

# Uncertainty Estimation for High-dimensional Nonparametric Forecasting

Nuwani Palihawadana

**Joint work with :** Rob Hyndman, Xiaoqian Wang

1 July 2025

# Nonparametric forecasting

$$y_t = f(\mathbf{a}_t, \mathbf{a}_{t-1}, \dots, \mathbf{a}_{t-s}, y_{t-1}, \dots, y_{t-k}) + \varepsilon_t,$$

- $y_t$  – observed response at time  $t$
- $f$  – arbitrary function
- $\mathbf{a}_t$  – vector of exogenous variables at time  $t$
- $\varepsilon_t$  – stationary error with mean zero and constant variance  $\sigma^2$ .

# Nonparametric forecasting

$$y_t = f(\mathbf{a}_t, \mathbf{a}_{t-1}, \dots, \mathbf{a}_{t-s}, y_{t-1}, \dots, y_{t-k}) + \varepsilon_t,$$

- $y_t$  – observed response at time  $t$
- $f$  – arbitrary function
- $\mathbf{a}_t$  – vector of exogenous variables at time  $t$
- $\varepsilon_t$  – stationary error with mean zero and constant variance  $\sigma^2$ .

We also define:

- $\mathbf{z}_t = (y_t, \mathbf{x}_t)$  – observation at time  $t$
- $\mathbf{x}_t = (\mathbf{a}_t, \mathbf{a}_{t-1}, \dots, \mathbf{a}_{t-s}, y_{t-1}, \dots, y_{t-k})$  – vector of all predictors at time  $t$ .
- $\hat{y}_t = \hat{f}(\mathbf{x}_t)$  – estimated model
- $e_t = y_t - \hat{y}_t$  – model residuals

# Sparse multiple index (SMI) model

$$y_i = \beta_0 + \sum_{j=1}^p g_j(\alpha_j^T \mathbf{x}_{ij}) + \sum_{k=1}^d f_k(w_{ik}) + \theta^T \mathbf{u}_i + \varepsilon_i, \quad i = 1, \dots, n,$$

- $y_i$  – univariate response
- $\mathbf{x}_{ij} \in \mathbb{R}^{\ell_j}, j = 1, \dots, p$  –  $p$  subsets of predictors entering indices
- $\alpha_j$  –  $\ell_j$ -dimensional vectors of index coefficients
- $g_j, f_k$  – smooth nonlinear functions
- Additional predictors :
  - ▶  $w_{ik}$  – nonlinear
  - ▶  $\mathbf{u}_i$  – linear

# Sparse multiple index (SMI) model

$$y_i = \beta_0 + \sum_{j=1}^p g_j(\alpha_j^T \mathbf{x}_{ij}) + \sum_{k=1}^d f_k(w_{ik}) + \theta^T \mathbf{u}_i + \varepsilon_i, \quad i = 1, \dots, n,$$

- $y_i$  – univariate response
- $\mathbf{x}_{ij} \in \mathbb{R}^{\ell_j}, j = 1, \dots, p$  –  $p$  subsets of predictors entering indices
- $\alpha_j$  –  $\ell_j$ -dimensional vectors of index coefficients
- $g_j, f_k$  – smooth nonlinear functions
- Additional predictors :
  - ▶  $w_{ik}$  – nonlinear
  - ▶  $\mathbf{u}_i$  – linear

Allow elements equal to zero in  
 $\alpha_j$  – "Sparse"

# Sparse multiple index (SMI) model

$$y_i = \beta_0 + \sum_{j=1}^p g_j(\alpha_j^T \mathbf{x}_{ij}) + \sum_{k=1}^d f_k(w_{ik}) + \theta^T \mathbf{u}_i + \varepsilon_i, \quad i = 1, \dots, n,$$

- $y_i$  – univariate response
- $\mathbf{x}_{ij} \in \mathbb{R}^{\ell_j}, j = 1, \dots, p$  –  $p$  subsets of predictors entering indices
- $\alpha_j$  –  $\ell_j$ -dimensional vectors of index coefficients
- $g_j, f_k$  – smooth nonlinear functions
- Additional predictors :
  - ▶  $w_{ik}$  – nonlinear
  - ▶  $\mathbf{u}_i$  – linear

Both "p" and the predictor grouping among indices are unknown.

