

A comparison of some conformal quantile regression methods

Matteo Sesia¹ and Emmanuel J. Candès^{1,2}

¹Department of Statistics, Stanford University

²Department of Mathematics, Stanford University

September 13, 2019

Abstract

We compare two recently proposed methods that combine ideas from conformal inference and quantile regression to produce locally adaptive and marginally valid prediction intervals under sample exchangeability (Romano et al., 2019 [1]; Kivaranovic et al., 2019 [2]). First, we prove that these two approaches are asymptotically efficient in large samples, under some additional assumptions. Then we compare them empirically on simulated and real data. Our results demonstrate that the method in Romano et al. (2019) typically yields tighter prediction intervals in finite samples. Finally, we discuss how to tune these procedures by fixing the relative proportions of observations used for training and conformalization.

1 Introduction

1.1 Background and motivation

Given a set of n points $\{(X_i, Y_i)\}_{i=1}^n$, with $Y_i \in \mathbb{R}$ and $X_i \in \mathbb{R}^d$, we consider the problem of constructing a prediction interval for a new point Y_{n+1} based on the observed value of X_{n+1} , assuming only that $\{(X_i, Y_i)\}_{i=1}^{n+1}$ are drawn exchangeably from some common distribution P_{XY} . There exist a vast selection of statistical and machine learning algorithms that can provide approximate answers to this question [3, 4]. However, the uncertainty in any of their predictions cannot be quantified without making strong assumptions and large-sample asymptotic approximations that may not be easily justifiable in applications. Conformal inference [5–13] addresses this problem by constructing an exact *marginal* prediction interval $\hat{C}_\alpha(X_{n+1})$ such that

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

while relying only on the exchangeability of the $n + 1$ points. This interval is said to be *marginal* because all variables in (1) are treated as random, including (X_{n+1}, Y_{n+1}) and the data used to train \hat{C} . Therefore, it is not guaranteed that the interval will cover Y_{n+1} conditional on a particular

observed value of X_{n+1} , or a fixed prediction model \hat{C} . Despite this limitation, conformal prediction intervals are attractive because their coverage is guaranteed on average regardless of the distribution of the data.

Ideally, prediction intervals should be as narrow as possible while maintaining coverage. Let us denote by $q_\alpha(x_{n+1})$ the α -th quantile of the conditional distribution of Y given $X_{n+1} = x_{n+1}$. Then a desirable oracle prediction interval would be

$$C_\alpha^{\text{oracle}}(X_{n+1}) = [q_{\alpha/2}(X_{n+1}), q_{1-\alpha/2}(X_{n+1})]. \quad (2)$$

By construction, this is the narrowest *symmetric* prediction interval that has valid coverage conditional on the value of X_{n+1} . Here, we say that a prediction interval is symmetric if Y_{n+1} is equally likely to be smaller or larger than predicted. Unfortunately, the oracle interval in (2) is unachievable in practice because we do not know $P_{Y|X}$. The goal of conformal quantile regression [1] is to form a practical prediction interval \hat{C}_α that estimates (2) as closely as possible while satisfying (1) exactly. In this work, we compare theoretically and empirically the method from [1] with a similar approach that was proposed independently in [2].

1.2 Conformal quantile regression

Throughout this paper, we follow the split-conformal approach to conformal inference [8, 11, 13] adopted in [1] and [2], since it is computationally feasible even with large data. The first step of the conformal quantile regression method in [1] is to split the data samples into two disjoint subsets, \mathcal{I}_1 and \mathcal{I}_2 . Lower and upper quantile regression functions, namely $\hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2} : \mathbb{R}^d \rightarrow \mathbb{R}$, are fitted on the observations in \mathcal{I}_1 . Any algorithm can be employed for this purpose; for example, one may rely on linear regression [14], neural networks [15] or random forests [16]. In any case, this algorithm is treated as a black box. The estimated quantile functions are used to compute a *conformity score* for each $i \in \mathcal{I}_2$:

$$E_i^{\text{CQR}} = \max \{ \hat{q}_{\alpha/2}(X_i) - Y_i, Y_i - \hat{q}_{1-\alpha/2}(X_i) \}. \quad (3)$$

Then, with $\hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2)$ defined as the $\lceil (1-\alpha)(|\mathcal{I}_2|+1) \rceil$ -th largest element of $\{E_i\}_{i \in \mathcal{I}_2}$, the conformal prediction interval for X_{n+1} is given by

$$\hat{C}_\alpha^{\text{CQR}}(X_{n+1}) = \left[\hat{q}_{\alpha/2}(X_{n+1}) - \hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2), \hat{q}_{1-\alpha/2}(X_{n+1}) + \hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2) \right] \quad (4)$$

This method is summarized in Algorithm 1, where it is denoted as CQR. It is shown in [1] that $\hat{C}_\alpha^{\text{CQR}}(X_{n+1})$ has marginal coverage at level $1-\alpha$.

The method described in [2] differs from CQR in the choice of the conformity scores, as outlined in Algorithm 1 as CQR-m. Instead of (3), one computes¹

$$E_i^{\text{CQR-m}} = \max \left\{ \frac{\hat{q}_{\alpha/2}(X_i) - Y_i}{\hat{q}_{1/2}(X_i) - \hat{q}_{\alpha/2}(X_i)}, \frac{Y_i - \hat{q}_{1-\alpha/2}(X_i)}{\hat{q}_{1-\alpha/2}(X_i) - \hat{q}_{1/2}(X_i)} \right\}, \quad (5)$$

¹Note that we present CQR-m with a slightly different notation than in [2] to facilitate the comparison.

Algorithm 1: Conformal quantile regression

Input:

data $\{(X_i, Y_i)\}_{i=1}^n$, covariates for new sample X_{n+1} ;
proportion of data for training $\gamma \in (0, 1)$;
quantile regression algorithm \hat{q} ;
conformalization method $\psi \in \{\text{CQR}, \text{CQR-m}, \text{CQR-r}\}$;
coverage level $\alpha \in (0, 1)$.

Procedure:

randomly split $\{1, \dots, n\}$ into $\mathcal{I}_1, \mathcal{I}_2$, of size $|\mathcal{I}_1| = \gamma n$, $|\mathcal{I}_2| = n - |\mathcal{I}_1|$;
fit the quantile regression functions $\hat{q}_{\alpha/2}$ and $\hat{q}_{1-\alpha/2}$ on $\{(X_i, Y_i)\}_{i \in \mathcal{I}_1}$;

if $\psi = \text{CQR}$ **then**

 compute the conformity scores E_i^{CQR} for each $i \in \mathcal{I}_2$, as in (3);
 compute $\hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2)$;
 compute the prediction interval $\hat{C}_\alpha(X_{n+1})$, as in (4).

else if $\psi = \text{CQR-m}$ **then**

 fit the median regression function $\hat{q}_{1/2}$ on $\{(X_i, Y_i)\}_{i \in \mathcal{I}_1}$;
 compute the conformity scores $E_i^{\text{CQR-m}}$ for each $i \in \mathcal{I}_2$, as in (5);
 compute $\hat{Q}_{1-\alpha}(E^{\text{CQR-m}}; \mathcal{I}_2)$;
 compute the prediction interval $\hat{C}_\alpha(X_{n+1})$, as in (6).

else if $\psi = \text{CQR-r}$ **then**

 compute the conformity scores $E_i^{\text{CQR-r}}$ for each $i \in \mathcal{I}_2$, as in (7);
 compute $\hat{Q}_{1-\alpha}(E^{\text{CQR-r}}; \mathcal{I}_2)$;
 compute the prediction interval $\hat{C}_\alpha(X_{n+1})$, as in (8).

Output:

A prediction interval $\hat{C}_\alpha(X_{n+1})$.

where $\hat{q}_{1/2}$ indicates an estimated median regression function obtained with the same black-box algorithm as $\hat{q}_{\alpha/2}$ and $\hat{q}_{1-\alpha/2}$. Then the conformal prediction interval for X_{n+1} is given by:

$$\begin{aligned} \hat{C}_\alpha^{\text{CQR-m}}(X_{n+1}) &= \left[\hat{q}_{\alpha/2}(X_{n+1}) - \hat{\Delta}_{\alpha, \text{lo}}^{\text{CQR-m}}, \hat{q}_{1-\alpha/2}(X_{n+1}) + \hat{\Delta}_{\alpha, \text{up}}^{\text{CQR-m}} \right] \\ \hat{\Delta}_{\alpha, \text{lo}}^{\text{CQR-m}} &= \hat{Q}_{1-\alpha}(E^{\text{CQR-m}}; \mathcal{I}_2) \left[\hat{q}_{1/2}(X_{n+1}) - \hat{q}_{\alpha/2}(X_{n+1}) \right], \\ \hat{\Delta}_{\alpha, \text{up}}^{\text{CQR-m}} &= \hat{Q}_{1-\alpha}(E^{\text{CQR-m}}; \mathcal{I}_2) \left[\hat{q}_{1-\alpha/2}(X_{n+1}) - \hat{q}_{1/2}(X_{n+1}) \right]. \end{aligned} \quad (6)$$

One can show that $\hat{C}_\alpha^{\text{CQR-m}}(X_{n+1})$ also has marginal coverage at level $1 - \alpha$ [2].

