



Department of Econometrics and Business Statistics

<http://monash.edu/business/ebs/research/publications>

Online conformal inference for multi-step time series forecasting

Xiaoqian Wang, Rob J Hyndman

May 2024

Working Paper no/yr

Online conformal inference for multi-step time series forecasting

Xiaoqian Wang

Department of Econometrics & Business Statistics

Monash University

Clayton VIC 3800

Australia

Email: xiaoqian.wang@monash.edu

Corresponding author

Rob J Hyndman

Department of Econometrics & Business Statistics

Monash University

Clayton VIC 3800

Australia

Email: rob.hyndman@monash.edu

12 May 2024

JEL classification: C10,C14,C32

Online conformal inference for multi-step time series forecasting

Abstract

A brief summary

Keywords: Conformal prediction; Distribution-free inference; Nonexchangeability; Valid prediction interval; Weighted quantile estimate.

1 Introduction

We consider a general sequential setting in which we observe a time series $\{y_t\}_{t \geq 1}$ generated by an unknown data generating process (DGP), which may depend on its own past, along with other exogenous predictors, $\mathbf{x}_t = (x_{1,t}, \dots, x_{p,t})'$, and their histories. The distribution of $\{(\mathbf{x}_t, y_t)\}_{t \geq 1} \subseteq \mathbb{R}^p \times \mathbb{R}$ is obviously allowed to vary over time in time series context. At each time point t , we aim to forecast H steps into the future, providing a *prediction set* (which is a prediction interval in this setting), $\hat{\mathcal{C}}_{t+h|t}$, for the realization y_{t+h} for each $h \in [H]$. The h -step-ahead forecast uses the previously observed data $\{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq t}$ along with the new information of the predictors $\{\mathbf{x}_{t+j}\}_{1 \leq j \leq h}$. Note that we can generate ex-ante forecasts by using forecasts of the predictors based on information available up to and including time t . Alternatively, ex-post forecasts are generated assuming that actual observations of the predictors from the forecast period are available. Given a nominal *miscoverage rate* $\alpha \in (0, 1)$ specified by the user, we expect the output $\hat{\mathcal{C}}_{t+h|t}$ to be a *valid* prediction interval so that y_{t+h} falls within the prediction interval $\hat{\mathcal{C}}_{t+h|t}$ at least $100(1 - \alpha)\%$ of the time.

- Throughout this paper, we use split conformal prediction.
- Uniform notation.

2 Related work

2.1 Conformal prediction for regression

In this section, we focus on the regression setting, which stands as one of the primary areas where conformal prediction methods have seen substantial development. Suppose we have n data points $Z_i = (X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i = 1, \dots, n$. The aim of conformal prediction is to construct a prediction interval $\hat{\mathcal{C}}_{n+1}^\alpha(X_{n+1})$ ensuring that the unseen response Y_{n+1} falls within $\hat{\mathcal{C}}_{n+1}^\alpha(X_{n+1})$ at least $100(1 - \alpha)\%$ of the time.

2.1.1 Split conformal prediction

Split conformal prediction (SCP, also called inductive conformal prediction, Papadopoulos et al. (2002); Vovk, Gammerman & Shafer (2005); Lei et al. (2018)), is a holdout method for building prediction intervals using a pre-trained model.

In regression setting, SCP randomly separates the available n data points into a *proper training set* \mathcal{D}_{tr} of size n_t and a *calibration set* \mathcal{D}_{cal} of size n_c . Given a regression model $\hat{\mu} : \mathbb{R}^d \rightarrow \mathbb{R}$ that is fitted on the training set and a score function \mathcal{S} , a *nonconformity score* $s_i = \mathcal{S}(X_i, Y_i)$, $i \in \mathcal{D}_{\text{cal}}$, is computed

on every data point in \mathcal{D}_{cal} to measure the nonconformity between the calibration's response values and the predicted values obtained from the fitted model $\hat{\mu}$. One commonly used nonconformity score function \mathcal{S} in regression is the absolute residual, i.e. $s_i = |Y_i - \hat{\mu}(X_i)|$. Then SCP computes the prediction interval for the test data Y_{n+1} using

$$\hat{\mathcal{C}}_{n+1}^\alpha(X_{n+1}) = \left\{ y \in \mathbb{R} : \mathcal{S}(X_{n+1}, y) \leq Q_{1-\alpha} \left(\sum_{i \in \mathcal{D}_{\text{cal}}} \frac{1}{n_c + 1} \cdot \delta_{s_i} + \frac{1}{n_c + 1} \cdot \delta_{+\infty} \right) \right\}, \quad (1)$$

where $Q_\tau(\cdot)$ denotes the τ -quantile of its argument, and δ_a denotes the point mass at a .

Theorem 2.1 (SCP, Vovk, Gammerman & Shafer (2005); Lei et al. (2018)). *Assume that the data points (X_i, Y_i) , $i = 1, \dots, n+1$, are i.i.d. (or more generally, exchangeable) from any distribution. For any score function \mathcal{S} , and any $\alpha \in (0, 1)$, the split conformal prediction interval defined in Equation 1 satisfies*

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

Moreover, if we assume additionally that the nonconformity scores on \mathcal{D}_{cal} are distinct with probability one, then we also have

$$\mathbb{P}\{Y_{n+1} \in \hat{\mathcal{C}}_{n+1}^\alpha(X_{n+1})\} < 1 - \alpha + \frac{1}{n_c + 1}.$$

2.1.2 Nonexchangeable conformal prediction

Barber et al. (2023) propose **nonexchangeable conformal prediction** (NexCP) that generalizes the SCP method to allow for some sources of nonexchangeability. For split conformal, the NexCP method can be considered as simply using weighted quantiles to obtain robust inference. Note that NexCP assumes the weights are fixed and data-independent. The intuition is that a higher weight should be assigned to a data point that is believed to originate from the same distribution as the test data.

