest variance—an intuitive result following from its transductive nature. Using aggregated conformal predictors greatly reduces the variance in PI performance across test sets as compared to the split conformal inference method. CI methods implemented with KDE have higher variance in coverage and average width than those implemented using the absolute residual as the nonconformity measure.

When considering aggregate and variable PI performance, both across data sets and replicates, cross-CI seems to provide the highest quality PIs. While implementing any number of folds appears adequate, using a relatively large number such as $K = 20$ appears to provide more stable results. Even so, cross-CI has a significantly smaller computational burden when compared to bootstrapping approaches and the full conformal inference method. Split CI has a potentially unacceptable degree of variance in its performance—induced by generating one random split in the data—while also having sub-optimal efficiency.

4.2.1 CASE STUDY: CONDITIONAL COVERAGE

In addition to the PI performance metrics discussed in previous sections, an additional area in which to evaluate PI methods is their conditional coverage. In the methods and results so far, coverage has been calculated and discussed in marginal terms (i.e., across the

entire test set). Conditional coverage is calculated separately across subsets of the feature space. Neither the bootstrap nor the set of conformal inference methods (in the simple implementations utilized here) guarantee conditional validity (Lei and Wasserman, 2014). Regardless, a practical consideration for constructing PIs is how well a $1 - \alpha$ PI method maintains nominal coverage when examining a particular subset of the feature space.

As a case study, the conditional coverage of three separate PI methods were calculated for combinations of features in the Boston Housing and Energy Efficiency data sets. The selected methods were the pivot bootstrap (500 resamples), the split CI method and the cross-CI method ($K = 20$). Table 11 displays the conditional coverage of these methods for each value of the binary variable “Charles River” in the Boston Housing data set. Values of one or zero correspond to whether the neighborhood borders the river or does not border the river, respectively. The means of the median home price corresponding to both cases are included in the table for reference. The conditional coverage of the bootstrap method exceeds 95% regardless of the value of the Charles River variable. However, the split CI and cross-CI predictors struggle to correctly predict median home price for neighborhood bordering the Charles River. The split CI method has a coverage rate of 79.0% for such homes, while cross-CI has a rate of 86.4%. Since homes in these neighbors tend to command a higher price, incorrect valuations of the home prices carry more cost.

Charles River	0	1
Count	919	81
Mean Home Price (\$1000’s)	22.236	30.868
Pivot Bootstrap (500 resamples)	0.9869	0.9506
Cross-Conformal (20 folds, no KDE)	0.9675	0.7901
Split Conformal (no KDE)	0.9445	0.8642

Table 11: Conditional coverage of select PI methods according to value of Charles River variable in Boston Housing data set.

A similar exercise is conducted for the Energy Efficiency data set. However, instead of examining conditional coverage given the value of a single feature, conditional coverage is computed across bins of the predicted target value: 0-10, 10-20, 20-30, 30-40, 40-50 (the chosen target variable for this data set is heating load, measured in British Thermal Units).

The results are summarized in Table 12. Again, the pivot bootstrap method maintains coverage of at least 95% in every bin of values for the ensemble’s predicted value. The coverages of the CI methods fall below 95% for the predicted heating loads between 20 and 40 BTUs. While deficiency is small for cross-CI, the split conformal inference method in particular struggles to predict this middle range of the target variable. When the neural network trained in conjunction with the split CI method predicts a heating load between 20 and 30 BTUs, the method correctly predicts the range of the true heating load in 80.6% of cases.

A further investigation into the performance of these intervals reveals a clear pattern into the incorrect predictions of the split CI PI method. As seen in Figure 7, the split CI PIs routinely over-predict the target value when the prediction is between 20 and 30

Predicted Heating Load (BTUs)	< 10	10 – 20	20 – 30	30 – 40	> 40
Pivot Bootstrap (500 resamples)	1.000	0.9892	0.9914	0.9563	0.9500
Cross-Conformal (20 folds, no KDE)	1.000	0.9871	0.9397	0.9301	0.9250
Split Conformal (no KDE)	1.000	0.9806	0.8060	0.8901	0.9500

Table 12: Conditional coverage of select PI methods according to predicted value of the target variable in the Energy Efficiency data set.

BTUs (incorrect predictions from pivot bootstrap method provided for reference). The lack of balance in the coverage of these intervals—that is, the relatively low left coverage rate compared to right coverage—suggests either some skewness in the distribution of the target \mathbf{y} , or the underlying neural network is not fully learning the relationship between the features \mathbf{X} and \mathbf{y} . Coverage balance is sometimes a consideration when constructing PIs, as providing an accurate lower or upper bound for an unknown outcome may be important in certain settings.

In general, the conservativeness of the pivot bootstrap method, as measured by its marginal coverage, appears to be a benefit when examining these conditional cases. The additional marginal coverage—resulting in relatively inefficient PIs—effectively provides the slack to maintain conditional coverage at a rate around $1 - \alpha$. The use of an aggregated conformal predictor, already noted for their slightly conservative marginal coverage, appears to be the best among the conformal predictors for optimizing conditional coverage. Alternate methods, such as Mondrian conformal predictors (Vovk et al., 2005), or “locally-weighted” CI (Lei et al., 2018), can be implemented to account for distributional differences in the covariate and response spaces, enabling more consistent coverage.

5. Conclusion

We investigated the relationship between various neural network architectures and PI performance. Additionally, we compared two common approaches for constructing PIs: bootstrapping and conformal inference. Overall, PI performance, especially in terms of efficiency, improves as the fit of the underlying neural network improves. Thus, using the ReLU activation function in the neural network—which generally results in a better fitting network as opposed to Sigmoid or Tanh when given equal amount of training time—and a sufficient number of layers (at least two) and nodes (at least 50) will result in better performing PIs. Analogously, ensuring network capacity is sufficient (three or more hidden layers and 3×3 convolutional kernels) when constructing CNNs, specifically in the number of layers and kernel size, will result in satisfactorily performing intervals.

However, as with point prediction, care should be taken to avoid extensive levels of over-fitting. Over-fitting can result from networks having more than sufficient capacity for the data set, or from overly-long training times. The efficiency of the pivot bootstrap PI method, in particular, appears to be prone to such issues since the effect of over-fitting is masked by the ensemble’s goodness-of-fit metrics (such as RMSE). This leads to an

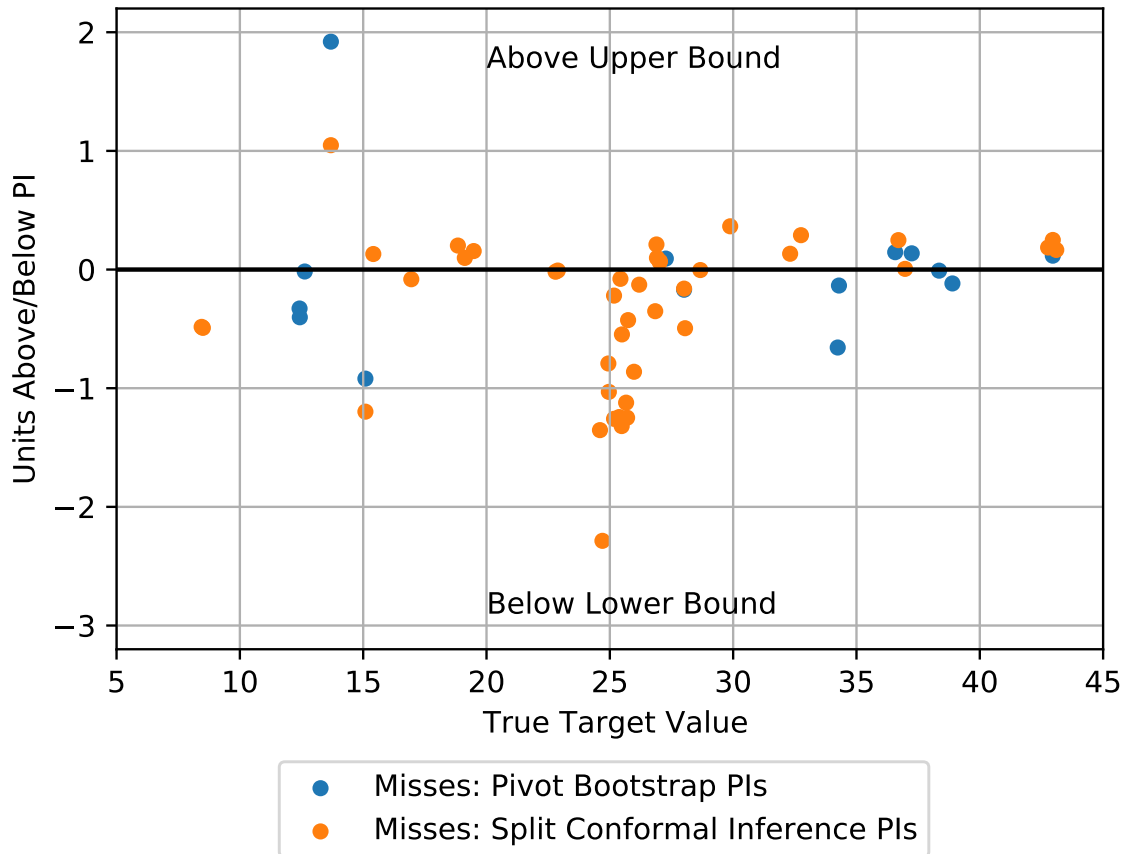


Figure 7: Incorrectly predicted test observations in the Energy Efficiency data set.

otherwise contradictory result; where the best fitting architectures can yield sub-optimal sets of PIs.

For practical implementations, cross-conformal inference is flexible enough to conform to the resources at hand. If resources are limited, or if the training of more than a few networks is not feasible, using the cross-conformal inference method with a small number of folds is likely sufficient. If the resources are not constrained, then using a larger number of folds (e.g., 20) appears to provide the best results, both in average metrics such as width and coverage, and in their variance across test sets.

Opportunities for future work along these lines of research include an extension into predicting categorical responses. The research herein focused solely on data sets having a continuous-valued target variable; in this way, accurate comparisons were made across multiple methods for metrics such as average width. Also interesting is further exploration of the cross-conformal inference method. This research was broad in its scope, seeking to evaluate and compare a suite of PI methods. As such, there was limited time to experiment with the different parameters of each PI method, such as the number of folds for cross-

conformal inference. The relationship between PI performance and the number of folds implemented is not clear from the results here.

Acknowledgments

We would like to thank the Air Force Research Laboratory (RHWA) for supporting this research.

References

- Allen Agresti and Brent A. Coull. Approximate is better than ‘exact’ for interval estimation of binomial proportions. *The American Statistician*, 52(2):119–126, 1998.
- Anastasios Angelopoulos, Stephen Bates, Jitendra Malik, and Michael I. Jordan. Uncertainty sets for image classifiers using conformal prediction. *ArXiv preprint arXiv:1704.07587v2*, 2020.
- Rushil Anirudh and Jayaraman Thiagarajan. Bootstrapping graph convolutional neural networks for autism spectrum disorder classification. *arXiv preprint arXiv:1704.07587v2*, 2017.
- Axel Berg, Magnus Oskarsson, and Mark O’Connor. Deep ordinal regression with label diversity. *arXiv preprint arXiv:2006.15864v1*, 2020.
- Lars Carlsson, Martin Eklund, and Ulf Norinder. Aggregated conformal prediction. In *Artificial Intelligence Applications and Innovations*, pages 231–240, Rhodes, Greece, 2014.
- George Casella and Roger Berger. *Statistical Inference (Second Edition)*. Cengage Learning, Belmont, CA, 2002.
- François Chollet. *Deep Learning with Python*. Manning, Shelter Island, NY, 2017.
- A. C. Davidson and D. V. Hinkley. *Bootstrap Method and Their Applications*. Cambridge University Press, New York, NY, 1997.
- Dheeru Dua and Casey Graff. Uci machine learning repository, 2021. URL <https://archive.ics.uci.edu/ml//index.php>.
- Bradley Efron and Robert Tibshirani. *An Introduction to the Bootstrap*. Chapman and Hall, London, 1993.
- Andrew Y. K. Foong, Yingzhen Li, Jose Miguel Hernandez-Lobato, and Robert E. Turner. In-between uncertainty in bayesian neural networks. *arXiv preprint arXiv:1906.11537*, 2019.
- D. A. S. Fraser and Irwin Guttman. Tolerance regions. *The Annals of Mathematical Statistics*, 27(1):162–179, 1956.

- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of the Thirty-Third International Conference on Machine Learning*, pages 1050–1059, New York, NY, 2016.
- Alex Gammerman, Vladimir Vovk, and Vladimir Vapnik. Learning by transduction. In *Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence*, pages 148–155, Madison, WI, 1998.
- Tianxiang Gao and Vladimir Jojic. Degrees of freedom in deep neural networks. *arXiv preprint arXiv:1603.09260*, 2016.
- James Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning with Applications in R*. Springer, New York, NY, 2013.
- James Gentle. *Computational Statistics*. Springer, New York, NY, 2009.
- Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *AISTATS*, 2011.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, Cambridge, MA, 2017.
- Jose Miguel Hernandez-Lobato and Ryan P. Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *Proceedings of the Thirty-Second International Conference on Machine Learning*, pages 1861–1869, Lille, France, 2015.
- Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861v1*, 2020.
- Kevin Jarrett, Koray Kavukcuoglu, Marc’Aurelio Ranzato, and Yann LeCun. What is the best multi-stage architecture for object recognition? In *2009 IEEE 12th International Conference on Computer Vision*, pages 2146–2153, 2009.
- Abbas Khosravi, Saeid Nahavandi, Dipti Srinivasan, and Rihanna Khosravi. Constructing optimal prediction intervals by using neural networks and bootstrap method. *IEEE Transactions on Neural Networks and Learning Systems*, 26(8):1810–1815, 2015.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980v9*, 2014.
- Jing Lei and Larry Wasserman. Distribution-free prediction bands for non-parametric regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):71–96, 2014.
- Jing Lei, James Robins, and Larry Wasserman. Efficient nonparametric conformal prediction regions. *arXiv preprint arXiv:1111.1418*, 2011.
- Jing Lei, Max G’Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.

- Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *arXiv preprint arXiv:1312.4400v3*, 2014.
- Henrik Linusson, Ulf Johansson, Henrik Boström, and Tuve Löfström. Efficiency comparison of unstable transductive and inductive conformal classifiers. In Lazaros Iliadis, Ilias Maglogiannis, Harris Papadopoulos, Spyros Sioutas, and Christos Makris, editors, *Artificial Intelligence Applications and Innovations*, pages 261–270. Springer Berlin Heidelberg, Berlin, Heidelberg, 2014.
- Eugene Ndiaye. Stable conformal prediction sets. In *Proceedings of the 39th International Conference on Machine Learning*, pages 16462–16479. PMLR, 17–23 Jul 2022.
- Ilya Nourtdinov, Thomas Melliush, and Volodya Vovk. Ridge regression confidence machine. In *Proceedings of the 18th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pages 385–392. PMLR, 2001.
- Harris Papadopoulos. Inductive conformal prediction: Theory and application to neural networks. In Paula Fritzsche, editor, *Tools in Artificial Intelligence*, pages 315–330. Intech, London, England, 2008.
- Georgios Papadopoulos, Peter J. Edwards, and Alan F. Murray. Confidence estimation methods for neural networks: A practical comparison. *IEEE Transactions on Neural Networks and Learning Systems*, 12(6):1278–1287, 2001.
- Murray Rosenblatt. Remarks on some nonparametric estimates of a density function. *Annals of Mathematical Statistics*, 27(3):832–837, 1956.
- Rasmus Rothe, Radu Tomofte, and Luc Van Gool. Deep expectation of real and apparent age from a single image without facial landmarks. *International Journal of Computer Vision*, 126(2–4):144–157, 2018.
- Craig Saunders, Alex Gammerman, and Vladimir Vovk. Transduction with confidence and credibility. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence*, pages 722–726, Stockholm, Sweden, 1999.
- Glen Schafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9:371–421, 2008.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Nitish Srivastava, Alex Krizhevsky, Geoffrey Hinton, Ilya Sutskever, , and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- Jake VanderPlas. Kernel density estimation in python, 2013. URL <http://jakevdp.github.io/blog/2013/12/01/kernel-density-estimation/>.
- Vladimir Vovk. Cross-conformal predictors. *Annals of Mathematics and Artificial Intelligence*, 74(1–2):9–28, 2015.

Vladimir Vovk, Alex Gammerman, and Glen Schafer. *Algorithmic Learning in a Random World*. Springer, New York, 2005.

Weidi Xie, J. Alison Noble, and Andrew Zisserman. Microscopy cell counting and detection with fully convolutional regression networks. *Computer Methods in Biomechanics and Biomedical Engineering: Imaging and Visualization*, 6(3):283–292, 2018.

Appendix A. PI Coverage Plots

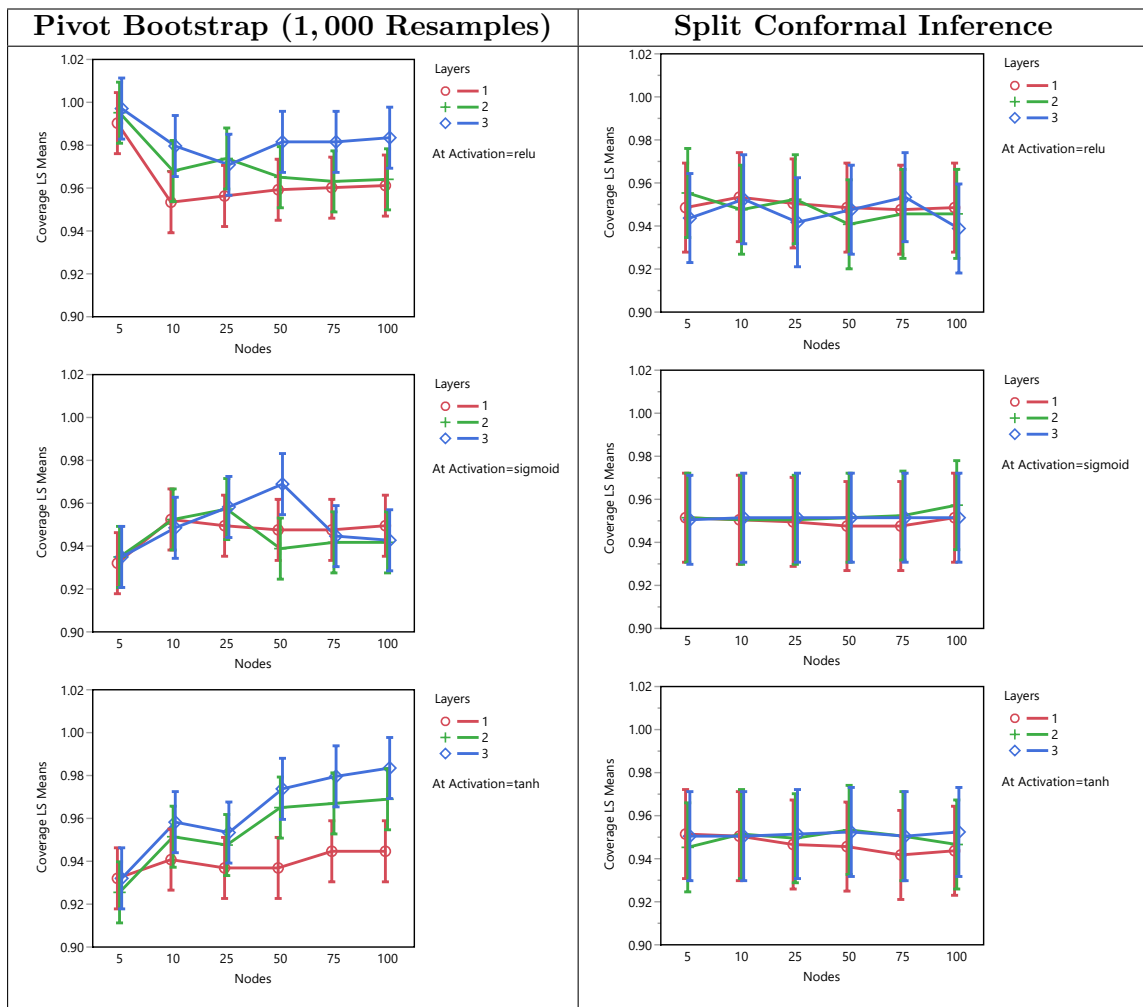


Figure 8: PI Coverage in the Concrete Strength Data Set.

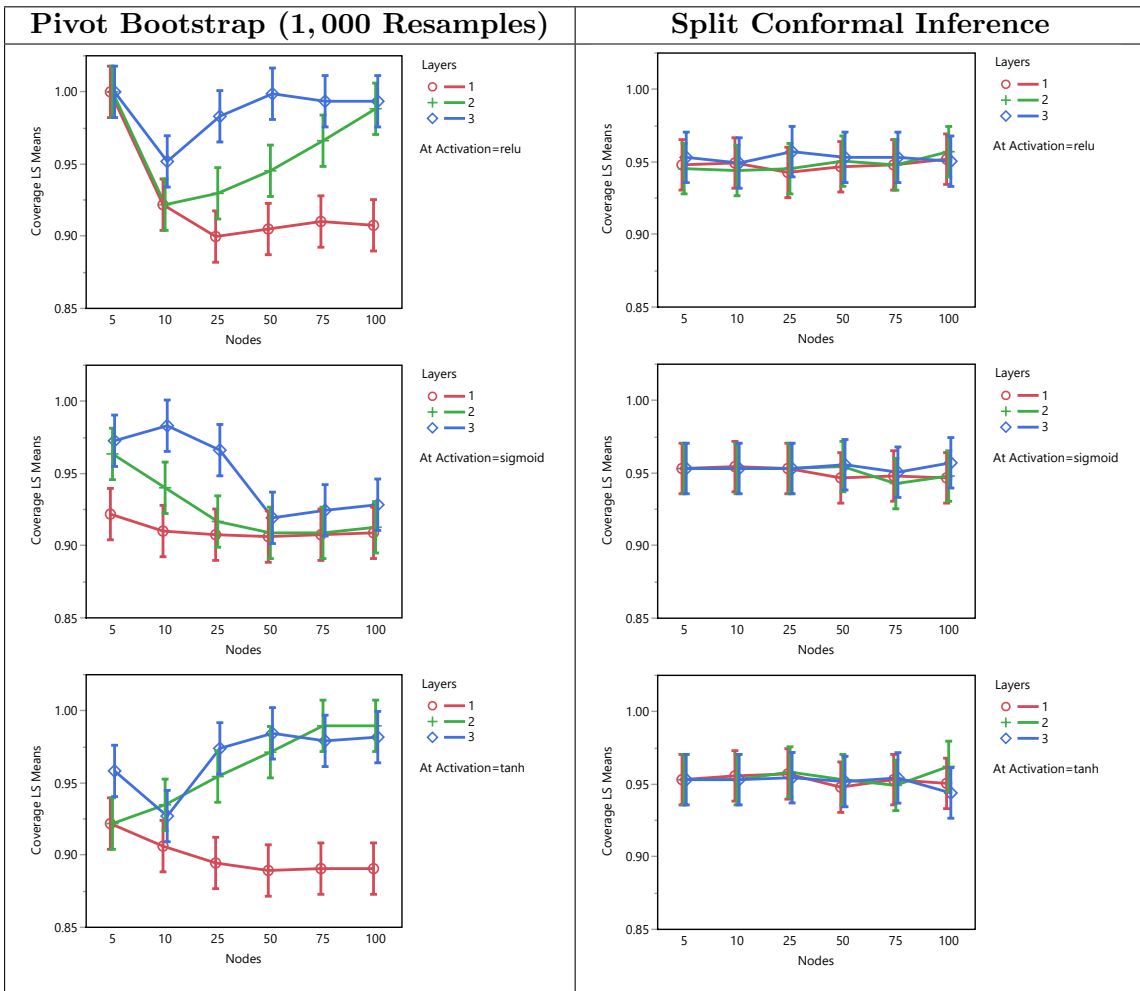


Figure 9: PI Coverage in the Energy Efficiency Data Set.

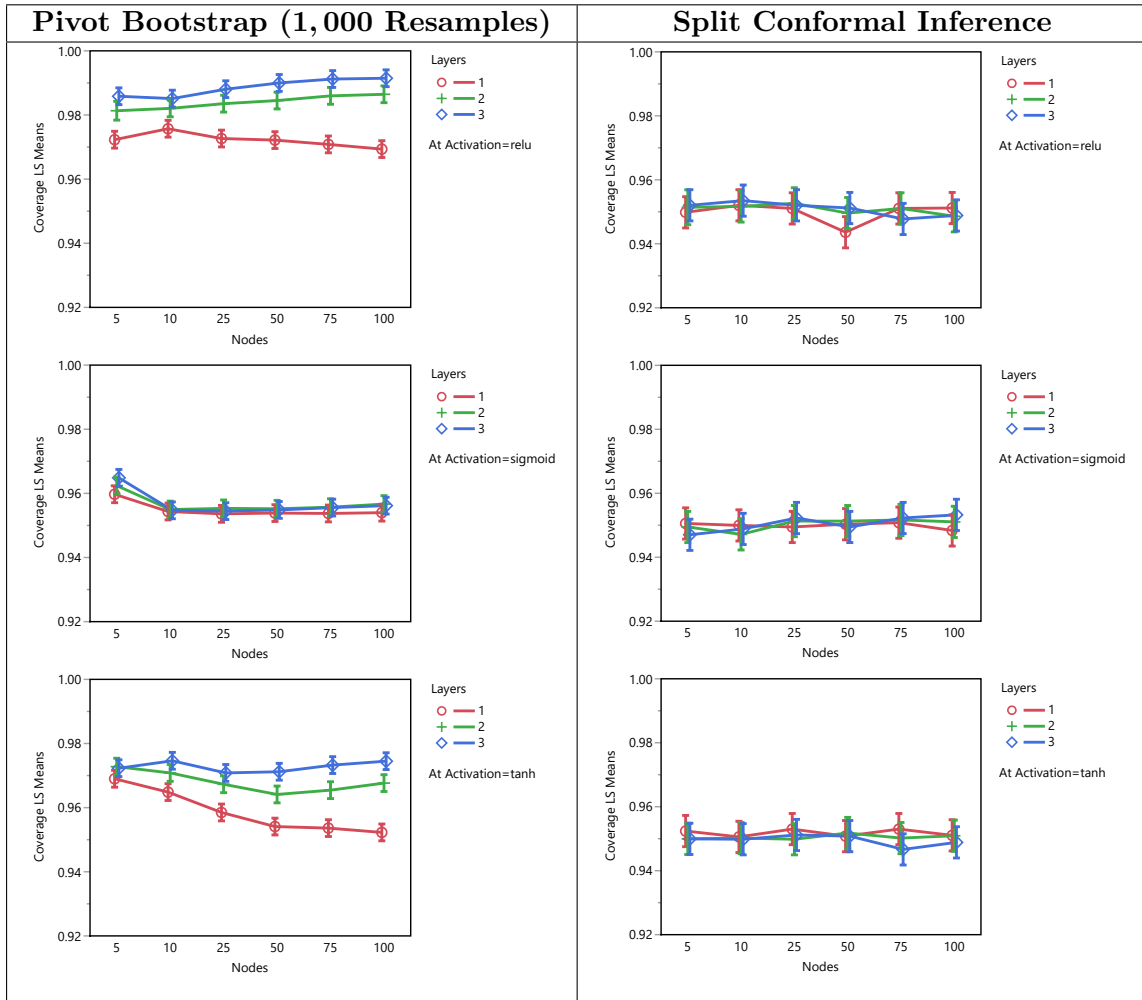


Figure 10: PI Coverage in the Kinematics Data Set.