We also find it interesting to consider a modified version of CQR-m that does not require estimating the regression median.² This third approach, listed in Algorithm 1 as CQR-r, is based on the following conformity scores:

$$E_i^{\text{CQR-r}} = \max \left\{ \frac{\hat{q}_{\alpha/2}(X_i) - Y_i}{\hat{q}_{1-\alpha/2}(X_i) - \hat{q}_{\alpha/2}(X_i)}, \frac{Y_i - \hat{q}_{1-\alpha/2}(X_i)}{\hat{q}_{1-\alpha/2}(X_i) - \hat{q}_{\alpha/2}(X_i)} \right\}. \quad (7)$$

²This was first suggested by Yaniv Romano through personal communication.

The CQR-r prediction intervals are

$$\begin{aligned}\hat{C}_\alpha^{\text{CQR-r}}(X_{n+1}) &= \left[\hat{q}_{\alpha/2}(X_{n+1}) - \hat{\Delta}_\alpha^{\text{CQR-r}}, \hat{q}_{1-\alpha/2}(X_{n+1}) + \hat{\Delta}_\alpha^{\text{CQR-r}} \right] \\ \hat{\Delta}_\alpha^{\text{CQR-r}} &= \hat{Q}_{1-\alpha}(E^{\text{CQR-r}}; \mathcal{I}_2) \left[\hat{q}_{1-\alpha/2}(X_{n+1}) - \hat{q}_{\alpha/2}(X_{n+1}) \right].\end{aligned}\tag{8}$$

It is easy to show that $\hat{C}_\alpha^{\text{CQR-r}}(X_{n+1})$ also attains marginal coverage at level $1 - \alpha$. A proof is omitted because it would be identical to those in [1] and [2]. CQR-r is similar in spirit to CQR-m, but it has a more direct interpretation: the conformity scores of CQR-r in (7) weight the distance of Y from the corresponding prediction interval by the inverse width of the interval. Therefore, the conformalization expands or contracts the black-box prediction bands proportionally to their width, instead of adding a constant shift as in CQR. Since it is not clear how the regression median $q_{1/2}$ should generally be related to the upper and lower α -quantiles of $P_{Y|X}$, we find this approach slightly more intuitive than CQR-m.

2 Theoretical analysis

We show that the output of Algorithm 1 converges to the oracle bands in (2) as n grows, if the black-box quantile regression estimates are consistent and a few additional assumptions hold. This can be established for any of the three alternative types of conformity scores discussed in this paper, which are therefore asymptotically equivalent in this sense.

Assumption 1 (i.i.d.). *The points $\{(X_i, Y_i)\}_{i=1}^{n+1}$ are drawn i.i.d. from some distribution P_{XY} .*

Assumption 2 (regularity). *The probability density of the conformity scores, either in (3), (5) or (7), depending on the conformalization method in Algorithm 1, is bounded away from zero in an open neighborhood of zero.*

Assumption 3 (consistency). *For simplicity, denote by n the size of the training data set \mathcal{I}_1 used to fit the quantile regression functions \hat{q} . Let X be a new observation independent of \mathcal{I}_1 . Then the assumption is that, for n large enough,*

$$\begin{aligned}\mathbb{P} \left[\mathbb{E} \left[\left(\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X) \right)^2 \mid \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2} \right] \leq \eta_n \right] &\geq 1 - \rho_n, \\ \mathbb{P} \left[\mathbb{E} \left[\left(\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X) \right)^2 \mid \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2} \right] \leq \eta_n \right] &\geq 1 - \rho_n,\end{aligned}$$

for some sequences $\eta_n = o(1)$ and $\rho_n = o(1)$, as $n \rightarrow \infty$.

Assumption 3 is similar to that used in [13] for mean regression estimators, and it is weaker than requiring $\hat{q}_{\alpha/2}(X) \xrightarrow{L^2} q_{\alpha/2}(X)$ and $\hat{q}_{1-\alpha/2}(X) \xrightarrow{L^2} q_{1-\alpha/2}(X)$ as $n \rightarrow \infty$, by Markov's inequality.

Theorem 1. *Under Assumptions 1–3, the conformal quantile regression bands \hat{C}_α obtained with Algorithm 1 satisfy*

$$L \left(\hat{C}_\alpha(X_{n+1}) \triangle C_\alpha^{\text{oracle}}(X_{n+1}) \right) = o_{\mathbb{P}}(1),$$

as $|\mathcal{I}_1|, |\mathcal{I}_2| \rightarrow \infty$. Here, $L(A)$ indicates the Lebesgue measure of the set A , and \triangle is the symmetric difference operator, i.e., $A \triangle B = (A \setminus B) \cup (B \setminus A)$.

The proof of Theorem 1 can be found in Appendix A and is inspired by that of Theorem 3.4 in [13], although the oracle and the conformalization methods considered here are different. Theorem 1 establishes a stronger form of statistical efficiency for conformal quantile regression compared to the result in [13], which assumes $Y = \mu(X) + \epsilon$, for some regression function μ , and homoscedastic noise ϵ . In general, the conformal prediction intervals described in [13] will not converge to those of our oracle if the noise is heteroscedastic, regardless of the consistency of the black-box regression estimator $\hat{\mu}$. By contrast, conformal quantile regression is efficient in the sense that, under Theorem 1, the prediction bands converge to those of the oracle, which are the narrowest possible bands achieving conditional coverage. Finally, the asymptotic consistency assumption may be verified theoretically for some specific algorithms under certain conditions, e.g. random forests [16]. In any case, our result provides some theoretical backing to conformal quantile regression even if the assumptions cannot be verified in practice.

As an immediate corollary of Theorem 1, note that it also follows that conformal quantile regression bands have asymptotic conditional coverage, which we define as in [13].

Definition 1 (Asymptotic conditional coverage). *We say that a sequence \hat{C}_n of random prediction bands has asymptotic conditional coverage at the level $1 - \alpha$ if there exists a sequence of random sets $\Lambda_n \subseteq \mathbb{R}^d$ such that $\mathbb{P}[X \in \Lambda_n] = 1 - o_{\mathbb{P}}(1)$ and*

$$\sup_{x \in \Lambda_n} \left| \mathbb{P} \left[Y \in \hat{C}_n(x) \mid X = x \right] - (1 - \alpha) \right| = o_{\mathbb{P}}(1).$$

Despite being asymptotically efficient under Assumptions 1–3, the three conformalization methods in Algorithm 1 typically perform differently with finite data, as discussed next.

3 Empirical comparison

The data and code used in this section are on <https://github.com/mnesia/cqr-comparison>.

3.1 Black-box quantile regression

In the following, we utilize two alternative black-box quantile regressors, implemented and trained as in [1]. The first procedure is based on quantile regression forests [16]. We fit 1000 trees and set the other tuning parameters equal to their default values. The second black box is a neural network [15] with three fully connected layers and ReLU non-linearities. We have chosen this design, which is slightly different from that in [2], because it leads to conformal prediction intervals that are tighter than those reported in [2]. If the estimated lower and upper quantiles overlap, which may sometimes occur, we swap them. The nominal level of the black boxes is tuned so that their empirical coverage, estimated by cross-validation, is approximately equal to $1 - 2\alpha$. We have observed that this heuristic generally leads to tighter conformal intervals compared to those obtained by directly requesting the black boxes to estimate $q_{\alpha/2}$ and $q_{1-\alpha/2}$; see Section 3.3 and Appendix B for empirical evidence. Throughout this section, we set $\alpha = 0.1$.

3.2 Experiments with artificial data

We begin by considering the same experiment based on artificial data as in [2]. We simulate $X \sim \text{Unif}([0, 1]^d)$, for $d = 100$, and $Y \in \mathbb{R}$ from:

$$Y = f(\beta'X) + \epsilon\sqrt{1 + (\beta'X)^2}, \tag{9}$$

where $f(x) = 2\sin(\pi x) + \pi x$, $\beta' = (1, 1, 1, 1, 1, 0, \dots, 0)$ and ϵ is independent standard Gaussian noise. Here, we have access to a natural benchmark: the oracle that knows $P_{Y|X}$ exactly. It follows from (9) that the expected width of the oracle prediction bands is:

$$\mathbb{E} [q_{1-\alpha/2}(X) - q_{\alpha/2}(X)] = 2 \mathbb{E} \left[\sqrt{1 + (\beta'X)^2} \right] Q_{1-\alpha/2}(\epsilon) \approx 8.91,$$

where $Q_\alpha(\epsilon)$ is the α -quantile of the standard Gaussian distribution.

The performances of CQR, CQR-m, and CQR-r are compared in Figure 1 as a function of the number of data points n . The proportion of observations used to train the black box is 3/4, as in [2]. The coverage and average width of the prediction bands is evaluated on an independent test set of size 20,000. The experiment is repeated for 100 independent realizations of the data and of the test set. The width and coverage of the conformal prediction bands approach those of the oracle as the sample size increases. This suggests that the estimated black-box quantiles may be approximately consistent. However, CQR typically produces narrower bands compared to the other methods.