Given weights $w_i \in [0, 1]$, $i \in \mathcal{D}_{\text{cal}}$, the prediction interval of the NexCP method for the test data Y_{n+1} is given by

$$\hat{\mathcal{C}}_{n+1}^\alpha(X_{n+1}) = \left\{ y \in \mathbb{R} : \mathcal{S}(X_{n+1}, y) \leq Q_{1-\alpha} \left(\sum_{i \in \mathcal{D}_{\text{cal}}} \tilde{w}_i \cdot \delta_{s_i} + \tilde{w}_{n+1} \cdot \delta_{+\infty} \right) \right\}, \quad (2)$$

where \tilde{w}_i and \tilde{w}_{n+1} are normalized weights calculated by

$$\tilde{w}_i = \frac{w_i}{\sum_{i \in \mathcal{D}_{\text{cal}}} w_i + 1}, \text{ for } i \in \mathcal{D}_{\text{cal}} \quad \text{and} \quad \tilde{w}_{n+1} = \frac{1}{\sum_{i \in \mathcal{D}_{\text{cal}}} w_i + 1}.$$

By placing different prespecified weights on data points, NexCP is able to deal with data that are not exchangeable and provide robustness against distribution drift. Barber et al. (2023) suggests using an exponential weighting scheme for time series data, where the weights decrease exponentially as data points come from the further in the past.

Theorem 2.2 (NexCP, Barber et al. (2023)). *Let $\mathcal{S}(Z)$ denote a vector of nonconformity scores for data points in the calibration and test sets, and $\mathcal{S}(Z^i)$ denote a vector of nonconformity scores after swapping the test point with the i th data point in the calibration set.*

For any score function \mathcal{S} , and any $\alpha \in (0, 1)$, the nonexchangeable split conformal prediction interval defined in Equation 2 satisfies

$$\mathbb{P}\{Y_{n+1} \in \hat{\mathcal{C}}_{n+1}^\alpha(X_{n+1})\} \geq 1 - \alpha - \sum_{i \in \mathcal{D}_{cal}} \tilde{w}_i \cdot d_{TV}(\mathcal{S}(Z), \mathcal{S}(Z^i)),$$

without the assumption of exchangeability of the data, where d_{TV} denotes the total variation distance between two distributions.

Moreover, if we assume additionally that the nonconformity scores on \mathcal{D}_{cal} are distinct with probability one, then the probability also has the upper bound:

$$\mathbb{P}\{Y_{n+1} \in \hat{\mathcal{C}}_{n+1}^\alpha(X_{n+1})\} < 1 - \alpha + \tilde{w}_{n+1} + \sum_{i \in \mathcal{D}_{cal}} \tilde{w}_i \cdot d_{TV}(\mathcal{S}(Z), \mathcal{S}(Z^i)).$$

2.1.3 Adaptive conformal prediction

The **adaptive conformal prediction** (ACP) proposed by Gibbs & Candès (2021) accounts for nonexchangeability by updating the quantile level in an online manner. Specifically, it treats α as a tunable parameter and estimates it recursively based on the historical performance. ACP can be used to deal with arbitrary online distribution shifts.

Similar to SCP, the initial step involves a random split on the observed data, yielding a training set for fitting a regression model and a withheld calibration set. ACP assumes that there exists an optimal value α_t^* to achieve the desired miscoverage rate α at each time t . To deal with cases where the data generating distribution is shifting over time, ACP recursively estimates the parameter α_t^* on the test points, using the updating equation

$$\alpha_{t+1} = \alpha_t + \gamma(\alpha - \text{err}_t), \quad (3)$$

rather than consistently using the target miscoverage rate α . Here, $\gamma > 0$ is a fixed step size parameter, α_1 is the initial estimate typically set as $\alpha_1 = \alpha$, and $\text{err}_t = \mathbb{1}\{Y_t \notin \hat{\mathcal{C}}_t^{\alpha_t}(X_t)\}$, where $\hat{\mathcal{C}}_t^{\alpha_t}(X_t)$ denotes the prediction set obtained using the $1 - \alpha_t$ quantile for the nonconformity scores available up to and including time t .

This update process adapts the estimation of α_t^* based on the historical frequency of miscoverage in the prediction sets. Specifically, it adjusts upwards (or downwards) the estimate of α_t^* if the prediction sets have shown over-coverage (or under-coverage) of the actual outcomes.

Selecting the parameter γ is pivotal yet challenging. Gibbs & Candès (2021) suggests setting γ in proportion to the degree of variation of the unknown α_t^* over time. Several strategies have been proposed to avoid the necessity of selecting γ . For example, Zaffran et al. (2022) use an adaptive aggregation of multiple ACPs with a set of candidate values for γ , determining weights based on their historical performance. Bastani et al. (2022) propose the multivalid prediction algorithm in which the prediction set is established by selecting a threshold from a sequence of candidate thresholds. However, Gibbs & Candès (2022) have empirically shown that both previous methods fail to promptly adapt to the local changes. To address this limitation, Gibbs & Candès (2022) propose adaptively tuning the step size parameter γ in an online setting, choosing an “optimal” value for γ from a candidate set of values by assessing their historical performance.

Theorem 2.3 (ACP, Gibbs & Candès (2021)). Assume that, with probability one, $Q_{\tau,t}$ is continuous and nondecreasing so that $Q_{0,t} = -\infty$ and $Q_{1,t} = \infty$. Then for any $\alpha \in (0, 1)$, $\gamma > 0$, and any $T \geq 1$ test points, the prediction sets given by ACP satisfy

$$\left| \frac{1}{T} \sum_{t=1}^T \text{err}_t - \alpha \right| \leq \frac{\max\{\alpha_1, 1 - \alpha_1\} + \gamma}{\gamma T}.$$

In particular, this means that the prediction intervals obtained by ACP yield long-run coverage, i.e. $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \text{err}_t = \alpha$.

Remark. Theorem 2.3 suggests that a larger value for γ generally results in less deviation from the target coverage. As there is no restriction on α_t and it can drift below 0 or above 1, a larger γ may lead to frequent output of null or infinite prediction sets in order to quickly adapt to the current miscoverage status.

2.1.4 Conformal PID control

Instead of iteratively updating the miscoverage rate α as in Gibbs & Candès (2021), Angelopoulos, Candès & Tibshirani (2024) draw inspiration from control theory and directly update the quantile estimate q_t in an online fashion to achieve long-run coverage. This method treats the system for producing prediction sets as a **proportional-integral-derivative** controller, later referred to as PID.

