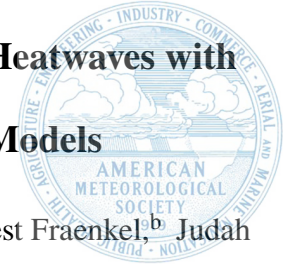


Sub-seasonal Prediction of Central European Summer Heatwaves with Linear and Random Forest Machine Learning Models



Elizabeth Weirich Benet,^a Maria Pyrina,^a Bernat Jiménez-Esteve,^a Ernest Fraenkel,^b Judah
Cohen,^{c,d} and Daniela I.V. Domeisen^{e,a}

^a *Institute for Atmospheric and Climate Science, ETH Zürich, Zürich, Switzerland*

^b *Department of Biological Engineering, MIT, Cambridge, USA*

^c *Atmospheric and Environmental Research (AER), Lexington, USA*

^d *Department of Civil and Environmental Engineering, MIT, Cambridge, USA*

^e *Institute of Earth Surface Dynamics, Université de Lausanne, Lausanne, Switzerland*

Corresponding author: Elizabeth Weirich Benet (weiriche@ethz.ch) and Daniela I.V. Domeisen (daniela.domeisen@env.ethz.ch)

1

Early Online Release: This preliminary version has been accepted for publication in *Artificial Intelligence for the Earth Systems*, may be fully cited, and has been assigned DOI 10.1175/AIES-D-22-0038.1. The final typeset copyedited article will replace the EOR at the above DOI when it is published.

ABSTRACT: Heatwaves are extreme near-surface temperature events that can have substantial impacts on ecosystems and society. Early Warning Systems help to reduce these impacts by helping communities prepare for hazardous climate-related events. However, state-of-the-art prediction systems can often not make accurate forecasts of heatwaves more than two weeks in advance, which are required for advance warnings. We therefore investigate the potential of statistical and machine learning methods to understand and predict central European summer heatwaves on timescales of several weeks. As a first step, we identify the most important regional atmospheric and surface predictors based on previous studies and supported by a correlation analysis: 2-m air temperature, 500-hPa geopotential, precipitation, and soil moisture in central Europe, as well as Mediterranean and North Atlantic sea surface temperatures, and the North Atlantic jet stream. Based on these predictors, we apply machine learning methods to forecast two targets: summer temperature anomalies and the probability of heatwaves for 1–6 weeks lead time at weekly resolution. For each of these two target variables, we use both a linear and a random forest model. The performance of these statistical models decays with lead time, as expected, but outperforms persistence and climatology at all lead times. For lead times longer than two weeks, our machine learning models compete with the ensemble mean of the European Centre for Medium-Range Weather Forecasts’ hindcast system. We thus show that machine learning can help improve sub-seasonal forecasts of summer temperature anomalies and heatwaves.

SIGNIFICANCE STATEMENT: Heatwaves (prolonged extremely warm temperatures) cause thousands of fatalities worldwide each year. These damaging events are becoming even more severe with climate change. This study aims to improve advance predictions of summer heatwaves in central Europe by using statistical and machine learning methods. Machine learning models are shown to compete with conventional physics-based models for forecasting heatwaves more than two weeks in advance. These early warnings can be used to activate effective and timely response plans targeting vulnerable communities and regions, thereby reducing the damage caused by heatwaves.

1. Introduction

A heatwave is an extended period of extremely hot weather relative to the expected local conditions at that time of the year. These high temperatures can cause substantial damage to human health, agriculture, infrastructure, and biodiversity (Barriopedro et al. 2011; Perkins 2015). Heatwaves are among the most dangerous natural hazards (Basu 2002; Lowe et al. 2011), having caused more than 166,000 deaths across the world between 1998 and 2017, including 70,000 fatalities during the 2003 European heatwave (Wallemacq et al. 2018). Summer heatwaves are associated with higher wet-bulb temperatures than winter heatwaves (Buzan and Huber 2020), resulting in higher mortality (Huynen et al. 2001). In addition, the probability of other natural disasters, such as wildfires, is higher during heatwaves (e.g., the 2020 Australian wildfires ignited amid a record-breaking heatwave (Deb et al. 2020)). Furthermore, climate change leads to more extreme hot weather (Barriopedro et al. 2011; Perkins 2015), and an increase in heatwave intensity, duration, and frequency (Ford et al. 2018; Perkins and Alexander 2013; Perkins-Kirkpatrick and Lewis 2020; Seneviratne et al. 2014).

Preparation for heatwaves is possible to a certain extent, for example through early warning systems (EWS) (Merz et al. 2020), which enable an effective and timely response targeting vulnerable populations and regions. For instance, EWS help to determine when crops will need more irrigation, when cooling centers must be set up, or when local hospitals must prepare for an additional number of patients (Bassil and Cole 2010). Moreover, measures for heatwave preparedness on the order of seasons to decades have to be taken by governments and municipalities (Casanueva et al. 2019; Kotharkar and Ghosh 2022). Hence, the time needed to prepare for heatwaves is

often beyond the timescales of medium-range weather forecasts (up to two weeks) (de Perez et al. 2018). While forecasts on seasonal timescales show potential, a skill gap between two weeks and seasonal timescales remains (Robertson et al. 2015; White et al. 2017). Alternative approaches must therefore be explored to extend the lead time of skillful forecasts to sub-seasonal timescales (two weeks to two months).

Central European heatwave predictability can be enhanced by a range of precursors, including both local and remote drivers linked to European temperatures via teleconnections. Heatwaves are generally associated with local persistent blocking anticyclones or upper-level ridges (Kautz et al. 2022; Suarez-Gutierrez et al. 2020). The atmospheric circulation associated with these persistent features exhibits predictability timescales of up to two weeks (Weyn et al. 2019; Zheng and Frederiksen 2007). In turn, the latitude and speed of the North Atlantic (NA) jet stream, which are influenced by the distribution of topography (JiménezEsteve and Domeisen 2022), affect the occurrence and location of these atmospheric features and, hence, central European heatwaves (Bladé et al. 2011; Oliveira et al. 2020). For instance, when the Summer East Atlantic (SEA) pattern (i.e., the second dominant mode of summer low-frequency variability in the Euro-Atlantic region) is in its positive phase, with low pressure west of the British Isles and high pressure to the east, the weather tends to be unusually warm over Europe (Wulff et al. 2017). The SEA pattern shows longer predictability timescales than local geopotential, on the order of 2–3 weeks (Vitart 2014; Zuo et al. 2016).

Cold sea surface temperature (SST) anomalies in the NA are also found to be present prior to the onset of the most extreme European heat waves since 1980 (Duchez et al. 2016) (e.g., anomalously cold NA SSTs were key to the development of the 2015 European heatwave (Mecking et al. 2019)). Moreover, northwestern Mediterranean (NWMED) SSTs are linked to temperatures over the European continent due to their proximity and large heat capacity, acting as a heat buffer for land temperatures (e.g., the 2003 European heatwave was connected to warm Mediterranean SSTs) (Black et al. 2004). Since SST anomalies tend to be highly persistent, in extratropical regions, weekly mean SST anomalies are associated with longer predictability of weeks to months (Hu et al. 2012; Kumar and Zhu 2018).

Furthermore, precipitation is linked to surface air temperature via several mechanisms, including changes in incoming solar radiation and surface sensible heat flux. Precipitation is characterized

by high-frequency variability and, thus, it is not expected to be predictable at lead times longer than a few weeks (Li and Robertson 2015; Wheeler et al. 2016). Precipitation directly influences soil moisture, which is another driver of summer heatwaves (Fischer et al. 2007). Dry soils (low soil moisture) and warming reinforce each other through a positive feedback effect (Kolstad et al. 2017): Moist soils mostly cool through latent heat flux to the atmosphere, while dry soils emit more sensible heat (Laguë et al. 2019) and hence heat up the atmosphere faster than moist soils. This warmer atmosphere, in turn, results in even more dryness, closing the positive feedback loop (Seneviratne et al. 2010). In addition, increased sensible heating can help maintain a blocking anticyclone over land (Miller et al. 2021). Consequently, dry springs are more likely to be followed by extremely high summertime temperatures (Mueller and Seneviratne 2012; Perkins 2015).

We here investigate whether the sub-seasonal forecasting accuracy of summer temperature anomalies and heatwaves in central Europe (CE) can be improved by using linear and random forest (RF) machine learning (ML) models based on these precursors. Other studies use ML and deep learning (DL) to forecast temperature and heatwaves, targeting timescales different from sub-seasonal (Khan et al. 2019; Kämäräinen et al. 2019; Pyrina et al. 2021) or North America instead of CE (Chattopadhyay et al. 2020; Miller et al. 2021; Sobhani et al. 2018; Vijverberg et al. 2020). Moreover, DL architectures successfully predict the onset of long-lasting extreme heatwaves in France two weeks in advance (Jacques-Dumas et al. 2022) and yield increased predictability with respect to the European Centre for Medium-Range Weather Forecasts (ECMWF) at lead times of 3–4 weeks (Lopez-Gomez et al. 2022), agreeing with the findings of the present study despite using a different set of predictors. Finally, additional predictors are identified in a related study by using explainable ML methods (van Straaten et al. 2022).

2. Methods

a. Data

1) PREDICTOR SELECTION

Seven atmospheric and surface predictors that are expected to be related to summer temperature and heatwaves in CE based on previous studies (Section 1) and a correlation analysis (Section 3b1) are selected: 2-m air *temperature*, 500-hPa *geopotential*, *precipitation*, *soil moisture*, the *SEA* index, *NWMED SST*, and *cold North Atlantic anomaly (CNA) SST*. Geopotential at the

500-hPa pressure level is used instead of sea-level pressure to avoid capturing the influence of high surface temperatures on the local low-level surface pressure (Suarez-Gutierrez et al. 2020). The following predictors were also evaluated but were not used, as they correlated only weakly with 2-m air *temperature*: deep soil moisture (28–289 cm underground), the Summer North Atlantic Oscillation (i.e., the first dominant mode of summer low-frequency variability in the Euro-Atlantic region), southeastern Mediterranean SST, Baltic SST, El Niño Southern Oscillation SST, the North Atlantic SST index by Ossó et al. (2020), and the Pacific-Caribbean Dipole index by Wulff et al. (2017). The seven selected predictors are considered in the extended summer season (MJJAS), during the time period between 1 May 1981 and 30 September 2018. Technical details about these predictors can be found in Table 1. Since both local predictors and remote teleconnections are included, their locations are shown in Fig. 1 and their latitude-longitude coordinates are provided in Table 2.

Calculation of the SEA index The changes in speed and location of the NA jet stream are included in our set of predictors through the *SEA* index. First, the *SEA* pattern is calculated via principal component analysis (PCA) (Storch and Zwiers 2003, chap. V), applied on the detrended 500-hPa geopotential height anomalies over the NA box for the summer season (JJA). The *SEA* index corresponds to the time-dependent coefficients (or PCA amplitudes) of the second PCA pattern (Wulff et al. 2017). Then, the daily *SEA* index is calculated for the extended summer season (MJJAS) by projecting the *SEA* pattern on the daily values of the 500-hPa geopotential height anomalies from May to September and the obtained time series are normalised ($\mu = 0$, $\sigma = 1$).

2) DATA PREPROCESSING PIPELINE

First, we select latitude-longitude boxes for each physical magnitude and take either the arithmetic mean of the area or perform PCA (Table 1) to obtain one-dimensional time series. The maximum overlap period for the selected predictors is chosen as 1 May 1981 to 30 September 2018 (38 summers). We then detrend each time series by subtracting the linear trend. Detrending the data removes linear long-term trends, which could be influenced by external climate forcing. Next, we compute the daily climatology (x_{clim}), defined as the mean over the full time period for a particular day of the year. We smooth the daily climatology by a centred 31-day rolling mean window. We then compute the anomalies with respect to climatology as: $x_{\text{anom}} = x - x_{\text{clim}}$. This way, also

Predictor	Physical magnitude (units)	Source (Space, Time Res.)	Level	Box	Method
Temperature	2-m air temperature (°C)	E-OBS (0.25°, daily)	2 m a.g.	CE	avg
Geopotential	geopotential (m ² s ⁻²)	ERA-Interim (2.5°, daily)	500 hPa	CE	avg
Precipitation	rainfall (mm)	E-OBS (0.25°, daily)	surface	CE	avg
Soil moisture	volumetric soil water layer (m ³ m ⁻³)	ERA5-Land (2.5°, daily)	0–28 cm u.g.	CE	avg
SEA index	geopotential (m ² s ⁻²)	ERA-Interim (2.5°, daily)	500 hPa	NA	PCA
NWMED SST	sea surface temperature (°C)	HadISST (1°, monthly)	sea level	NWMED	avg
CNAA SST	sea surface temperature (°C)	HadISST (1°, monthly)	sea level	CNAA	avg

TABLE 1. **Properties of the predictors.** For each predictor, the name of the corresponding variable (physical magnitude) as labeled in the dataset (source) is presented. We also indicate the temporal and spatial resolution at which each variable was downloaded, the extracted vertical level, the selected spatial location, and the method used to convert the three-dimensional time-latitude-longitude space into a one-dimensional time series. The soil moisture (0–28 cm u.g.) is calculated as the average over the first two layers (layer one: 0–7 cm u.g. and layer two: 7–28 cm u.g.). The monthly sea surface temperature (SST) predictors are interpolated to daily time resolution. Notation: Summer East Atlantic (SEA), northwestern Mediterranean (NWMED), cold North Atlantic anomaly (CNAA), above ground (a.g.), and underground (u.g.).

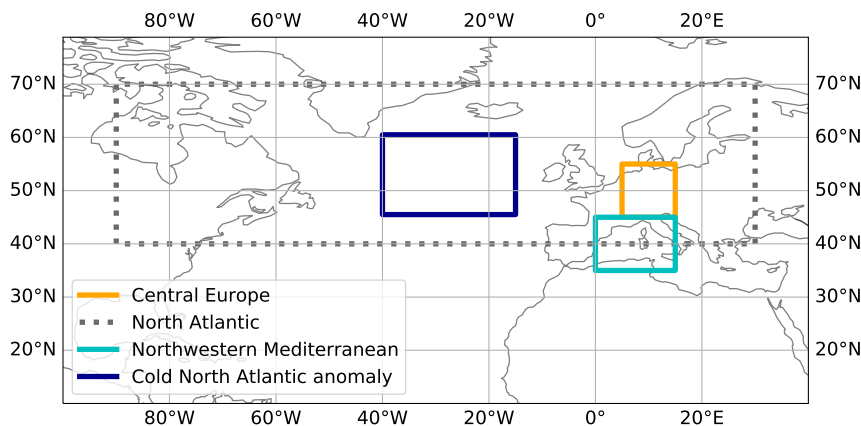


FIG. 1. **Location of latitude-longitude boxes.** Used to define the location of the predictors shown in Table 1. The latitude-longitude coordinates of the boxes are shown in Table 2.

periodic changes due to seasonality are removed. Afterwards, to reduce the noise caused by natural variability, which might lead to overfitted statistical models, these anomalies are smoothed out via a 7-day centred rolling mean. Then, we standardize the predictors: $x_{\text{std anom}} = \frac{x_{\text{anom}}}{x_{\text{std}}}$, where $x_{\text{std anom}}$ are the standardized anomalies and x_{std} the standard deviation of the distribution of each predictor. Furthermore, for each of the six prediction lead times (1–6 weeks), the predictors are provided to

Box	Latitude	Longitude
Central Europe (CE)	45°N–55°N	5°E–15°E
North Atlantic (NA)	40°N–70°N	90°W–30°E
Northwestern Mediterranean (NWMED)	35°N–45°N	0°–15°E
Cold North Atlantic anomaly (CNAА) (Duchez et al. 2016)	45°N–60°N	15°W–40°W

TABLE 2. **Coordinates of latitude-longitude boxes.** The boxes correspond to the location of the predictors of Table 1 as seen in Fig. 1.

the ML models for the four weeks before initialization. For example, for a forecast at two weeks lead time (meaning that we are using a statistical model initialized two weeks before the target week for which we make the forecast), the *precipitation* from two, three, four, and five weeks before the target week is used as a predictor by the ML models. Finally, since we want to investigate the predictability of summer temperature, the extended summer months (MJJAS) are selected.