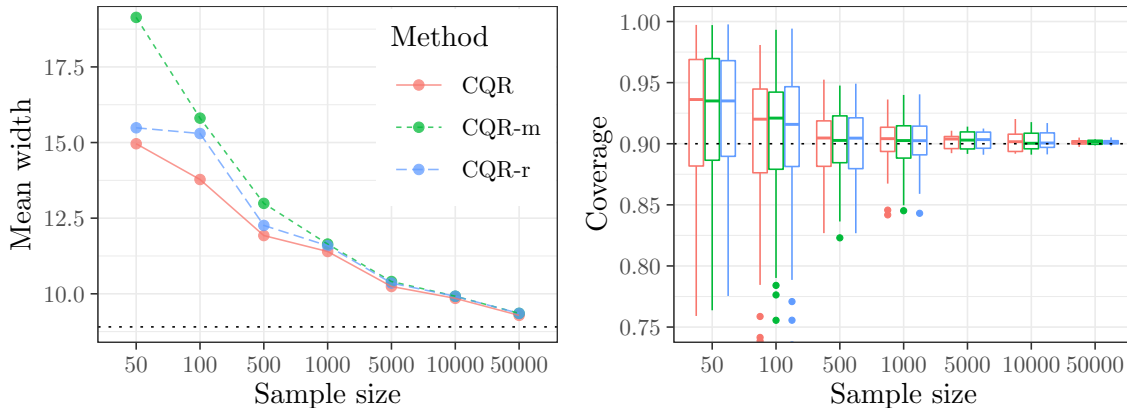


Figure 1: Conformal prediction bands obtained with different conformalization methods on artificial data, as a function of the sample size. The black dotted line on the left indicates the width of the oracle predictions. The black dotted line on the right indicates the nominal coverage level (90%).

3.3 Experiments with real data

We now apply Algorithm 1 on the same data analyzed in [1] and [2].³ Some details about these data sets and information on the corresponding sources are summarized in Table 1. For all data

³We have excluded the X-ray data in [2] because we are unsure of how to replicate the pre-processing.

sets except *homes*, we randomly hold out 20% of the samples for testing. Then we divide the remaining observations into two disjoint sets, \mathcal{I}_1 and \mathcal{I}_2 , to train the black box and conformalize the prediction bands, respectively. The response variables Y are standardized as in [1] and [2]. We explore different values of the fraction of samples used for training: $|\mathcal{I}_1| = \gamma n$, with $\gamma \in \{0.1, 0.25, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.98\}$. We are interested in this comparison because different values are used in [1] and [2]: $\gamma = 0.5$ and $\gamma = 0.75$, respectively. In the case of the *homes* data, we follow in the footsteps of [2]: first, we randomly hold out 3613 test samples; then, we train the black box on 15,000 samples and conformalize on the remaining 3000. All experiments are repeated 10 times, starting from the data splitting.

Name	Description	n	d	Source
bike	bike sharing	10886	18	[17]
bio	physicochemical properties of protein tertiary structures	45730	9	[18]
blog	blog feedback	52397	280	[19]
community	community and crime	1994	100	[20]
concrete	concrete compressive strength	1030	8	[21]
facebook 1	facebook comment volume	40948	53	[22]
facebook 2	facebook comment volume	81311	53	[22]
homes	sale prices of homes in King County, Washington	21613	19	[23]
meps 19	medical expenditure panel survey	15785	139	[24]
meps 20	medical expenditure panel survey	17541	139	[25]
meps 21	medical expenditure panel survey	15656	139	[26]
star	Tennessee’s student-teacher achievement ratio	2161	39	[27]

Table 1: Data sets for our empirical analysis, with numbers of samples (n) and features (d).

The test-set performances of CQR, CQR-m, and CQR-r are summarized in Tables 2 and 3. These quantities correspond to the best choice of black box and the optimal value of the hyperparameter γ , defined separately for each algorithm. The CQR method consistently produces the narrowest valid prediction bands, while CQR-m and CQR-r are often comparable.

The performances obtained with different black boxes and values of γ are reported in Figure 2 for the *community* data, and in Appendix B for the other data sets. The results are shown as a function of γ , which affects the average width of the prediction intervals as well as their variability. If γ is small, the prediction intervals are not sufficiently adaptive because the black box cannot estimate the regression quantiles accurately. Larger values of γ may lead to tighter predictions on average, but at the cost of increased variability in the conditional coverage. In fact, the conditional coverage for new observations given the data may be lower than the expected marginal level, especially when γ is very close to one and the sample size is not very large. The empirical results in Figure 2 and Appendix B suggest that fixing $\gamma \in [0.7, 0.9]$ achieves a reasonable compromise for all data sets analyzed in this paper. This observation is also consistent with the choice in [2].

The CQR-m method sometimes produces very wide intervals because the denominator in (5) can be close to zero (we added a small constant to prevent overflowing). An example is visible in the second plot in Figure 2, where some of the CQR-m prediction intervals based on a random forest black box are extremely large (hence the discontinuous vertical axis in Figure 2) when $\gamma \geq 0.9$. The CQR-r method is less susceptible to this problem because the denominator in the conformity scores in (7) is larger.

Dataset	Width		
	CQR	CQR-r	CQR-m
bike	0.503 (0.024)	0.520 (0.024)	0.521 (0.023)
bio	0.995 (0.037)	1.048 (0.049)	1.114 (0.019)
blog	1.269 (0.040)	1.462 (0.148)	1.351 (0.109)
community	1.461 (0.116)	1.548 (0.066)	1.617 (0.063)
concrete	0.378 (0.056)	0.387 (0.063)	0.384 (0.059)
facebook-1	1.117 (0.048)	1.188 (0.043)	1.164 (0.127)
facebook-2	1.110 (0.051)	1.172 (0.066)	1.116 (0.066)
homes	0.477 (0.013)	0.491 (0.013)	0.492 (0.013)
meps-19	2.300 (0.164)	2.349 (0.175)	2.442 (0.364)
meps-20	2.309 (0.121)	2.467 (0.313)	2.467 (0.168)
meps-21	2.201 (0.076)	2.273 (0.119)	2.343 (0.337)
star	0.179 (0.006)	0.180 (0.010)	0.181 (0.006)

Table 2: Average width (and standard deviation) of the conformal quantile regression prediction bands obtained with different conformalization methods on the data sets listed in Table 1. The corresponding coverage is reported in Table 3. The smallest value on each row is written in bold.

Dataset	Coverage		
	CQR	CQR-r	CQR-m
bike	0.899 (0.012)	0.901 (0.012)	0.900 (0.012)
bio	0.891 (0.012)	0.893 (0.016)	0.895 (0.008)
blog	0.905 (0.003)	0.901 (0.007)	0.903 (0.004)
community	0.896 (0.025)	0.899 (0.025)	0.902 (0.017)
concrete	0.875 (0.061)	0.879 (0.061)	0.877 (0.059)
facebook-1	0.901 (0.006)	0.898 (0.004)	0.902 (0.002)
facebook-2	0.900 (0.003)	0.900 (0.002)	0.900 (0.002)
homes	0.902 (0.009)	0.904 (0.009)	0.904 (0.009)
meps-19	0.902 (0.008)	0.902 (0.007)	0.900 (0.011)
meps-20	0.897 (0.004)	0.898 (0.004)	0.899 (0.006)
meps-21	0.899 (0.008)	0.898 (0.009)	0.897 (0.009)
star	0.905 (0.024)	0.904 (0.024)	0.903 (0.020)

Table 3: Average coverage (and standard deviation) of the prediction bands in Table 2.

4 Conclusion

Early work on conformal prediction focused on estimating a mean regression function for $Y | X$ and building a fixed-width band around it [5–7, 13]. Even though this strategy produces valid marginal prediction intervals regardless of $P_{Y|X}$, it is clearly designed with a homoscedastic regression model in mind and it may lead to be unnecessarily wide intervals in other cases. Locally-adaptive conformal prediction [9, 10, 12, 13] goes a step beyond this model by weighting the residuals according to a local estimate of their variance. Conformal quantile regression [1] goes further by observing that the estimation of the regression mean is unnecessary if the ultimate goal is to build prediction intervals. This approach has already been shown to outperform earlier methods in practice [1].