The iteration of the PID method is given by

$$q_{t+1} = \underbrace{q_t + \eta(\text{err}_t - \alpha)}_P + \underbrace{r_t \left(\sum_{i=1}^t (\text{err}_i - \alpha) \right)}_I + \underbrace{g'_t}_D.$$

The PID method integrates three modules, namely quantile tracking (P control), error integration (I control), and scorecasting (D control).

The P control module updates the quantile iteratively with a constant learning rate $\eta > 0$. The underlying intuition is similar to that of ACP: it increases (or decreases) the quantile estimate if the prediction set at time t miscovered (or covered) the corresponding realization. ACP can be considered as a special case of the P control, while the P control has the ability to prevent the generation of null or infinite prediction sets after a sequence of miscoverage events.

Theorem 2.4 (The P control, Angelopoulos, Candès & Tibshirani (2024)). *Assume that the nonconformity scores are bounded within $[-b, b]$, for $0 < b < \infty$. Then for any $\alpha \in (0, 1)$, $\eta > 0$, and any $T \geq 1$ the P control iteration satisfies*

$$\left| \frac{1}{T} \sum_{t=1}^T \text{err}_t - \alpha \right| \leq \frac{b + \eta}{\eta T}.$$

In particular, this means that the prediction intervals obtained by the P control iteration yield long-run coverage, i.e. $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \text{err}_t = \alpha$.

The I control involves the sum of historical coverage errors $\sum_{i=1}^t (\text{err}_i - \alpha)$ in a saturation function r_t when updating the quantile estimate, further stabilizing the coverage.

Theorem 2.5 (The I control, Angelopoulos, Candès & Tibshirani (2024)). *Assume that the nonconformity scores are bounded within $[-b, b]$, for $b > 0$, and that the saturation function r_t satisfies*

$$x \geq c \cdot g(t) \implies r_t(x) \geq b, \quad \text{and} \quad x \leq -c \cdot g(t) \implies r_t(x) \leq -b, \quad (4)$$

for positive constants b, c and an admissible function g . Then for any $\alpha \in (0, 1)$, and any $T \geq 1$ the I control iteration satisfies

$$|\frac{1}{T} \sum_{t=1}^T \text{err}_t - \alpha| \leq \frac{c \cdot g(T) + 1}{T}.$$

In particular, this means that the prediction intervals obtained by the I control iteration yield long-run coverage, i.e. $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \text{err}_t = \alpha$.

Finally, the D control g'_t is the forecast of q_{t+1} produced by a scorecaster (forecasting model) fitted using the nonconformity scores available up to and including time t . Instead of reacting to the past miscoverage events, this module looks forward and identifies the leftover signal not captured by the regression model $\hat{\mu}$. However, it may introduce variability and result in wider prediction sets if the scorecaster is aggressive or if there is not much leftover signal in the nonconformity scores.

2.2 Conformal prediction for time series

- Brief literature review of conformal prediction methods or applications on time series data.
- Multi-step conformal prediction.

3 Multi-step conformal prediction for time series

We now consider multi-step time series forecasting problems. In the following sections, we first introduce the setup for time series forecasting problems, and then generalize the existing conformal prediction methods described in Section 2.1 to deal with multi-step time series forecasting. Finally, we explore the properties of forecast errors for optimal multi-step forecasts, and then propose the **autocorrelated multi-step conformal prediction** (AcMCP) method to account for the serial correlation of multi-step forecast errors.

3.1 Setup

Let $z_t = (\mathbf{x}_t, y_t)$ denote the data point (including the response and possibly predictors) at time t . Suppose that, at each time t , we have a forecasting model \hat{f}_t trained using the historical data $z_{1:t}$. Throughout the paper, we assume that the predictors are known into the future. In this way, we perform ex-post forecasting and there is no additional uncertainty introduced from forecasting the exogenous predictors. Using the forecasting model \hat{f}_t , we are able to produce H -step point forecasts, $\{\hat{y}_{t+h|t}\}_{h \in [H]}$, using the future values for the predictors. The task is to employ conformal inference to build H -step prediction intervals, $\{\hat{\mathcal{C}}_{t+h|t}^\alpha(z_{1:t}, \mathbf{x}_{t+1:h})\}_{h \in [H]}$, at the target coverage level $1 - \alpha$. For brevity, we will use $\hat{\mathcal{C}}_{t+h|t}^\alpha$ to denote the h -step-ahead prediction interval produced using the information up to time t .

Sequential split. In time series context, it is infeasible to perform random splitting in split conformal due to the temporal dependency present in the data. Instead, throughout the conformal prediction methods in this section, we use a sequential split to preserve the temporal structure. For example, the t available data points, $z_{1:t}$, are sequentially split into two consecutive sets, a proper training set $\mathcal{D}_{\text{tr}} \subset \{1, \dots, t_r\}$ and a calibration set $\mathcal{D}_{\text{cal}} \subset \{t_r + 1, \dots, t\}$, where $t - t_r \gg h$.

Online learning. Here we consider a generic online learning framework to adapt to all conformal prediction methods we will discuss in subsequent sections. The framework adopts a standard rolling window evaluation strategy. Let the length of the training set be denoted as t_r and the calibration set length as t_c . The online learning framework consists of the following steps.

1. Initialization. Train a forecasting model on the initial proper training set $z_{t-t_r+1:t}$, setting $t = t_r$. Then generate H -step point forecasts $\{\hat{y}_{t_r+h|t_r}\}_{h \in [H]}$ and compute the corresponding nonconformity scores $\{s_{t_r+h|t_r} = \mathcal{S}(z_{1:t_r}, y_{t_r+h})\}_{h \in [H]}$ based on the true values H time steps ahead, i.e. $\{y_{t_r+h}\}_{h \in [H]}$.