3) HEATWAVE INDEX DEFINITIONS

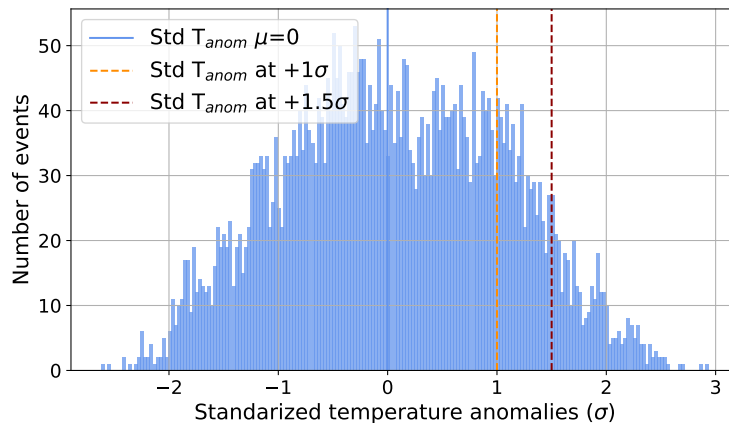


FIG. 2. **Histogram of temperature anomalies averaged over CE for the definition of heatwave indices.** The blue bars correspond to the standardized ($\mu = 0$, $\sigma = 1$) temperature anomalies. The data is smoothed by a 7-day running mean (Section 2a2). The vertical blue line marks the mean ($\mu = 0$) of the distribution. The stippled orange (red) line marks $+1$ ($+1.5$) standard deviations (σ) from the mean and is used to define heatwaves.

We define weekly heatwaves via a binary index: one for a heatwave week and zero, otherwise. While there is no universal definition for heatwaves and a range of different indices are found across the literature, percentile-based definitions are widely used (Perkins and Alexander 2013; Perkins 2015; Perkins-Kirkpatrick and Lewis 2020; Spensberger et al. 2020). We use two different

heatwave definitions, thereby defining two independent classification problems: $+1\sigma$ for high and $+1.5\sigma$ for extremely high temperature anomalies (Fig. 2). The $+1\sigma$ weekly heatwave index is defined as one for the weekly mean temperature anomalies above one standard deviation (σ) (i.e., to the right of the orange line in Fig. 2) and zero, otherwise. Analogously, the $+1.5\sigma$ weekly heatwave index is defined as one for the weekly mean temperature anomalies above 1.5 standard deviations (i.e., to the right of the red line in Fig. 2) and zero, otherwise. The number of heatwave and no-heatwave samples can be found in Table 3.

Weekly heatwave index	$+1\sigma$	$+1.5\sigma$
Absolute number of heatwave events	1,121	430
Absolute number of no-heatwave events	4,813	5,504
Percentage of heatwaves	18.89%	7.25%

TABLE 3. **Class imbalance.** Class distribution of the 5,934 samples in the extended summer (MJJAS) and the 1981–2018 time period.

b. Lead time

We forecast at 1–6 weeks lead time. The statistical models are trained separately for each lead time and do not learn from each other. For instance, the two-weeks-lead-time forecast does not receive the one-week-lead-time forecast as an additional input. Moreover, since our data is averaged via a seven-day rolling mean (Section 2a2), weeks are labeled by their central day. A one-week-lead-time prediction leaves no gap between the days used to calculate the one-week-lag predictors and the days used to determine the target. For instance, the one-week-lead-time forecast run on June 4th (average over June 1st–June 7th) forecasts June 11th (average over June 8th–June 14th). Similarly, a lead time of two weeks leaves a gap of seven unused days.

c. Machine learning models

For our study, we choose statistical models at the two extremes of the bias-variance tradeoff (Mehta et al. 2019). (1) The simpler linear models are prone to have high bias, meaning that the model will match the training set less closely. These models have a higher potential for underfitting. Linear models, however, have low variance, meaning that the predictions of the model do

not fluctuate much with a change of dataset. Overall, these models are focused on the larger trends rather than on the complicated patterns of the training set. (2) By contrast, the more complex decision trees (DTs) are likely to overfit the data, but also to capture most of the relevant patterns. They tend to have high variance, but low bias. To mitigate the risk of DTs overfitting, we use RFs instead.

Two statistical models from each of these two families (1 and 2) are used for the regression and classification forecasts: ridge regressor (RR), ridge classifier (RC), random forest regressor (RFR), and random forest classifier (RFC). Moreover, the final forecasts by each model are the average of an ensemble of these ML models trained on slightly different samples (Section 2h).

1) LINEAR MODELS

Linear regression models forecast the target time series $\mathbf{y} = (y_t)$ as a linear combination of N predictor time series $\mathbf{x}_n = (x_{n,t})$:

$$\hat{\mathbf{y}}(\boldsymbol{\omega}, \mathbf{X}) = \omega_0 + \omega_1 \mathbf{x}_1 + \dots + \omega_N \mathbf{x}_N \quad (1)$$

where ω_0 is the intercept, ω_n ($0 < n \leq N$) are the regression coefficients, and $t \in [1, T]$ is the time step. The coefficients are chosen to minimize the residual sum of squares between the forecast ($\hat{\mathbf{y}}$) and the observed target (\mathbf{y}): $\min_{\boldsymbol{\omega}} \|\hat{\mathbf{y}} - \mathbf{y}\|$. Linear classification models first convert binary targets to $\{-1, 1\}$ and then treat the problem as a regression task. The forecast class corresponds to the sign of the regressors forecast. We use Ridge regularization to control excessively fluctuating functions by adding an additional penalty term in the error function, such that the coefficients do not take extreme values (Hastie et al. 2009, chap. 3). Ridge shrinks the predictor coefficients based on the L2-norm ($\|\boldsymbol{\omega}\|_2 = \sqrt{\sum_{n=1}^N \omega_n^2}$). The loss function for minimization then becomes $\|\hat{\mathbf{y}} - \mathbf{y}\| + \alpha \|\boldsymbol{\omega}\|_2^2$, where the complexity parameter α is a hyper-parameter which controls the amount of shrinkage.

2) RANDOM FORESTS

A DT makes a recursive partition of the input space into rectangles, by selecting the predictor and the respective cutting point that discriminate best at each node. The resulting leaves correspond to a specific forecast value (regression) or to a probability of belonging to the positive class (binary classification). However, DTs have two key disadvantages: (1) Trees usually have high variance

due to their greedy split process, which implies that a small change in training data can result in significantly different splits. (2) Since the tree estimate is not smooth, DTs may not be appropriate when the underlying function is smooth (Khan et al. 2019). A more accurate and robust statistical model can be constructed by creating a random ensemble of DTs whose averaged prediction is more accurate than that of any individual tree. RFs use two sources of randomness while training: bagging and feature randomness (Breiman 2001). (1) Bagging (or bootstrap aggregation) consists in selecting a random subset of the training set with replacement –meaning that individual data points can be chosen more than once– to train each individual tree. (2) When splitting a node in a classical DT, all features are considered and the one that provides the greatest separation between observations is selected. In contrast, each individual tree in a RF can pick only from a random subset of features (Hastie et al. 2009, chap. 15). Finally, the mean or majority-vote forecast of all the regression or classification trees in the forest is selected as the final result, respectively. RFs are chosen over other tree-based algorithms since they are more interpretable (Rudin 2019) than gradient boosting and less prone to overfit than single DTs.

d. Hyper-parameter optimization

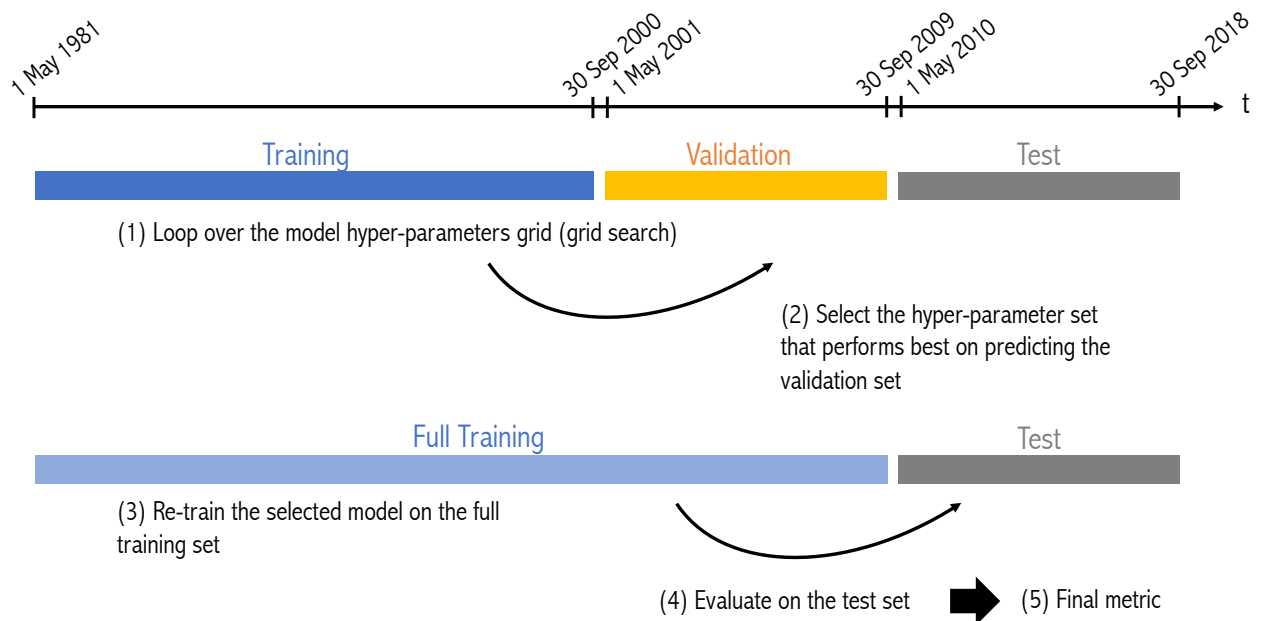


FIG. 3. Schematic of the training-validation-test split

We split the available data into a training period (1 May 1981 – 30 September 2000), a validation period (1 May 2001 – 30 September 2009), and a testing period (1 May 2010 – 30 September 2018) (Fig. 3). The validation period is used to optimize the statistical model’s hyper-parameters for each lead time. After the hyper-parameter optimization, the model is re-trained on the full training period (1 May 1981 – 30 September 2009), which is the combination of the validation and the training period. A nested cross-validation (CV) scheme is also implemented (Appendix, Fig. B1).

For the RFs, we use an exhaustive grid-search hyper-parameter optimization including all possible combinations (750) of the following parameters: number of trees in the forest $\in \{50, 100, 200, 400, 600\}$, maximum tree depth $\in 5\text{--}14$, and a range of 15 values centered around the full training set’s length T_{ft} divided by 100 in steps of $T_{\text{ft}}/500$ for the minimum number of samples per leaf. The minimum number of samples for splitting a node is set to the minimum number of samples per leaf multiplied by a factor of two and, for classification, the class weight is set to *balanced*. For the two linear models, the complexity parameter α is selected from the range $[0, 1]$ in steps of 0.05. The reference metrics for optimization are the root mean-square error (RMSE) for regression and the Brier score (BS) for classification (Section 2e). The selected hyper-parameters are shown in the Appendix (Table C1).

e. Metrics for the evaluation of forecasting performance

1) REGRESSION METRICS

For regression, two different metrics are considered: RMSE and Pearson correlation. The RMSE evaluates how far away the forecast ($\hat{\mathbf{y}}$) and the ground truth (\mathbf{y}) time series are from each other and is defined as:

$$\text{RMSE}(\hat{\mathbf{y}}, \mathbf{y}) = \sqrt{\text{MSE}(\hat{\mathbf{y}}, \mathbf{y})} = \sqrt{\frac{1}{T} \sum_{t=1}^T (\hat{y}_t - y_t)^2} \quad (2)$$

for T the number of time steps (sample size).

The Pearson correlation measures to what extent the curve follows the changes and is given by:

$$\text{Corr}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{\sum_{t=1}^T (\hat{y}_t - \bar{\hat{\mathbf{y}}})(y_t - \bar{\mathbf{y}})}{\sqrt{\sum_{t=1}^T (\hat{y}_t - \bar{\hat{\mathbf{y}}})^2} \sqrt{\sum_{t=1}^T (y_t - \bar{\mathbf{y}})^2}} \quad (3)$$

for $\bar{z} = \frac{1}{T} \sum_{t=1}^T z_t$ the mean over all time steps.

2) CLASSIFICATION METRICS

For classification, the BS and the Receiver Operating Characteristic (ROC) Area Under Curve (AUC) are used to evaluate the probabilistic forecast. The BS is the mean squared error of the probability forecasts (i.e., Eq. 2 squared), considering that an observation is $y_t = 1$ if the event occurs and $y_t = 0$ if the event does not occur at time t . Since individual probabilistic forecasts and observations are bounded by zero and one, the BS can only take values in the range $[0,1]$ (Wilks 2019, chap. 9).

The ROC is the true positive rate (TPR) as a function of the false positive rate (FPR) (Bradley 1997). The TPR (or Recall) is defined as the proportion of positive data points that are correctly considered positive, with respect to all positive data points. The TPR is given by $TP / (FN+TP)$ for true positives (TPs) and false negatives (FNs). The FPR (or False Alarm) is defined as the proportion of negative data points that are mistakenly considered positive, with respect to all negative data points. The FPR is calculated as $FP / (FP+TN)$ for false positives (FPs) and true negatives (TNs) (see Table 4 for the definition of TP, FP, FN, and TN).

		Actual value (y)	
		Positive (1)	Negative (0)
Forecast value (\hat{y})	Positive (1)	TP	FP
	Negative (0)	FN	TN

TABLE 4. **Confusion matrix.** The positive class corresponds to a heatwave and the negative class to no heatwave. For a sensible model, the principal diagonal values must be high and the off-diagonal values must be low (Bradley 1997).

Moreover, the performance of the binary classification is assessed via the FPR-to-TPR ratio, extremal dependence index (EDI), and frequency bias (B). The EDI is used to evaluate forecasts of rare binary events and is calculated as (Ferro and Stephenson 2011):

$$EDI = \frac{\ln(FPR) - \ln(TPR)}{\ln(FPR) + \ln(TPR)} \quad (4)$$

This score is ill-defined if any of the four cells in the confusion matrix (Table 4) equals zero, since $\ln(0)$ or a division by zero yield infinity. However, such models can still be interpreted by adding an infinitely small number (pseudo-count) to those cells containing zeros (Wunderlich et al. 2019).

The frequency bias is the ratio of the number of positive-class forecasts to the number of positive-class observations:

$$B = \frac{TP + FP}{TP + FN} \quad (5)$$

Unbiased forecasts exhibit $B = 1$, indicating that the event is forecast the same number of times as observed (Wilks 2019, chap. 9).

We define a *useful* probabilistic forecast as having $BS < 0.25$ (Steyerberg et al. 2010) and $ROC\ AUC > 0.5$ (Bradley 1997). We consider a binary forecast *useful* if $FPR/TPR < 1$ and $EDI > 0$ (Wilks 2019, chap. 9). In addition, B should be as close to one as possible.

f. Calibration of the classification forecasts

Good forecasts should not only be accurate (as measured by ROC AUC, EDI and the FPR-to-TPR ratio) but also well-calibrated (as measured by BS and B) (Jolliffe and Stephenson 2005), meaning that the sub-sample relative frequency should be exactly equal to the forecast probability in each sub-sample (Wilks 2019, chap. 9). For example, if a model forecasts 100 positive-class events (e.g., heatwave weeks), each with a probability of 80%, we expect 80 of the events to be correctly classified (i.e., to actually be a heatwave).

1) PLATT SCALING FOR THE PROBABILISTIC FORECASTS

Unlike accuracy, reliability can be improved in a post-processing step by calibrating the probabilistic forecasts (Jolliffe and Stephenson 2005). The linear ML models already predict calibrated probabilities and do not need an additional calibration step. We use Platt scaling to re-calibrate the probabilistic forecasts by the RFs. Platt scaling consists in projecting the (ill-calibrated) probabilities predicted by the ML models onto the right probability distribution using a logistic regression model (Smola et al. 2000, chap. 5). The RFs are trained on the training set and calibrated on the validation set to determine the parameters of the logistic regression. The calibrated RF models are

then used to predict the test set. These datasets correspond to the ones defined in Fig. 3. Since the logistic function is monotonic, the calibration via Platt scaling does not change the ordering of the samples, and, consequently, the ROC AUC score remains the same. Instead, the BS is considerably reduced after the calibration step.

2) PROBABILITY THRESHOLD MOVING FOR THE BINARY FORECASTS

Forecasting the two weekly summer heatwave indices defined in Section 2a3 ($+1\sigma$ and $+1.5\sigma$) results in imbalanced classification problems (Table 3). A binary classifier trained on these imbalanced data will learn to always forecast the negative class, leading to a trivial and ill-calibrated statistical model. Balancing the data before the training or moving the probability threshold are two potential solutions to this problem. Random undersampling and oversampling methods have been explored to balance the training data (Lemaitre et al. 2017). However, these methods are not used for the final version of the statistical models since, in this particular case, they result in over-forecasting heatwaves.