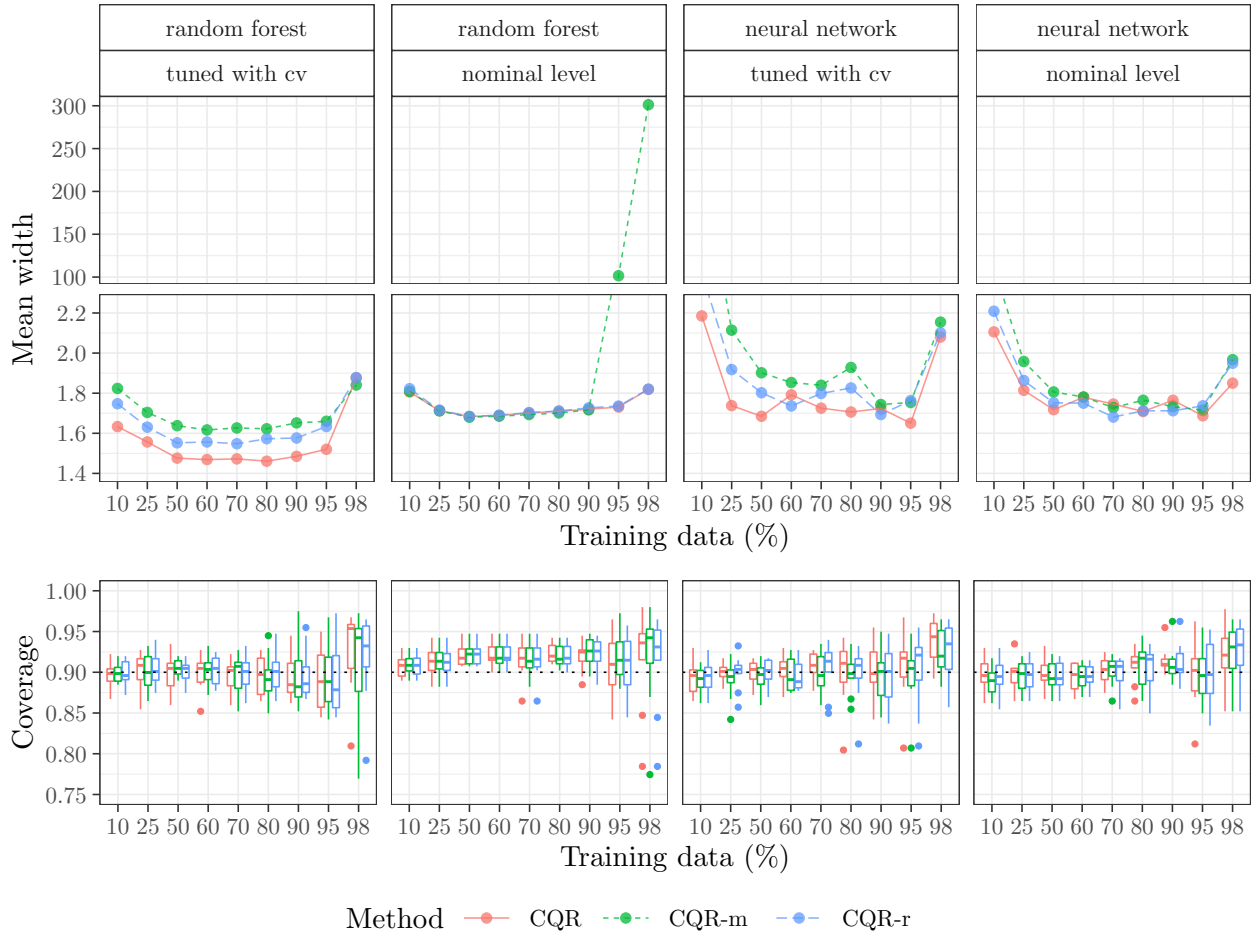


Figure 2: Conformal prediction bands obtained with different black boxes and conformalization methods on the *community* data set from Table 1. The dotted black line in the lower plots indicates the nominal coverage level (90%). A different black box is considered in each column. The vertical axis in the upper panels is discontinuous to facilitate the visualization of values on different scales.

In this paper, we have strengthened the case for conformal quantile regression by proving that it is asymptotically efficient in large samples, if the quantile regression estimates are consistent. The empirical comparison of three alternative conformity scores has shown that those proposed in [1] are preferable because they typically lead to shorter prediction intervals in practice. Even though we have only explicitly considered symmetric intervals for simplicity, it is straightforward to generalize these methods to asymmetric intervals and conformity scores [1]. Finally, we have highlighted a bias-variance tradeoff in the choice of the proportion of data points used to train the black-box quantile regressors. Our empirical results show that it is usually better to invest more of the available data (between 70% and 90%, indicatively) to train the black-box than to conformalize the predictions. We hope that these results will be helpful to practitioners and may inspire others to develop even more powerful variations of conformal quantile regression.

Acknowledgements

M. S. and E. C. are supported by the National Science Foundation under grant DMS 1712800. E. C. is also supported by the Army Research Office under grant W911NF-17-1-0304. We thank Yaniv Romano for helpful discussions, during which he suggested the CQR-r method.

References

- [1] Y. Romano, E. Patterson, and E. J. Candès, “Conformalized quantile regression,” *arXiv preprint arXiv:1905.03222*, 2019.
- [2] D. Kivaranovic, K. D. Johnson, and H. Leeb, “Adaptive, distribution-free prediction intervals for deep neural networks,” *arXiv preprint arXiv:1905.10634*, 2019.
- [3] G. Papadopoulos, P. J. Edwards, and A. F. Murray, “Confidence estimation methods for neural networks: A practical comparison,” *IEEE Transactions on Neural Networks*, vol. 12, no. 6, pp. 1278–1287, 2001.
- [4] S. Wager, T. Hastie, and B. Efron, “Confidence intervals for random forests: The jackknife and the infinitesimal jackknife,” *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1625–1651, 2014.
- [5] V. Vovk, A. Gammerman, and C. Saunders, “Machine-learning applications of algorithmic randomness,” in *Proceedings of the Sixteenth International Conference on Machine Learning, ICML '99*, (San Francisco, CA, USA), pp. 444–453, Morgan Kaufmann Publishers Inc., 1999.
- [6] V. Vovk, A. Gammerman, and G. Shafer, *Algorithmic Learning in a Random World*. Berlin, Heidelberg: Springer-Verlag, 2005.
- [7] V. Vovk, I. Nouretdinov, A. Gammerman, *et al.*, “On-line predictive linear regression,” *The Annals of Statistics*, vol. 37, no. 3, pp. 1566–1590, 2009.
- [8] H. Papadopoulos, K. Proedrou, V. Vovk, and A. Gammerman, “Inductive confidence machines for regression,” in *European Conference on Machine Learning*, pp. 345–356, Springer, 2002.
- [9] H. Papadopoulos, V. Vovk, and A. Gammerman, “Conformal prediction with neural networks,” in *19th IEEE International Conference on Tools with Artificial Intelligence (ICTAI 2007)*, vol. 2, pp. 388–395, Oct 2007.
- [10] H. Papadopoulos, A. Gammerman, and V. Vovk, “Normalized nonconformity measures for regression conformal prediction,” in *Proceedings of the 26th IASTED International Conference on Artificial Intelligence and Applications, AIA '08*, (Anaheim, CA, USA), pp. 64–69, ACTA Press, 2008.
- [11] H. Papadopoulos, “Inductive conformal prediction: Theory and application to neural networks,” in *Tools in Artificial Intelligence* (P. Fritzsche, ed.), ch. 18, Rijeka: IntechOpen, 2008.
- [12] H. Papadopoulos, V. Vovk, and A. Gammerman, “Regression conformal prediction with nearest neighbours,” *Journal of Artificial Intelligence Research*, vol. 40, pp. 815–840, Jan. 2011.

- [13] J. Lei, M. Gsell, A. Rinaldo, R. J. Tibshirani, and L. Wasserman, “Distribution-free predictive inference for regression,” *Journal of the American Statistical Association*, vol. 113, no. 523, pp. 1094–1111, 2018.
- [14] R. Koenker and G. Bassett Jr, “Regression quantiles,” *Econometrica: Journal of the Econometric Society*, pp. 33–50, 1978.
- [15] J. W. Taylor, “A quantile regression neural network approach to estimating the conditional density of multiperiod returns,” *Journal of Forecasting*, vol. 19, no. 4, pp. 299–311, 2000.
- [16] N. Meinshausen, “Quantile regression forests,” *Journal of Machine Learning Research*, vol. 7, no. Jun, pp. 983–999, 2006.
- [17] “Bike sharing dataset data set.” <https://archive.ics.uci.edu/ml/datasets/bike+sharing+dataset>. Accessed: July, 2019.
- [18] “Physicochemical properties of protein tertiary structure data set.” <https://archive.ics.uci.edu/ml/datasets/Physicochemical+Properties+of+Protein+Tertiary+Structure>. Accessed: July, 2019.
- [19] “BlogFeedback data set.” <https://archive.ics.uci.edu/ml/datasets/BlogFeedback>. Accessed: July, 2019.
- [20] “Communities and crime data set.” <http://archive.ics.uci.edu/ml/datasets/communities+and+crime>. Accessed: July, 2019.
- [21] “Concrete compressive strength data set.” <http://archive.ics.uci.edu/ml/datasets/concrete+compressive+strength>. Accessed: July, 2019.
- [22] “Facebook comment volume data set.” <https://archive.ics.uci.edu/ml/datasets/Facebook+Comment+Volume+Dataset>. Accessed: July, 2019.
- [23] “House sales in King County, USA.” <https://www.kaggle.com/harlfoxem/housesalesprediction/metadata>. Accessed: August, 2019.
- [24] “Medical expenditure panel survey, panel 19.” https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-181. Accessed: July, 2019.
- [25] “Medical expenditure panel survey, panel 20.” https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-181. Accessed: July, 2019.
- [26] “Medical expenditure panel survey, panel 21.” https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-192. Accessed: July, 2019.
- [27] C. Achilles, H. P. Bain, F. Bellott, J. Boyd-Zaharias, J. Finn, J. Folger, J. Johnston, and E. Word, “Tennessee’s student teacher achievement ratio (STAR) project,” 2008. Accessed: July, 2019.
- [28] A. W. v. d. Vaart, *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press, 1998.