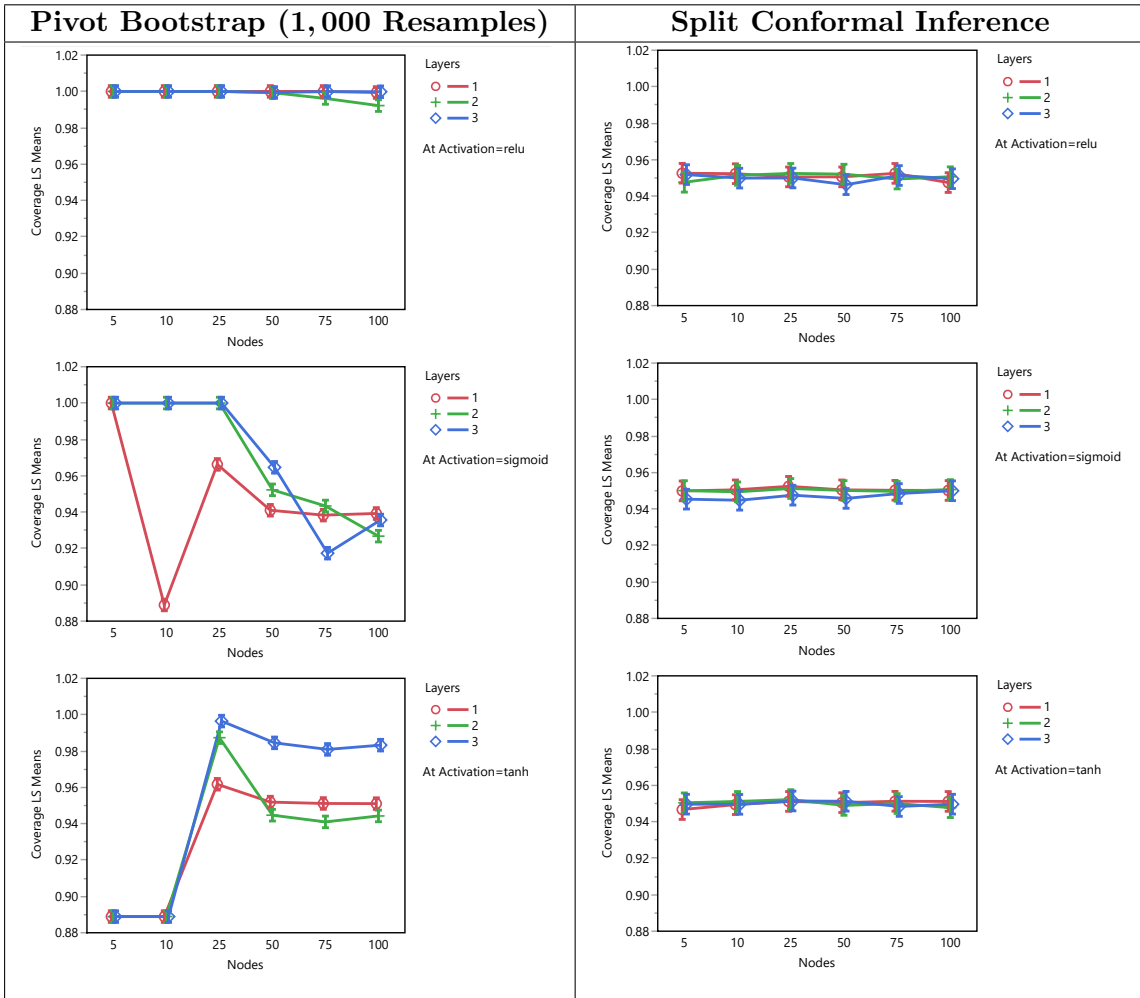


Figure 11: PI Coverage in the Naval Propulsion Data Set.

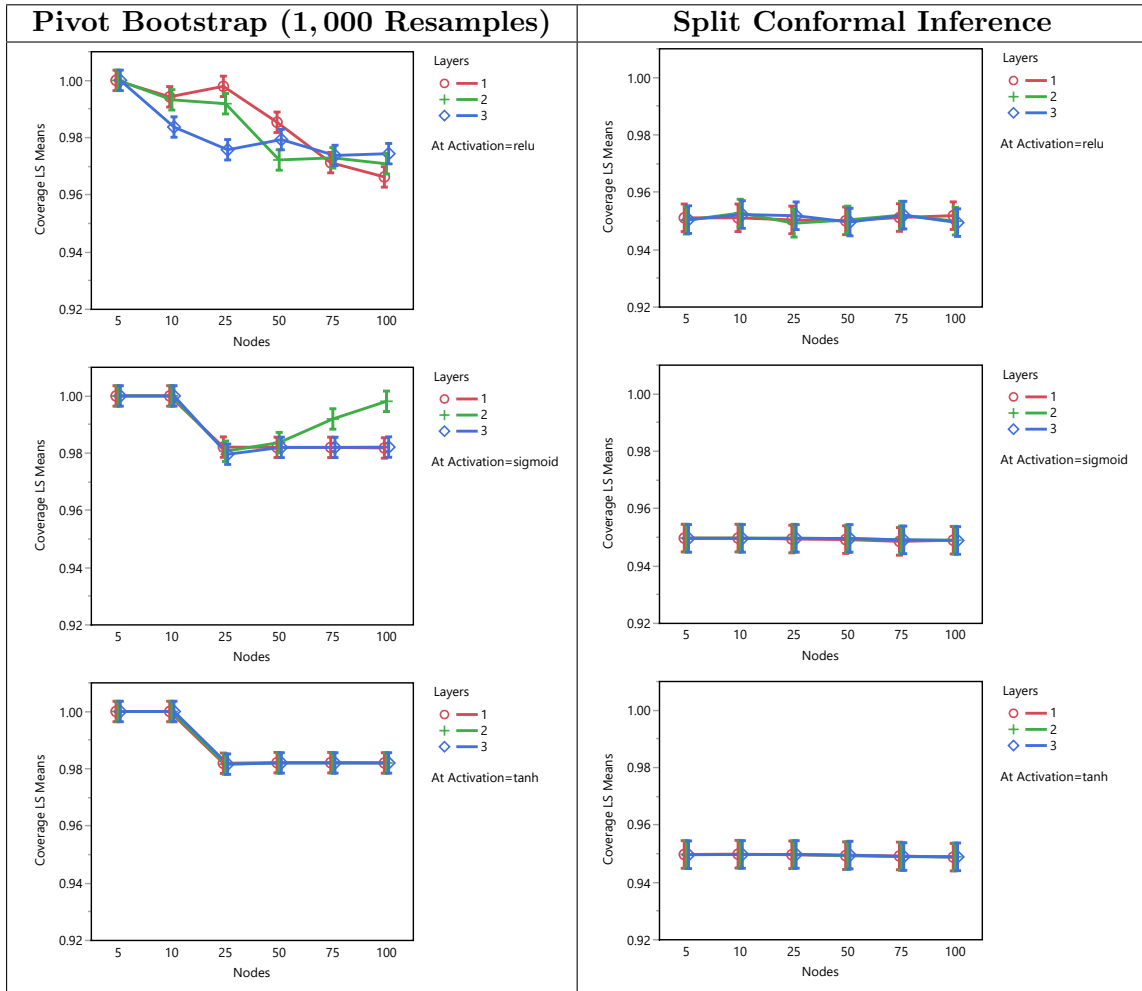


Figure 12: PI Coverage in the Power Plant Data Set.

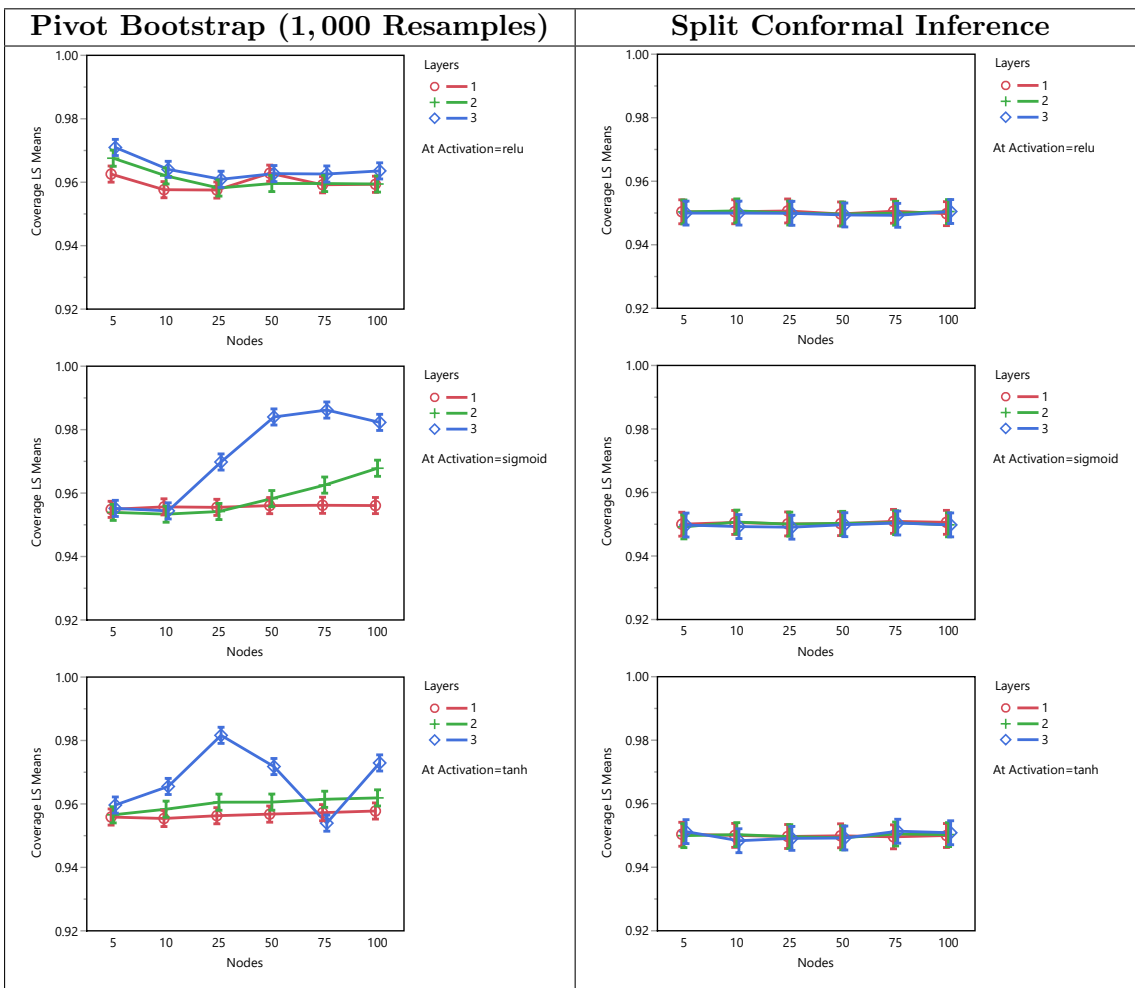


Figure 13: PI Coverage in the Protein Structure Data Set.

