



## ***Interactive comment on “Improving snowfall representation in climate simulations via statistical models informed by air temperature and total precipitation” by Flavio Maria Emanuele Pons and Davide Faranda***

**Flavio Maria Emanuele Pons and Davide Faranda**

flavio.pons@lsce.ipsl.fr

Received and published: 14 March 2021

The manuscript by Pons and Farada assesses the performance of several snowfall separation methods to reproduce simulated snowfall in the ERA5 reanalysis on a European scale by taking into account simulated near-surface air temperature and total precipitation at daily resolution. The two best-performing methods are in a second stage applied to bias-adjusted output of the IPSL-WRF regional climate model (historical period) to obtain a bias-adjusted estimate of simulated snowfall in the RCM. The

C1

evaluation reveals a satisfying representation of the PDF of the daily ERA5 reference snowfall amount in the historical period by the bias-adjusted and separated IPSL-WRF simulation. Overall, the paper fits well into the journal's scope. Data and methods are for most parts clearly introduced and explained. The presentation of the results has some weaknesses but is still acceptable. The major drawback of the work, however, is the unclear relevance of the work for a broader audience and for RCM snowfall bias-adjustment. Essentially, the authors search for a method to emulate the ERA5 micro-physics scheme that simulates the actual snowfall flux in the reanalysis model taking into account simulated near-surface temperature and simulated total precipitation only. The two best performing methods are then applied to a different model (IPSL-WRF) to separate snowfall from total precipitation after bias-adjustment of simulated temperature and precipitation. Results look satisfying, but there is

Q: (1) no evaluation of the ERA5 snowfall flux (which is the basic reference in the entire work, and the entire analysis is geared towards a reproduction of ERA5-simulated snowfall flux; the paper frequently uses the term "observed" for ERA5 snowfall flux, although it is essentially a simulated flux probably subject to systematic biases)

A: We agree that the word "observation" is incorrectly used to describe ERA5. We corrected all sentences containing such inaccuracy, and we will instead use the term "reanalysis" or the expression "reference dataset".

It is indeed possible that ERA5 presents some biases compared to observations, even though it is not in the scope of our paper to evaluate the accuracy of ERA5 with respect to direct measurement.

In general, the bias correction of climate projection models with respect to observations or reanalysis is a well established practice. Reanalysis datasets such as ERA5, ERA-interim or NCEP are often considered as reference datasets in this context, even if it is known and accepted that they have limitations, as these are generally balanced by the advantages.

C2

Q: (2) no analysis to what extent the satisfying results of the application of the method to the RCM are specific for the chosen RCM and the bias-adjustment method of temperature and precipitation that was carried out beforehand (a different RCM might, even after bias-adjustment, have a completely different multivariate structure of daily temperature and precipitation, at least a structure that is different to ERA5, and the method might not hold in these cases)

A: We agree about the fact that our results cannot be generalized to different types of BC used to adjust the RCM. We remark that BC is a computationally expensive and very time consuming operation, and very rarely one can try and compare several different types of BC in a climate study. In our specific case, we considered a model bias adjusted with univariate CDFt, as mentioned at lines 154-7. This method has been widely applied and validated, and it has been used to prepare the datasets constituting the CORDEX-Adjust project, from which we downloaded the already bias adjusted output.

Our method requires the use of an effective BC method for temperature and precipitation beforehand, in the same way it would require well calibrated measuring stations if we were dealing with in-situ observations. Unfortunately, while official measuring stations are regulated by WMO standards, there is not a BC method considered a universal standard. Exploring several BC methods in this study, their multivariate performance and its impact on catching the microphysics over the reference period would be a very heavy task which goes beyond our objectives, and it would produce an extremely large amount of supplementary data.

As a further consideration, we remark that all of the methods representing snowfall with the same philosophy require the knowledge of temperature and precipitation, so the same objection should be true for all the empirical methods already existing in the literature and cited or even put to the test in the present paper.