A Proofs

Proof of Theorem 1. We begin by considering the case of CQR. For ease of notation and without loss of generality, assume that we have $2n$ data points and $n_1 = n_2 = n$. Then, we can equivalently rewrite Assumption 3 as follows:

$$\begin{aligned} \mathbb{P} \left[\mathbb{E} \left[(\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X))^2 \mid \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2} \right] \leq \frac{\eta_n}{2} \right] &\geq 1 - \frac{\rho_n}{2}, \\ \mathbb{P} \left[\mathbb{E} \left[(\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X))^2 \mid \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2} \right] \leq \frac{\eta_n}{2} \right] &\geq 1 - \frac{\rho_n}{2}, \end{aligned}$$

for n large enough, $X \perp\!\!\!\perp \mathcal{I}_1$ and for some sequences $\eta_n = o(1)$ and $\rho_n = o(1)$, as $n \rightarrow \infty$.

Recall that the conformal quantile regression prediction band is defined as:

$$\hat{C}_\alpha^{\text{CQR}}(X_{2n+1}) = \left[\hat{q}_{\alpha/2}(X_{2n+1}) - \hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2), \hat{q}_{1-\alpha/2}(X_{2n+1}) + \hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2) \right],$$

while the oracle band is:

$$C_\alpha^{\text{oracle}}(X_{2n+1}) = [q_{\alpha/2}(X_{2n+1}), q_{1-\alpha/2}(X_{2n+1})].$$

It suffices to show:

$$\begin{aligned} |\hat{q}_{\alpha/2}(X_{2n+1}) - \hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2) - q_{\alpha/2}(X_{2n+1})| &= o_{\mathbb{P}}(1), \\ |\hat{q}_{1-\alpha/2}(X_{2n+1}) + \hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2) - q_{1-\alpha/2}(X_{2n+1})| &= o_{\mathbb{P}}(1). \end{aligned}$$

We will proceed in two steps, proving:

- (i) $|\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X)| = o_{\mathbb{P}}(1)$ and $|\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X)| = o_{\mathbb{P}}(1)$, for $X \perp\!\!\!\perp \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2}$;
- (ii) $|\hat{Q}_{1-\alpha}(E^{\text{CQR}}; \mathcal{I}_2)| = o_{\mathbb{P}}(1)$.

Then the proof will be completed by the triangle inequality.

- (i) Define the random sets

$$B_{n,\text{up}} = \{x : |\hat{q}_{1-\alpha/2}(x) - q_{1-\alpha/2}(x)| \geq \eta_n^{1/3}\}, \quad B_{n,\text{lo}} = \{x : |\hat{q}_{\alpha/2}(x) - q_{\alpha/2}(x)| \geq \eta_n^{1/3}\},$$

and $B_n = B_{n,\text{up}} \cup B_{n,\text{lo}}$. We can prove that for a new $X \perp\!\!\!\perp \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2}$,

$$\mathbb{P} [X \in B_n \mid \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2}] \leq \eta_n^{1/3} + \rho_n. \quad (10)$$

In fact, in the event

$$\left\{ \mathbb{E} \left[(\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X))^2 \right] \leq \frac{\eta_n}{2} \right\} \cap \left\{ \mathbb{E} \left[(\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X))^2 \right] \leq \frac{\eta_n}{2} \right\}, \quad (11)$$

we have:

$$\begin{aligned}
& \mathbb{P} [X \in B_n \mid \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2}] \\
&= \mathbb{P} [X \in B_{n,\text{lo}} \cup B_{n,\text{up}} \mid \hat{q}_{\alpha/2}, \hat{q}_{1-\alpha/2}] \\
&\leq \mathbb{P} [X \in B_{n,\text{lo}} \mid \hat{q}_{\alpha/2}] + \mathbb{P} [X \in B_{n,\text{up}} \mid \hat{q}_{1-\alpha/2}] \\
&= \mathbb{P} \left[|\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X)| \geq \eta_n^{1/3} \mid \hat{q}_{\alpha/2} \right] + \mathbb{P} \left[|\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X)| \geq \eta_n^{1/3} \mid \hat{q}_{1-\alpha/2} \right] \\
&= \mathbb{P} \left[|\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X)|^2 \geq \eta_n^{2/3} \mid \hat{q}_{\alpha/2} \right] + \mathbb{P} \left[|\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X)|^2 \geq \eta_n^{2/3} \mid \hat{q}_{1-\alpha/2} \right] \\
&\leq \eta_n^{-2/3} \mathbb{E} \left[(\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X))^2 \right] + \eta_n^{-2/3} \mathbb{E} \left[(\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X))^2 \right] \\
&\leq \eta_n^{1/3}.
\end{aligned}$$

Equation (10) follows because the event in (11) has probability at least $1 - \rho_n$, by Assumption 3. This implies:

$$|\hat{q}_{1-\alpha/2}(X) - q_{1-\alpha/2}(X)| = o_{\mathbb{P}}(1), \quad |\hat{q}_{\alpha/2}(X) - q_{\alpha/2}(X)| = o_{\mathbb{P}}(1).$$

(ii) With B_n defined as above, consider the following further partition of the data in \mathcal{I}_2 :

$$\mathcal{I}_{2,a} = \{i \in \{n+1, \dots, 2n\} : X_i \in B_n^c\}, \quad \mathcal{I}_{2,b} = \{i \in \{n+1, \dots, 2n\} : X_i \in B_n\}.$$

By definition, $\mathcal{I}_2 = \{n+1, \dots, 2n\} = \mathcal{I}_{2,a} \cup \mathcal{I}_{2,b}$. Since B_n only depends on the data in \mathcal{I}_1 , it is independent of (X_i, Y_i) for all $i \in \mathcal{I}_2$. Therefore, the size of $\mathcal{I}_{2,b}$ conditional on the data in \mathcal{I}_1 can be bounded using Hoeffding's inequality, thanks to Assumption 1. We already know that the probability that any $i \in \{n+1, \dots, 2n\}$ belongs to $\mathcal{I}_{2,b}$ is smaller than $\eta_n^{1/3}$. In particular, conditional on $\hat{q}_{\alpha/2}$ and $\hat{q}_{1-\alpha/2}$,

$$\mathbb{P} \left[|\mathcal{I}_{2,b}| \geq n\eta_n^{1/3} + t \right] \leq \mathbb{P} \left[\frac{1}{n} \sum_{i=n+1}^{2n} \mathbb{1}[X_i \in B_n] \geq \mathbb{P}[X_i \in B_n] + \frac{t}{n} \right] \leq \exp \left(-\frac{2t^2}{n} \right).$$

Choosing $t = c\sqrt{n \log n}$ leads to $|\mathcal{I}_{2,b}| = o_{\mathbb{P}}(n)$ because

$$\mathbb{P} \left[|\mathcal{I}_{2,b}| \geq n\eta_n^{1/3} + c\sqrt{n \log n} \right] \leq n^{-2c^2}.$$

Now, define $\tilde{E}_i^{\text{CQR}} = \max\{q_{\alpha/2}(X_{n+1}) - Y, Y - q_{1-\alpha/2}(X_{n+1})\}$ for any $i \in \mathcal{I}_2$. By definition of E_i^{CQR} , for all $i \in \mathcal{I}_{2,a}$,

$$\begin{aligned}
E_i^{\text{CQR}} &= \max\{\hat{q}_{\alpha/2}(X_{n+1}) - Y, Y - \hat{q}_{1-\alpha/2}(X_{n+1})\} \\
&= \max\{\hat{q}_{\alpha/2}(X_{n+1}) - q_{\alpha/2}(X_{n+1}) + q_{\alpha/2}(X_{n+1}) - Y, \\
&\quad Y - q_{1-\alpha/2}(X_{n+1}) + q_{1-\alpha/2}(X_{n+1}) - \hat{q}_{1-\alpha/2}(X_{n+1})\} \\
&\leq \max\{\eta_n^{1/3} + q_{\alpha/2}(X_{n+1}) - Y, Y - q_{1-\alpha/2}(X_{n+1}) + \eta_n^{1/3}\} \\
&= \eta_n^{1/3} + \max\{q_{\alpha/2}(X_{n+1}) - Y, Y - q_{1-\alpha/2}(X_{n+1})\} \\
&= \eta_n^{1/3} + \tilde{E}_i^{\text{CQR}}.
\end{aligned}$$

Proceeding similarly, one can also show that $E_i^{\text{CQR}} \geq \tilde{E}_i^{\text{CQR}} - \eta_n^{1/3}$. Hence, for all $i \in \mathcal{I}_{2,a}$,

$$\left| E_i^{\text{CQR}} - \tilde{E}_i^{\text{CQR}} \right| \leq \eta_n^{1/3}.$$

Therefore, all empirical quantiles of E_i^{CQR} and \tilde{E}_i^{CQR} , for $i \in \mathcal{I}_{2,a}$, are within $\eta_n^{1/3}$.