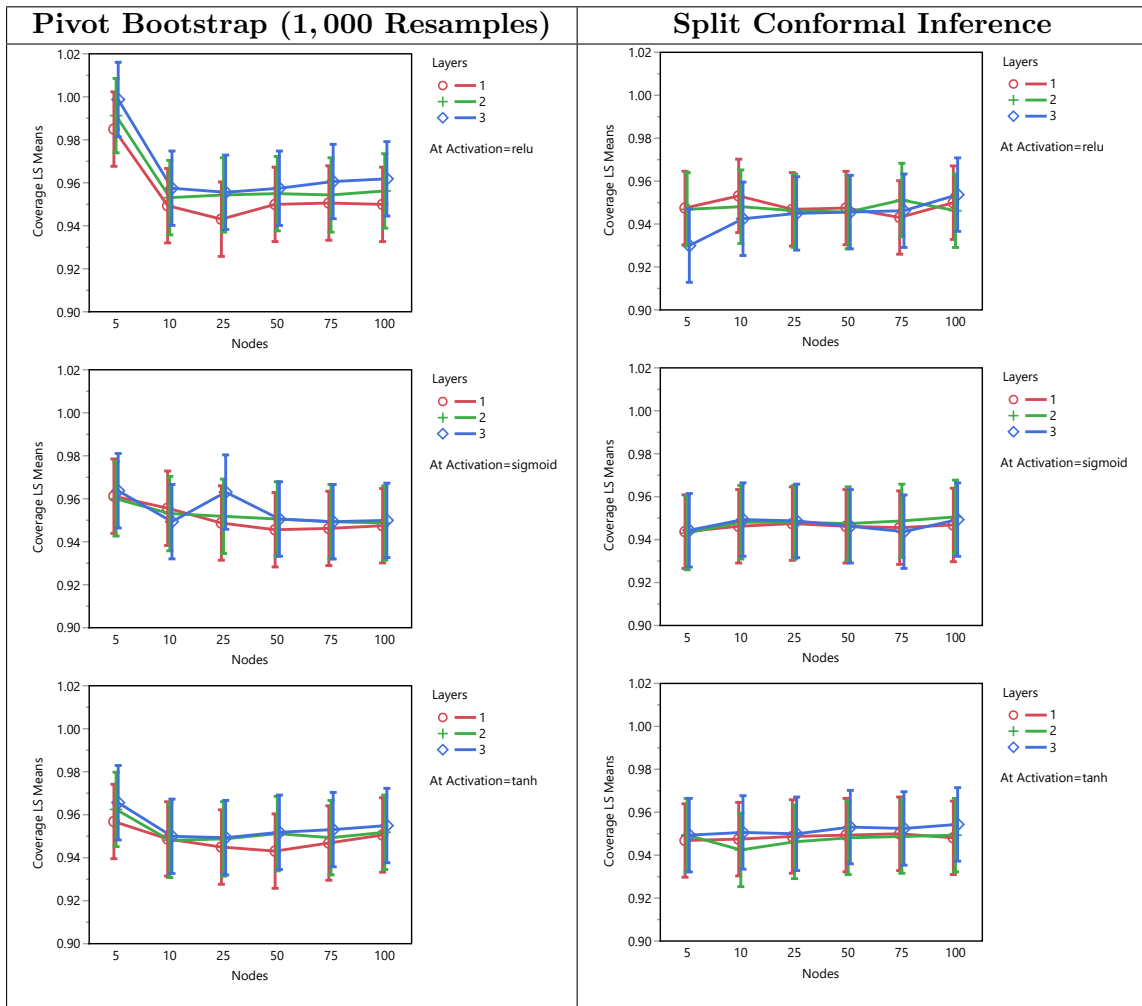


Figure 14: PI Coverage in the Wine Quality Data Set.

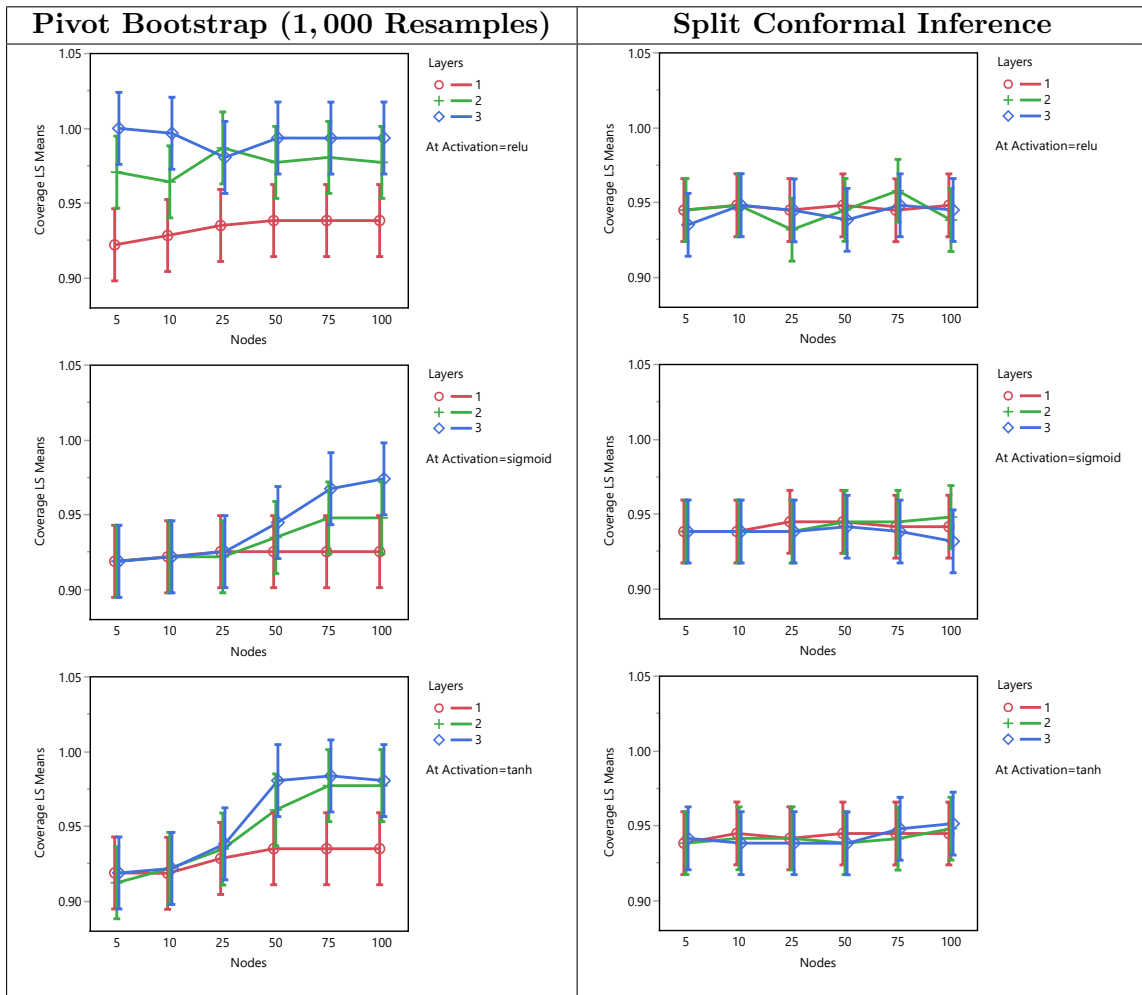


Figure 15: PI Coverage in the Yacht Hydrodynamics Data Set.

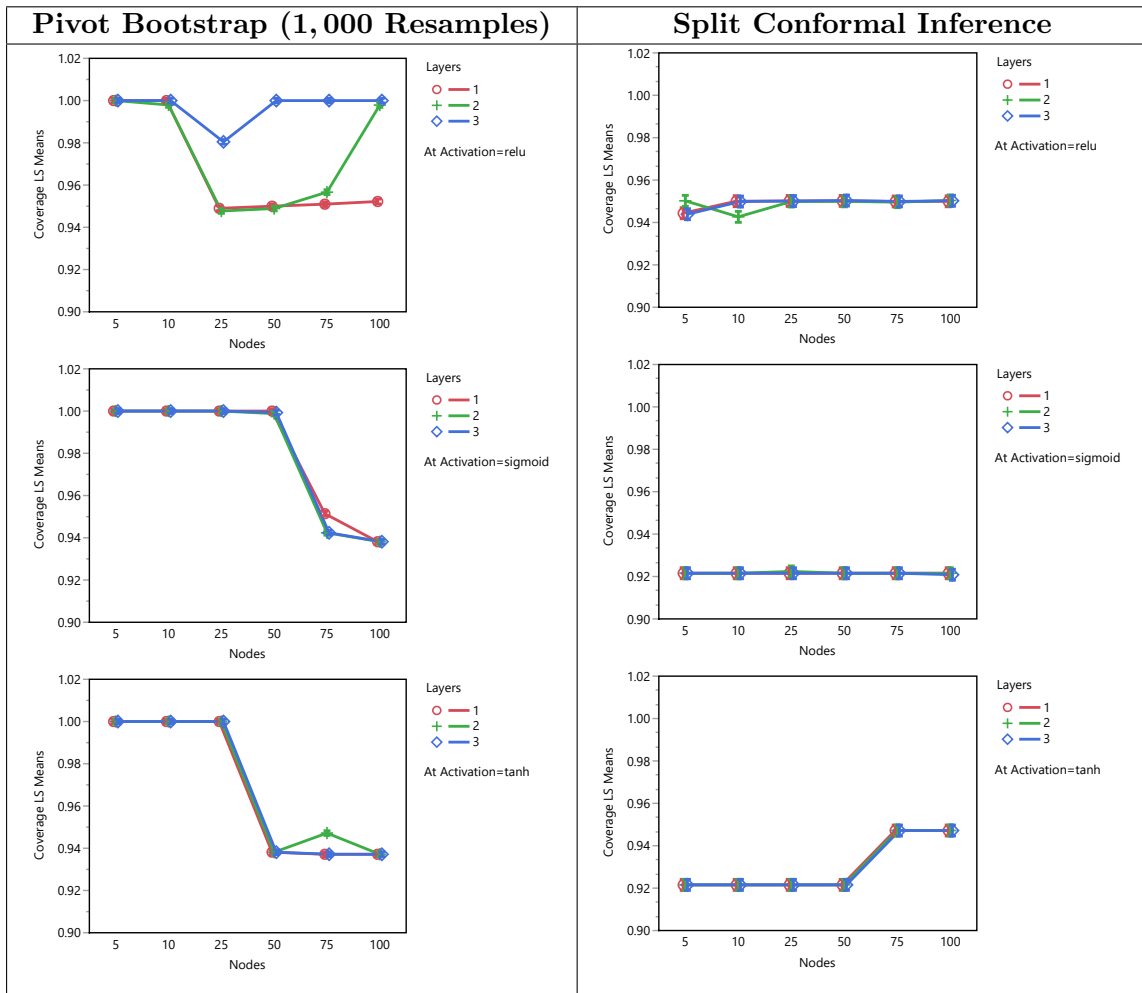


Figure 16: PI Coverage in the Year Prediction Data Set.

Appendix B. PI Average Width Plots

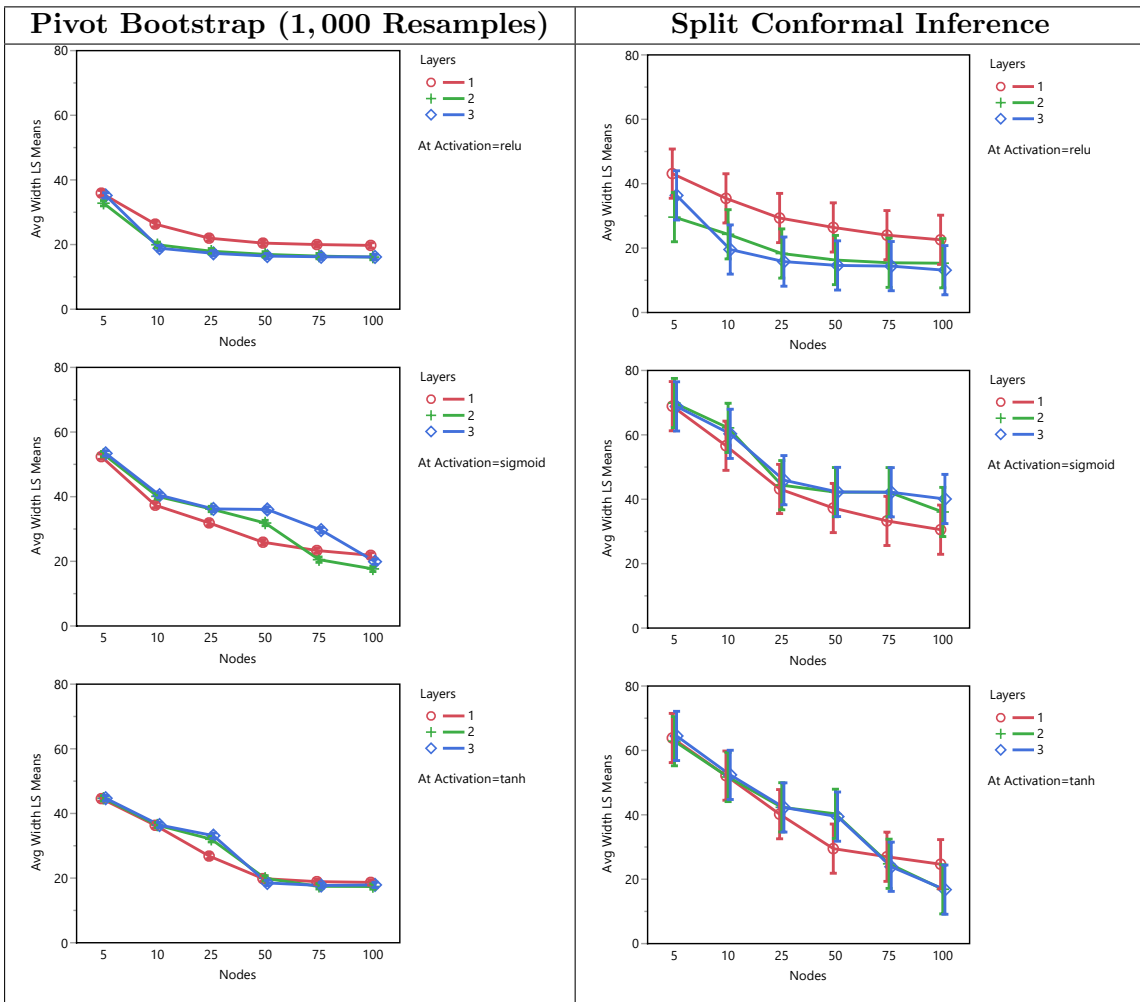


Figure 17: PI Coverage in the Boston Housing Data Set.