Instead, for this study, the data imbalance is accounted for by adjusting the probability threshold: The (non-calibrated) classification models output a probability for each validation sample to belong to the positive class. Then, the probability threshold between zero and one that corresponds to no frequency bias (i.e., $B = 1$) on the validation set is selected to binarize the output (Wilks 2019, chap. 9). To avoid a strong dependency on the distribution of the validation set, an internal cross-validation scheme is used for selecting the probability threshold. Thirty validation sets of nine randomly selected non-consecutive years belonging to the full training set (1981–2009) are constructed. The remaining 20 years are used for training. The threshold that minimises the deviation from the mean frequency bias of the 30 validation sets from one is selected.

g. Reference forecasts

We compare our statistical models to the climatology, persistence, and ECMWF hindcast forecasts:

(i) *Climatology* For regression, temperature anomalies with respect to climatology are forecast. Thus, the climatology forecast is zero for all times per definition. For classification, the climatology forecast is the mode class for each day of the year. Since, in our dataset, the negative class

predominates strongly over the positive class, the climatology forecast is found to always predict the negative class (no heatwave).

(ii) *Persistence* Persistence forecasts predict that the future weather condition will be the same as the present condition. In practice, the persistence forecast is defined as keeping the value from initialization time until verification time. For instance, for the regression forecast at two weeks lead time, the persistence is the temperature anomaly two weeks before verification time.

(iii) *ECMWF* Early warnings are issued by the operational ECMWF sub-seasonal prediction system, using 51 ensemble members and information beyond the ensemble mean. However, these forecasts are currently only available for the years 2015–2022. Therefore, in order to evaluate our ML models' skill for the full test period (2010–2018), we compare to ECMWF sub-seasonal hindcast system's ensemble mean instead. This hindcast system is initialized twice a week and provides 20-year hindcasts with 11 ensemble members integrated over 46 days. The hindcasts used here cover the period 2000–2019 and use the model version of the Integrated Forecasting System cycle 47r1 (Haiden et al. 2019).

The mean daily 2m-air temperature is downloaded at a spatial resolution of $1^\circ \times 1^\circ$ and the arithmetic mean of the area over CE (as defined in Fig. 1) is calculated. Then, the temperature anomalies are calculated by removing the lead-time-dependent climatology at each initialization, calculated by the 20-year mean of the 11-member ensemble started on the same day and month for each year of the reference period (2000–2019). For instance, if a hindcast was initialized on May 31st, the lead time dependent climatology corresponding to that hindcast is calculated as the mean of the 11-member ensemble initialized on May 31st and averaged over the 20-year reference period (2000–2019) separately for each of the 46 days. After the calculation of the temperature anomalies, a 7-day rolling mean is applied for each initialization. In this way, we end up with 40 days per initialization, with each day being the centre of the 7-day rolling mean. For instance, the first day predicted by the initialization on May 31st will be June 4th (average over June 1st–June 7th).

Removing different climatologies for individual dynamical models and reanalysis or observational datasets is standard practice, as the climatological normals are slightly different across datasets (IPCC 2013, chap. 9). Moreover, in the case of sub-seasonal forecasting, calculating anomalies with respect to a lead-time dependent climatology is expected to remove systematic biases which are

lead-time dependent (Manzanas 2020; Molteni et al. 2011). However, the methodology followed for the calculation of the dynamical model’s climatology can influence the forecast’s skill (Manrique-Suñén et al. 2020).

h. Ensembles and uncertainty estimation

For both ECMWF and the ML models, the final forecast is calculated as the mean forecast by an ensemble of K models:

$$\mu(\hat{\mathbf{Y}}) = \frac{1}{K} \sum_{k=1}^K \hat{\mathbf{y}}_k \quad (6)$$

with $\hat{\mathbf{y}}_k$ the time series prediction by each ensemble member. Then, the M metrics ψ_m defined in Section 2e for the final forecast are calculated as $\psi_m(\mu(\hat{\mathbf{Y}}), \mathbf{y})$, for $m = 1, \dots, M$. To quantify the uncertainty of these metrics, the M metrics are calculated with respect to the ground truth (\mathbf{y}) for each ensemble member ($\psi_{m,k} = \psi_m(\hat{\mathbf{y}}_k, \mathbf{y})$). Then, for each metric m , the unbiased standard deviation of the ensemble ($\sigma_m(\hat{\mathbf{Y}})$) is used to represent the uncertainty of the final forecast’s metrics:

$$\sigma_m(\hat{\mathbf{Y}}) = \sqrt{\frac{1}{K-1} \sum_{k=1}^K (\psi_{m,k} - \mu(\psi_m))^2} \quad (7)$$

for $\mu(\psi_m) = \frac{1}{K} \sum_{k=1}^K \psi_{m,k}$ the mean metric m of all models in the ensemble.

For ECMWF, the considered ensemble consists of $K = 11$ sub-seasonal hindcasts. For both the linear and RF models, block bootstrapping is used to create an ensemble. Bootstrapping consists of randomly drawing samples with replacement from the full training dataset (as defined in Section 2d), with each sample having the same size as the original training dataset. Bootstrap resampling generally results in $\approx 37\%$ of the observations not being selected. This resampling procedure is repeated $K = 600$ times, producing K bootstrap training datasets used to train K ML models (Hastie et al. 2009, chap. 7). However, standard bootstrapping fails to represent the statistics of dependent data, like time series. Block bootstrapping overcomes this limitation by resampling independent chunks of continuous observations instead of single dependent ones (Kunsch 1989). Therefore, under the assumption of inter-annual independency of summers, we apply block bootstrapping

with a block size of one year, which means that the smallest unit considered for resampling is one year instead of one day.

3. Results and discussion

a. Forecasts

1) REGRESSION FORECASTS

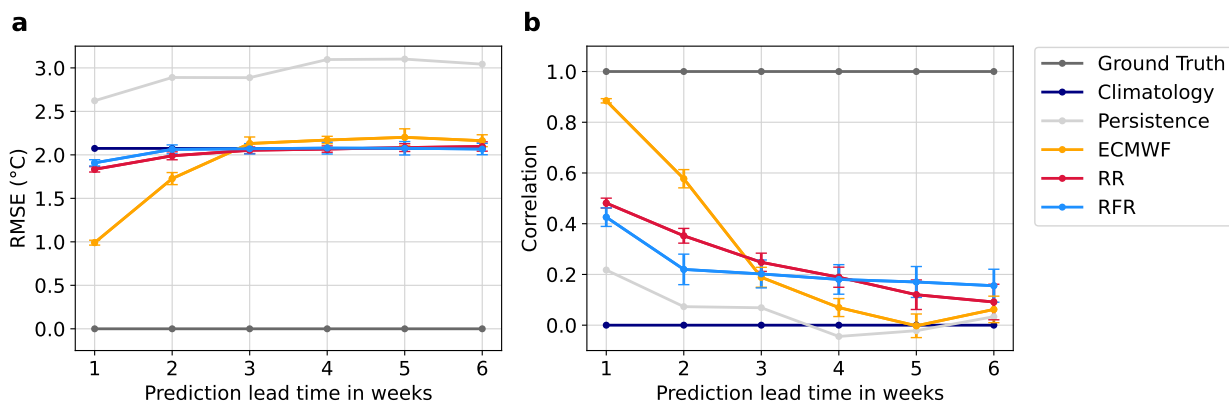


FIG. 4. **Performance of the regression models for six different lead times.** (a) RMSE and (b) correlation for the regression forecasts. An accurate forecast is characterized by a low RMSE and a high correlation. The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble.

In Figure 4, the regression forecasts by two different ML models (RR and RFR) at six different lead times (1–6 weeks) are compared to three reference forecasts: climatology, persistence, and ECMWF. The analogous results for nested CV are shown in the Appendix (Fig. B2).

As can be observed in Fig. 4, all metrics are best for a lead time of one week. The uncertainty in the forecasts by most models, which is represented by the error bars, increases with lead time. The RR’s performance decays linearly with increasing lead time, with a correlation that ranges from 0.48 for one week lead time to 0.09 for six weeks lead time. The RF’s correlation decreases overall from one to six weeks lead time (from 0.43 to 0.16) but remains noticeably constant for lead times longer than two weeks. The evolution of the RMSE is similar, but with the difference that it saturates when reaching the RMSE value that corresponds to the climatology forecast. The RMSE for the best statistical model at each lead time ranges between 1.83 for one week lead time and 2.07 at six weeks lead time.

The linear ML model outperforms the RF in terms of correlation at short lead times (up to three weeks), but the RF model provides a better forecast at long lead times (5–6 weeks). Both ML models outperform the persistence forecast at all lead times. However, the climatology forecast has a relatively low RMSE, being a comparatively good guess at long lead times, when forecasting becomes difficult. For lead times longer than two weeks, the RMSEs of the ML models saturate at the climatology’s RMSE and the ensemble mean of ECMWF’s hindcast has a worse RMSE than the climatology forecast. Still, the climatology forecast does not correlate with the ground truth and the ML and ECMWF models outperform climatology at all lead times in terms of correlation, since these models always correlate positively with the ground truth. While ECMWF provides highly skilled forecasts in terms of correlation and RMSE for one and two weeks lead time, the skill decreases fast with increasing lead time; for lead times of three weeks and longer, the ML models forecast the temperature anomalies more accurately than the ensemble mean of ECMWF’s hindcast.

The ML models generally pick up the sign of the anomalies but their sharpness, which refers to the ability of a probabilistic forecast to spread away from the climatological average (Gneiting et al. 2007), is lower than the one from ECMWF and extreme values are not well-captured (Appendix, Fig. A1). For longer lead times, all models exhibit low sharpness in their forecasts, tending to the climatology forecast. In the case of the ML models, this tendency towards climatology can be a consequence of the loss function. The loss functions for the RR and the RFR models are the linear least squares function and the mean squared error, respectively. Both metrics measure the distance between the forecast and the target curves. Since forecasting anomalies accurately becomes more difficult with increasing lead time, a statistical model that is trained to minimise the error will tend to forecast the mean of the distribution of possible outcomes, becoming smoother and losing sharpness compared to the observations (Rasp and Thuerey 2021). ML models trained to optimize alternative loss functions, like in the study by Lopez-Gomez et al. (2022), would be worth exploring.

2) CLASSIFICATION FORECASTS

The classification models output a probability for each sample in the test set to belong to the positive class (i.e., for a week to be classified as a heatwave week). These probabilities are calibrated

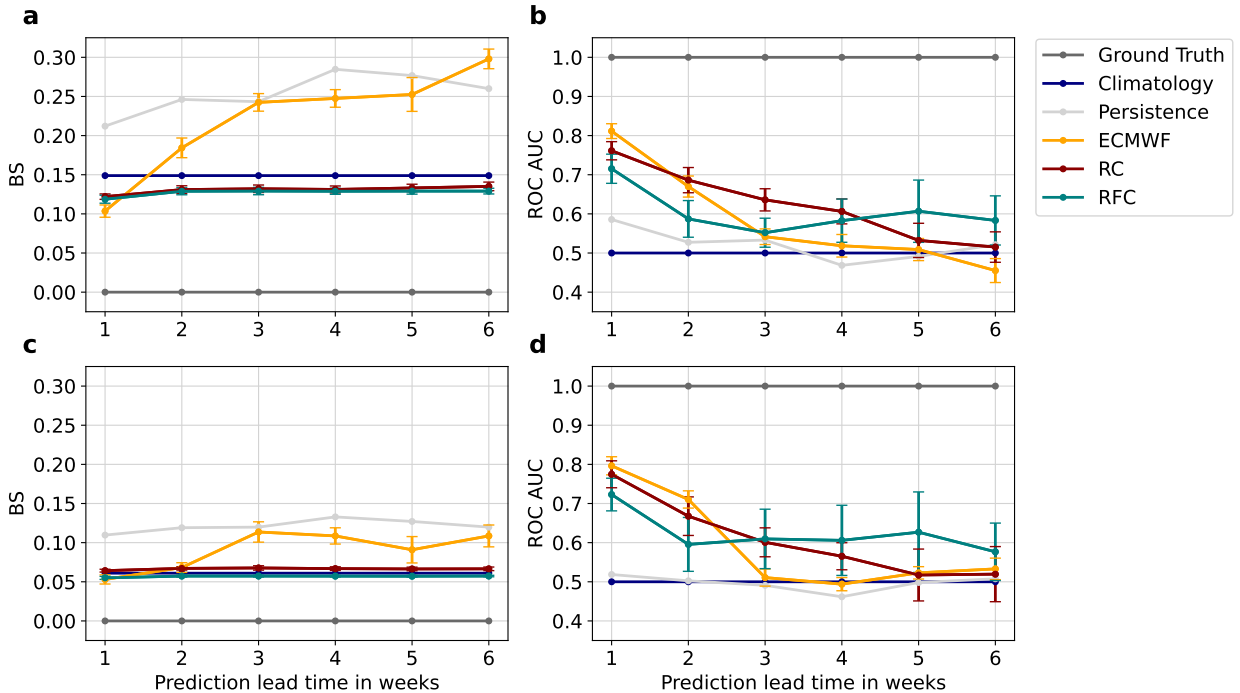


FIG. 5. **Performance of the probabilistic classification models for six different lead times.** BS and ROC AUC for the $+1\sigma$ (a&b) and $+1.5\sigma$ (c&d) weekly heatwave indices. An accurate probabilistic classification forecast is characterized by a low BS and a high ROC AUC. A no-skill probabilistic classification forecast is represented by a BS of 1 and a ROC AUC of 0.5 (as indicated by the climatology). The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble.

to obtain the probabilistic forecast for the RFC model and kept unchanged for the RC model. For both classifiers, the non-calibrated probabilities are binarized via a probability threshold, meaning that a zero (no heatwave) or a one (heatwave) is assigned to each sample in the test set (Section 2f). In Figure 5, the probabilistic classification forecasts by two ML models (RC and RFC) at six different lead times (1–6 weeks) are compared to the three reference forecasts. In Figure 6, the performance of the binary classification is shown. The analogous results for nested CV are shown in the Appendix (Figs. B3 and B4). Two different heatwave indices are used: $+1\sigma$ for high and $+1.5\sigma$ for extremely high temperature anomalies (Section 2a3).

For the probabilistic forecasts, the linear models have a higher ROC AUC than the RFCs for short lead times (up to four weeks for the $+1\sigma$ heatwave index and up to two weeks for the $+1.5\sigma$ heatwave index). However, the RFCs' ROC AUC remains more constant than the linear models'

ROC AUC across lead times, outperforming the linear models for longer lead times (Figs. 5b&d). Moreover, the probabilistic forecasts by both classification ML models outperform persistence and climatology at all lead times and the ensemble mean of ECMWF's hindcast for lead times longer than two weeks, except for the $+1.5\sigma$ forecast at lead times of 5–6 weeks by the RC model. Overall, the forecast uncertainties by all models increase with lead time, resulting in overlapping error bars. These patterns are analogous to the ones observed for the regression forecast (Fig. 4b). In terms of BS, both statistical models present a smaller loss than the ensemble mean of ECMWF's hindcast at lead times of two weeks and higher (Figs. 5a&c). As for regression, the climatology shows a constant Brier loss, which is comparable to the BS of the ML models. The probabilistic forecasts by both statistical models (taking the uncertainty into account) are *useful* at each of the considered lead times (1–6 weeks), except for the RC model at 5–6 weeks lead time, where the uncertainty bars overlap with the no-skill ROC AUC score. Meant by *useful* is $BS < 0.25$ and $ROC\ AUC > 0.5$. It is remarkable that non-null skill by the RFC model is present at these long lead times.

Moreover, in terms of Brier loss, extremely high temperature anomalies ($+1.5\sigma$) are easier to forecast than high temperature anomalies ($+1\sigma$), which agrees with the findings of Wulff and Domeisen (2019). The performance of the ensemble mean of ECMWF's hindcast in predicting extremely high temperature anomalies ($+1.5\sigma$) drops drastically between two and three weeks lead time and remains constant for lead times longer than three weeks. In contrast, ECMWF's classification skill when forecasting high temperature anomalies ($+1\sigma$) decays close to linearly with lead time. The probabilistic RFC is slightly more skilled in capturing extremes than the probabilistic linear model: the RFC forecasts extremely high temperature anomalies ($+1.5\sigma$) more accurately than high temperature anomalies ($+1\sigma$) compared to the linear model. This difference in skill is possibly due to non-linear effects driving extreme temperature which the RFC is able to capture but the linear model is not.