Overlapping of predictors among indices is not allowed.

# Benchmarks

- Nonparametric additive model with backward elimination (Backward):
  - ▶ No linear combinations (indices)
  - ▶ Fully additive

# Benchmarks

- Nonparametric additive model with backward elimination (Backward):
  - ▶ No linear combinations (indices)
  - ▶ Fully additive
  
- Groupwise Additive Index Model (GAIM):
  - ▶ Predefined predictor groups
  - ▶ No overlapping predictors among groups

# Benchmarks

- Nonparametric additive model with backward elimination (Backward):
  - ▶ No linear combinations (indices)
  - ▶ Fully additive
- Groupwise Additive Index Model (GAIM):
  - ▶ Predefined predictor groups
  - ▶ No overlapping predictors among groups
- Projection Pursuit Regression model (PPR):
  - ▶ All predictors enter all indices

# Forecast uncertainty

- Uncertainty of a forecast → **Prediction Interval (PI)**

# Forecast uncertainty

- Uncertainty of a forecast → **Prediction Interval (PI)**
- Theoretical  $100(1 - \alpha)\%$  prediction interval:

$$\hat{y}_{t+h|t} \pm z_{\alpha/2} \times \hat{\sigma}_h,$$

where

- ▶  $y$  – time series  $y_1, \dots, y_T$
- ▶  $\hat{y}_{t+h|t}$  –  $h$ -step-ahead point forecast for  $y_{t+h}$  given observations up to  $t$
- ▶  $z_{\alpha/2}$  –  $\alpha/2$  quantile of standard normal distribution
- ▶  $\hat{\sigma}_h$  – estimate of std. deviation of  $h$ -step forecast distribution

# Forecast uncertainty

- Uncertainty of a forecast → **Prediction Interval (PI)**
- Theoretical  $100(1 - \alpha)\%$  prediction interval:

$$\hat{y}_{t+h|t} \pm z_{\alpha/2} \times \hat{\sigma}_h,$$

where

- ▶  $y$  – time series  $y_1, \dots, y_T$
- ▶  $\hat{y}_{t+h|t}$  –  $h$ -step-ahead point forecast for  $y_{t+h}$  given observations up to  $t$
- ▶  $z_{\alpha/2}$  –  $\alpha/2$  quantile of standard normal distribution
- ▶  $\hat{\sigma}_h$  – estimate of std. deviation of  $h$ -step forecast distribution

- Main issue:
  - ▶ Difficult to analytically calculate  $h$ -step forecast variances for  $h > 1$

# Block bootstrap

- Resampling from empirical distribution of historical model residuals  
→ Bootstrapping