Let F_n and \tilde{F}_n denote the empirical distributions of E_i^{CQR} and \tilde{E}_i^{CQR} for $i \in \mathcal{I}_2$, respectively. Define also $F_{n,a}$ and $\tilde{F}_{n,a}$ as the corresponding empirical distributions when i is restricted to $\mathcal{I}_{2,a}$. For n large enough, one can assume without loss of generality that $|\mathcal{I}_{2,b}|/|\mathcal{I}_{2,a}| \leq \alpha$ because $|\mathcal{I}_{2,b}| = o_{\mathbb{P}}(n)$. Then we can show that

$$F_{n,a}^{-1} \left(1 - \frac{n\alpha}{|\mathcal{I}_{2,a}|} \right) \leq F_n^{-1}(1 - \alpha) \leq F_{n,a}^{-1} \left(1 - \frac{n\alpha - |\mathcal{I}_{2,b}|}{|\mathcal{I}_{2,a}|} \right). \quad (12)$$

To prove the second inequality in (12), note that if all the E_i^{CQR} , for $i \in \mathcal{I}_{2,b}$, are in the upper α -quantile of F_n , then

$$F_{n,a}^{-1} \left(1 - \frac{n\alpha - |\mathcal{I}_{2,b}|}{n} \frac{n}{|\mathcal{I}_{2,a}|} \right) = F_n^{-1}(1 - \alpha).$$

However, in general the quantiles of $F_{n,a}$ will be larger and

$$F_{n,a}^{-1} \left(1 - \frac{n\alpha - |\mathcal{I}_{2,b}|}{|\mathcal{I}_{2,a}|} \right) \geq F_n^{-1}(1 - \alpha).$$

To prove the first inequality in (12), note that if all the E_i^{CQR} for $i \in \mathcal{I}_{2,b}$ are in the lower $1 - \alpha$ quantile of F_n ,

$$F_{n,a}^{-1} \left(1 - \frac{n\alpha}{|\mathcal{I}_{2,a}|} \right) = F_n^{-1}(1 - \alpha).$$

However, in general the quantiles of $F_{n,a}$ will be smaller and

$$F_{n,a}^{-1} \left(1 - \frac{n\alpha}{|\mathcal{I}_{2,a}|} \right) \leq F_n^{-1}(1 - \alpha).$$

This completes the proof of (12). By combining this with the previous result that all empirical quantiles of E_i^{CQR} and \tilde{E}_i^{CQR} , for $i \in \mathcal{I}_{2,a}$, are within $\eta_n^{1/3}$, we obtain:

$$\tilde{F}_{n,a}^{-1} \left(1 - \frac{n\alpha}{|\mathcal{I}_{2,a}|} \right) - \eta_n^{1/3} \leq F_n^{-1}(1 - \alpha) \leq \tilde{F}_{n,a}^{-1} \left(1 - \frac{n\alpha - |\mathcal{I}_{2,b}|}{|\mathcal{I}_{2,a}|} \right) + \eta_n^{1/3}.$$

Recall that we have defined $\hat{Q}_{1-\alpha}(E; \mathcal{I}_2) = F_n^{-1}(1 - \alpha_n)$, where $\alpha_n = \alpha - (1 - \alpha)/n$. Therefore,

$$\tilde{F}_{n,a}^{-1} \left(1 - \frac{n\alpha_n}{|\mathcal{I}_{2,a}|} \right) - \eta_n^{1/3} \leq \hat{Q}_{1-\alpha}(E; \mathcal{I}_2) \leq \tilde{F}_{n,a}^{-1} \left(1 - \frac{n\alpha_n - |\mathcal{I}_{2,b}|}{|\mathcal{I}_{2,a}|} \right) + \eta_n^{1/3}.$$

Note that $Q_{1-\alpha}(\tilde{E}^{\text{CQR}}) = 0$. This follows immediately from its definition:

$$\begin{aligned} \mathbb{P} \left[\tilde{E}^{\text{CQR}} \leq 0 \right] &= \mathbb{P} \left[\max\{q_{\alpha/2}(X_{n+1}) - Y, Y - q_{1-\alpha/2}(X_{n+1})\} \leq 0 \right] \\ &= \mathbb{P} \left[Y \in [q_{\alpha/2}(X_{n+1}), q_{1-\alpha/2}(X_{n+1})] \right] = 1 - \alpha. \end{aligned}$$

Then we know from Assumptions 1-2 and classical asymptotic theory [28, Chapter 21] that $\tilde{F}_{n,a}^{-1}(1 - \alpha)$ is within $o_{\mathbb{P}}(1)$ of the upper α population quantile of \tilde{E}_i^{CQR} , which we denote by $Q_{1-\alpha}(\tilde{E}^{\text{CQR}})$. Therefore, since both $n\alpha_n/|\mathcal{I}_{2,a}|$ and $(n\alpha_n - |\mathcal{I}_{2,b}|)/|\mathcal{I}_{2,a}|$ are $\alpha + o_{\mathbb{P}}(1)$, we have

$$\left| \hat{Q}_{1-\alpha}(E; \mathcal{I}_2) - Q_{1-\alpha}(\tilde{E}^{\text{CQR}}) \right| = o_{\mathbb{P}}(1).$$

This completes the proof in the case of CQR. The proof for CQR-m and CQR-r are analogous. □

B Supplementary figures

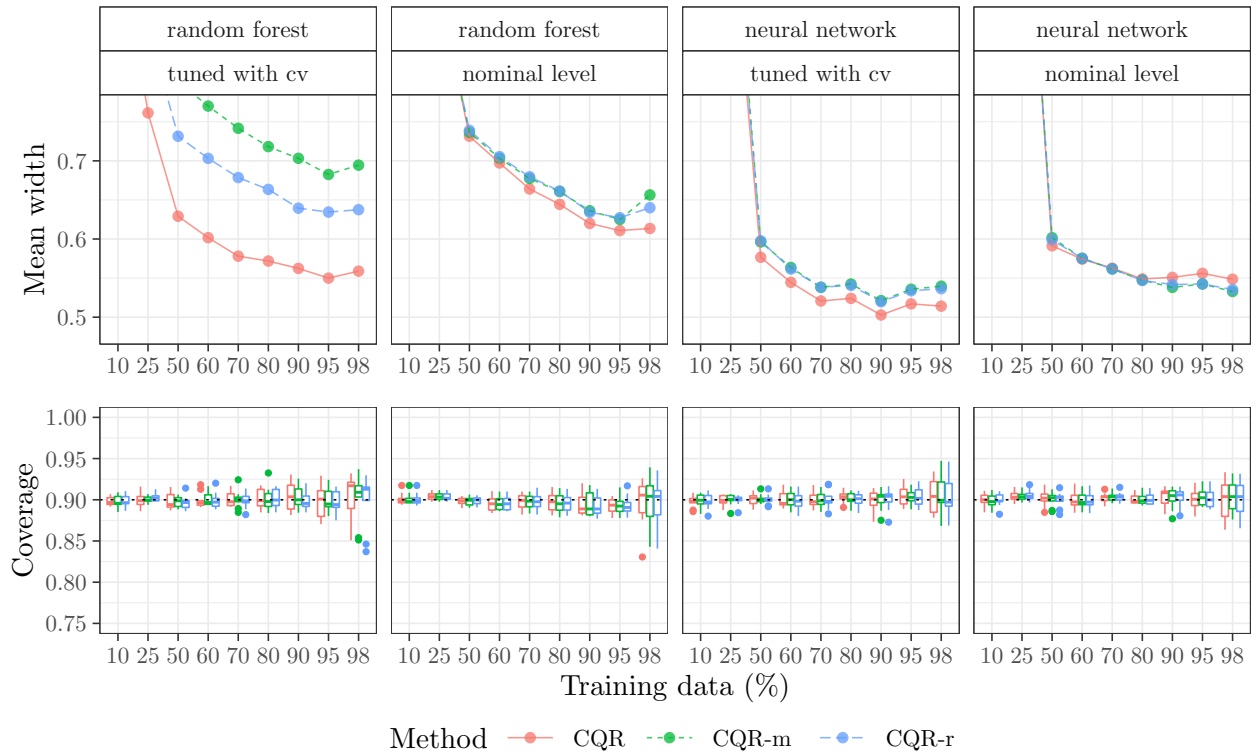


Figure 3: Results on the *bike* data. Other details as in Figure 2. The plots are truncated from above on the vertical axis to focus on the most interesting region. In the second through fourth plots on top, the curves corresponding to CQR-m and CQR-r are almost overlapping.