2. Recurring procedure. Roll the training set forward by one data point by setting $t \rightarrow t + 1$. Then repeat the step 1 until the nonconformity scores on the entire initial calibration set, $\{s_{t+h|t}\}_{t_r \leq t \leq t_r+t_c-h}$ for $h \in [H]$, are computed.
3. Quantile estimation and prediction interval calculation. Use nonconformity scores obtained from the calibration set to perform quantile estimation and compute H -step prediction intervals on the test set.
4. Online updating. Recursively roll the training set and calibration set forward by one data point to update the nonconformity scores for calibration, and then repeat the step 3 until prediction intervals for the entire test set are obtained, i.e., $\{\hat{\mathcal{C}}_{t+h|t}^\alpha\}_{t_r+t_c \leq t \leq T-h}$ for $h \in [H]$, where $T - t_r - t_c$ is the length of the test set. Therefore, our goal is to achieve long-run coverage in time.

For simplicity, so far we have only presented the nonconformity score defined as the (signed) forecast error

$$s_{t+h|t} = \mathcal{S}(z_{1:t}, y_{t+h}) := y_{t+h} - \hat{f}_t(z_{1:t}, \mathbf{x}_{t+1:h}) = y_{t+h} - \hat{y}_{t+h|t},$$

which is the most commonly used accuracy measure in the context of time series forecasting. We also note that the online learning setting can also be easily adjusted to work with expanding windows for the training and calibration sets.

Remark. With sequential splitting, multiple H -step forecasts and their respective nonconformity scores can be computed on the calibration set. These nonconformity scores have diverse forecast horizons, ranging from 1 to H , i.e., the number of periods between the forecast origin and the time at which nonconformity scores are evaluated. Thus, we can not uniformly treat these nonconformity scores and generate H -step prediction intervals of identical width.

3.2 Related methods extensions to multi-step forecasting

One of the key properties of optimal forecast errors is that the variance of the forecast error $e_{t+h|t}$ is non-decreasing in h (Diebold & Lopez 1996; Patton & Timmermann 2007). Therefore, instead of uniformly treating H -step nonconformity scores and generating H -step prediction intervals of identical width, we consider a setting wherein a separate conformal prediction procedure is applied for each $h \in [H]$ in an online manner.

3.2.1 Online multi-step split conformal prediction

We introduce online **multi-step split conformal prediction** (MSCP) as a generalization of SCP to recursively update all H -step prediction intervals over time. Specifically, for each $h \in [H]$, we consider the following simple online update to construct prediction intervals on the test set:

$$\hat{\mathcal{C}}_{t+h|t}^\alpha = \left\{ y \in \mathbb{R} : s_{t+h|t}^y \leq Q_{1-\alpha} \left(\sum_{i=t-t_c+1}^t \frac{1}{t_c+1} \cdot \delta_{s_{i|i-h}} + \frac{1}{t_c+1} \cdot \delta_{+\infty} \right) \right\}, \quad (5)$$

where $s_{t+h|t}^y := \mathcal{S}(z_{1:t}, y)$ denote the h -step-ahead nonconformity score calculated at time t using a hypothesized test observation y .

- Quantile estimation.
- Theorem?

3.2.2 Online multi-step weighted conformal prediction

The online **multi-step weighted conformal prediction** (MWCP) method adapts the NexCP method to the online setting. MWCP uses weighted quantile estimate for constructing prediction intervals,

contrasting with the MSCP definitions where all nonconformity scores for calibration are implicitly assigned equal weight.

We choose fixed weights $w_i = b^{t+1-i}$, $b \in (0, 1)$ and $i = t - t_c + 1, \dots, t$, for nonconformity scores on the corresponding calibration set. In this setting, weights decay exponentially as the nonconformity scores get older, akin to the rationale behind the exponential smoothing method in time series forecasting. Then for each $h \in [H]$, MWCP consider the online update for h -step-ahead prediction interval:

$$\hat{\mathcal{C}}_{t+h|t}^\alpha = \left\{ y \in \mathbb{R} : s_{t+h|t}^y \leq Q_{1-\alpha} \left(\sum_{i=t-t_c+1}^t \tilde{w}_i \cdot \delta_{s_{i|t-h}} + \tilde{w}_{t+1} \cdot \delta_{+\infty} \right) \right\},$$

where \tilde{w}_i and \tilde{w}_{t+1} are normalized weights given by

$$\tilde{w}_i = \frac{w_i}{\sum_{i=t-t_c+1}^t w_i + 1}, \text{ for } i \in \{t - t_c + 1, \dots, t\} \quad \text{and} \quad \tilde{w}_{t+1} = \frac{1}{\sum_{i=t-t_c+1}^t w_i + 1}.$$

Throughout the experiments in the paper, we will take $b = 0.99$.

- Weighted quantile estimation.
- Theorem?

3.2.3 Multi-step adaptive conformal prediction

In the online learning framework outlined in Section 3.1, we extend the ACP method to address multi-step time series forecasting, introducing the **multi-step adaptive conformal prediction** (MACP) method. Specifically, for each $h \in [H]$, we iteratively estimate α^* (treated as a tunable parameter) using the update equation

$$\alpha_{t+h|t} := \alpha_{t+h-1|t-1} + \gamma (\alpha - \text{err}_{t|t-h}), \quad (6)$$

and compute the h -step-ahead prediction interval using Equation 5 by setting $\alpha = \alpha_{t+h|t}$. Here, $\gamma > 0$ denotes a fixed step size parameter, $\alpha_{h+1|1}$ denotes the initial estimate typically set to α , and $\text{err}_{t|t-h}$ denotes the miscoverage event $\text{err}_{t|t-h} = \mathbb{1} \{y_t \notin \hat{\mathcal{C}}_{t|t-h}^{\alpha_{t|t-h}}\}$.

Equation 5 indicates that the correction to the estimation of α at time $t + h$ is determined by the historical miscoverage frequency up to time t . At each iteration, we raise the estimate of α used for quantile estimation at time $t + h$ if $\hat{\mathcal{C}}_{t|t-h}^{\alpha_{t|t-h}}$ covered y_t , whereas we lower the estimate if $\hat{\mathcal{C}}_{t|t-h}^{\alpha_{t|t-h}}$ miscovered y_t . Thus the miscoverage event has a delayed impact on the estimation of α over h periods, indicating that the correction of the α estimate becomes less prompt with increasing values of h . Particularly, Equation 6 reduces to the update for ACP as given by Equation 3 for $h = 1$.