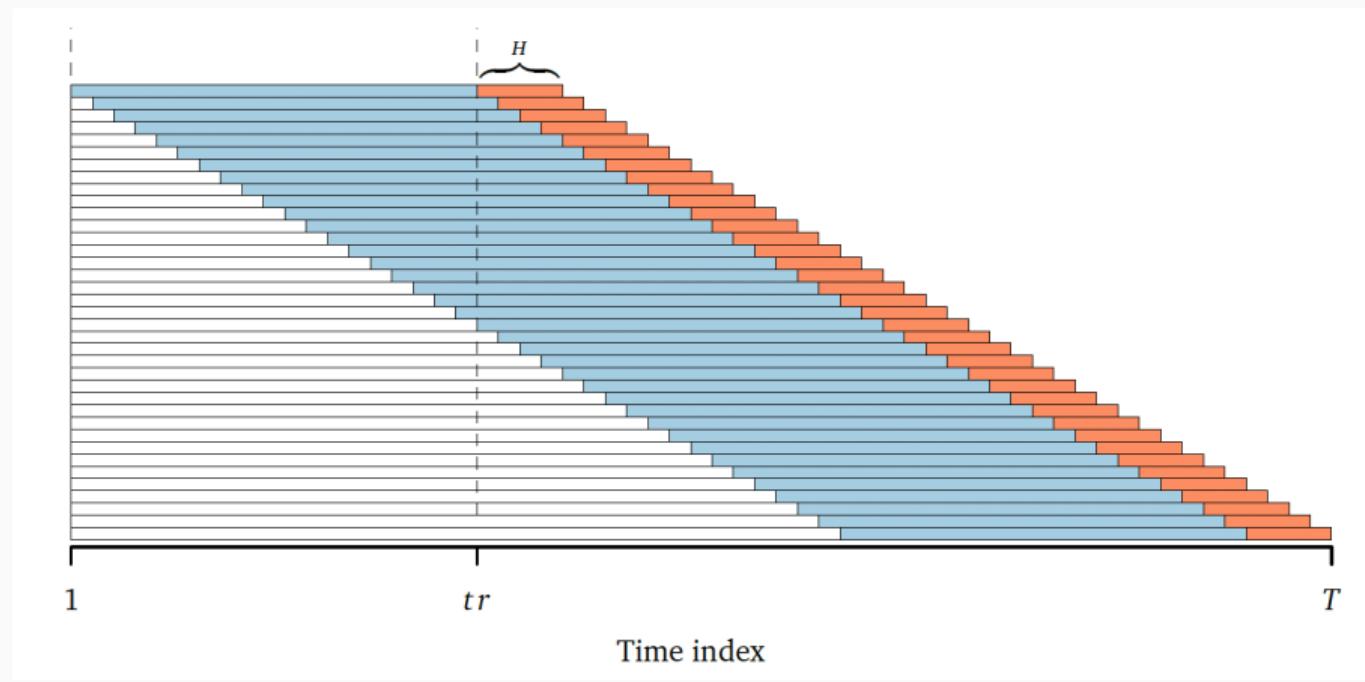
# Block bootstrap

- Resampling from empirical distribution of historical model residuals  
→ Bootstrapping
- Randomly resample blocks from the historical model residuals, and join together → **Block Bootstrapping**
- Retains serial correlation in the data