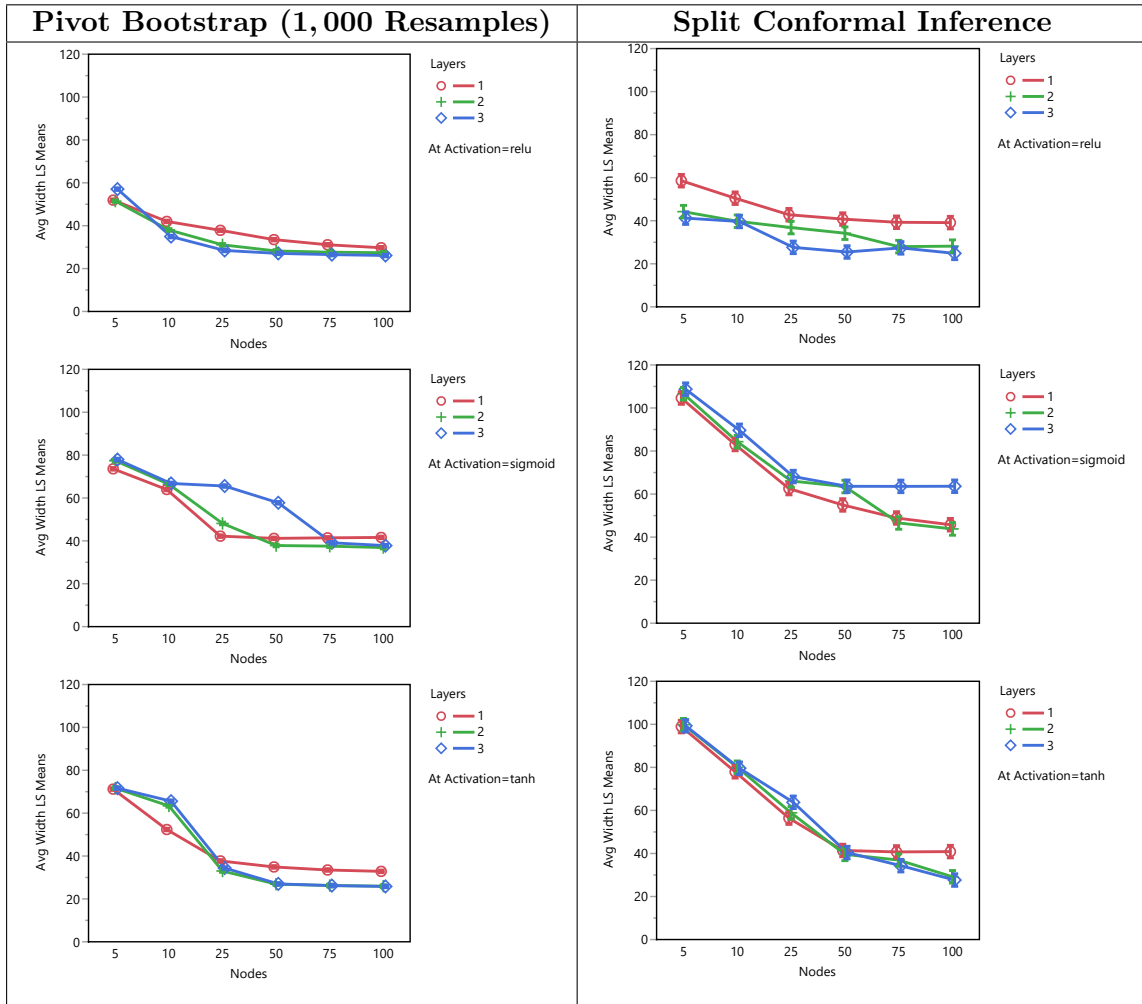


Figure 18: PI Average Width in the Concrete Strength Data Set.

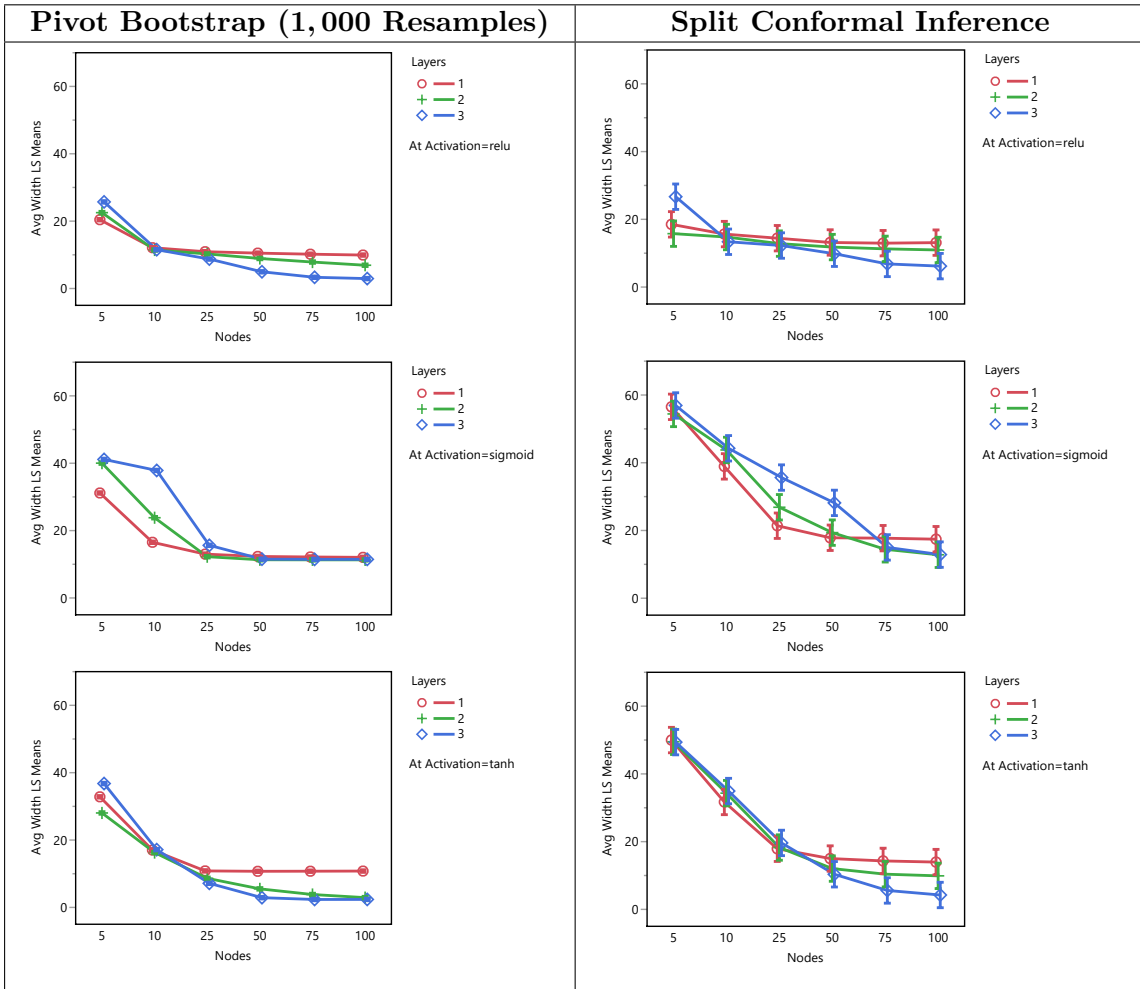


Figure 19: PI Average Width in the Energy Efficiency Data Set.

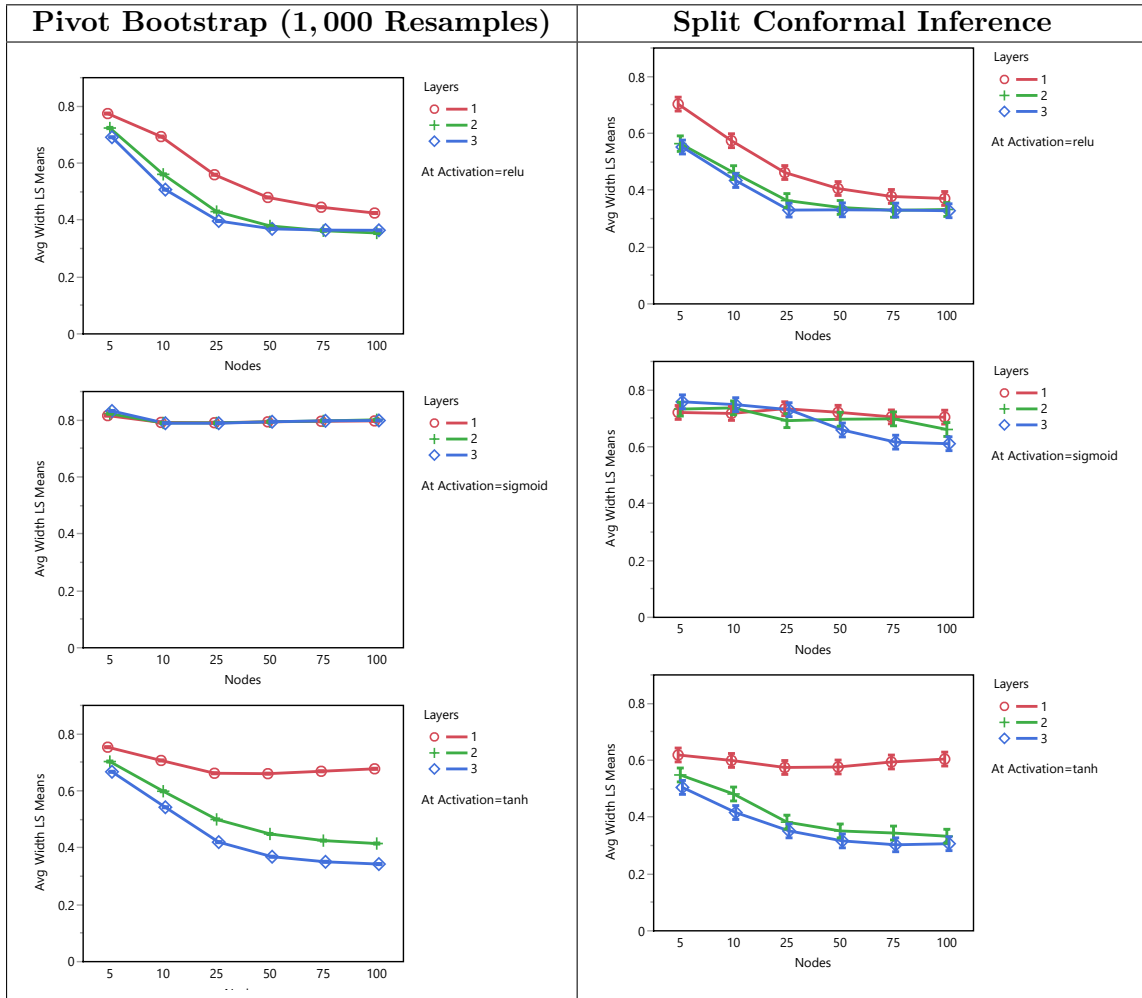


Figure 20: PI Average Width in the Kinematics Data Set.