For the binary classification, the overall skill of the statistical models is poorer than for the probabilistic classification. As the lead time increases, the two statistical models and the ensemble mean of ECMWF's hindcast predict fewer weekly heatwave events and the TPR decreases with lead time (Figs. 6b&d). Moreover, despite moving the probability threshold to forecast an unbiased validation set (Section 2f2), the binary forecasts of the test set by the statistical models (in particular, for the $+1.5\sigma$ heatwave index) are considerably biased compared to the predictions by the ensemble

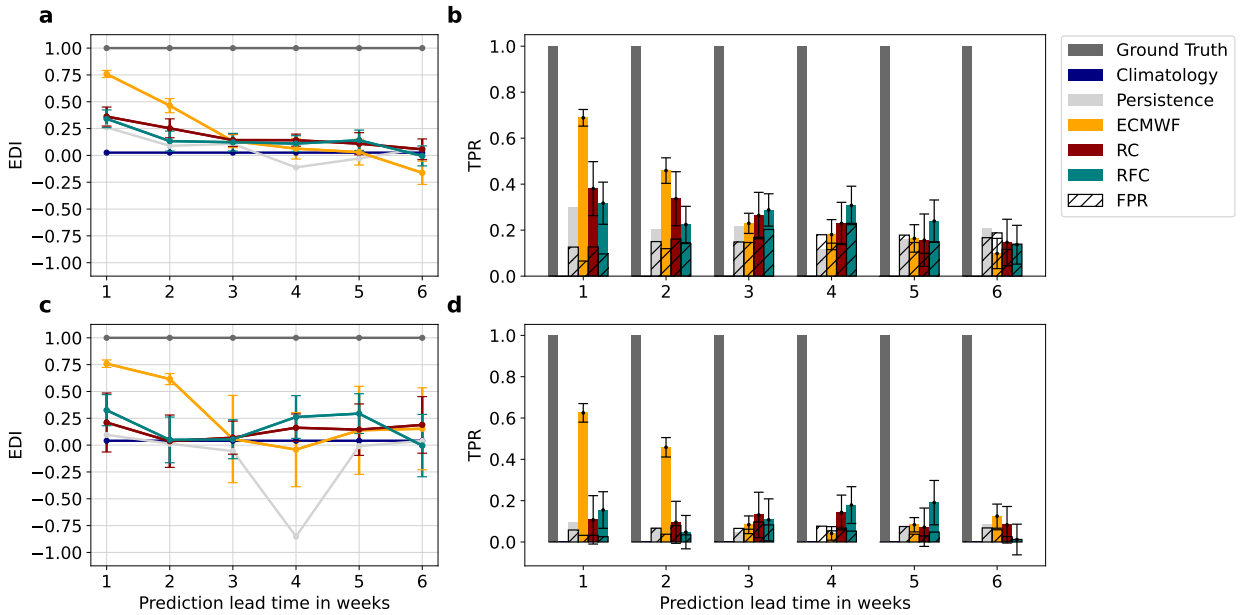


FIG. 6. **Performance of the binary classification models for six different lead times.** (a) EDI and (b) TPR (coloured bars) and FPR (stippled bars) for the $+1\sigma$ weekly heatwave index. (c) and (d) are the corresponding forecasts for the $+1.5\sigma$ weekly heatwave index. An accurate binary classification forecast is characterized by a high EDI, a high TPR, and a low FPR. The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble. Since the climatology forecast predicts only zeros (no heatwave), both its TPR and FPR are equal to zero at all lead times (Figs. b&d). Moreover, at a lead time of four weeks, there is no overlapping between the $+1.5\sigma$ heatwave events in the ground truth and persistence forecast, resulting in zero hits ($TP = 0$). Therefore, the EDI is not defined for the persistence forecast at this lead time and the pseudo-count correction yields a considerably lower value for the EDI compared to the persistence forecast at the other lead times (Fig. c). This is an artifact of the limited sample size and does not appear in nested CV (Appendix, Fig. B4c).

mean of ECMWF’s hindcast (Table 5). *Useful* binary forecasts by at least one of the statistical models (taking the uncertainty into account) are found at 1–5 weeks lead time for the $+1\sigma$ heatwave index and at lead times of one, four, and five weeks for the $+1.5\sigma$ heatwave index, where *useful* is defined as $FPR/TPR < 1$ and $EDI > 0$.

Finally, the RFC tends to overfit the training set considerably, with ROC AUCs and EDIs above 0.99 at all considered lead times (1–6 weeks). The hyper-parameters chosen during the grid search

Heatwave index	Model	1 week	2 weeks	3 weeks	4 weeks	5 weeks	6 weeks
$+1\sigma$	RC	1.11 ± 0.37	1.26 ± 0.47	1.23 ± 0.49	1.03 ± 0.46	0.72 ± 0.58	0.81 ± 0.57
	RFC	0.87 ± 0.29	1.03 ± 0.31	1.45 ± 0.36	1.62 ± 0.44	1.09 ± 0.46	0.93 ± 0.43
	ECMWF	1.05 ± 0.04	1.11 ± 0.10	1.03 ± 0.11	0.97 ± 0.14	0.97 ± 0.18	1.13 ± 0.12
$+1.5\sigma$	RC	0.61 ± 0.71	1.32 ± 0.95	1.62 ± 1.23	1.18 ± 1.07	0.52 ± 1.13	0.49 ± 0.92
	RFC	0.55 ± 0.42	0.58 ± 0.58	1.38 ± 0.81	0.99 ± 0.59	0.93 ± 0.75	0.20 ± 0.63
	ECMWF	1.12 ± 0.08	1.04 ± 0.14	1.04 ± 0.22	0.88 ± 0.18	0.67 ± 0.31	1.04 ± 0.27

TABLE 5. **Frequency bias** of the ensemble mean forecasts of each of the two classification targets in the test period (2010–2018) by the two ML models (RC and RFC) and ECMWF’s hindcast. A well-calibrated model should have $B = 1$. For $B < 1$, the forecast underestimates the total number of heatwave events and for $B > 1$, the events are overestimated. Biases of the ensemble mean forecasts above 1.5 or below 0.5 are bold.

for the RFC correspond to the deepest possible trees and the smallest possible leaves (Appendix, Table C1).

b. Predictor importance

The relevance of each of the seven predictors for forecasting summer temperature anomalies is investigated by performing a linear correlation analysis and examining which predictors were predominantly used by each ML model.

1) LINEAR CORRELATION ANALYSIS

In Figure 7, the linear correlations between the *temperature* and the predictors in the extended summer season (MJJAS) are shown for six different time lags (1–6 weeks). At short time lags, the *temperature* shows a strong autocorrelation. The *geopotential* has an even stronger positive correlation to the *temperature*, indicating that during anticyclonic conditions higher temperatures than normal are expected. In contrast, *precipitation*, *soil moisture*, and the *SEA* index correlate negatively with *temperature* at short time lags. *Precipitation* is associated with cyclones, cloudy conditions, and lower surface air temperatures. Moreover, dryness (low *soil moisture*) and high *temperature* reinforce each other (Section 1). The correlations with the atmospheric predictors (*temperature*, *geopotential*, *precipitation*, and *SEA*) decay fast. In addition, the linear correlation with *soil moisture* becomes non-significant for lead times of two weeks and longer. In contrast, the SST predictors show a more constant linear correlation over time and dominate on timescales

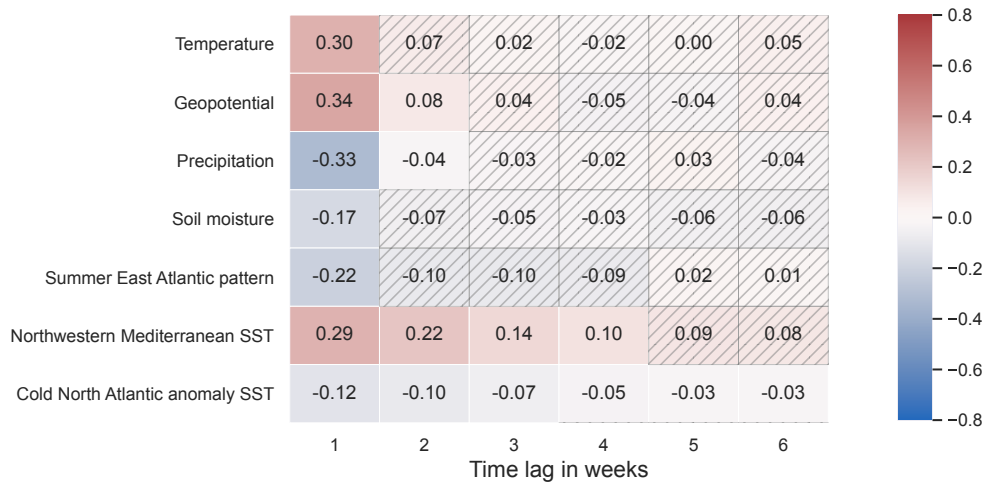


FIG. 7. **Lagged linear correlations between the predictors and the *temperature*** in the extended summer season (MJJAS) at weekly time resolution. Hatched cells correspond to non-significant linear Pearson correlation coefficients at 5% significance level.

longer than a week, since they are more persistent. While the *NWMED SST* correlates positively with the *temperature* over CE, the *CNAA SST* correlates negatively with both.

2) RELEVANCE OF LAGGED PREDICTORS FOR THE MACHINE LEARNING MODELS

Each of the seven predictors is provided to the ML models at four time lags, building a set of 28 lagged predictors for each lead time (Section 2a2). The relevance of a lagged predictor for each ML model is given by the absolute value of its correlation coefficient for the linear models and its feature importance for the RF models. Here, the impurity-based feature (or Gini) importance for a predictor X_i is computed by the sum of all impurity decrease measures of all nodes in the forest at which a split based on X_i has been conducted, normalized by the number of trees (Menze et al. 2009; Nembrini et al. 2018). These values are shown in Tables D1 and D2 for the linear models (RR and RC, respectively) and in Tables D3 and D4 for the RFs (RFR and RFC, respectively) in the Appendix.

In general, predictors at short lags are more useful to the statistical models. Also, the longer the forecast's lead time, the higher the relative contribution from SST becomes. The location of the most important SST region is lead-time dependent: the *NWMED SST* dominates for short lead times (up to two weeks) and the *CNAA SST* prevails for longer lead times (3–6 weeks). The *CNAA*

SST's dominance at long lead times is consistent with the linear correlation shown in Fig. 7, which remains significant for *CNAA SSTs* at the longest lead times.

When forecasting the $+1\sigma$ and the $+1.5\sigma$ heatwave indices, the overall set of relevant lagged predictors is similar, with two exceptions: First, the *SST* is used more to forecast high temperature anomalies ($+1\sigma$) compared to extremely high temperature anomalies ($+1.5\sigma$). Second, the RFC model relies more on *soil moisture* to forecast extremely high temperature anomalies ($+1.5\sigma$) compared to high temperature anomalies ($+1\sigma$), coinciding with the findings by Lopez-Gomez et al. (2022). The different importances of the *SST* and *soil moisture* for forecasting the two heatwave indices could be due to the positive feedback between temperature and soil moisture (Section 1) being more pronounced for extremely high compared to high temperature anomalies. Nevertheless, we can find more marked differences between the two families of statistical models:

(i) *Linear models* For the linear models, *SSTs* dominate at all lead times. In particular, the *CNAA SST* is the most relevant predictor for the RR model at lead times of 2–6 weeks. Nonetheless, the *temperature* is a useful predictor for the RR model at short lead times (1–3 weeks) as well. At a lead time of one week, also the *precipitation* and *soil moisture* contribute to the regression forecast. In contrast, these three lagged predictors are not used by the RC model, which relies almost exclusively on *SSTs*. Therefore, the prediction skill of the ML models incorporating only the *NWMED* and *CNAA SST* predictors has been tested additionally (Appendix, Figs. E1–E3). The regression models have poorer prediction skill when using *SST*-based predictors only. The RC probabilistic classification model benefits from including *SST*-only predictors at lead times of 4–6 weeks for $+1.5\sigma$, indicating that the *SSTs* are the most important predictors for these forecasts (Appendix, Table D2) and the other predictors only increase the model's complexity. Overall, poorer prediction skill is observed for the binary classification models that use only *SST* predictors, especially for the $+1.5\sigma$ prediction.

(ii) *RF models* For the RF models, *temperature*, *geopotential*, *precipitation*, the *SEA* index, and *NWMED SST* at short lags are the most important predictors at short lead times (one week) and *SSTs* are found to dominate for longer lead times (2–6 weeks). In addition, *soil moisture* and the *SEA* index are useful at lead times of 3–6 and 1–5 weeks, respectively. At lead times longer than one week, these two predictors have no significant linear correlation with the *temperature* (Fig. 7) and are used by the RF models but not by the linear models. A plausible explanation

for this phenomenon is the presence of highly non-linear links between *temperature* and *soil moisture*, and *temperature* and the *SEA* index. The physical mechanism behind the non-linear link between *temperature* and *soil moisture* can be the positive feedback described in Section 1 as well as threshold behavior. For example, over transitional wet/dry regimes, soil moisture exhibits large variability and therefore air temperature can be altered by up to 6–7K, while typical soil moisture variations can impact air temperature by up to 1.1–1.3K (Schwingshackl et al. 2017). The SEA pattern and its relation to enhanced summer temperature anomalies resemble the one of air temperature and the summer North Atlantic Oscillation (Folland et al. 2009). The anomalous subsidence associated with the positive geopotential center of the SEA pattern over CE causes a reduction of cloud cover and thus increased solar radiation and surface sensible heating. Increased sensible heating can help maintain the anticyclone over land, contribute to further dryness of the soil, and thus lead to a positive feedback loop with increasing temperatures. These two non-linear links between *temperature* and *soil moisture*, and *temperature* and the *SEA* index (including *soil moisture*) would explain the enhanced skill of the RF models compared to the linear models at lead times higher than four weeks (Section 3a).

4. Limitations and downstream tasks

In this section, the current limitations are discussed and further research ideas to improve the forecasts are suggested: (1) alternative statistical models, (2) approaches to overcome the limitations due to the small sample size, and (3) non-operational statistical models.

(1) The statistical models used in our study belong to the field of classical ML. The complex nature of climate data (e.g., non-linear dependencies between predictors, autocorrelation, and unobserved predictors) poses important challenges to traditional ML models. As discussed in Section 1, DL is also being used for extreme weather forecasting. DL can capture more complex relationships between predictors and target, and might therefore be better suited to describe the mechanisms behind heatwaves, which most likely include non-linear processes. In addition, classical ML approaches benefit from domain-specific hand-crafted features to account for dependencies in time or space but rarely exploit spatio-temporal dependencies exhaustively. In contrast, DL can automatically extract abstract spatio-temporal features (Reichstein et al. 2019). Yet, DL models require larger datasets than the ones used for this study and were therefore not used.

(2) One of the main limitations of this study is the size of the dataset. The initial dataset is considerably larger, but precious information gets lost when taking the average over latitude-longitude boxes. It might be interesting to explore the effect of using several smaller sub-boxes instead of one large box. Additional columns could be added to the dataset, such as a box label or its latitude-longitude coordinates. Also, the currently used boxes are rectangular and their coordinates are chosen based on our physical understanding and the correlation to the target. This could be refined by letting an algorithm select sub-regions of different shapes for each predictor based on the correlation of each grid cell to the target (Vijverberg et al. 2020) or even including the spatial information of the predictors (van Straaten et al. 2022). While lower-dimensional statistical models like RR and RC might not be able to distinguish between distinct mechanisms acting in different regions, RFs are expected to benefit from additional gridded observational data.

(3) The proposed ML models use input data at daily resolution and make weekly predictions. Therefore, to provide the predictions by these models operationally, there is a need for input data updates with at least weekly frequency. Since this high frequency of updates is not available for the data from gridded observations used in this study, the proposed ML models cannot be used operationally. ERA5 reanalysis data, which provides preliminary product updates every 5 days (Hersbach et al. 2020), could be explored as an alternative input.

5. Conclusions

To conclude, we summarize the improvements on sub-seasonal central European temperature anomalies and heatwave prediction by the chosen ML models: The performance of the linear and RF models decays with lead time but outperforms persistence and climatology at all lead times. ECMWF yields accurate forecasts for 1–2 weeks lead time but our ML models compete with the ensemble mean of ECMWF’s hindcast at lead times longer than two weeks. While the linear models perform better for shorter lead times (1–3 weeks), the RFs take over at lead times longer than four weeks.

The statistical regression forecast of summer temperature is better than a random prediction in forecasting the sign of the anomalies at all considered lead times (1–6 weeks) and outperforms the ensemble mean of ECMWF’s hindcast at long lead times (3–6 weeks). However, extreme values are poorly captured. For the classification problem, both statistical models yield a *useful*

probabilistic forecast (meaning $BS < 0.25$ and $ROC\ AUC > 0.5$) for each of the considered lead times (1–6 weeks), except for the RC model at 5–6 weeks lead time. It is remarkable that non-null skill by the RFC model is present at these long lead times. The binary forecast by at least one of the statistical models is *useful* (meaning $FPR/TPR < 1$ and $EDI > 0$) at 1–5 weeks lead time for the $+1\sigma$ heatwave index and at lead times of one, four, and five weeks for the $+1.5\sigma$ heatwave index (Section 3a).

At short lead times (1 week), the following variables are found to be the best predictors of summer temperature anomalies and heatwaves in CE: local 2-m air *temperature*, 500-hPa *geopotential*, *precipitation*, and *NWMED SST*. At longer lead times (2–6 weeks), *NWMED* and *CNAA SST* are the most relevant predictors. Moreover, the *SEA* index and *soil moisture* have a linear link with *temperature* at one week lead time and a possible non-linear link at longer lead times (Section 3b).