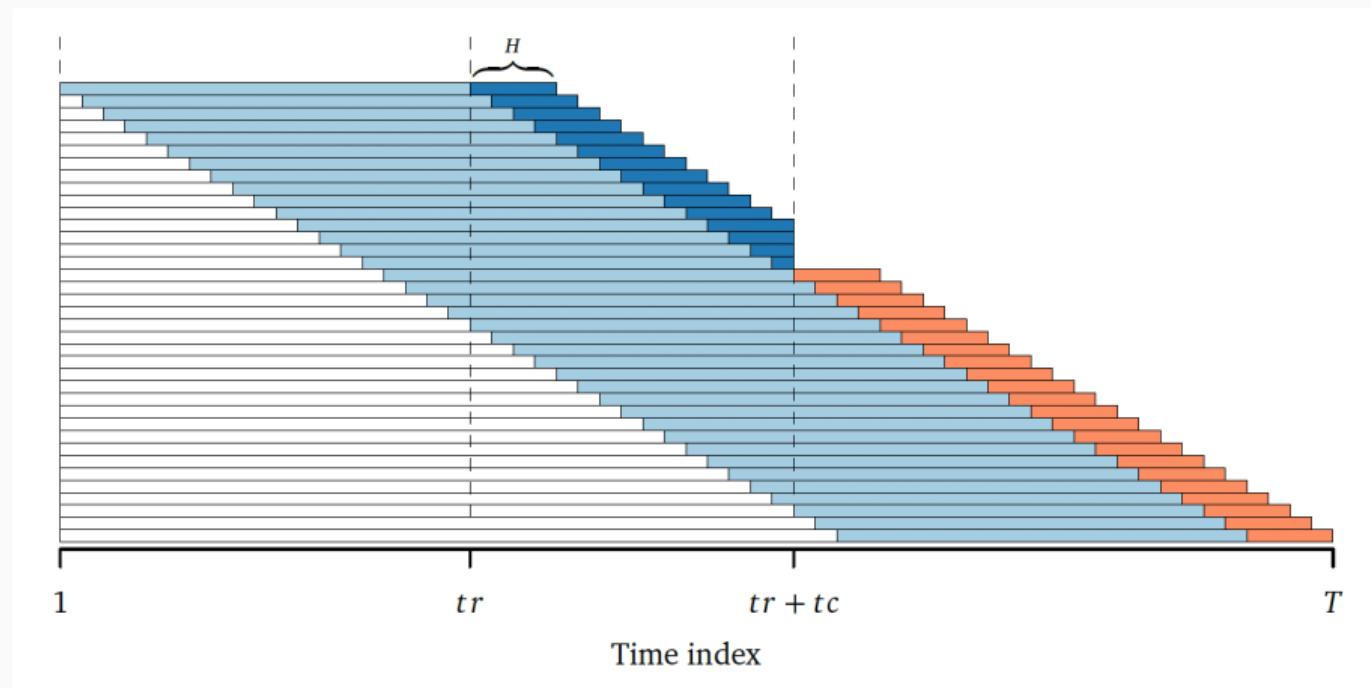
# Block bootstrap

- Resampling from empirical distribution of historical model residuals  
→ Bootstrapping
- Randomly resample blocks from the historical model residuals, and join together → **Block Bootstrapping**
- Retains serial correlation in the data
- **block length:**
  - ▶ Long enough to capture autocorrelation patterns
  - ▶ Short enough to construct sufficient number of blocks