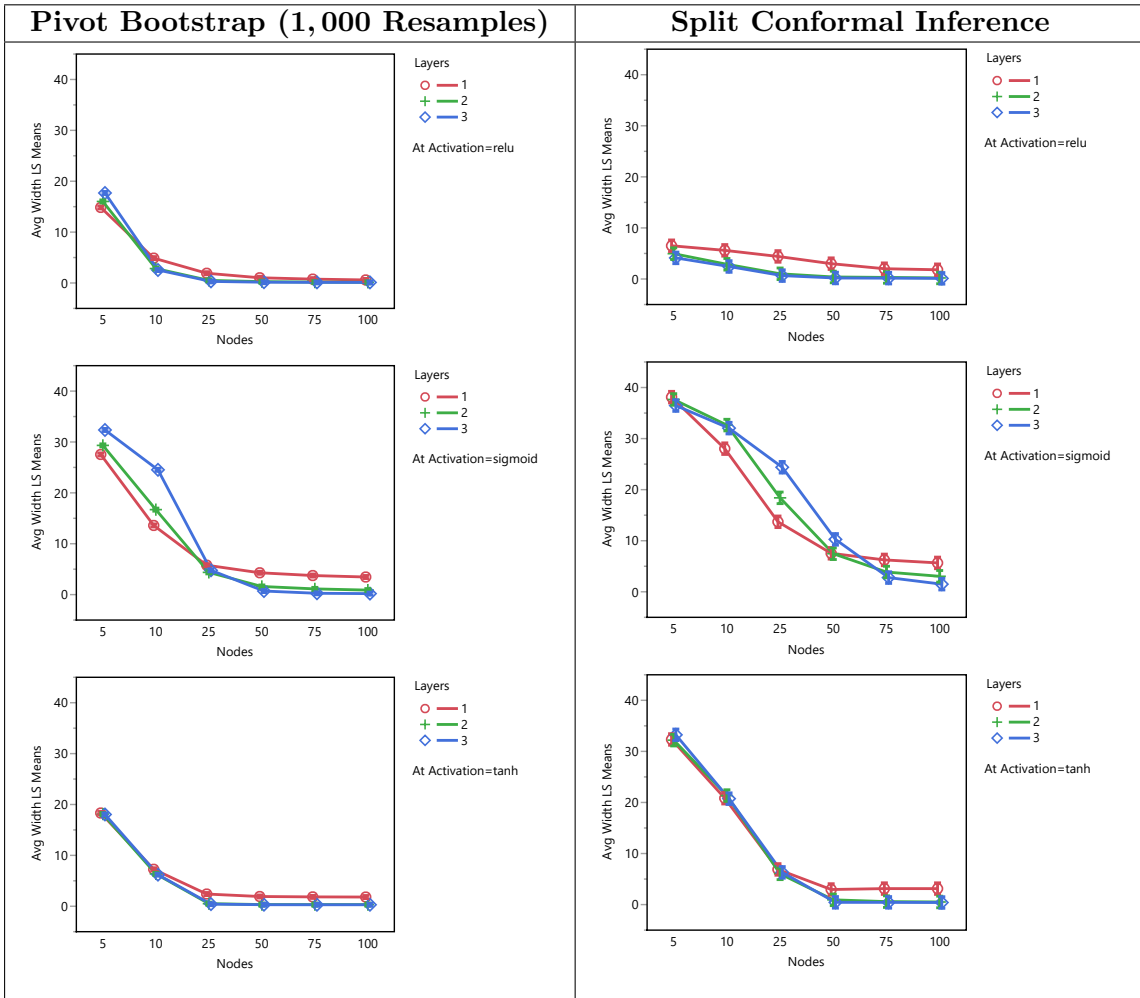


Figure 21: PI Average Width in the Naval Propulsion Data Set.

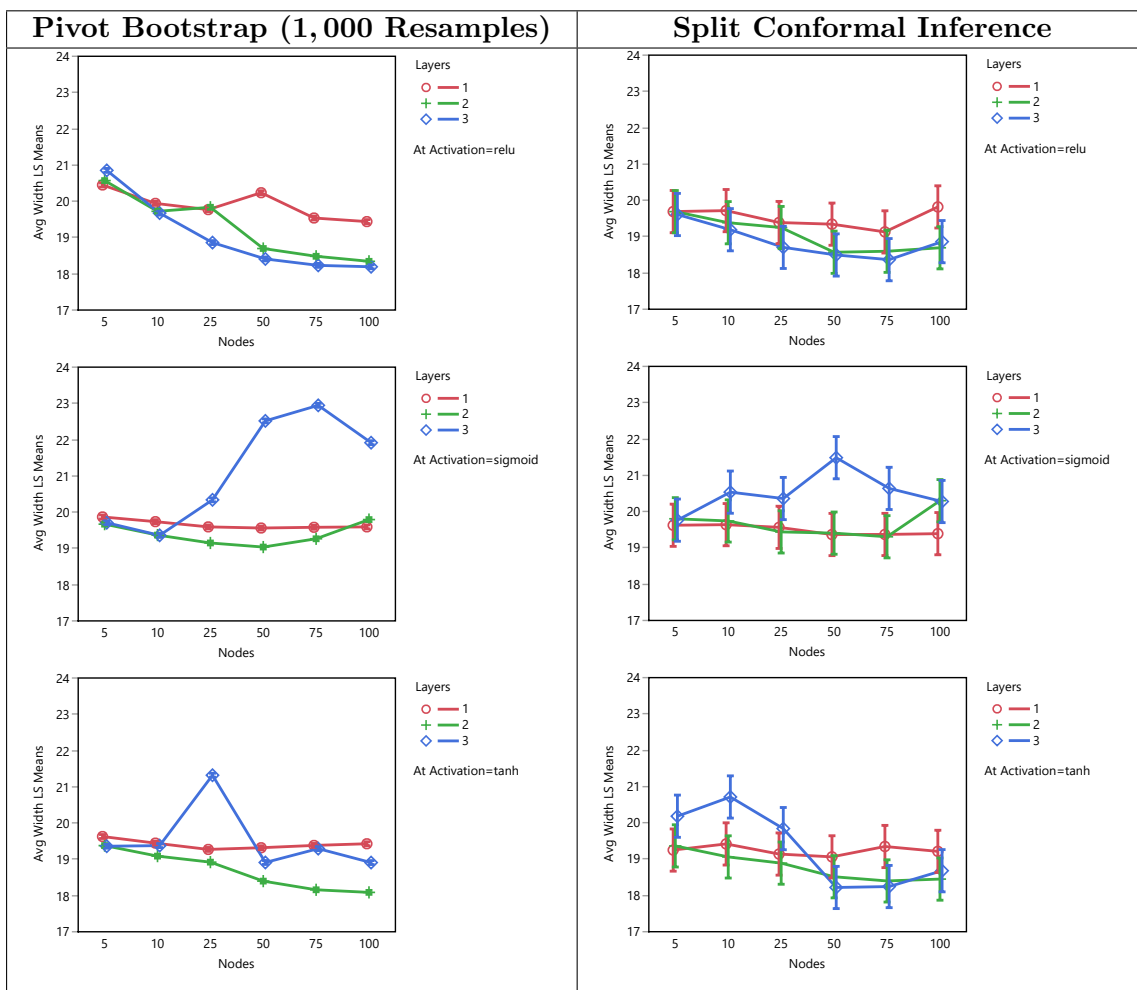


Figure 22: PI Average Width in the Protein Structure Data Set.

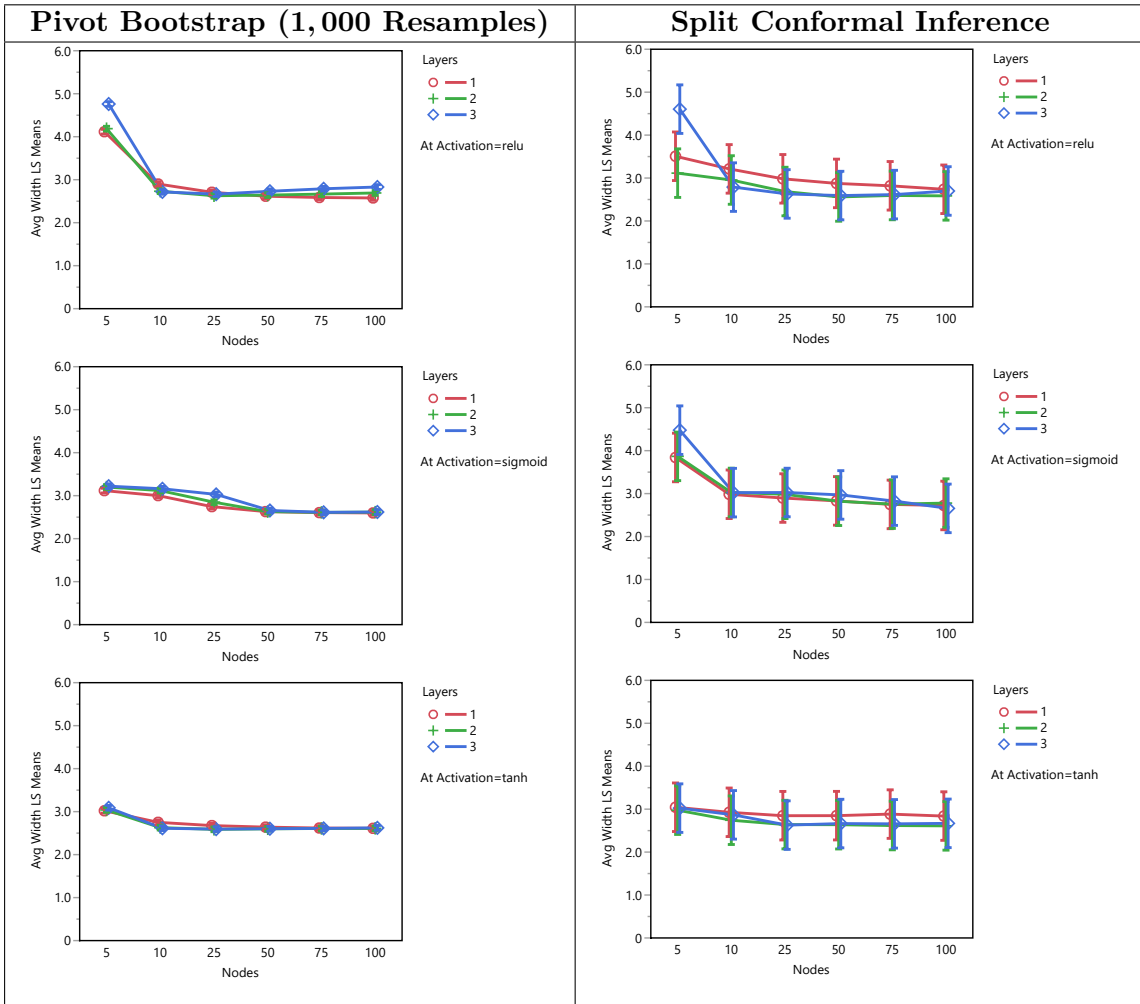


Figure 23: PI Average Width in the Wine Quality Data Set.

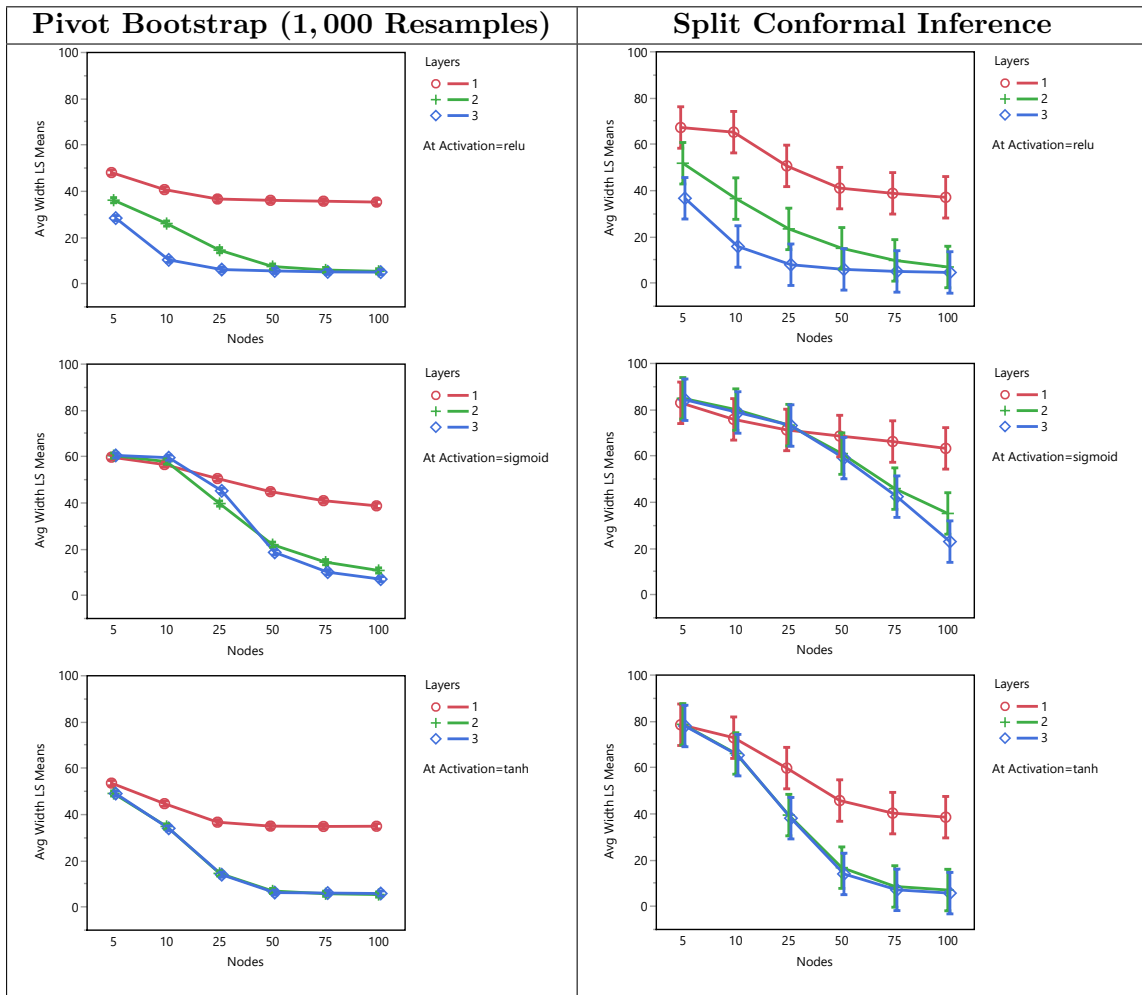


Figure 24: PI Average Width in the Yacht Hydrodynamics Data Set.