In summary, even though our ML models cannot currently be used operationally, these statistical models seem to capture a signal that the ensemble mean of ECMWF's hindcast is not capturing. ML models can, therefore, help extend the forecasting lead time of summer temperature anomalies and heatwaves to sub-seasonal scales, and are a promising direction for further research in sub-seasonal forecasting. Nevertheless, making better forecasts is not enough. Forecasts acquire value through their ability to influence the decisions made by their users (Murphy 1993). As discussed in the Introduction (Section 1), EWS involve not only forecasting the heatwave event but also triggering effective and timely response plans that target vulnerable populations and regions. This second step must also be successfully implemented to reduce the impact of such damaging events (Merz et al. 2020; White et al. 2021).

Acknowledgments. This project has received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (project "HEATforecast", Grant agreement No. 847456). Support from the Swiss National Science Foundation through projects PP00P2_170523 and PP00P2_198896 to M. P. and D. D. is gratefully acknowledged. J. C. is supported by the US National Science Foundation grants AGS-1657748, PLR-1901352, and ARCSS-2115068. We thank O. Wulff for downloading a part of the data used for this study. Also, we appreciate the recommendations from M. Murphy and A. Baldus Benet. We acknowledge L. Grunwald for improving the English in the manuscript. Finally, we thank the three anonymous reviewers for providing insightful comments which improved our work. The authors declare no conflicts of interests.

Data availability statement. We have made the *Python* code used to perform the calculations and generate the figures publicly available on GitHub.¹ The RR and RC functions belong to the *linear model*, and the RFR and the RFC functions belong to the *ensemble* modules from *Sklearn*, respectively (Pedregosa et al. 2011). The Pearson linear correlation test uses the TIGRAMITE code by J. Runge, which is publicly available (Runge et al. 2019).² We acknowledge the E-OBS dataset from the EU-FP6 project UERRA³ and the Copernicus Climate Change Service, and the data providers in the ECA&D project (Cornes et al. 2018).⁴ The ERA-Interim (Dee et al. 2011) and ERA5-Land (Muñoz-Sabater et al. 2021) data are provided by ECMWF.⁵ The HadISST data (version 1.1) are provided by the Met Office Hadley Centre⁶ (Rayner et al. 2003). The ECMWF S2S data are publicly accessible.⁷

¹www.github.com/bethweirich/hwai.git

²www.github.com/jakobrunge/tigramite

³www.uerra.eu

⁴www.ecad.eu

⁵www.ecmwf.int

⁶www.metoffice.gov.uk/hadobs

⁷apps.ecmwf.int/datasets/data/s2s

APPENDIX A

Regression forecasts' time series

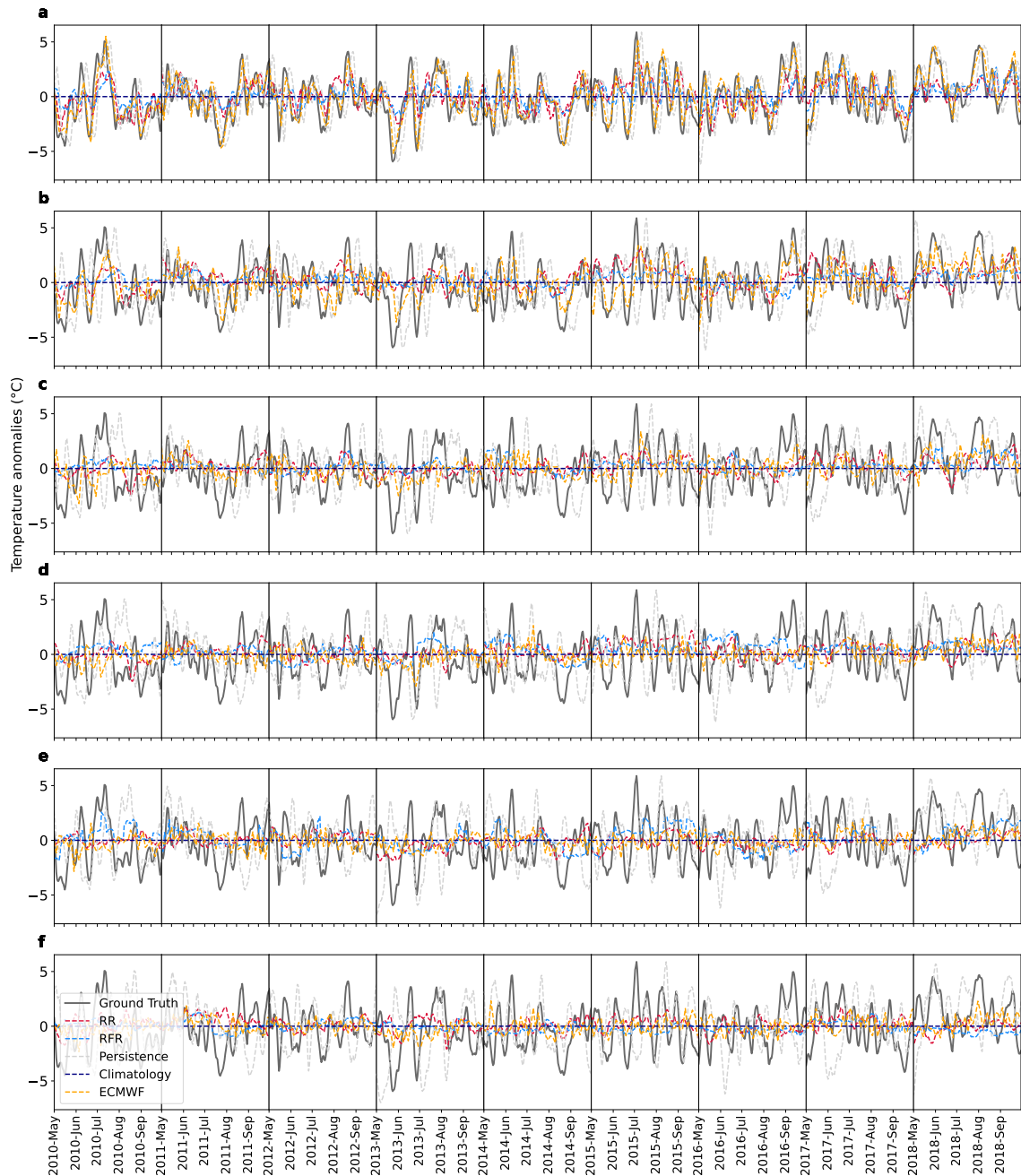


FIG. A1. **Regression time series.** The ground truth time series, the reference forecasts, and the predictions by the ML regression models of the temperature anomalies are shown for the nine summers in the test time period (2010–2018). Figs. a–f correspond to lead times 1–6, respectively.

Nested Cross-Validation

To assess the robustness of our ML models, a CV scheme is implemented. In CV, the model is trained on different data subsets, which reduces overfitting and results in a better generalisation. Moreover, CV removes the dependency on an arbitrarily-selected test set (i.e., on decadal climate variability), making the metrics more robust (Vabalas et al. 2019). Here, a nested CV scheme with five outer and two inner splits is used (Fig. B1). The main benefit of nested CV compared to other CV schemes is that the statistical model is trained and tested on the full dataset while maintaining the independence of the test set, making this method well-suited for a limited sample size.

Nested CV is generally not used for time series data since consecutive time steps are strongly correlated. However, since the correlation between the considered predictors decays after a maximum of a few months and only summer data points are selected for this study, summers belonging to different years can be considered independent. To avoid a strong correlation between the sets at the splitting points, the data is split during the winter months.



FIG. B1. **Nested cross-validation scheme.** $N = 5$ different test sets are predicted by the statistical models and the metrics with respect to the ground truth are calculated for each test set. The final metrics are obtained by averaging the metrics for the five test sets. The uncertainties of these metrics are estimated via the standard deviation of these 5-member ensembles. This figure is adopted from Vabalas et al. (2019).

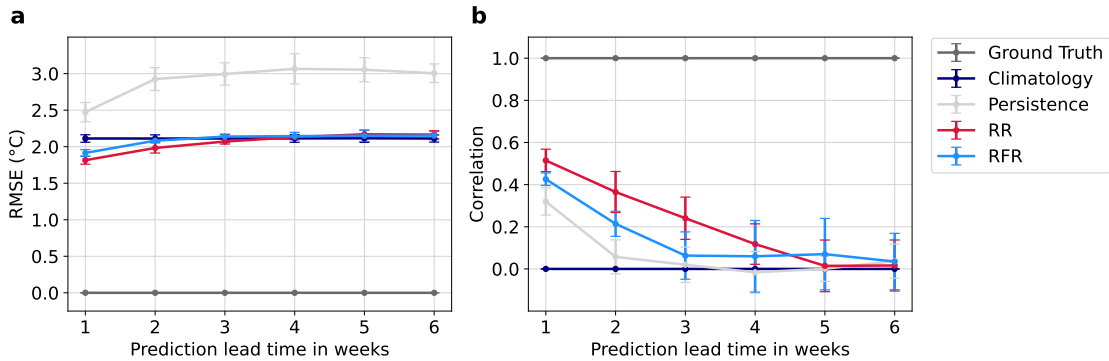


FIG. B2. **Performance of the regression models for six different lead times with nested CV.** (a) RMSE and (b) correlation for the regression forecasts. An accurate forecast is characterized by a low RMSE and a high correlation. The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble.

The metrics obtained with nested CV (Figs. B2, B3, and B4) are similar, although smoother, compared to the results without CV (Figs. 4, 5, and 6), except for the binary classification by the RC model (Fig. B4c). The linear models also show higher skill than the RF models for lead times up to three weeks and the RFs outperform the linear models at 5–6 weeks lead time. While the skill of the ML models at short lead times (up to three weeks) is similar with and without CV, the models in nested CV perform slightly worse for longer lead times. Moreover, the uncertainty of the ML models is higher with nested CV. Therefore, while at least two ML models outperform persistence and climatology on average for all lead times, the error bars overlap with the reference forecasts for lead times of three weeks and longer. A comparison to the ECMWF forecast can not be included for nested CV, because the dynamical model is not available during the full test period used for this CV scheme (1981–2018).

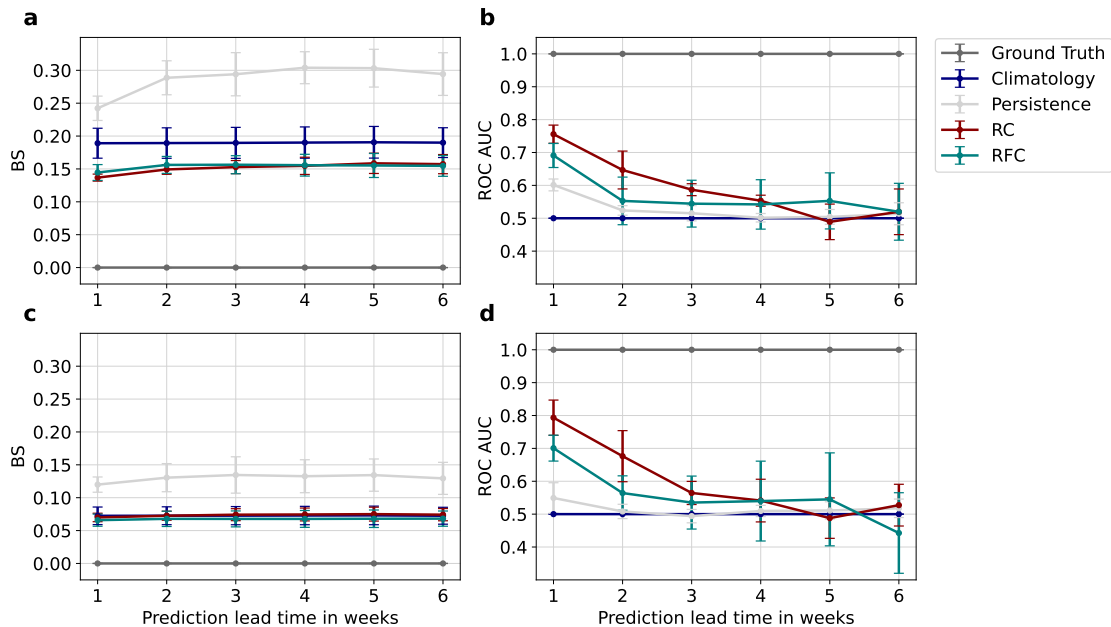


FIG. B3. Performance of the probabilistic classification models for six different lead times with nested CV. BS and ROC AUC for the $+1\sigma$ (a&b) and $+1.5\sigma$ (c&d) weekly heatwave indices. An accurate probabilistic classification forecast is characterized by a low BS and a high ROC AUC. A no-skill probabilistic classification forecast is represented by a BS of 1 and a ROC AUC of 0.5 (as indicated by the climatology). The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble.

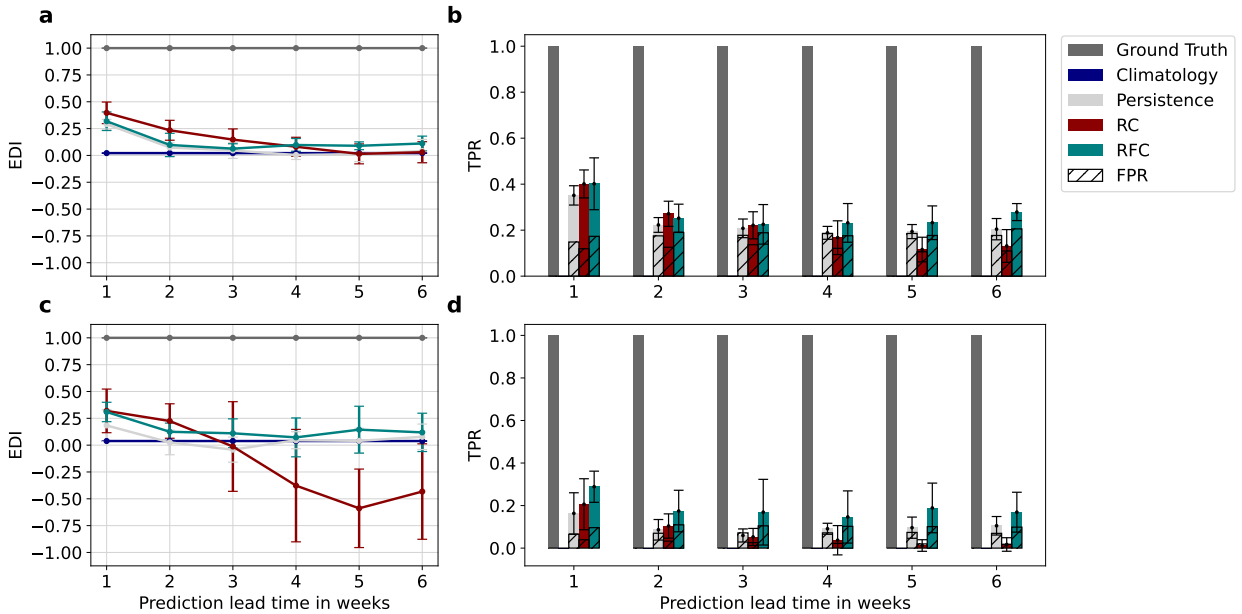


FIG. B4. **Performance of the binary classification models for six different lead times with nested CV.** (a) EDI and (b) TPR (coloured bars) and FPR (stippled bars) for the $+1\sigma$ weekly heatwave index. (c) and (d) are the corresponding forecasts for the $+1.5\sigma$ weekly heatwave index. An accurate binary classification forecast is characterized by a high EDI, a high TPR, and a low FPR. The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble. Since the climatology forecast predicts only zeros (no heatwave), both its TPR and FPR are equal to zero at all lead times (Figs. b&d).

APPENDIX C

Hyper-parameters

Target	Lead time (weeks)	α	Number of estimators	Min. samples/leaf	Max. depth
Temperature anomalies	1	1.0	100	20	5
	2	0.0	200	116	8
	3	1.0	100	52	5
	4	1.0	50	4	5
	5	1.0	200	12	5
	6	0.0	400	100	5
+1σ heatwave index	1	1.0	600	4	14
	2	0.95	400	4	14
	3	1.0	400	4	14
	4	0.0	600	4	14
	5	1.0	600	4	14
	6	1.0	600	4	14
+1.5σ heatwave index	1	1.0	600	4	14
	2	0.75	400	4	14
	3	1.0	600	4	14
	4	1.0	600	4	14
	5	1.0	600	4	14
	6	1.0	600	4	14

TABLE C1. **Optimized hyper-parameters.** Linear (α) and RF (number of estimators, minimum samples per leaf, and maximum depth) hyper-parameters for three targets and six lead times.