It should be noted that the update equation $\alpha_{t+1|t-h+1} := \alpha_{t|t-h} + \gamma (\alpha - \text{err}_{t|t-h})$ is not utilized in this context. This arises from its limitation that the available information at time t is insufficient to derive the α estimate for forecasting h steps ahead.

In our experiments we will use $\gamma = 0.01$.

- Theorem?

3.2.4 Multi-step conformal PID control

We propose **multi-step conformal PID control** method (referred to as MPID hereafter), a generalization of the PID method to deal with multi-step forecasting.

For each individual forecast horizon $h \in [H]$, the iteration of the h -step-ahead quantile estimate is given by

$$q_{t+h|t} = \underbrace{q_{t+h-1|t-1} + \eta (\text{err}_{t|t-h} - \alpha)}_P + r_t \left(\underbrace{\sum_{i=h+1}^t (\text{err}_{i|i-h} - \alpha)}_I \right) + \underbrace{\hat{s}_{t+h|t}}_D,$$

where as before, $\eta > 0$ is a constant learning rate, and r_t is a saturation function that adheres to the following conditions

$$x \geq c \cdot g(t-h) \implies r_t(x) \geq b, \quad \text{and} \quad x \leq -c \cdot g(t-h) \implies r_t(x) \leq -b, \quad (7)$$

for constant $b, c > 0$, and an admissible function g that is sublinear, nonnegative, nondecreasing. Here, we define $\hat{s}_{t+h|t}$ as the h -step-ahead forecast of the nonconformity score (defined as the forecast error), produced by any suitable scorecaster (forecasting model) trained using the h -step-ahead nonconformity scores available up to and including time t . With this updating equation, we can obtain all required h -step-ahead forecast intervals using information available up to time t .

The P control in MPID shows a delayed correction of the α estimate for a length of h periods. The I control accounts for the cumulative historical coverage errors associated with h -step-ahead prediction intervals during the update process, thereby enhancing the stability of the interval coverage. Moreover, the D control performs h -step-ahead forecasting, which tends to result in increased forecast variance for larger forecast horizon h .

- Theorem?

3.3 Autocorrelated multi-step conformal prediction

In the PID method proposed by Angelopoulos, Candès & Tibshirani (2024), a notable feature is the inclusion of a scorecaster, a model trained on the score sequence, to forecast the next score. The rationale behind it is to residualize out any leftover signal in the score distribution not captured by the base forecasting model, such as trend and seasonality. However, in the context of time series forecasting, good forecasts are essential for making good decisions. We naturally expect to use a good forecasting model and ensure there is no useful signal in forecast errors (defined as nonconformity scores in this paper). If the forecasts are not optimal, the forecasting model should be improved to enhance its performance. Hence, we typically assume the use of a good forecasting model, and therefore, relying on another model to predict forecast errors to capture leftover information is not a commonly used solution. Moreover, the inclusion of a scorecaster often only introduces variance to the quantile estimate, resulting in wider prediction intervals.

On the other hand, in our general setup outlined in Section 1 and Section 3.1, the DGP of a time series may depend on its own past, along with other exogenous predictors and their histories. Consequently, the h -step-ahead forecast errors $e_{t+h|t}$ may depend on the forecast errors from the past $h-1$ steps, i.e. $e_{t+1|t}, \dots, e_{t+h-1|t}$, and forecast errors may accumulate over the forecast horizon. However, no conformal prediction methods have taken this potential dependence into account in their methodological construction.

In this section, we will explore the properties of multi-step forecast errors and propose a novel conformal prediction method that considers the autocorrelations of multi-step forecast errors.

3.3.1 Properties of multi-step forecast errors

We assume that a time series $\{y_t\}_{t \geq 1}$ is generated by a general non-stationary autoregressive process given by:

$$y_t = f_t(y_{(t-d):(t-1)}, \mathbf{x}_{(t-k):t}) + \epsilon_t, \quad (8)$$

where f_t is considered a non-linear function in d lagged values of y_t ($y_{(t-d):(t-1)}$) and the current value along with the preceding k values of the exogenous predictors ($\mathbf{x}_{(t-k):t}$), and ϵ_t is white noise.

It is well-established in the forecasting literature that, for optimal h -step-ahead forecasts, the sequence of forecast errors is *at most* an $\text{MA}(h-1)$ process (Harvey, Leybourne & Newbold 1997; Diebold 2017). We now present the property under the assumption of a non-stationary autoregressive DGP and provide its proof in Section B.1 based on the proof of Proposition 3.2 that we will introduce later.

Proposition 3.1 (MA($h-1$) process for h -step-ahead optimal forecast errors). *Let $\{y_t\}_{t \geq 1}$ be a time series generated by a general non-stationary autoregressive process as given in Equation 8. Assume that the exogenous predictors are known into the future if applicable. The forecast errors of optimal h -step-ahead forecasts follow an approximate MA($h-1$) process:*

$$e_{t+h|t} = c + \epsilon_{t+h} + \theta_1 \epsilon_{t+h-1} + \cdots + \theta_{h-1} \epsilon_{t+1},$$

where $c = 0$, motivated by the property that optimal forecasts are unbiased.

We proceed by exploring the autocorrelations of multi-step forecast errors for optimal forecasts.

Proposition 3.2 (Autocorrelations of multi-step optimal forecast errors). *Let $\{y_t\}_{t \geq 1}$ be a time series generated by a general non-stationary autoregressive process as given in Equation 8. Assume that the exogenous predictors are known into the future if applicable. The forecast errors of optimal h -step-ahead forecasts are at most an approximate AR($h-1$) process given by:*

$$e_{t+h|t} = c + \epsilon_{t+h} + \phi_1 e_{t+h-1|t} + \cdots + \phi_{h-1} e_{t+1|t},$$

where $e_{t+h|t}$ is the h -step-ahead forecast error with variance non-decreasing in h , and the intercept $c = 0$, given the property that optimal forecasts are unbiased.