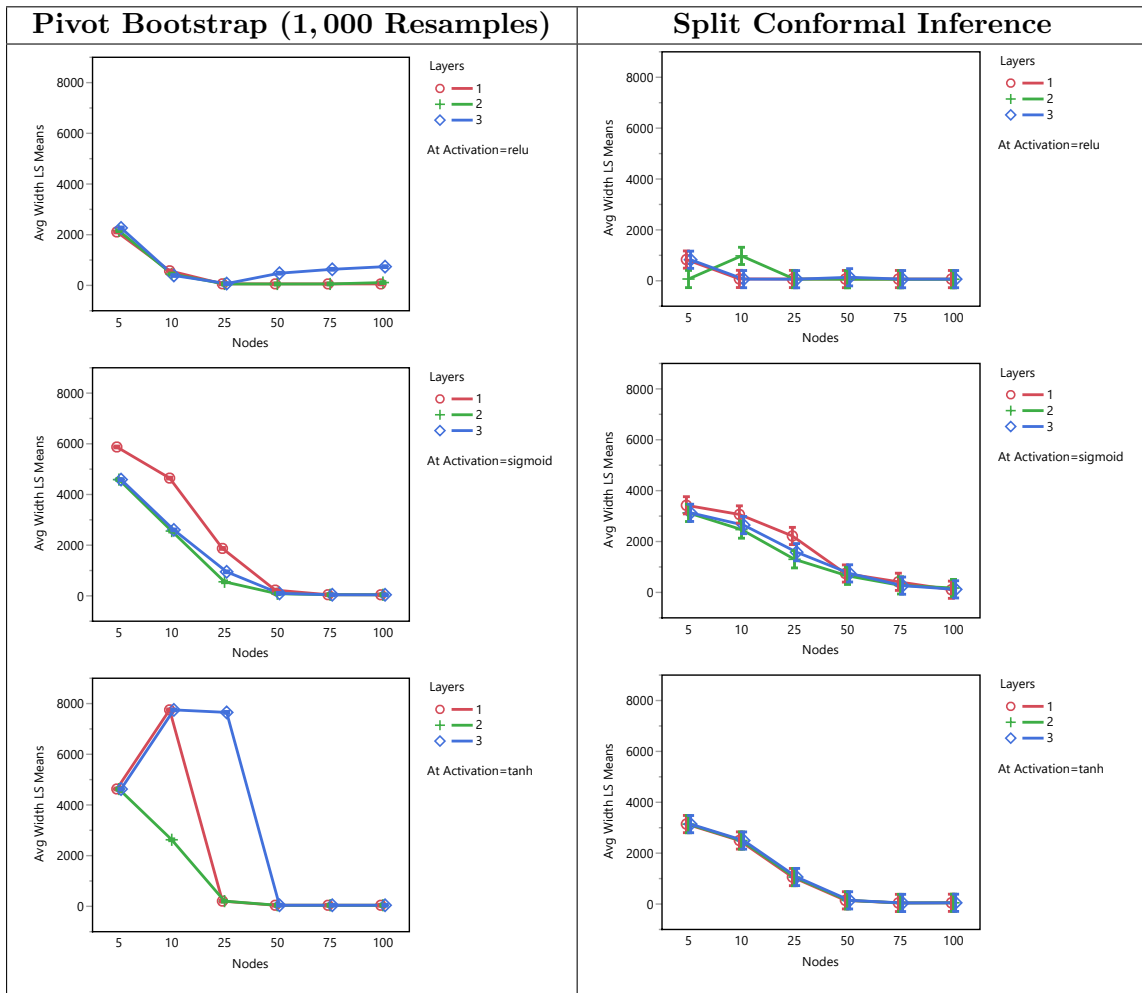


Figure 25: PI Average Width in the Year Prediction Data Set.

Appendix C. Bootstrapped PI Coverage vs RMSE Plots

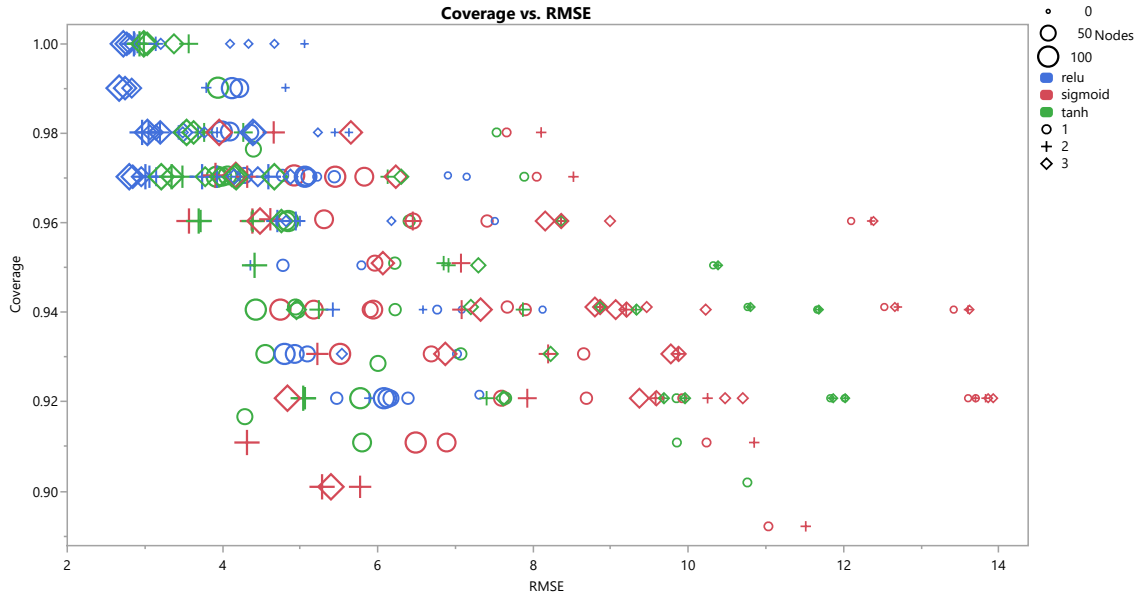


Figure 26: PI Coverage versus RMSE for the Boston Housing Data Set.

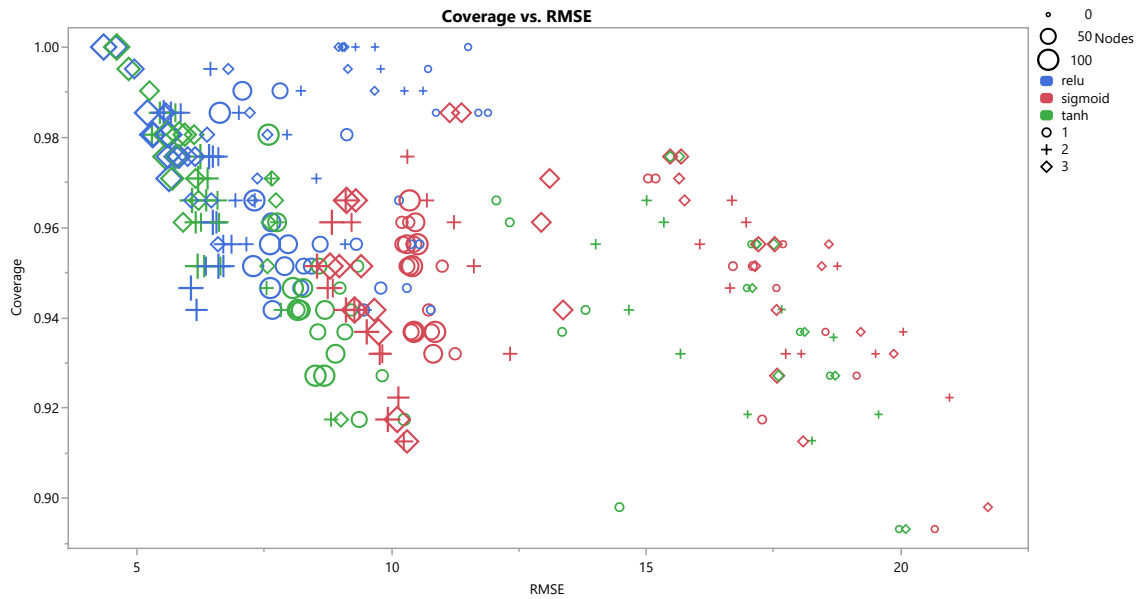


Figure 27: PI Coverage versus RMSE for the Concrete Strength Data Set.

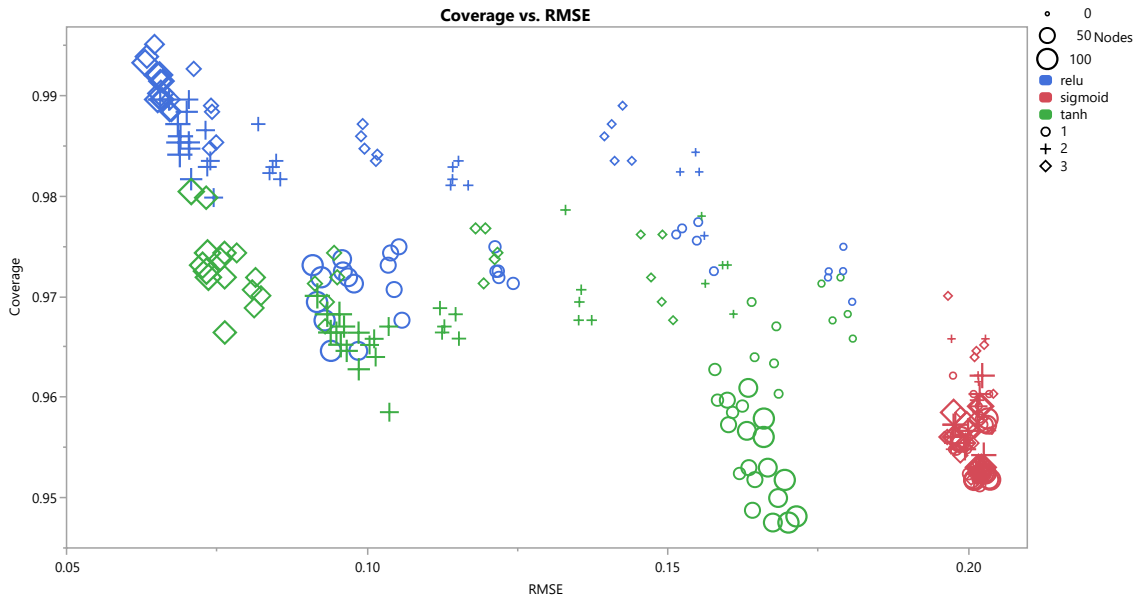


Figure 28: PI Coverage versus RMSE for the Kinematics Data Set.

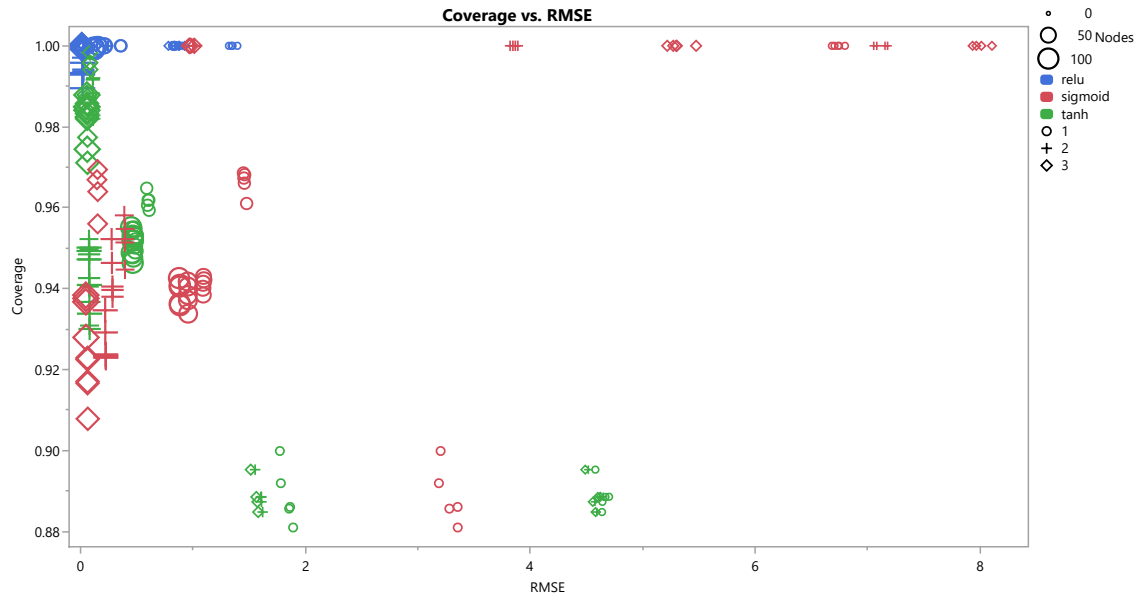


Figure 29: PI Coverage versus RMSE for the Naval Propulsion Data Set.