APPENDIX D

Correlation coefficients and feature importances

Lead time		1 week	2 weeks	3 weeks	4 weeks	5 weeks	6 weeks
Predictor	Lag (weeks)						
Temperature	1	0.47	-	-	-	-	-
	2	-0.4	-0.3	-	-	-	-
	3	-0.23	-0.51	-0.42	-	-	-
	4	0.05	0.02	-0.07	-0.11	-	-
	5	-	0.26	0.35	0.31	0.26	-
	6	-	-	0.2	0.32	0.29	0.31
	7	-	-	-	-0.28	-0.22	-0.14
	8	-	-	-	-	-0.14	-0.08
	9	-	-	-	-	-	-0.07
Geopotential	1	0.07	-	-	-	-	-
	2	0.21	0.21	-	-	-	-
	3	0.14	0.33	0.25	-	-	-
	4	-0.22	-0.17	-0.14	-0.13	-	-
	5	-	-0.3	-0.38	-0.37	-0.4	-
	6	-	-	-0.18	-0.34	-0.31	-0.32
	7	-	-	-	0.29	0.15	0.08
	8	-	-	-	-	0.25	0.18
	9	-	-	-	-	-	0.15
Precipitation	1	-0.66	-	-	-	-	-
	2	0.07	0.22	-	-	-	-
	3	0.21	0.27	0.3	-	-	-
	4	-0.03	0.02	0.04	-0.01	-	-
	5	-	-0.05	-0.05	0.02	-0.04	-
	6	-	-	-0.1	-0.01	0.04	-0.05
	7	-	-	-	0.08	0.17	0.13
	8	-	-	-	-	0.2	0.28
	9	-	-	-	-	-	0.33
Soil moisture	1	0.94	-	-	-	-	-
	2	-0.65	-0.08	-	-	-	-
	3	-0.24	-0.28	-0.39	-	-	-
	4	0.04	0.08	-0.04	-0.32	-	-
	5	-	0.03	0.14	-0.02	-0.27	-
	6	-	-	0.08	0	-0.05	-0.17
	7	-	-	-	0.19	-0.06	-0.06
	8	-	-	-	-	0.18	-0.11
	9	-	-	-	-	-	0.03
SEA	1	-0.06	-	-	-	-	-
	2	-0.01	-0.04	-	-	-	-
	3	-0.14	-0.12	-0.13	-	-	-
	4	-0.11	-0.14	-0.14	-0.17	-	-
	5	-	0.17	0.2	0.24	0.18	-
	6	-	-	0.03	0.08	0.13	0.14
	7	-	-	-	0.01	0.04	0
	8	-	-	-	-	0.04	0.04
	9	-	-	-	-	-	-0.1
NWMED SST	1	2.1	-	-	-	-	-
	2	-1.67	3.05	-	-	-	-
	3	-0.2	-3.31	1.99	-	-	-
	4	0.31	0.4	-2.37	1.35	-	-
	5	-	0.46	0.12	-2.5	0.46	-
	6	-	-	0.69	1.52	-1.09	-0.35
	7	-	-	-	-0.02	1.45	0.98
	8	-	-	-	-	-0.56	-0.23
	9	-	-	-	-	-	-0.26
CNAASST	1	-1.74	-	-	-	-	-
	2	1.8	-3.24	-	-	-	-
	3	0.36	3.67	-3.27	-	-	-
	4	-0.39	0.47	3.25	-4.15	-	-
	5	-	-1	2.04	7.83	-0.97	-
	6	-	-	-2.16	-4.93	2.34	1.38
	7	-	-	-	1.08	-3.27	-3.73
	8	-	-	-	-	1.74	3.05
	9	-	-	-	-	-	-0.76

TABLE D1. Regression coefficients for a single RR model trained on the full training set. Coefficients with absolute values above 0.5 are bold.

Lead time		1 week		2 weeks		3 weeks		4 weeks		5 weeks		6 weeks	
Target		+1 σ	+1.5 σ	+1 σ	+1.5 σ	+1 σ	+1.5 σ	+1 σ	+1.5 σ	+1 σ	+1.5 σ	+1 σ	+1.5 σ
Predictor	Lag (weeks)												
Temperature	1	0.16	0.09	-	-	-	-	-	-	-	-	-	-
	2	-0.13	-0.06	-0.1	-0.03	-	-	-	-	-	-	-	-
	3	-0.05	-0.08	-0.13	-0.12	-0.11	-0.09	-	-	-	-	-	-
	4	-0.06	-0.03	-0.07	-0.04	-0.1	-0.06	-0.11	-0.07	-	-	-	-
	5	-	-	0.06	0.05	0.07	0.06	0.05	0.05	0.04	0.04	-	-
	6	-	-	-	-	0.07	0.04	0.09	0.07	0.07	0.06	0.07	0.06
	7	-	-	-	-	-	-	-0.03	-0.09	-0.01	-0.08	0	-0.07
	8	-	-	-	-	-	-	-	-	-0.01	-0.02	0.03	0.01
	9	-	-	-	-	-	-	-	-	-	-	-0.09	-0.08
Geopotential	1	-0.02	-0.04	-	-	-	-	-	-	-	-	-	-
	2	0.09	0.06	0.08	0.05	-	-	-	-	-	-	-	-
	3	0.02	0.07	0.08	0.1	0.06	0.09	-	-	-	-	-	-
	4	0.01	-0.01	0.02	-0.01	0.04	0.01	0.04	0.01	-	-	-	-
	5	-	-	-0.05	-0.03	-0.06	-0.02	-0.04	-0.02	-0.07	-0.03	-	-
	6	-	-	-	-	-0.04	-0.04	-0.06	-0.06	-0.04	-0.06	-0.04	-0.05
	7	-	-	-	-	-	-	0.03	0.05	-0.02	0.03	0.03	-0.04
	8	-	-	-	-	-	-	-	-	0.06	0.04	0.04	0.03
	9	-	-	-	-	-	-	-	-	-	-	0.12	0.08
Precipitation	1	-0.19	-0.1	-	-	-	-	-	-	-	-	-	-
	2	-0.01	-0.03	0.04	0.01	-	-	-	-	-	-	-	-
	3	0	0	0.02	0.02	0.03	0.04	-	-	-	-	-	-
	4	-0.01	0	-0.02	0	-0.01	0.01	-0.01	0.01	-	-	-	-
	5	-	-	-0.02	0	-0.02	-0.02	-0.01	-0.01	0	-0.01	-	-
	6	-	-	-	-	-0.02	-0.02	-0.01	-0.01	0.01	0	-0.02	-0.02
	7	-	-	-	-	-	-	0.03	0	0.07	0.02	0.05	0.01
	8	-	-	-	-	-	-	-	-	0.08	0.03	0.09	0.03
	9	-	-	-	-	-	-	-	-	-	-	0.15	0.07
Soil moisture	1	0.29	0.16	-	-	-	-	-	-	-	-	-	-
	2	-0.17	0	0	0.08	-	-	-	-	-	-	-	-
	3	-0.01	-0.05	-0.02	-0.06	-0.04	-0.02	-	-	-	-	-	-
	4	-0.02	-0.05	0.03	-0.05	0	-0.05	-0.04	-0.07	-	-	-	-
	5	-	-	-0.01	0.02	0.01	0.07	0	0.05	-0.06	-0.02	-	-
	6	-	-	-	-	0.03	-0.02	0	0	0	-0.01	-0.01	0
	7	-	-	-	-	-	-	0.02	0	-0.08	-0.04	-0.08	-0.04
	8	-	-	-	-	-	-	-	-	0.08	0.04	0.04	0.04
	9	-	-	-	-	-	-	-	-	-	-	-0.06	-0.04
SEA	1	-0.07	-0.03	-	-	-	-	-	-	-	-	-	-
	2	-0.03	-0.01	-0.03	-0.01	-	-	-	-	-	-	-	-
	3	-0.07	-0.04	-0.05	-0.03	-0.05	-0.03	-	-	-	-	-	-
	4	-0.06	-0.03	-0.07	-0.03	-0.06	-0.03	-0.06	-0.03	-	-	-	-
	5	-	-	0.05	0.02	0.05	0.03	0.06	0.03	0.04	0.02	-	-
	6	-	-	-	-	0.03	0.02	0.04	0.03	0.06	0.03	0.06	0.04
	7	-	-	-	-	-	-	0	-0.01	0.01	-0.01	0	-0.02
	8	-	-	-	-	-	-	-	-	0.01	0.02	0.02	0.02
	9	-	-	-	-	-	-	-	-	-	-	-0.02	-0.03
NWMED SST	1	0.66	0.37	-	-	-	-	-	-	-	-	-	-
	2	-0.71	-0.29	0.7	0.47	-	-	-	-	-	-	-	-
	3	0.25	0.01	-0.66	-0.54	0.46	0.25	-	-	-	-	-	-
	4	-0.04	0.01	-0.02	0.14	-0.39	-0.23	0.49	0.25	-	-	-	-
	5	-	-	0.15	0.02	-0.32	-0.11	-0.9	-0.43	0.16	0.03	-	-
	6	-	-	-	-	0.38	0.15	0.41	0.21	-0.39	0.03	-0.09	0.08
	7	-	-	-	-	-	-	0.11	0.02	0.34	-0.12	0.15	-0.08
	8	-	-	-	-	-	-	-	-	-0.03	0.12	0.01	-0.03
	9	-	-	-	-	-	-	-	-	-	-	-0.02	0.08
CNA A SST	1	-0.18	0	-	-	-	-	-	-	-	-	-	-
	2	0.54	0.09	-0.45	-0.24	-	-	-	-	-	-	-	-
	3	-0.29	-0.05	0.4	0.18	-0.67	-0.42	-	-	-	-	-	-
	4	0.02	-0.01	0.25	0.19	0.25	0.23	-1.55	-0.73	-	-	-	-
	5	-	-	-0.16	-0.12	1.17	0.58	2.8	1.3	-0.52	-0.18	-	-
	6	-	-	-	-	-0.75	-0.4	-1.53	-0.67	0.97	0.18	0.12	-0.12
	7	-	-	-	-	-	-	0.27	0.09	-0.66	0.03	-0.25	0.11
	8	-	-	-	-	-	-	-	-	0.2	-0.05	0.21	0.08
	9	-	-	-	-	-	-	-	-	-	-	-0.06	-0.08

TABLE D2. Regression coefficients for a single RC model trained on the full training set. Coefficients with absolute values above 0.5 are bold.

Lead time		1 week	2 weeks	3 weeks	4 weeks	5 weeks	6 weeks
Predictor	Lag (weeks)						
Temperature	1	0.02	-	-	-	-	-
	2	0.01	0.03	-	-	-	-
	3	0.01	0.02	0.01	-	-	-
	4	0.01	0.05	0.03	0.01	-	-
	5	-	0.01	0	0.01	0.01	-
	6	-	-	0.01	0.02	0.02	0.02
	7	-	-	-	0.03	0.03	0.01
	8	-	-	-	-	0.01	0.01
	9	-	-	-	-	-	0.01
Geopotential	1	0.23	-	-	-	-	-
	2	0.01	0.01	-	-	-	-
	3	0.01	0	0	-	-	-
	4	0.01	0.01	0.01	0	-	-
	5	-	0	0.01	0.01	0.01	-
	6	-	-	0	0.01	0.01	0
	7	-	-	-	0.02	0.02	0.01
	8	-	-	-	-	0.01	0
	9	-	-	-	-	-	0.01
Precipitation	1	0.18	-	-	-	-	-
	2	0.03	0.01	-	-	-	-
	3	0.01	0.01	0.01	-	-	-
	4	0	0	0	0.01	-	-
	5	-	0	0	0.01	0.01	-
	6	-	-	0.01	0.02	0.02	0.02
	7	-	-	-	0.01	0	0.01
	8	-	-	-	-	0.01	0.01
	9	-	-	-	-	-	0.02
Soil moisture	1	0.01	-	-	-	-	-
	2	0.01	0.02	-	-	-	-
	3	0.01	0.02	0.02	-	-	-
	4	0.02	0.01	0.02	0.02	-	-
	5	-	0.04	0.05	0.04	0.05	-
	6	-	-	0.05	0.05	0.05	0.06
	7	-	-	-	0.01	0.01	0.01
	8	-	-	-	-	0.02	0.04
	9	-	-	-	-	-	0.03
SEA	1	0.07	-	-	-	-	-
	2	0.01	0.03	-	-	-	-
	3	0.01	0.01	0.03	-	-	-
	4	0.01	0.02	0.01	0.02	-	-
	5	-	0.06	0.08	0.06	0.05	-
	6	-	-	0.04	0.02	0.02	0.04
	7	-	-	-	0.03	0.03	0.04
	8	-	-	-	-	0.01	0.02
	9	-	-	-	-	-	0.01
NWMED SST	1	0.21	-	-	-	-	-
	2	0.01	0.35	-	-	-	-
	3	0.03	0.05	0.13	-	-	-
	4	0.01	0.03	0.03	0.07	-	-
	5	-	0.01	0.04	0.04	0.05	-
	6	-	-	0.06	0.04	0.05	0.05
	7	-	-	-	0.12	0.1	0.07
	8	-	-	-	-	0.04	0.04
	9	-	-	-	-	-	0.05
CNAASST	1	0.02	-	-	-	-	-
	2	0.02	0.1	-	-	-	-
	3	0.01	0.01	0.12	-	-	-
	4	0.02	0.03	0.03	0.06	-	-
	5	-	0.09	0.07	0.1	0.13	-
	6	-	-	0.12	0.15	0.16	0.23
	7	-	-	-	0.03	0.02	0.01
	8	-	-	-	-	0.07	0.02
	9	-	-	-	-	-	0.16

TABLE D3. Predictor importances for a single RFR model trained on the full training set. Values above 0.04 are bold.

Lead time		1 week		2 weeks		3 weeks		4 weeks		5 weeks		6 weeks	
Target		$+1\sigma$	$+1.5\sigma$	$+1\sigma$	$+1.5\sigma$	$+1\sigma$	$+1.5\sigma$	$+1\sigma$	$+1.5\sigma$	$+1\sigma$	$+1.5\sigma$	$+1\sigma$	$+1.5\sigma$
Predictor	Lag (weeks)												
Temperature	1	0.06	0.08	-	-	-	-	-	-	-	-	-	-
	2	0.02	0.02	0.03	0.02	-	-	-	-	-	-	-	-
	3	0.03	0.03	0.02	0.03	0.03	0.02	-	-	-	-	-	-
	4	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	-	-	-	-
	5	-	-	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	-	-
	6	-	-	-	-	0.03	0.02	0.03	0.02	0.03	0.02	0.03	0.02
	7	-	-	-	-	-	-	0.03	0.04	0.03	0.04	0.03	0.04
	8	-	-	-	-	-	-	-	-	0.03	0.03	0.03	0.03
	9	-	-	-	-	-	-	-	-	-	-	0.03	0.03
Geopotential	1	0.06	0.06	-	-	-	-	-	-	-	-	-	-
	2	0.02	0.02	0.03	0.03	-	-	-	-	-	-	-	-
	3	0.02	0.02	0.02	0.02	0.02	0.02	-	-	-	-	-	-
	4	0.02	0.02	0.03	0.03	0.03	0.02	0.03	0.02	-	-	-	-
	5	-	-	0.03	0.03	0.03	0.03	0.03	0.02	0.03	0.02	-	-
	6	-	-	-	-	0.03	0.02	0.03	0.02	0.02	0.02	0.03	0.03
	7	-	-	-	-	-	-	0.03	0.03	0.03	0.03	0.03	0.03
	8	-	-	-	-	-	-	-	-	0.03	0.03	0.03	0.02
	9	-	-	-	-	-	-	-	-	-	-	0.02	0.03
Precipitation	1	0.07	0.06	-	-	-	-	-	-	-	-	-	-
	2	0.02	0.02	0.03	0.03	-	-	-	-	-	-	-	-
	3	0.02	0.02	0.02	0.02	0.02	0.03	-	-	-	-	-	-
	4	0.02	0.02	0.02	0.02	0.02	0.02	0.02	0.03	-	-	-	-
	5	-	-	0.03	0.03	0.02	0.03	0.02	0.02	0.02	0.02	-	-
	6	-	-	-	-	0.03	0.02	0.03	0.02	0.03	0.02	0.03	0.02
	7	-	-	-	-	-	-	0.03	0.02	0.02	0.02	0.02	0.02
	8	-	-	-	-	-	-	-	-	0.02	0.03	0.02	0.03
	9	-	-	-	-	-	-	-	-	-	-	0.03	0.03
Soil moisture	1	0.03	0.03	-	-	-	-	-	-	-	-	-	-
	2	0.03	0.04	0.03	0.04	-	-	-	-	-	-	-	-
	3	0.03	0.02	0.03	0.03	0.03	0.03	-	-	-	-	-	-
	4	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	-	-	-	-
	5	-	-	0.04	0.04	0.04	0.04	0.04	0.05	0.04	0.05	-	-
	6	-	-	-	-	0.04	0.04	0.04	0.03	0.04	0.03	0.04	0.03
	7	-	-	-	-	-	-	0.03	0.03	0.03	0.03	0.03	0.03
	8	-	-	-	-	-	-	-	-	0.03	0.03	0.04	0.03
	9	-	-	-	-	-	-	-	-	-	-	0.03	0.03
SEA	1	0.05	0.06	-	-	-	-	-	-	-	-	-	-
	2	0.03	0.03	0.04	0.04	-	-	-	-	-	-	-	-
	3	0.03	0.04	0.04	0.05	0.04	0.05	-	-	-	-	-	-
	4	0.03	0.03	0.03	0.04	0.03	0.04	0.04	0.04	-	-	-	-
	5	-	-	0.04	0.04	0.04	0.04	0.03	0.04	0.04	0.04	-	-
	6	-	-	-	-	0.04	0.04	0.03	0.03	0.03	0.04	0.03	0.03
	7	-	-	-	-	-	-	0.03	0.03	0.03	0.02	0.03	0.03
	8	-	-	-	-	-	-	-	-	0.03	0.03	0.03	0.03
	9	-	-	-	-	-	-	-	-	-	-	0.03	0.04
NWMED SST	1	0.06	0.08	-	-	-	-	-	-	-	-	-	-
	2	0.04	0.04	0.06	0.07	-	-	-	-	-	-	-	-
	3	0.04	0.04	0.05	0.05	0.05	0.06	-	-	-	-	-	-
	4	0.03	0.03	0.04	0.04	0.05	0.04	0.05	0.04	-	-	-	-
	5	-	-	0.05	0.04	0.04	0.04	0.04	0.04	0.05	0.04	-	-
	6	-	-	-	-	0.05	0.05	0.04	0.05	0.04	0.05	0.04	0.05
	7	-	-	-	-	-	-	0.04	0.05	0.04	0.05	0.04	0.05
	8	-	-	-	-	-	-	-	-	0.05	0.04	0.04	0.04
	9	-	-	-	-	-	-	-	-	-	-	0.05	0.06
CNAASST	1	0.04	0.03	-	-	-	-	-	-	-	-	-	-
	2	0.04	0.03	0.06	0.04	-	-	-	-	-	-	-	-
	3	0.04	0.03	0.04	0.04	0.05	0.04	-	-	-	-	-	-
	4	0.04	0.03	0.05	0.04	0.05	0.04	0.05	0.04	-	-	-	-
	5	-	-	0.06	0.05	0.06	0.06	0.06	0.06	0.07	0.06	-	-
	6	-	-	-	-	0.06	0.06	0.06	0.06	0.06	0.06	0.07	0.05
	7	-	-	-	-	-	-	0.05	0.06	0.05	0.06	0.05	0.05
	8	-	-	-	-	-	-	-	-	0.05	0.05	0.05	0.05
	9	-	-	-	-	-	-	-	-	-	-	0.06	0.05