# Block bootstrap



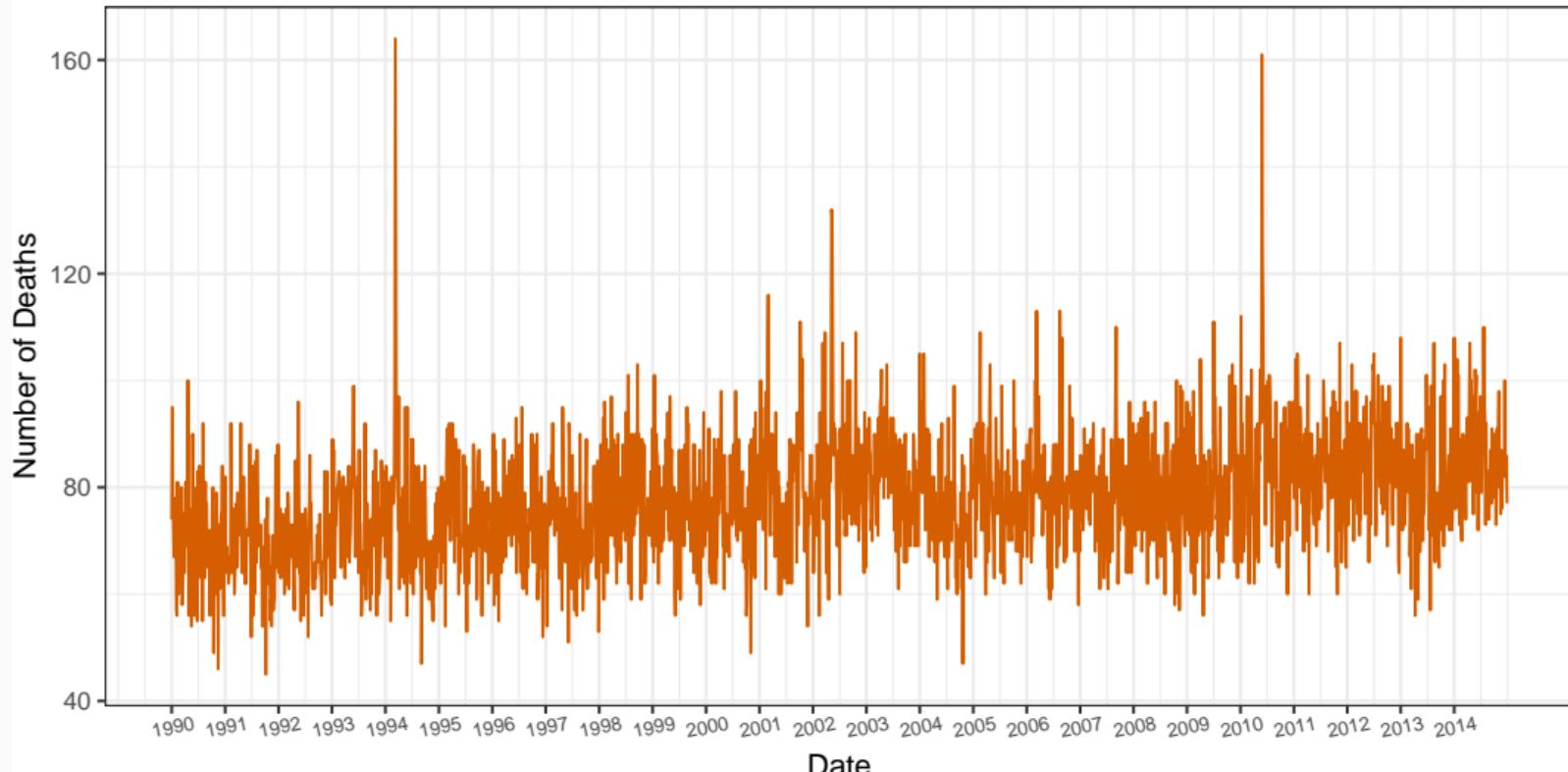
# Conformal prediction



# Conformal bootstrap

# Forecasting heat exposure-related daily mortality

Daily Deaths in Summer – Montreal, Canada



# Forecasting heat exposure-related daily mortality

## Data

- **Response:** Daily deaths in Summer
  - 1990 to 2014 – Montreal, Canada
- **Index Variables:**
  - ▶ Death lags
  - ▶ Max temperature lags
  - ▶ Min temperature lags
  - ▶ Vapor pressure lags
- **Nonlinear:** DOS (day of the season),  
Year