Overall, we consider it impossible to evaluate several BC methods and their multivariate

C3

impact as a part of our study, as this step alone would constitute a completely different (and probably larger) paper. In a similar way, also testing a model ensemble would be beyond the scopes of this paper, and would still not be exhaustive. For example, if we considered the entire CORDEX ensemble, this would still not make the results directly transferable to CMIP5 models.

We will add a “Limitations” paragraph to the Conclusions section concerning issues with datasets and BC methodology:

“Limitations

We also clarify some of the limitations of our analysis. The nature of climate datasets makes multiple comparisons among models and BC techniques very demanding in terms of data storage and computational time. For this reason, we limited our analysis to one reanalysis dataset (ERA5), one marginal bias correction technique (CDF-t), and one climate projection model (IPSL\$\\\_WRF).

We do not consider the choice of ERA5 problematic with respect to other gridded datasets that could be observational (e.g. E-OBSv20) or other reanalysis (e.g. NCEP/NCAR): while the actual values could change between datasets, we do not foresee this affecting directly the performance of the methodology we presented in terms of improvement of raw simulations respect the chosen reference dataset.

On the other hand, the choice of the BC may influence the outcome of our modelling procedure. The CDF-t is applied marginally to each variable, so that there is no guarantee that the inter-variable correlations are correctly reproduced in the target climate simulation. Indeed, \\cite{meyer2019effects} showed that applying multivariate as opposed to univariate BC produces significant changes in estimated snow accumulation, stressing the importance of modelling the interdependence between precipitation and air temperature in hydrological studies focused on snowy areas. The choice of the BC, in general, should be tuned on the trade-off between complexity and need for controlling specific features, in this case inter-variable correlation. In our case, we con-

C4

sidered a climate dataset prepared in the context of the CORDEX-Adjust project, which is made available already adjusted with respect to ERA5 using marginal CDF-t. Our results show an improvement in snowfall representation even relying on marginal BC; however, we stress that the methodology should be validated again if used on datasets prepared with different BC techniques, to assess whether this difference affects the predicting performance of the model.

On the same note, we remark that prediction accuracy may vary across different climate models, due to the different physical approximations and parameterizations, which are likely to affect the relationship between near-surface temperature and precipitation. Due to these differences, even other RGMs from the EURO-CORDEX project may exhibit variability in the performance of the snowfall reconstruction. This holds true for all statistical models cited in Section \ref{int1}, as it is rarely the case that snowfall reconstruction techniques are tested over an ensemble of different climate models. Once more, we underline the importance of assessing the performance of the chosen methodology to approximate snow (or compare several of them) by validating it on the historical period of the available models in reference to the available reanalysis/observation dataset."

Q: (3) no discussion of potential problems with inter-variable dependencies even after bias-adjustment of an RCM (-> see, for instance, Meyer et al., HESS, 2019 Effects of univariate and multivariate bias correction on hydrological impact projections in alpine catchments)

A: We thank the reviewer for this suggestion, we included this in the Limitation section mentioned above.

Q: (4) no indication if the identified methods will also produce robust snowfall estimates in a future climate change scenario (which is, as far as I can guess, the basic motivation of the entire work -> a possibility to investigate such an applicability would be to split the ERA5 period into "cold" and "warm" years and to calibrate on the cold and validate

C5

on the warm sample)

A: We thank the reviewer for suggesting such an interesting and yet feasible addition to our validation procedure. We performed the suggested experiment and summarized the results in Fig\_warmcold.pdf (see attachments). We will add the following paragraph to the manuscript:

"\Subsection{Robustness to climate change}

As an additional element to evaluate the performance of the identified methods, we assess if they can produce robust snowfall estimates in a climate change scenario. In order to do so, we repeat the validation procedure described in Section \ref{design}; after ordering the ERA5 dataset based on the annual DJF average temperature, we take the coldest 25% as the train set and the warmest 25% as the test set. We run this procedure for the two best performing models, the segmented linear regression and the cubic spline regression.