CONSTRUCTING PREDICTION INTERVALS WITH NEURAL NETWORKS

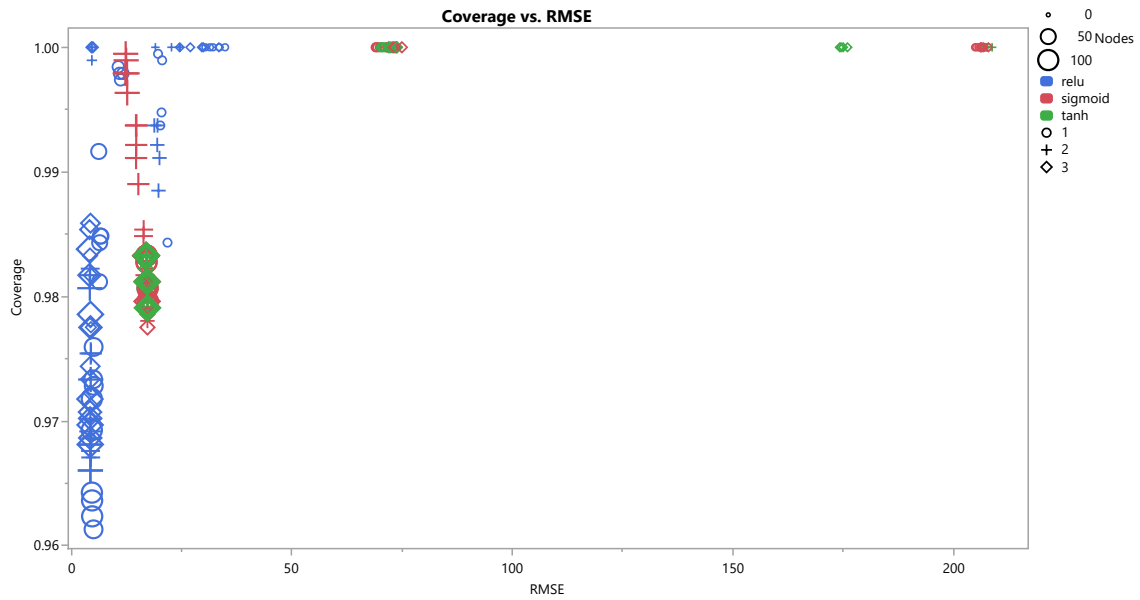


Figure 30: PI Coverage versus RMSE for the Power Plant Data Set.

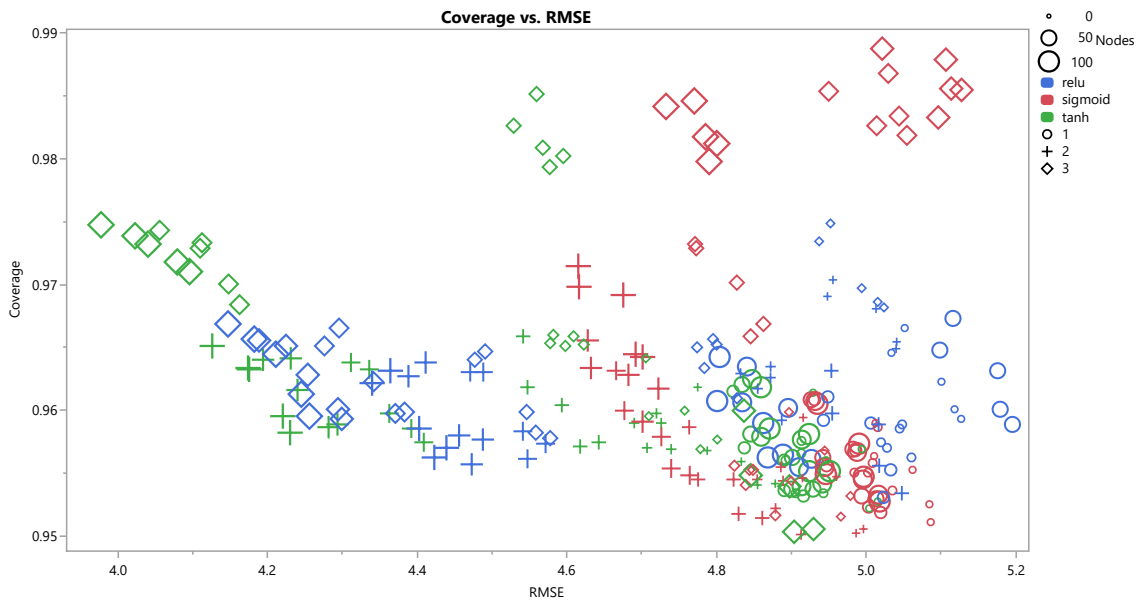


Figure 31: PI Coverage versus RMSE for the Protein Structure Data Set.

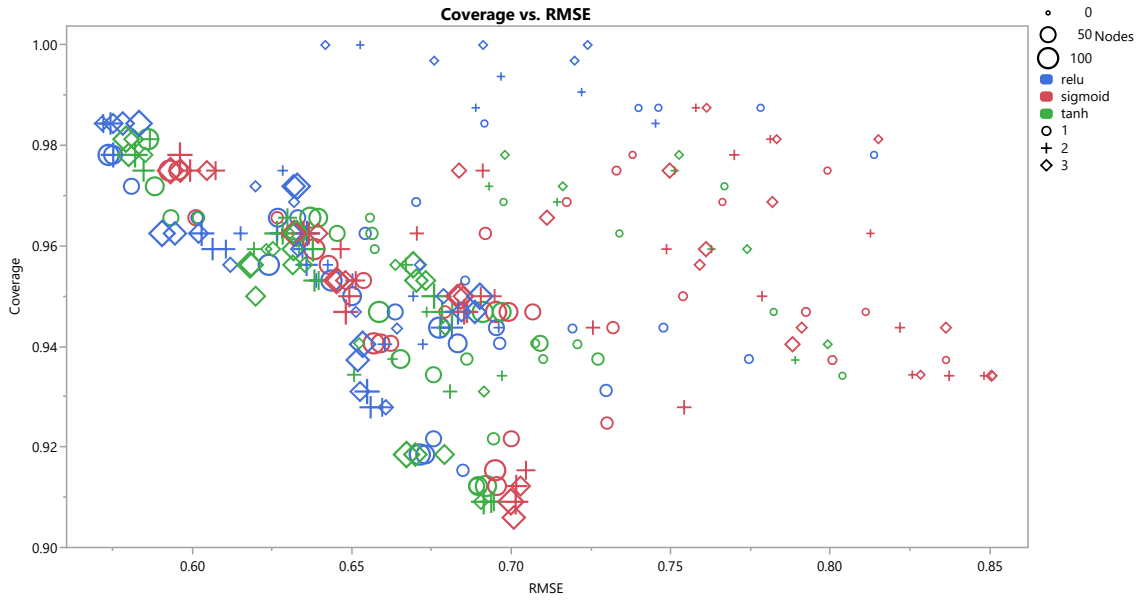


Figure 32: PI Coverage versus RMSE for the Wine Quality Data Set.

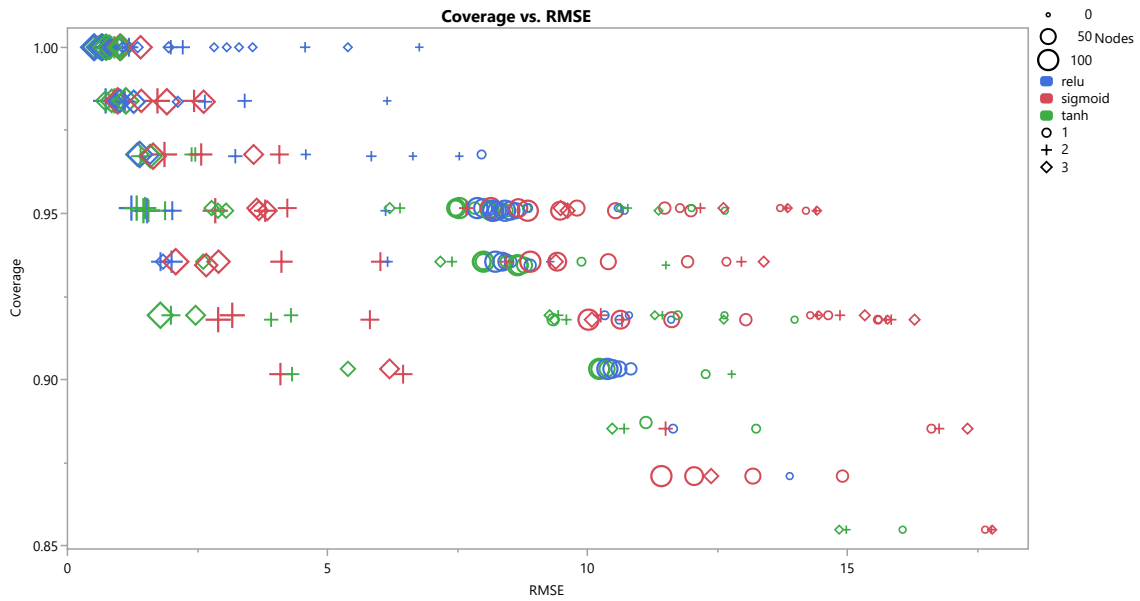


Figure 33: PI Coverage versus RMSE for the Yacht Hydrodynamics Data Set.

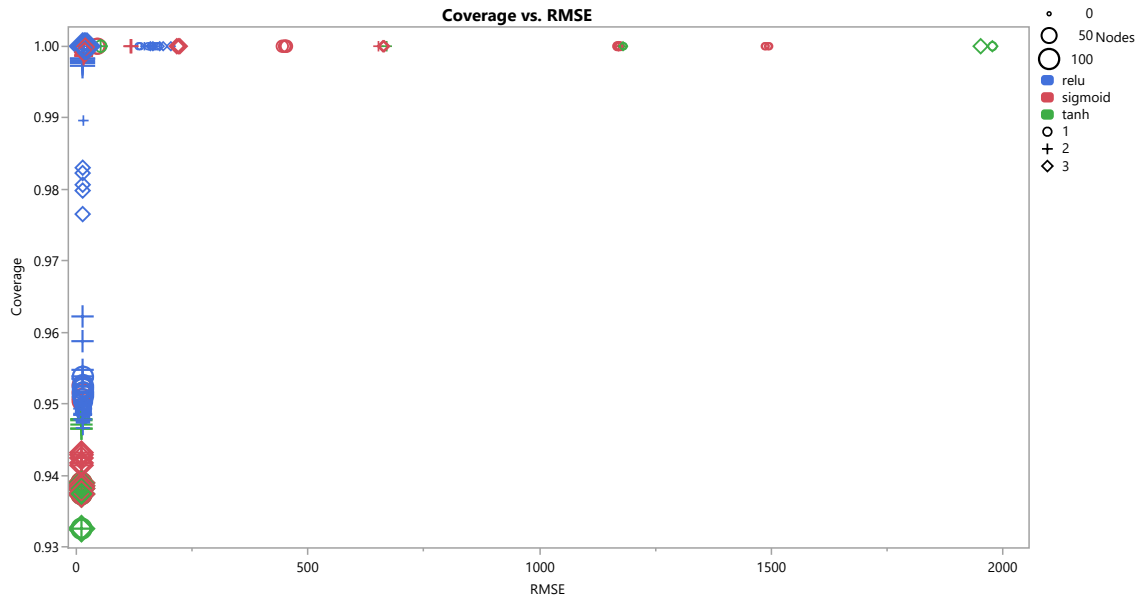


Figure 34: PI Coverage versus RMSE for the Year Prediction Data Set.

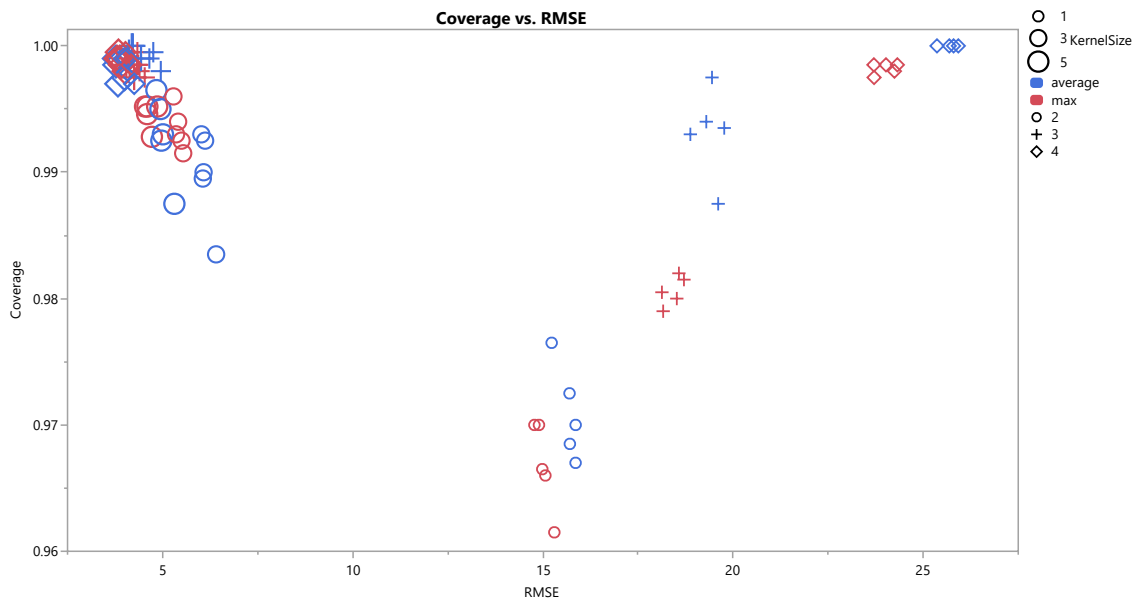


Figure 35: PI Coverage versus RMSE for the RotNIST Data Set.