Proof. Let $\hat{y}_{t+h|t}$ be the optimal h -step-ahead point forecast, and $e_{t+h|t}$ be the h -step-ahead forecast error. Denote $\mathbf{u}_{t+h} = \mathbf{x}_{(t-k+h):(t+h)}$. Then we have

$$\begin{aligned} \hat{y}_{t+1|t} &= f_{t+1}(y_t, \dots, y_{t-d+1}, \mathbf{u}_{t+1}), \text{ for } h = 1, \\ \hat{y}_{t+h|t} &= f_{t+h}(\hat{y}_{t+h-1|t}, \dots, \hat{y}_{t+1|t}, y_t, \dots, y_{t+h-d}, \mathbf{u}_{t+h}), \text{ for } 1 < h \leq d. \text{ and} \\ \hat{y}_{t+h|t} &= f_{t+h}(\hat{y}_{t+h-1|t}, \dots, \hat{y}_{t+h-d|t}, \mathbf{u}_{t+h}), \text{ for } h > d. \end{aligned}$$

For $h = 1$, we simply have $e_{t+1|t} = \epsilon_{t+1}$.

For $1 < h \leq d$, applying the first order Taylor series expansion, we can write

$$\begin{aligned} y_{t+h} &= f_{t+h}(y_{t+h-1}, \dots, y_{t+h-d}, \mathbf{u}_{t+h}) + \epsilon_{t+h} \\ &= f_{t+h}(\hat{y}_{t+h-1|t} + e_{t+h-1|t}, \dots, \hat{y}_{t+1|t} + e_{t+1|t}, y_t, \dots, y_{t+h-d}, \mathbf{u}_{t+h}) + \epsilon_{t+h} \\ &\approx_{\text{te}} f_{t+h}(\mathbf{a}) + Df_{t+h}(\mathbf{a})(\mathbf{v} - \mathbf{a}) + \epsilon_{t+h} \\ &= f_{t+h}(\hat{y}_{t+h-1|t}, \dots, \hat{y}_{t+1|t}, y_t, \dots, y_{t+h-d}, \mathbf{u}_{t+h}) \\ &\quad + e_{t+h-1|t} \frac{\partial f_{t+h}(\mathbf{a})}{\partial v_1} + \cdots + e_{t+2|t} \frac{\partial f_{t+h}(\mathbf{a})}{\partial v_{h-2}} + e_{t+1|t} \frac{\partial f_{t+h}(\mathbf{a})}{\partial v_{h-1}} + \epsilon_{t+h} \\ &= \hat{y}_{t+h|t} + e_{t+h|t}, \end{aligned}$$

where $\mathbf{v} = (y_{t+h-1}, \dots, y_{t+h-d}, \mathbf{u}_{t+h})$, $\mathbf{a} = (\hat{y}_{t+h-1|t}, \dots, \hat{y}_{t+1|t}, y_t, \dots, y_{t+h-d}, \mathbf{u}_{t+h})$, $Df_{t+h}(\mathbf{a})$ denotes the matrix of partial derivative of $f_{t+h}(\mathbf{v})$ at $\mathbf{v} = \mathbf{a}$, and $\frac{\partial}{\partial v_i}$ denotes the partial derivative with respect to the i th component in f_{t+h} .

Similarly, for $h > d$, we can write

$$\begin{aligned} y_{t+h} &= f_{t+h}(y_{t+h-1}, \dots, y_{t+h-d}, \mathbf{u}_{t+h}) + \epsilon_{t+h} \\ &= f_{t+h}(\hat{y}_{t+h-1|t} + e_{t+h-1|t}, \dots, \hat{y}_{t+h-d|t} + e_{t+h-d|t}, \mathbf{u}_{t+h}) + \epsilon_{t+h} \\ &\approx f_{t+h}(\mathbf{a}) + Df_{t+h}(\mathbf{a})(\mathbf{v} - \mathbf{a}) + \epsilon_{t+h} \\ &= f_{t+h}(\hat{y}_{t+h-1|t}, \dots, \hat{y}_{t+h-d|t}, \mathbf{u}_{t+h}) \\ &\quad + e_{t+h-1|t} \frac{\partial f_{t+h}(\mathbf{a})}{\partial v_1} + e_{t+h-d|t} \frac{\partial f_{t+h}(\mathbf{a})}{\partial v_d} + \epsilon_{t+h} \\ &= \hat{y}_{t+h|t} + e_{t+h|t}, \end{aligned}$$

Therefore, the forecast errors of optimal h -step-ahead forecasts are at most an approximate $AR(h-1)$ process. \square

Proposition 3.2 can be viewed as an extension of Proposition 3.1. It suggests that the h -step ahead forecast error, $e_{t+h|t}$, is serially correlated with the forecast errors from at most the past $h-1$ steps, i.e., $e_{t+1|t}, \dots, e_{t+h-1|t}$. However, we note that the autocorrelation among errors associated with optimal forecasts can not be used to improve forecasting performance, as it does not incorporate any information available when the forecast was made. It is reasonable because if we could forecast the forecast error, we could improve the forecast, indicating that the initial forecast couldn't have been optimal.

The proof of Proposition 3.2 suggests that, if f_t is a linear autoregressive model, then the AR coefficients are the linear coefficients of the optimal forecasting model. However, when f_t takes on a more complex non-linear structure, the AR coefficients become complicated functions of observed data and unobserved model coefficients.

3.3.2 The AcMCP method

Inspired by the properties of multi-step forecast errors discussed in Section 3.3.1, we now propose the **autocorrelated multi-step conformal prediction** (AcMCP) method. Unlike extensions of existing conformal prediction methods that treat multi-step forecasting as independent events (see Section 3.2), the AcMCP method integrates the autocorrelations inherent in multi-step forecast errors, thereby making the output multi-step prediction intervals more logically structured.

The AcMCP method updates the quantile estimate q_t in an online setting to achieve the goal of long-run coverage. Specifically, the iteration of the h -step-ahead quantile estimate is given by

$$q_{t+h|t} = q_{t+h-1|t-1} + \eta(\text{err}_{t|t-h} - \alpha) + r_t \left(\sum_{i=h+1}^t w_i (\text{err}_{i|i-h} - \alpha) \right) + \tilde{e}_{t+h|t},$$

for $h \in [H]$. Obviously, the AcMCP method can be viewed as a further extension of the MPID method. Nevertheless, AcMCP diverges from MPID with several innovations and differences.