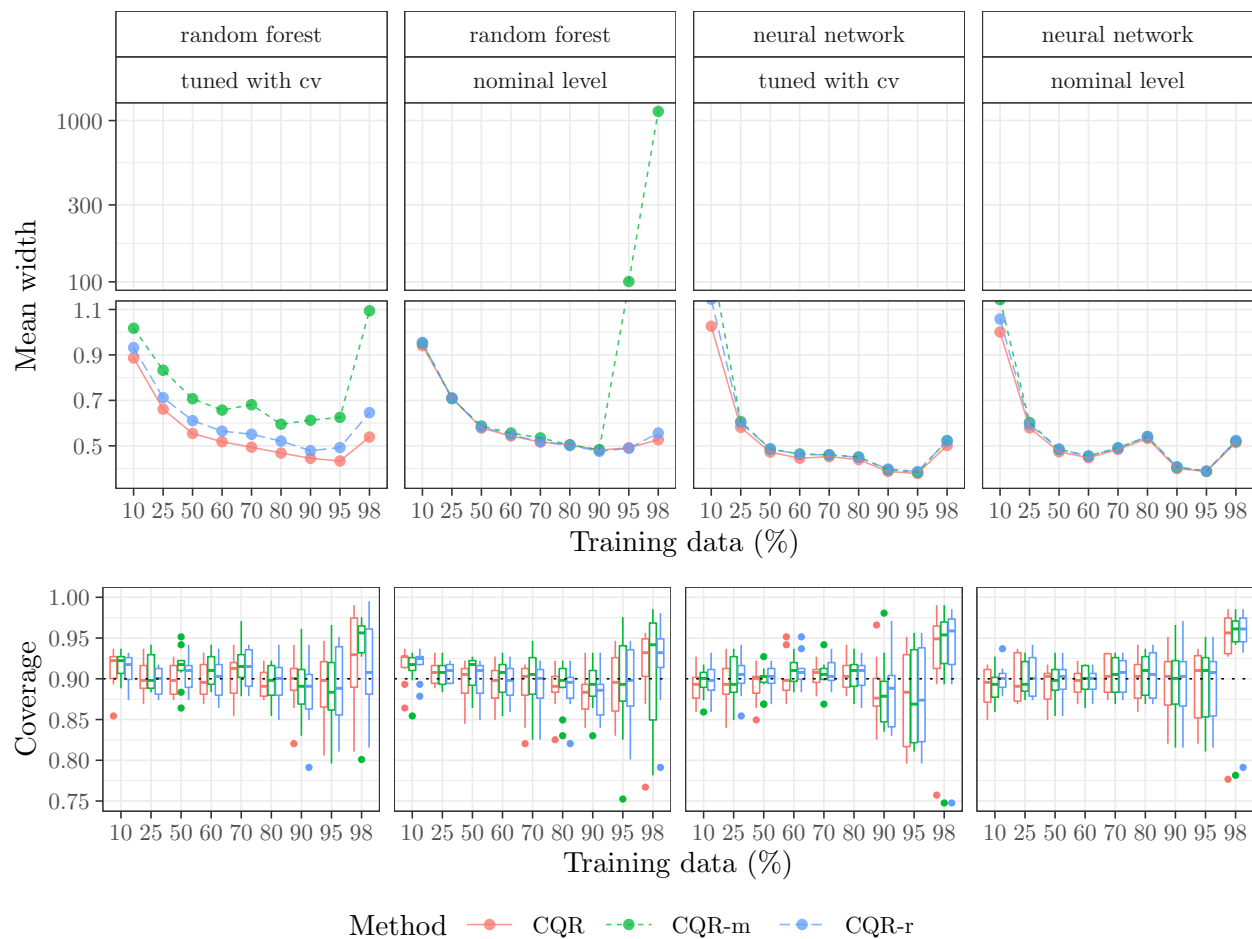


Figure 4: Results on the *concrete* data. Other details as in Figure 3. In the second through fourth plots on top, the three curves are mostly overlapping. The vertical axis in the upper panels is discontinuous to facilitate the visualization of values on different scales.

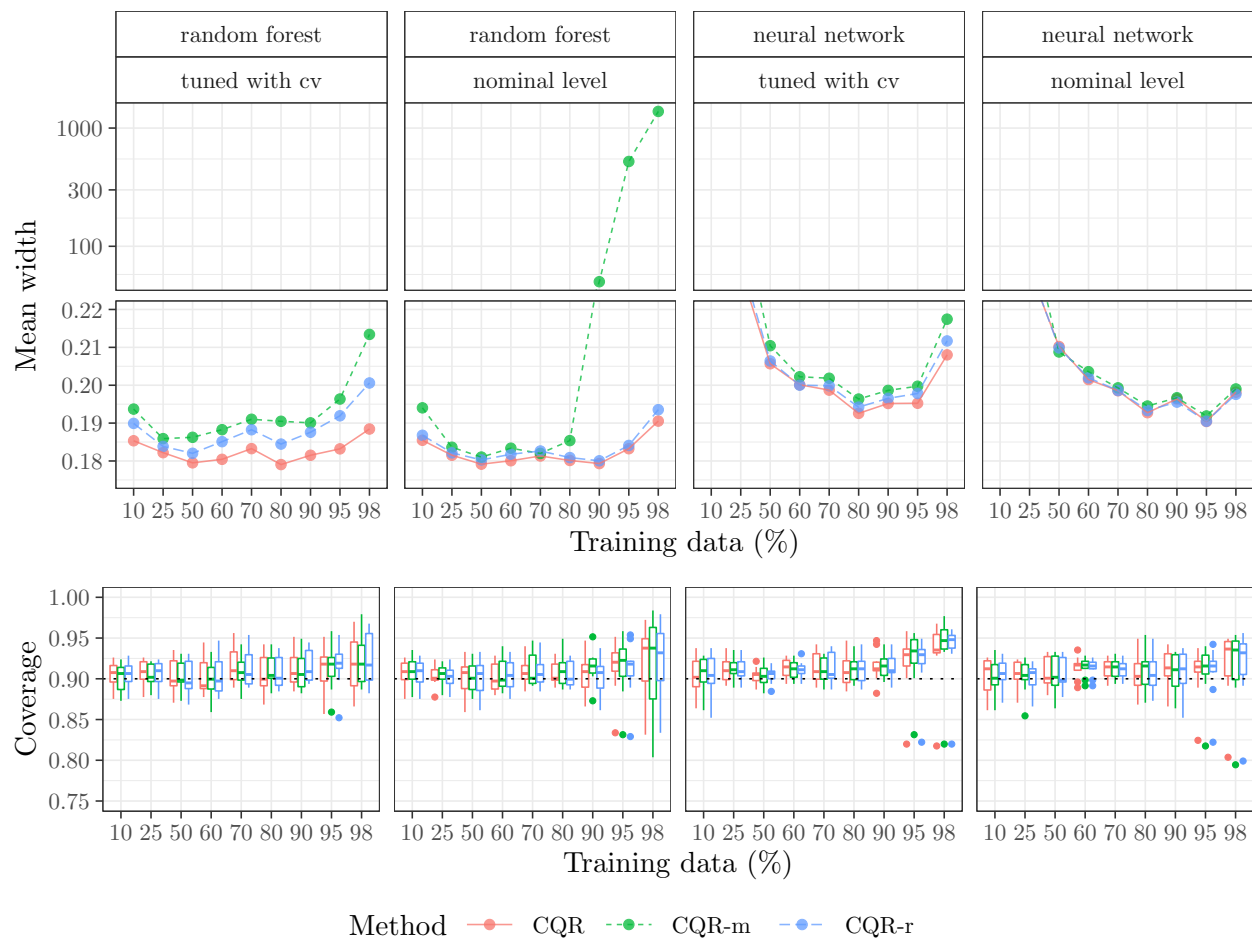


Figure 5: Results on the *star* data. Other details as in Figure 3. In the second through fourth plots on top, the three curves are mostly overlapping. The vertical axis in the upper panels is discontinuous to facilitate the visualization of values on different scales.

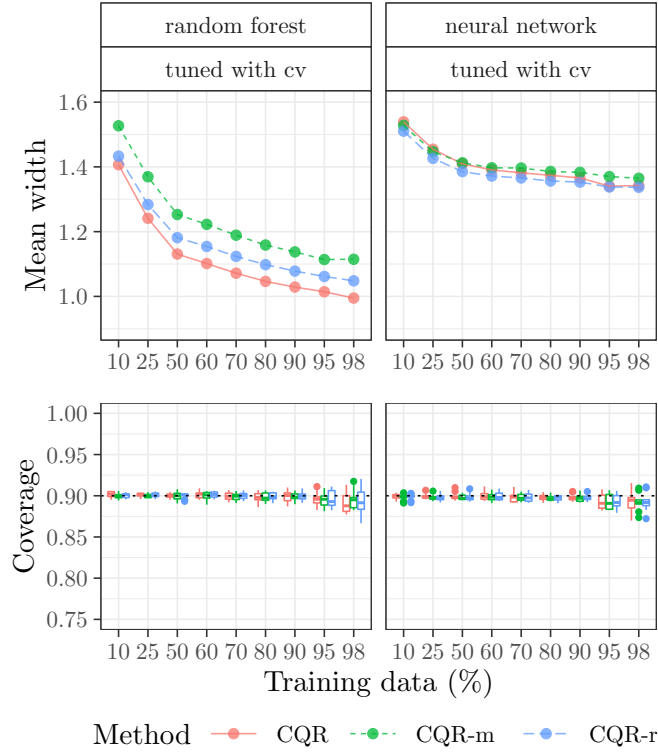


Figure 6: Results on the *bio* data. Other details as in Figure 3.