TABLE D4. Predictor importances for a single RFC model trained on the full training set. Values above 0.04 are bold.

APPENDIX E

Only-SST runs

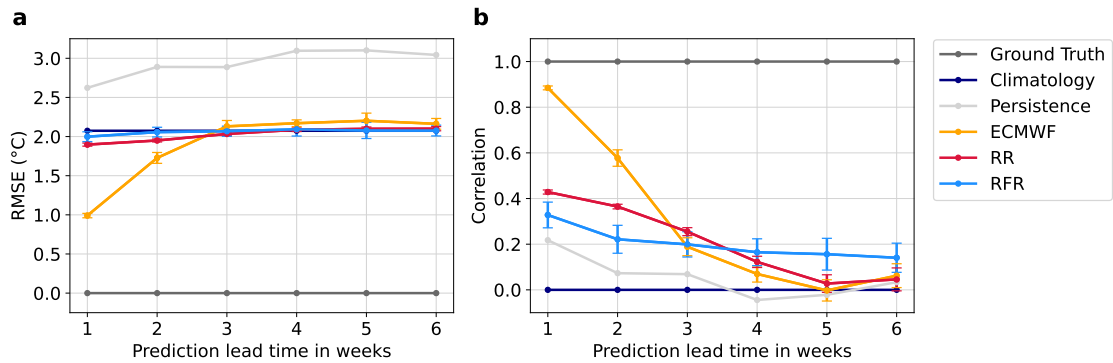


FIG. E1. Performance of the regression models for six different lead times using only the *NWMED* and *CNAA SST* predictors. (a) RMSE and (b) correlation for the regression forecasts. An accurate forecast is characterized by a low RMSE and a high correlation. The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble.

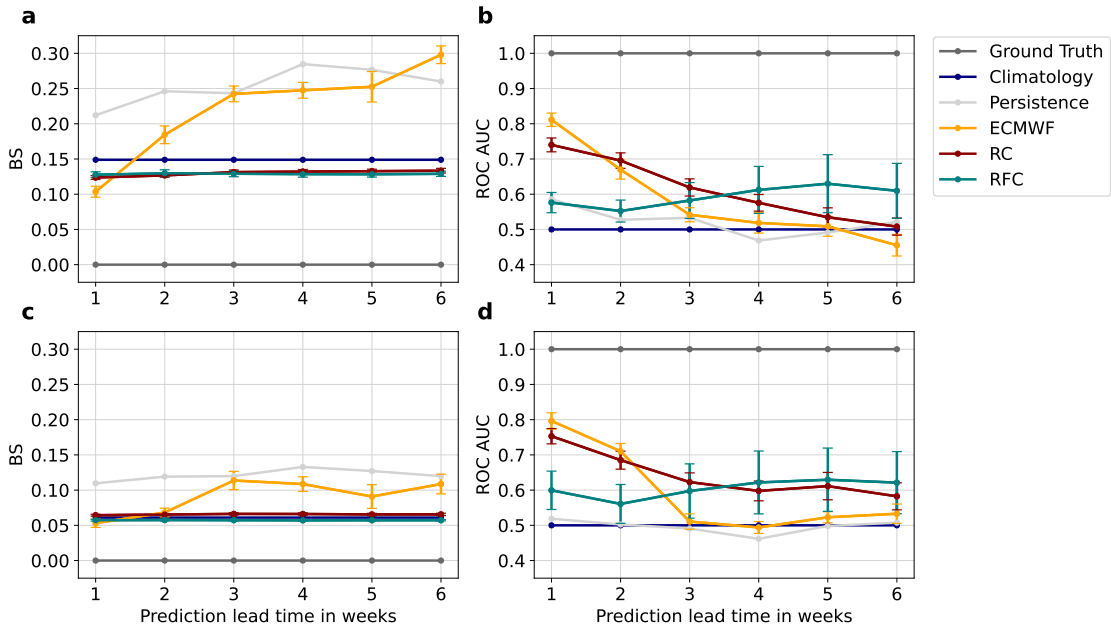


FIG. E2. Performance of the probabilistic classification models for six different lead times using only the *NWMED* and *CNA SST* predictors. BS and ROC AUC for the $+1\sigma$ (a&b) and $+1.5\sigma$ (c&d) weekly heatwave indices. An accurate probabilistic classification forecast is characterized by a low BS and a high ROC AUC. A no-skill probabilistic classification forecast is represented by a BS of 1 and a ROC AUC of 0.5 (as indicated by the climatology). The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble.

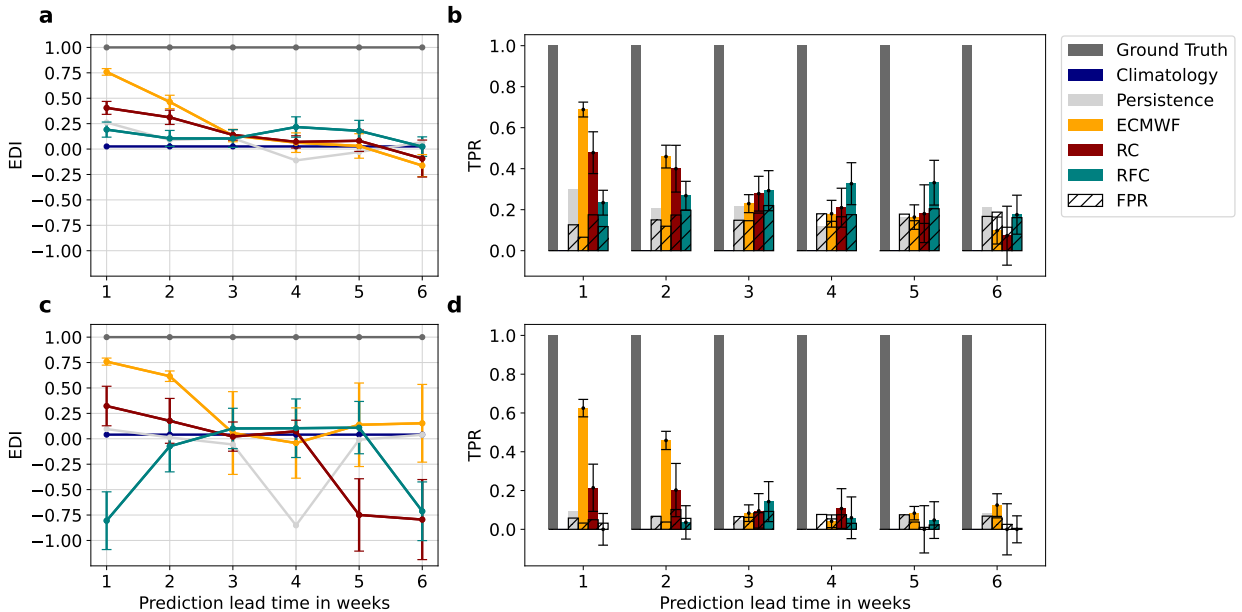


FIG. E3. **Performance of the binary classification models for six different lead times using only the *NWMED* and *CNAASST* predictors.** (a) EDI and (b) TPR (coloured bars) and FPR (stippled bars) for the $+1\sigma$ weekly heatwave index. (c) and (d) are the corresponding forecasts for the $+1.5\sigma$ weekly heatwave index. An accurate binary classification forecast is characterized by a high EDI, a high TPR, and a low FPR. The error bars show the uncertainty of each forecast estimated via the standard deviation of the ensemble. Since the climatology forecast predicts only zeros (no heatwave), both its TPR and FPR are equal to zero at all lead times (Figs. b&d).

References

- Barriopedro, D., E. M. Fischer, J. Luterbacher, R. M. Trigo, and R. Garcia-Herrera, 2011: The hot summer of 2010: Redrawing the temperature record map of europe. *Science*, **332**, 220–224, <https://doi.org/10.1126/science.1201224>.
- Bassil, K., and D. Cole, 2010: Effectiveness of public health interventions in reducing morbidity and mortality during heat episodes: a structured review. *International Journal of Environmental Research and Public Health*, **7**, 991–1001, <https://doi.org/10.3390/ijerph7030991>.
- Basu, R., 2002: Relation between elevated ambient temperature and mortality: a review of the epidemiologic evidence. *Epidemiologic Reviews*, **24**, 190–202, <https://doi.org/10.1093/epirev/mxf007>.
- Black, E., M. Blackburn, G. Harrison, B. Hoskins, and J. Methven, 2004: Factors contributing to the summer 2003 european heatwave. *Weather*, **59**, 217–223, <https://doi.org/10.1256/wea.74.04>.
- Bladé, I., B. Liebmann, D. Fortuny, and G. J. van Oldenborgh, 2011: Observed and simulated impacts of the summer nao in europe: Implications for projected drying in the mediterranean region. *Climate Dynamics*, **39**, 709–727, <https://doi.org/10.1007/s00382-011-1195-x>.
- Bradley, A. P., 1997: The use of the area under the roc curve in the evaluation of machine learning algorithms. *Pattern Recognition*, **30**, 1145–1159, [https://doi.org/10.1016/s0031-3203\(96\)00142-2](https://doi.org/10.1016/s0031-3203(96)00142-2).
- Breiman, L., 2001: Random forests. *Machine Learning*, **45**, 5–32, <https://doi.org/10.1023/a:1010933404324>.
- Buzan, J. R., and M. Huber, 2020: Moist heat stress on a hotter earth. *Annual Review of Earth and Planetary Sciences*, **48**, <https://doi.org/10.1146/annurev-earth-053018-060100>.
- Casanueva, A., and Coauthors, 2019: Overview of existing heat-health warning systems in europe. *International Journal of Environmental Research and Public Health*, **16**, <https://doi.org/10.3390/ijerph16152657>.

- Chattopadhyay, A., E. Nabizadeh, and P. Hassanzadeh, 2020: Analog forecasting of extreme causing weather patterns using deep learning. *Journal of Advances in Modeling Earth Systems*, **12**, <https://doi.org/10.1029/2019ms001958>.
- Cornes, R. C., G. van der Schrier, E. J. M. van den Besselaar, and P. D. Jones, 2018: An ensemble version of the e-obs temperature and precipitation data sets. *Journal of Geophysical Research: Atmospheres*, **123**, 9391–9409, <https://doi.org/10.1029/2017jd028200>.
- de Perez, E. C., and Coauthors, 2018: Global predictability of temperature extremes. *Environmental Research Letters*, **13**, 1748–9318, <https://doi.org/10.1088/1748-9326/aab94a>.
- Deb, P., H. Moradkhani, P. Abbaszadeh, A. S. Kiem, J. Engström, D. Keellings, and A. Sharma, 2020: Causes of the widespread 2019/2020 Australian bushfire season. *Earth's Future*, **8**, 2328–4277, <https://doi.org/10.1029/2020ef001671>.
- Dee, D. P., and Coauthors, 2011: The era-interim reanalysis: configuration and performance of the data assimilation system. *Quarterly Journal of the Royal Meteorological Society*, **137**, 553–597, <https://doi.org/10.1002/qj.828>.
- Duchez, A., and Coauthors, 2016: Drivers of exceptionally cold North Atlantic ocean temperatures and their link to the 2015 European heat wave. *Environmental Research Letters*, **11**, <https://doi.org/10.1088/1748-9326/11/7/074004>.
- Ferro, C. A. T., and D. B. Stephenson, 2011: Extremal dependence indices: Improved verification measures for deterministic forecasts of rare binary events. *Weather and Forecasting*, **26**, 699–713, <https://doi.org/10.1175/waf-d-10-05030.1>.
- Fischer, E. M., S. I. Seneviratne, P. L. Vidale, D. Lüthi, and C. Schär, 2007: Soil moisture-atmosphere interactions during the 2003 European summer heat wave. *J. Climate*, **20**, 5081–5099, <https://doi.org/10.1175/jcli4288.1>.
- Folland, C. K., J. Knight, H. W. Linderholm, D. Fereday, S. Ineson, and J. W. Hurrell, 2009: The summer North Atlantic oscillation: Past, present, and future. *Journal of Climate*, **22**, 1082–1103, <https://doi.org/10.1175/2008jcli2459.1>.

- Ford, T. W., P. A. Dirmeyer, and D. O. Benson, 2018: Evaluation of heat wave forecasts seamlessly across subseasonal timescales. *Npj Climate and Atmospheric Science*, **1**, <https://doi.org/10.1038/s41612-018-0027-7>.
- Gneiting, T., F. Balabdaoui, and A. E. Raftery, 2007: Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **69**, 243–268, <https://doi.org/10.1111/j.1467-9868.2007.00587.x>.
- Haiden, T., M. Janousek, F. Vitart, Z. B. Bouallegue, L. Ferranti, F. Prates, and D. Richardson, 2019: Technical memorandum: Evaluation of ecmwf forecasts, including the 2019 upgrade. 10.21957/mlvapkke, URL <https://www.ecmwf.int/node/19277>.
- Hastie, T., R. Tibshirani, and J. Friedman, 2009: *The Elements of Statistical learning: Data mining, inference, and Prediction*. 2nd ed., Springer, 61–68, 249–254, and 587–588 pp.
- Hersbach, H., and Coauthors, 2020: The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, **146**, 1999–2049, <https://doi.org/10.1002/qj.3803>.
- Hu, Z.-Z., A. Kumar, B. Huang, W. Wang, J. Zhu, and C. Wen, 2012: Prediction skill of monthly sst in the north atlantic ocean in ncep climate forecast system version 2. *Climate Dynamics*, **40**, 2745–2759, <https://doi.org/10.1007/s00382-012-1431-z>.
- Huynen, M. M., P. Martens, D. Schram, M. P. Weijnenberg, and A. E. Kunst, 2001: The impact of heat waves and cold spells on mortality rates in the dutch population. *Environmental Health Perspectives*, **109**, 463–470, <https://doi.org/10.1289/ehp.01109463>.
- IPCC, 2013: *Climate Change 2013 - The Physical Science Basis Working Group I Contribution to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, Vol. Chapter 9: Evaluation of Climate Models. Cambridge University Press, 768 pp.
- Jacques-Dumas, V., F. Ragone, P. Borgnat, P. Abry, and F. Bouchet, 2022: Deep learning-based extreme heatwave forecast. *Frontiers in Climate*, **4**, <https://doi.org/10.3389/fclim.2022.789641>.
- JiménezEsteve, B., and D. I. Domeisen, 2022: The role of atmospheric dynamics and largescale topography in driving heatwaves. *Quarterly Journal of the Royal Meteorological Society*, **148**, 2344–2367, <https://doi.org/10.1002/qj.4306>.