First, we assign different weights $w_i \in [0, 1]$ to historical coverage errors with the intuition that a higher weight should be given to a historical coverage error obtained closer to the present forecast period. Such a setting is particularly suitable for time series forecasting problems. In line with MWCP we choose normalized exponentially decaying weights, i.e., $w_i = \frac{b^{t-i}}{\sum_{i=h+1}^t b^{t-i}}$, for $b \in (0, 1)$. Throughout the experiments in the paper, we will take $b = 0.99$.

Additionally, $\tilde{e}_{t+h|t}$ is a forecast combination of two simple models: one being an $\text{MA}(h-1)$ model trained on h -step-ahead forecast errors available up to and including time t (i.e. $e_{1+h|1}, \dots, e_{t|t-h}$), and the other an $\text{AR}(h-1)$ model (with respect to h instead of t) trained by regressing $e_{t+h|t}$ on forecast errors from past steps (i.e. $e_{t+h-1|t}, \dots, e_{t+1|t}$). Thus, we perform multi-step conformal prediction recursively, contrasting with the independent approach employed in MPID. Moreover, the inclusion of $\tilde{e}_{t+h|t}$ is not intended to forecast the nonconformity scores (i.e., forecast errors in this paper), but rather to incorporate autocorrelations present in multi-step forecast errors within the resulting multi-step prediction intervals. As previously explained, in the context of time series forecasting, we typically assume the use of a good base forecasting model, making it unnecessary to train an additional model to predict forecast errors in order to capture leftover information. If the forecasts are not optimal, the base forecasting model should be improved to enhance its performance.

Proposition 3.3. *Let $\{s_{t+h|t}\}_{t \in \mathbb{N}}$ be any sequence of numbers in $[-b, b]$ for any $h \in [H]$, where $b > 0$, and may be infinite. Assume that r_t is a saturation function obeying Equation 7, for an admissible function g . Then the iteration $q_{t+h|t} = r_t(\sum_{i=h+1}^t (\text{err}_{i|i-h} - \alpha))$ satisfies*

$$\left| \frac{1}{T-h} \sum_{t=h+1}^T (\text{err}_{t|t-h} - \alpha) \right| \leq \frac{c \cdot g(T-h) + h}{T-h}, \quad (9)$$

for any $T \geq h+1$, where $c > 0$ is the constant in Equation 7.

In particular, this means the prediction intervals obtained by the iteration yield long-run coverage, i.e., $\lim_{T \rightarrow \infty} \frac{1}{T-h} \sum_{t=h+1}^T \text{err}_{t|t-h} = \alpha$.

Proof. Let $E_T = \sum_{t=h+1}^T (\text{err}_{t|t-h} - \alpha)$. The inequality given by Equation 9 can be expressed as $|E_T| \leq c \cdot g(T-h) + h$. We will prove one side of the absolute inequality, specifically $E_T \leq c \cdot g(T-h) + h$, with the other side following analogously. We proceed with the proof using induction.

For $T = h+1, \dots, 2h$, $E_T = \sum_{t=h+1}^T (\text{err}_{t|t-h} - \alpha) \leq (T-h) - (T-h)\alpha \leq T-h \leq h \leq c \cdot g(T-h) + h$ as $c > 0$, $h \geq 1$, g is nonnegative, and $\text{err}_{t|t-h} \leq 1$. Thus, Equation 9 holds for $T = h+1, \dots, 2h$.

Now, assuming Equation 9 is true up to T . We partition the argument into $h+1$ cases:

$$\begin{cases} cg(T-h) + h - 1 < E_T \leq cg(T-h) + h, & \dots \text{ case (1)} \\ cg(T-h) + h - 2 < E_T \leq cg(T-h) + h - 1, & \dots \text{ case (2)} \\ \dots & \\ cg(T-h) < E_T \leq cg(T-h) + 1, & \dots \text{ case (h)} \\ E_T \leq cg(T-h). & \dots \text{ case (h+1)} \end{cases}$$

In case (1), we observe that $E_T > cg(T-h) + h - 1 > cg(T-h)$, implying $q_{T+h|T} = r_t(E_T) \geq b$ according to Equation 7. Thus, $s_{T+h|T} \leq q_{T+h|T}$ and $\text{err}_{T+h|T} = 0$. Furthermore, we have $E_{T-1} = E_T - (\text{err}_{T|T-h} - \alpha) > cg(T-h) + h - 2 > cg(T-h-1)$ as g is nondecreasing. This implies $q_{T+h-1|T-1} = r_t(E_{T-1}) \geq b$, hence $s_{T+h-1|T-1} \leq q_{T+h-1|T-1}$ and $\text{err}_{T+h-1|T-1} = 0$. Similarly, $E_{T-2} = E_{T-1} - (\text{err}_{T-1|T-h-1} - \alpha) > cg(T-h) + h - 3 > cg(T-h-2)$, thus $\text{err}_{T+h-2|T-2} = 0$. This process iterates, leading to $\text{err}_{T+h|T} = \text{err}_{T+h-1|T-1} = \dots = \text{err}_{T+1|T-h+1} = 0$. Consequently,

$$E_{T+h} = E_T + \sum_{t=T+1}^{T+h} (\text{err}_{t|t-h} - \alpha) \leq cg(T-h) + h - h\alpha \leq cg(T) + h,$$

which is the desired result at $T+h$.

In case (2), we observe that $E_T > cg(T-h) + h - 2 > cg(T-h)$, thus $s_{T+h|T} \leq q_{T+h|T}$ and $\text{err}_{T+h|T} = 0$. Moving forward, we have $\text{err}_{T+h|T} = \text{err}_{T+h-1|T-1} = \dots = \text{err}_{T+2|T-h+2} = 0$. Along with $\text{err}_{T+1|T-h+1} \leq 1$, this means that

$$E_{T+h} = E_T + \sum_{t=T+1}^{T+h} (\text{err}_{t|t-h} - \alpha) \leq cg(T-h) + h - 1 + 1 - h\alpha \leq cg(T) + h,$$

which again gives the desired result at $T+h$.

Similarly, in cases (3)-(h), we can always get the desired result at $T+h$.