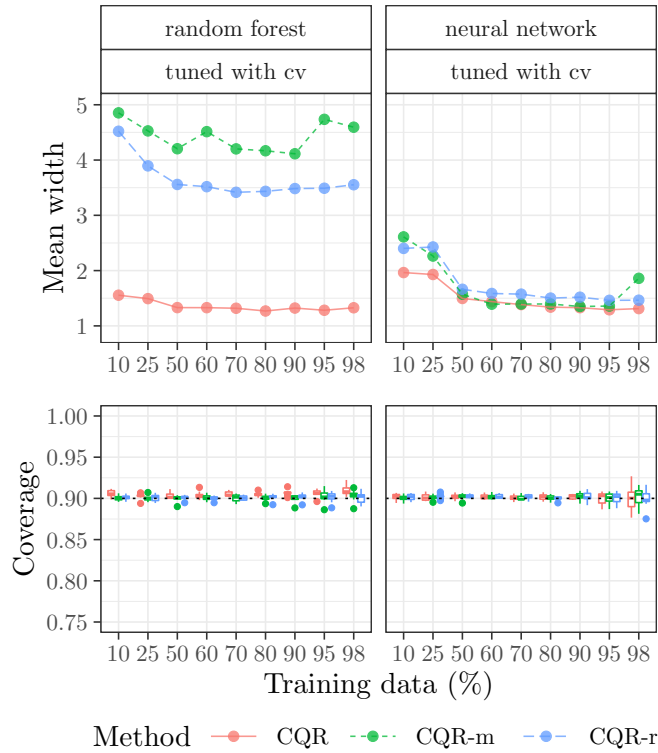


Figure 7: Results on the *blog* data. Other details as in Figure 3.

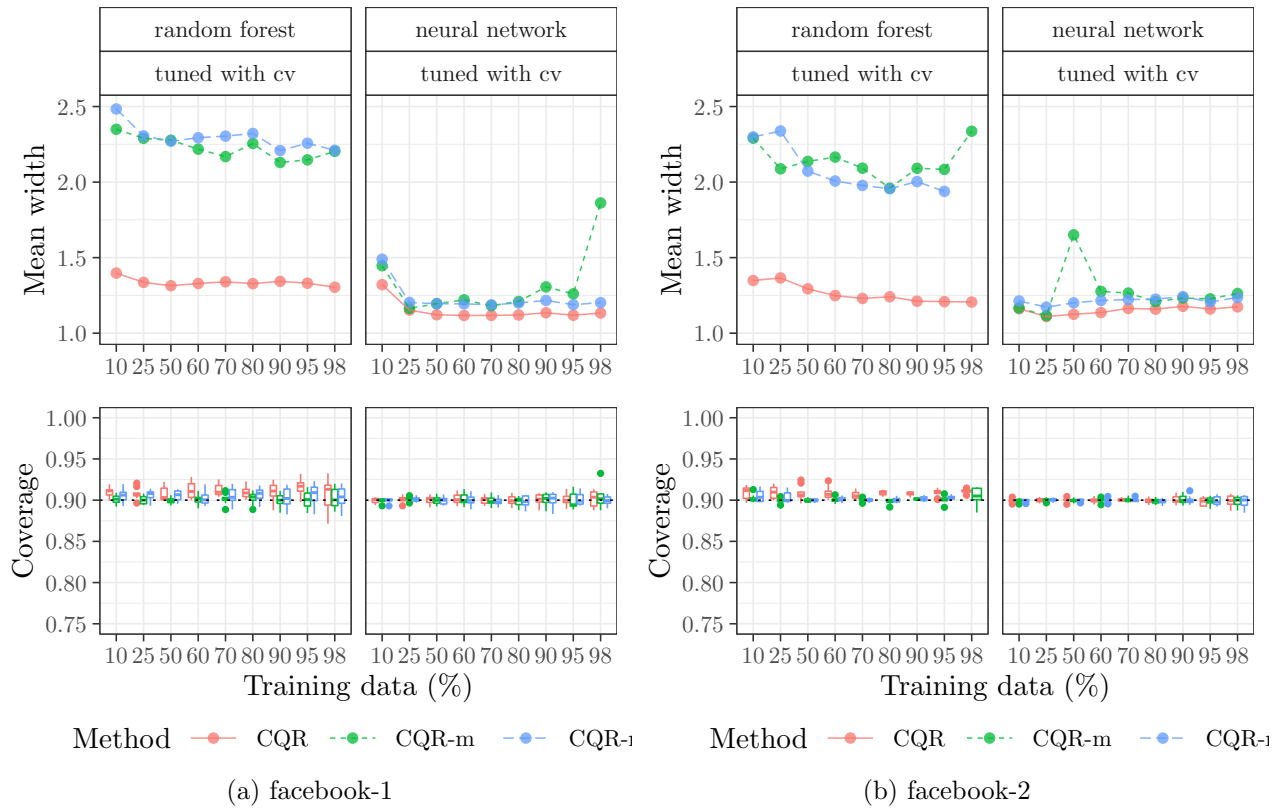


Figure 8: Results on the *facebook* data. Other details as in Figure 3.

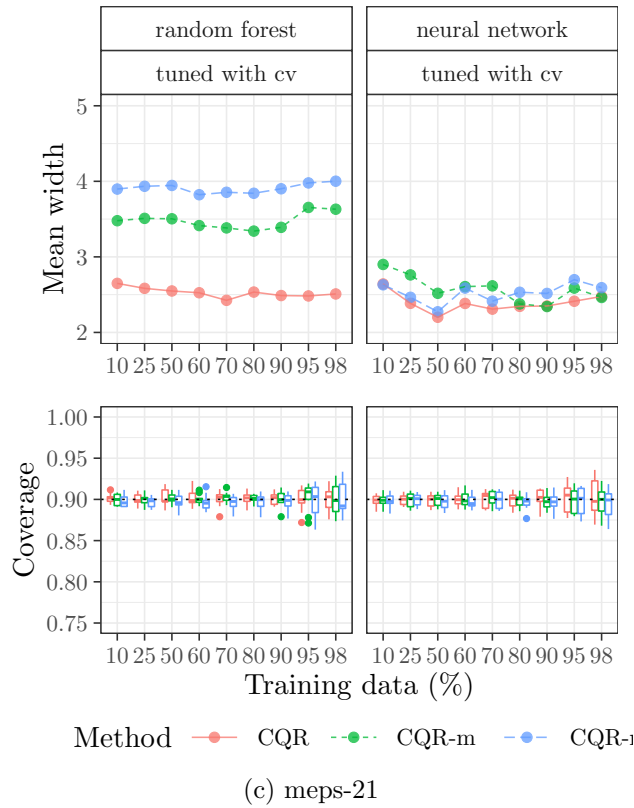
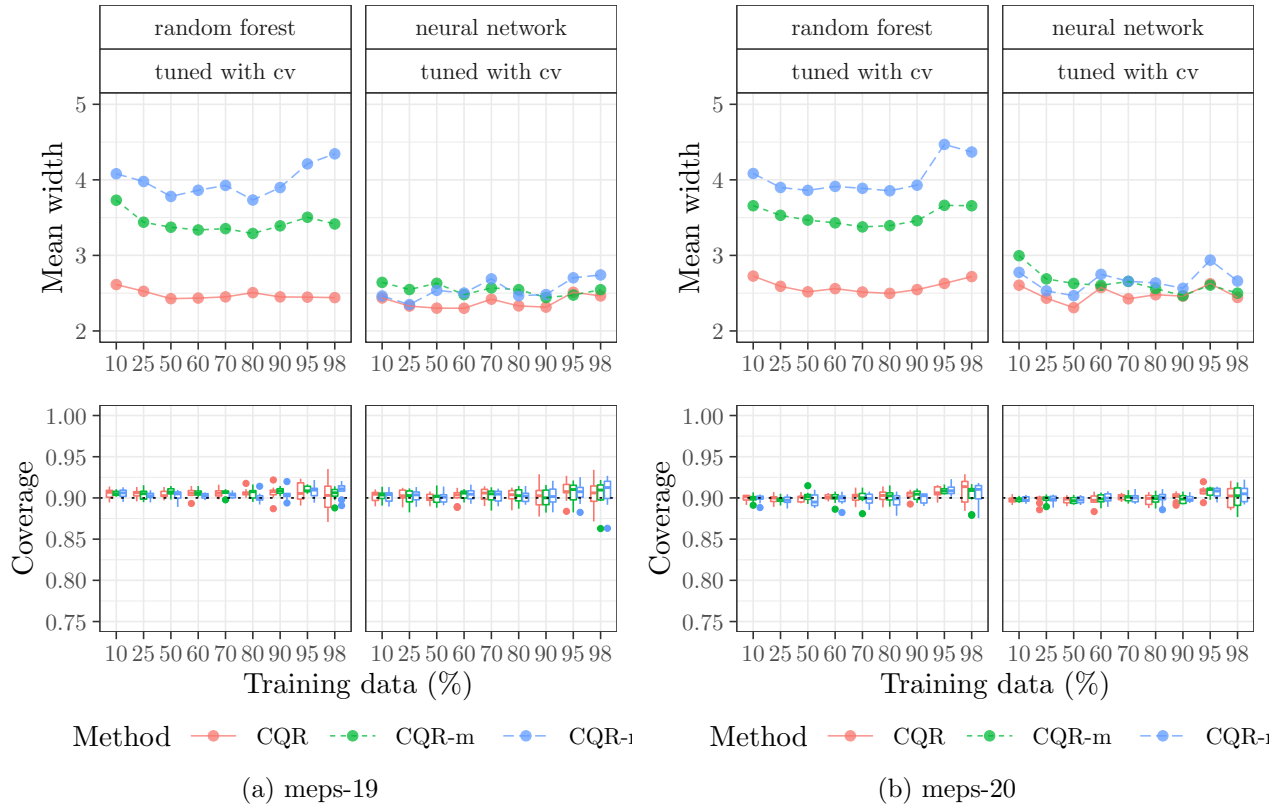


Figure 9: Results on the *meps* data. Other details as in Figure 3.