# Forecasting heat exposure-related daily mortality

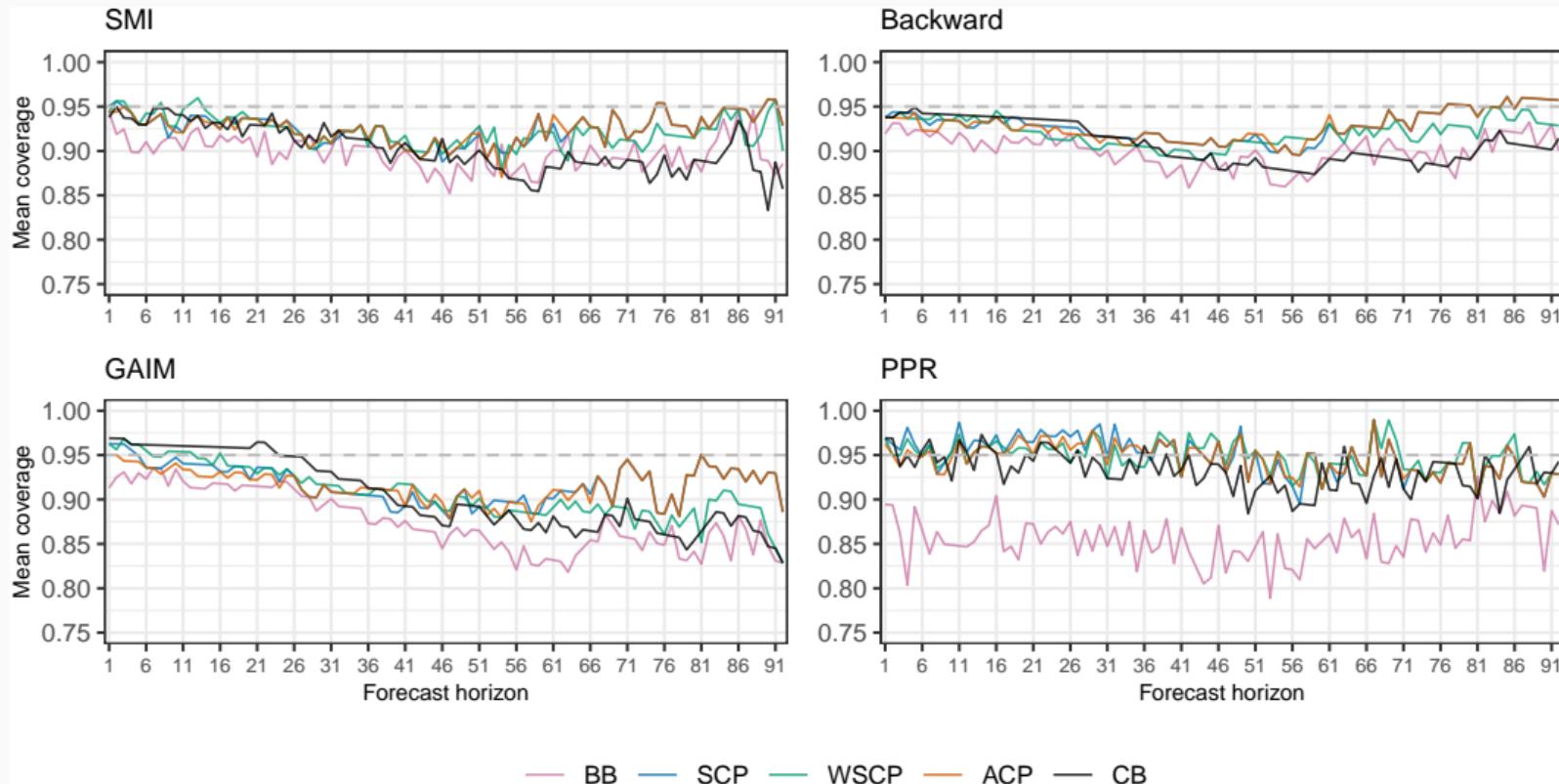
## Data

- **Response:** Daily deaths in Summer
  - 1990 to 2014 – Montreal, Canada
- **Index Variables:**
  - ▶ Death lags
  - ▶ Max temperature lags
  - ▶ Min temperature lags
  - ▶ Vapor pressure lags
- **Nonlinear:** DOS (day of the season), Year