- Jolliffe, I. T., and D. B. Stephenson, 2005: Comments on discussion of verification concepts in forecast verification: A practitioners guide in atmospheric science. *Weather and Forecasting*, **20**, 796–800, <https://doi.org/10.1175/waf877.1>.
- Kautz, L.-A., O. Martius, S. Pfahl, J. G. Pinto, A. M. Ramos, P. M. Sousa, and T. Woollings, 2022: Atmospheric blocking and weather extremes over the euro-atlantic sector a review. *Weather and Climate Dynamics*, **3**, 305–336, <https://doi.org/10.5194/wcd-3-305-2022>.
- Khan, N., S. Shahid, L. Juneng, K. Ahmed, T. Ismail, and N. Nawaz, 2019: Prediction of heat waves in pakistan using quantile regression forests. *Atmospheric Research*, **221**, 1–11, <https://doi.org/10.1016/j.atmosres.2019.01.024>.
- Kolstad, E. W., E. A. Barnes, and S. P. Sobolowski, 2017: Quantifying the role of land-atmosphere feedbacks in mediating near-surface temperature persistence. *Quarterly Journal of the Royal Meteorological Society*, **143**, 1620–1631, <https://doi.org/10.1002/qj.3033>.
- Kotharkar, R., and A. Ghosh, 2022: Progress in extreme heat management and warning systems: A systematic review of heat-health action plans (1995-2020). *Sustainable Cities and Society*, **76**, <https://doi.org/10.1016/j.scs.2021.103487>.
- Kumar, A., and J. Zhu, 2018: Spatial variability in seasonal prediction skill of ssts: Inherent predictability or forecast errors? *Journal of Climate*, **31**, 613–621, <https://doi.org/10.1175/jcli-d-17-0279.1>.
- Kunsch, H. R., 1989: The jackknife and the bootstrap for general stationary observations. *The Annals of Statistics*, **17**, <https://doi.org/10.1214/aos/1176347265>.
- Kämäräinen, M., P. Uotila, A. Y. Karpechko, O. Hyvärinen, I. Lehtonen, and J. Räisänen, 2019: Statistical learning methods as a basis for skillful seasonal temperature forecasts in europe. *J. Climate*, **32**, 5363–5379, <https://doi.org/10.1175/jcli-d-18-0765.1>.
- Laguë, M. M., G. B. Bonan, and A. L. S. Swann, 2019: Separating the impact of individual land surface properties on the terrestrial surface energy budget in both the coupled and uncoupled land-atmosphere system. *Journal of Climate*, **32**, 5725–5744, <https://doi.org/10.1175/jcli-d-18-0812.1>.

- Lemaitre, G., F. Nogueira, and C. K. Aridas, 2017: Imbalanced-learn: a python toolbox to tackle the curse of imbalanced datasets in machine learning. *Journal of Machine Learning Research*, **18**, 1–5, <https://doi.org/10.1162/jmlr.2017.18.1>, <https://arxiv.org/abs/1609.06570>.
- Li, S., and A. W. Robertson, 2015: Evaluation of submonthly precipitation forecast skill from global ensemble prediction systems. *Monthly Weather Review*, **143**, 2871–2889, <https://doi.org/10.1175/mwr-d-14-00277.1>.
- Lopez-Gomez, I., A. McGovern, S. Agrawal, and J. Hickey, 2022: Global extreme heat forecasting using neural weather models. *arXiv*, <https://doi.org/10.48550/ARXIV.2205.10972>.
- Lowe, D., K. L. Ebi, and B. Forsberg, 2011: Heatwave early warning systems and adaptation advice to reduce human health consequences of heatwaves. *International Journal of Environmental Research and Public Health*, **8**, 4623–4648, <https://doi.org/10.3390/ijerph8124623>.
- Manrique-Suñén, A., N. Gonzalez-Reviriego, V. Torralba, N. Cortesi, and F. J. Doblas-Reyes, 2020: Choices in the verification of s2s forecasts and their implications for climate services. *Monthly Weather Review*, **148**, 3995–4008, <https://doi.org/10.1175/mwr-d-20-0067.1>.
- Manzanas, R., 2020: Assessment of model drifts in seasonal forecasting: Sensitivity to ensemble size and implications for bias correction. *Journal of Advances in Modeling Earth Systems*, **12**, <https://doi.org/10.1029/2019ms001751>.
- Mecking, J. V., S. S. Drijfhout, J. J.-M. Hirschi, and A. T. Blaker, 2019: Ocean and atmosphere influence on the 2015 european heatwave. *Environmental Research Letters*, **14**, <https://doi.org/10.1088/1748-9326/ab4d33>.
- Mehta, P., M. Bukov, C.-H. Wang, A. G. Day, C. Richardson, C. K. Fisher, and D. J. Schwab, 2019: A high-bias, low-variance introduction to machine learning for physicists. *Physics Reports*, **810**, 1–124, <https://doi.org/10.1016/j.physrep.2019.03.001>.
- Menze, B. H., B. M. Kelm, R. Masuch, U. Himmelreich, P. Bachert, W. Petrich, and F. A. Hamprecht, 2009: A comparison of random forest and its gini importance with standard chemometric methods for the feature selection and classification of spectral data. *BMC Bioinformatics*, **10**, 213, <https://doi.org/10.1186/1471-2105-10-213>.

- Merz, B., and Coauthors, 2020: Impact forecasting to support emergency management of natural hazards. *Reviews of Geophysics*, **58**, 8755–1209, <https://doi.org/10.1029/2020rg000704>.
- Miller, D. E., Z. Wang, B. Li, D. S. Harnos, and T. Ford, 2021: Skillful subseasonal prediction of united states extreme warm days and standardized precipitation index in boreal summer. *Journal of Climate, American Meteorological Society*, **34**, 5887–5898, <https://doi.org/10.1175/jcli-d-20-0878.1>.
- Molteni, F., T. Stockdale, and M. Balmaseda, 2011: The new ecmwf seasonal forecast system (system 4). *ECMWF Technical Memoranda*, **656**, 35, <https://doi.org/10.21957/4nery093i>.
- Mueller, B., and S. I. Seneviratne, 2012: Hot days induced by precipitation deficits at the global scale. *Proceedings of the National Academy of Sciences*, **109**, 12 398–12 403, <https://doi.org/10.1073/pnas.1204330109>.
- Murphy, A. H., 1993: What is a good forecast? an essay on the nature of goodness in weather forecasting. *Wea. Forecasting*, **8**, 281–293, [https://doi.org/10.1175/1520-0434\(1993\)008<0281:wiagfa>2.0.co;2](https://doi.org/10.1175/1520-0434(1993)008<0281:wiagfa>2.0.co;2).
- Muñoz-Sabater, J., and Coauthors, 2021: Era5-land: a state-of-the-art global reanalysis dataset for land applications. *Earth System Science Data*, **13**, 4349–4383, <https://doi.org/10.5194/essd-13-4349-2021>.
- Nembrini, S., I. R. König, and M. N. Wright, 2018: The revival of the gini importance? *Bioinformatics*, **34**, 3711–3718, <https://doi.org/10.1093/bioinformatics/bty373>, URL <https://repository.publisso.de/resource/fri:6411640/data>.
- Oliveira, J. C., E. Zorita, V. Koul, T. Ludwig, and J. Baehr, 2020: Forecast opportunities for european summer climate ensemble predictions using self-organising maps. *Proceedings of the 10th International Conference on Climate Informatics*, 67–71, <https://doi.org/10.1145/3429309.3429319>.
- Ossó, A., R. Sutton, L. Shaffrey, and B. Dong, 2020: Development, amplification, and decay of atlantic/european summer weather patterns linked to spring north atlantic sea surface temperatures. *J. Climate*, **33**, 5939–5951, <https://doi.org/10.1175/JCLI-D-19-0613.1>.

- Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, **12**, 2825–2830, <https://doi.org/arXiv:1201.0490>, URL <https://scikit-learn.org/stable/>.
- Perkins, S. E., 2015: A review on the scientific understanding of heatwaves -their measurement, driving mechanisms, and changes at the global scale. *Atmospheric Research*, **164–165**, 242–267, <https://doi.org/10.1016/j.atmosres.2015.05.014>.
- Perkins, S. E., and L. V. Alexander, 2013: On the measurement of heat waves. *J. Climate*, **26**, 4500–4517, <https://doi.org/10.1175/jcli-d-12-00383.1>.
- Perkins-Kirkpatrick, S. E., and S. C. Lewis, 2020: Increasing trends in regional heatwaves. *Nature Communications*, **11**, <https://doi.org/10.1038/s41467-020-16970-7>.
- Pyrina, M., M. Nonnenmacher, S. Wagner, and E. Zorita, 2021: Statistical seasonal prediction of european summer mean temperature using observational, reanalysis and satellite data. *Weather and Forecasting*, **36**, <https://doi.org/10.1175/waf-d-20-0235.1>.
- Rasp, S., and N. Thuerey, 2021: Datadriven mediumrange weather prediction with a resnet pretrained on climate simulations: A new model for weatherbench. *Journal of Advances in Modeling Earth Systems*, **13**, 1942–2466, <https://doi.org/10.1029/2020ms002405>.
- Rayner, N. A., D. E. Parker, E. B. Horton, C. K. Folland, L. V. Alexander, D. P. Rowell, E. C. Kent, and A. Kaplan, 2003: Global analyses of sea surface temperature, sea ice, and night marine air temperature since the late nineteenth century. *Journal of Geophysical Research*, **108**, 148–227, <https://doi.org/10.1029/2002jd002670>.
- Reichstein, M., G. Camps-Valls, B. Stevens, M. Jung, J. Denzler, and N. Carvalhais, 2019: Deep learning and process understanding for data-driven earth system science. *Nature*, **566**, 195–204, <https://doi.org/10.1038/s41586-019-0912-1>.
- Robertson, A. W., A. Kumar, M. Peña, and F. Vitart, 2015: Improving and promoting sub-seasonal to seasonal prediction. *Bull. Amer. Meteor. Soc.*, **96**, ES49–ES53, <https://doi.org/10.1175/bams-d-14-00139.1>.

- Rudin, C., 2019: Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, **1**, 206–215, <https://doi.org/10.1038/s42256-019-0048-x>.
- Runge, J., P. Nowack, M. Kretschmer, S. Flaxman, and D. Sejdinovic, 2019: Detecting and quantifying causal associations in large non-linear time series datasets. *Science Advances*, **5**, eaau4996, <https://doi.org/10.1126/sciadv.aau4996>.
- Schwingshackl, C., M. Hirschi, and S. I. Seneviratne, 2017: Quantifying spatiotemporal variations of soil moisture control on surface energy balance and near-surface air temperature. *Journal of Climate*, **30**, 7105–7124, <https://doi.org/10.1175/jcli-d-16-0727.1>.
- Seneviratne, S. I., T. Corti, E. L. Davin, M. Hirschi, E. B. Jaeger, I. Lehner, B. Orlowsky, and A. J. Teuling, 2010: Investigating soil moisture-climate interactions in a changing climate: a review. *Earth-Science Reviews*, **99**, 125–161, <https://doi.org/10.1016/j.earscirev.2010.02.004>.
- Seneviratne, S. I., M. G. Donat, B. Mueller, and L. V. Alexander, 2014: No pause in the increase of hot temperature extremes. *Nature Climate Change*, **4**, 161–163, <https://doi.org/10.1038/nclimate2145>.
- Smola, A. J., P. L. Bartlett, B. Schölkopf, and D. Schuurmans, 2000: *Probabilities for SVM Machines*, 61–75. *Advances in Large Margin Classifiers*, The MIT Press.
- Sobhani, N., D. del Vento, and A. Fanfarillo, 2018: Long-lead forecast of heatwaves in the eastern united states using artificial intelligence. *Proceedings of the Amer. Geophysical Union, Fall Meeting 2018*.
- Spensberger, C., and Coauthors, 2020: Dynamics of concurrent and sequential central european and scandinavian heatwaves. *Quarterly Journal of the Royal Meteorological Society*, **146**, 2998–3013, <https://doi.org/10.1002/qj.3822>.
- Steyerberg, E. W., A. J. Vickers, N. R. Cook, T. Gerds, M. Gonen, N. Obuchowski, M. J. Pencina, and M. W. Kattan, 2010: Assessing the performance of prediction models. *Epidemiology*, **21**, 128–138, <https://doi.org/10.1097/ede.0b013e3181c30fb2>.
- Storch, H. V., and F. W. Zwiers, 2003: *Statistical analysis in climate research*. Cambridge University Press, 293–299 pp.

- Suarez-Gutierrez, L., W. A. Mueller, C. Li, and J. Marotzke, 2020: Dynamical and thermodynamical drivers of variability in european summer heat extremes. *Climate Dynamics*, **54**, 4351–4366, <https://doi.org/10.1007/s00382-020-05233-2>.
- Vabalas, A., E. Gowen, E. Poliakoff, and A. J. Casson, 2019: Machine learning algorithm validation with a limited sample size. *PLOS ONE*, **14**, <https://doi.org/10.1371/journal.pone.0224365>.
- van Straaten, C., K. Whan, D. Coumou, B. van den Hurk, and M. Schmeits, 2022: Using explainable machine learning forecasts to discover sub-seasonal drivers of high summer temperatures in western and central europe. *Mon. Wea. Rev.*, <https://doi.org/10.1175/mwr-d-21-0201.1>.
- Vijverberg, S., M. Schmeits, K. van der Wiel, and D. Coumou, 2020: Subseasonal statistical forecasts of eastern u.s. hot temperature events. *Mon. Wea. Rev.*, **148**, 4799–4822, <https://doi.org/10.1175/mwr-d-19-0409.1>.
- Vitart, F., 2014: Evolution of ecmwf sub-seasonal forecast skill scores. *Quarterly Journal of the Royal Meteorological Society*, **140**, 1889–1899, <https://doi.org/10.1002/qj.2256>.
- Wallemacq, P., R. Below, and D. McClean, 2018: Economic losses, poverty and disasters (1998–2017). URL <https://www.undrr.org/publication/economic-losses-poverty-disasters-1998-2017>, 1–9 pp.
- Weyn, J. A., D. R. Durran, and R. Caruana, 2019: Can machines learn to predict weather? using deep learning to predict gridded 500hpa geopotential height from historical weather data. *Journal of Advances in Modeling Earth Systems*, **11**, 2680–2693, <https://doi.org/10.1029/2019ms001705>.
- Wheeler, M. C., H. Zhu, A. H. Sobel, D. Hudson, and F. Vitart, 2016: Seamless precipitation prediction skill comparison between two global models. *Quarterly Journal of the Royal Meteorological Society*, **143**, 374–383, <https://doi.org/10.1002/qj.2928>.
- White, C. J., and Coauthors, 2017: Potential applications of subseasonal-to-seasonal (s2s) predictions. *Meteorological Applications*, **24**, 315–325, <https://doi.org/10.1002/met.1654>.
- White, C. J., and Coauthors, 2021: Advances in the application and utility of subseasonal-to-seasonal predictions. *Bull. Amer. Meteor. Soc.*, **aop**, 1–57, <https://doi.org/10.1175/bams-d-20-0224.1>.

- Wilks, D. S., 2019: *Statistical Methods in the Atmospheric Sciences*. 4th ed., Elsevier, 379 and 386–388 (Chapter 9) pp.
- Wulff, C. O., and D. I. V. Domeisen, 2019: Higher subseasonal predictability of extreme hot european summer temperatures as compared to average summers. *Geophysical Research Letters*, **46**, 11 520–11 529, <https://doi.org/10.1029/2019gl084314>.
- Wulff, C. O., R. J. Greatbatch, D. I. V. Domeisen, G. Gollan, and F. Hansen, 2017: Tropical forcing of the summer east atlantic pattern. *Geophysical Research Letters*, **44**, 94–8276, <https://doi.org/10.1002/2017gl075493>.
- Wunderlich, R. F., Y.-P. Lin, J. Anthony, and J. R. Petway, 2019: Two alternative evaluation metrics to replace the true skill statistic in the assessment of species distribution models. *Nature Conservation*, **35**, 97–116, <https://doi.org/10.3897/natureconservation.35.33918>.
- Zheng, X., and C. S. Frederiksen, 2007: Statistical prediction of seasonal mean southern hemisphere 500-hpa geopotential heights. *Journal of Climate*, **20**, 2791–2809, <https://doi.org/10.1175/jcli4180.1>.
- Zuo, J., H.-L. Ren, J. Wu, Y. Nie, and Q. Li, 2016: Subseasonal variability and predictability of the arctic oscillation/north atlantic oscillation in bcc_agcm2.2. *Dynamics of Atmospheres and Oceans*, **75**, 33–45, <https://doi.org/10.1016/j.dynatmoce.2016.05.002>.