In case (h+1), noticing $E_T \leq cg(T-h)$, and simply using $\text{err}_{T+h-i|T-i} \leq 1$ for $i = 0, \dots, h-1$, we have

$$E_{T+h} = E_T + \sum_{t=T+1}^{T+h} (\text{err}_{t|t-h} - \alpha) \leq cg(T-h) + h - h\alpha \leq cg(T) + h.$$

Therefore, we can deduce the desired outcome at any $T \geq h+1$. This completes the proof for the first part of Proposition 3.3.

Regarding the second part, $g(t-h)/(t-h) \rightarrow 0$ as $t \rightarrow \infty$ due to the sublinearity of the admissible function g . Hence the second part holds trivially. \square

- Remark on Proposition 3.3.
- Choosing the learning rate? Analysis of the parameter effect.

4 Simulations

- Evaluated on two properties: validity and efficiency.
- Evaluation and visualization.
 - an easy-to-interpret visualisation combining validity and efficiency: plot average/median interval length against mean coverage for different methods or different parameter values.
 - instead of plotting forecasts and prediction intervals, plot errors and quantile estimate.

5 Applications

- COVID-19 case counts
- Stock price
- Electricity demand
- Hierarchical time series (tourism data)?

6 Discussion

- Large forecast horizon
- Evaluation of multi-step intervals
- Use the forecasts of predictors, introducing more variability

References

- Angelopoulos, AN, EJ Candès & RJ Tibshirani (2024). Conformal PID control for time series prediction. *Advances in Neural Information Processing Systems* **36**.
- Barber, RF, EJ Candès, A Ramdas & RJ Tibshirani (1, 2023). Conformal prediction beyond exchangeability. *The Annals of Statistics* **51**(2).
- Bastani, O, V Gupta, C Jung, G Noarov, R Ramalingam & A Roth (2022). Practical adversarial multivalid conformal prediction. *Advances in Neural Information Processing Systems* **35**, 29362–29373.
- Diebold, FX (2017). *Forecasting*. Department of Economics, University of Pennsylvania. <http://www.ssc.upenn.edu/~fdiebold/Textbooks.html>.
- Diebold, FX & JA Lopez (1996). “Forecast evaluation and combination”. In: *Statistical Methods in Finance*. Vol. 14. Handbook of Statistics. Elsevier, pp.241–268.
- Gibbs, I & E Candès (2021). Adaptive conformal inference under distribution shift. *Advances in Neural Information Processing Systems* **34**, 1660–1672.
- Gibbs, I & E Candès (2022). Conformal inference for online prediction with arbitrary distribution shifts. *arXiv preprint arXiv:2208.08401*.
- Harvey, D, S Leybourne & P Newbold (1997). Testing the equality of prediction mean squared errors. *International Journal of forecasting* **13**(2), 281–291.
- Lei, J, M G’Sell, A Rinaldo, RJ Tibshirani & L Wasserman (2018). Distribution-free predictive inference for regression. *Journal of the American Statistical Association* **113**(523), 1094–1111.
- Papadopoulos, H, K Proedrou, V Vovk & A Gammerman (2002). Inductive confidence machines for regression. In: *Machine learning: ECML 2002: 13th European conference on machine learning Helsinki, Finland, August 19–23, 2002 proceedings 13*. Springer, pp.345–356.
- Patton, AJ & A Timmermann (2007). Properties of optimal forecasts under asymmetric loss and nonlinearity. *Journal of Econometrics* **140**(2), 884–918.
- Sommer, B (2023). “Forecasting and decision-making for empty container repositioning”. PhD thesis.
- Vovk, V, A Gammerman & G Shafer (2005). *Algorithmic Learning in a Random World*. Springer-Verlag. <http://dx.doi.org/10.1007/b106715>.
- Zaffran, M, O Féron, Y Goude, J Josse & A Dieuleveut (2022). Adaptive conformal predictions for time series. In: *International Conference on Machine Learning*. PMLR, pp.25834–25866.

Appendix A Weighted quantile estimation

- median-unbiased estimates

Appendix B Proof

B.1 Proof of Proposition 3.1

Proof. Here, we give the proof of Proposition 3.1 based on Proposition 3.2, and the idea is motivated by Sommer (2023).

For 1-step-ahead optimal forecast, we have

$$y_{t+1} = f_{t+1}(y_{(t-d+1):t}, \mathbf{x}_{(t-k+1):(t+1)}) + \epsilon_{t+1},$$

so $e_{t+1|t} = \epsilon_{t+1}$.

Based on Proposition 3.2, we have

$$\begin{aligned} e_{t+2|t} &= \epsilon_{t+2} + \phi_1^{(2)} e_{t+1|t} \\ e_{t+3|t} &= \epsilon_{t+3} + \phi_1^{(3)} e_{t+2|t} + \phi_2^{(3)} e_{t+1|t} \\ &\dots \\ e_{t+d|t} &= \epsilon_{t+d} + \phi_1^{(d)} e_{t+d-1|t} + \dots + \phi_{d-1}^{(d)} e_{t+1|t} \\ e_{t+d+1|t} &= \epsilon_{t+d+1} + \phi_1^{(d+1)} e_{t+d|t} + \dots + \phi_{d-1}^{(d+1)} e_{t+2|t} + \phi_d^{(d+1)} e_{t+1|t} \\ &\dots \\ e_{t+H|t} &= \epsilon_{t+H} + \phi_1^{(H)} e_{t+H-1|t} + \dots + \phi_{d-1}^{(H)} e_{t+H-d+1|t} + \phi_d^{(H)} e_{t+H-d|t}, \quad H > d + 1. \end{aligned}$$

Substituting all equations above into the following equation, we can obtain

$$e_{t+h|t} = \epsilon_{t+h} + \sum_{i=1}^{h-1} \theta_i \epsilon_{t+h-i}, \text{ for each } h \in [H],$$

where θ_i is a complex function derived from the AR coefficients of all $\text{AR}(j-1)$ models, for $j = 1, 2, \dots, h-1$. So we conclude that the h -step-ahead forecast error sequence $\{e_{t+h|t}\}$ follows an $\text{MA}(h-1)$ process. \square