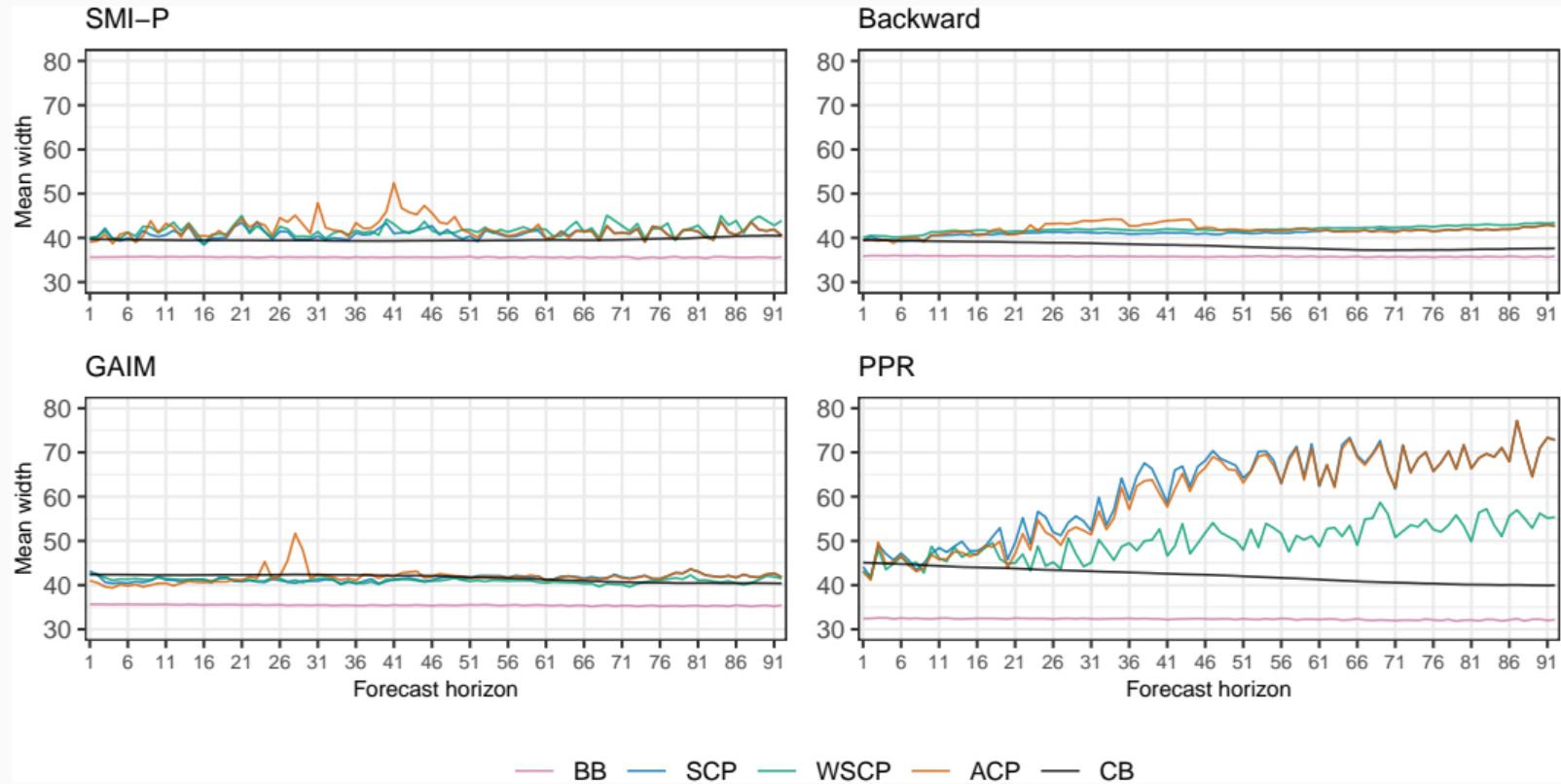
## Data split

- **Training Set:** 1990 to 2007
- **Validation Set:** 2008
- **Test Set:** 2009 to 2014

# Mean coverage



# Mean width



# Conclusion

## i Summary of Results (work-in-progress):

- **Block Bootstrap** – Under-coverage; too narrow
- **Conformal Prediction** – Better achieves a target coverage, with acceptable sharpness



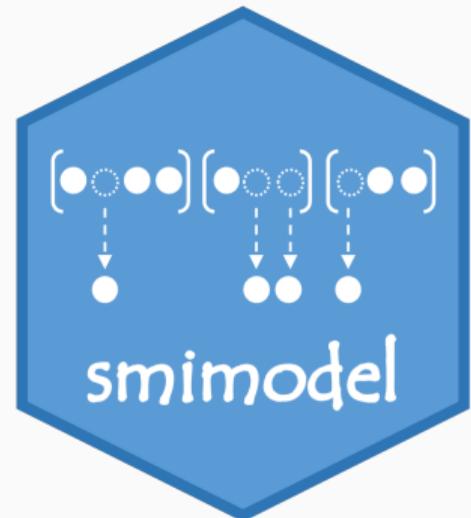
## Limitations:

- Test set is not long enough for larger forecast horizons
- Hyper-parameter choices

# R Package - smimodel

- Block bootstrap
  - ▶ **bb\_cvforecast()**
- Conformal bootstrap
  - ▶ **cb\_cvforecast()**

[github.com/nuwani-palihawadana/smimodel](https://github.com/nuwani-palihawadana/smimodel)



## Find me :

-  [nuwanipalihawadana.netlify.app](https://nuwanipalihawadana.netlify.app)
-  [in/nuwani-palihawadana](https://in/nuwani-palihawadana)
-  [@nuwani-palihawadana](https://github.com/nuwani-palihawadana)
-  [nuwani.kodikarapalihawadana@monash.edu](mailto:nuwani.kodikarapalihawadana@monash.edu)