Fig. \ref{warmcold} shows the model performance metrics in analogy with Figures \ref{boxplots} and \ref{boxcor}. Panel (a) and (b) display the map of the event-to-event correlation coefficient, showing overall higher value than for the random train and test sets. The two models also perform much more similarly in terms of correlation than in the random sets case, as it can be seen from the boxplots in panel (c); the performance in terms of RMSE and MAE is also comparable (panel (d)), as it was in the overall validation presented before. Overall, assuming that separating cold and warm year can be a proxy of climate change to assess model performance, the two technique perform very similarly to the general case in terms of forecasting error, without any visible improvement or accuracy decrease. However, we observe an improvement in the correlation between predicted and true forecasting values: we argue that this effect is likely due to precipitation patterns in years characterized by extreme temperatures in the historical period, and it should not be expected to happen under future climate change."

C6

Q: (5) no reference to differences in spatial resolution of the models employed and the fact the subgrid orography can actually have a considerable influence on simulated snowfall (or, the other way round, neglecting subgrid scale orographic variability in model bias-adjustment could result in false derived snowfall sums)

A: Concerning the difference in spatial resolution of the adopted models, all the dataset we consider have the same lon-lat grid with  $0.25^\circ$  resolution, so there are no differences in resolution to be considered in the data we used.

It is plausible that the performance of the algorithm could change if we considered datasets with a resolution sensibly different from the one chosen here, for example 1 km or 100 km. However, comparable methods, including the ones cited here, are applied to anything from a single station time series to gridded datasets without necessarily exploring all possible scales.

The reviewer also underlines that “subgrid orography can actually have a considerable influence on simulated snowfall” and that “neglecting subgrid scale orographic variability in model bias-adjustment could result in false derived snowfall sums”. We are aware of the limitation of neglecting subgrid scales, but this is something we always have to live with when dealing with climate simulation models. Considering the lack of scale separation in the atmosphere, the correct description of any phenomenon would benefit from including more fine scales, but sometimes this is not possible. Indeed, Frei et al. 2018 underline that the choice of more complex functional form (which we replace with segmented and spline regression) instead of the binary threshold separation is made precisely to the purpose to approximate subgrid effects.

Q: (6) no analysis of a calibrated threshold within the "naive" STM method (which I assume could yield even better results than the two best identified methods, as even the performance with a fixed  $2^\circ\text{C}$  threshold is very close to the two best-performing methods)

A: Indeed, we chose not to try a finer calibration of the threshold for the STM model.

C7

This is because, dealing with a high number of grid points, such calibration would still require an explorative technique such as a breakpoint search. However, the STM model with threshold determined in such a way, would correspond to the case of our spline regression, but constraining the number of thresholds at 1 and using 0th-order splines. Since admitting up to two thresholds and cubic splines is not more complicated or sensibly more time demanding, we did not think it is worth to add these constraints: where the optimal model would be a binary threshold, our algorithm can still reproduce that feature while being more versatile if more complex parameterizations are needed.

We do not agree that the STM result is “very close” to the best performing model: it is somehow halfway in terms of average values (both for error measures and correlation) but showing a high variability. This is likely due to the fixed threshold, but as we explained above, it makes no sense from a practical point of view to test this method with more complex thresholds. In fact, such a method is used (for example, in Frei et al. 2018; Bai et al 2019) in the literature taking  $2^\circ\text{C}$  as an accepted, overall well working typical threshold. As we mentioned in the manuscript at lines 205 and following, a sensitivity analysis over Europe has been conducted, for example, in Faranda 2020, finding that thresholds varying between 0 and  $2.5^\circ\text{C}$  produce rather comparable results.

Q: (7) no analysis of the importance of variations on the sub-daily scale which might be important for daily snowfall sums.

A: We agree that including small scales in space and time would improve the representation of snow, but we stress again that it is not typical to deal with sub-daily datasets of long term climate projections. Even gridded observations datasets, such as E-OBS, are provided at the daily frequency, making such an evaluation de facto impossible (see e.g. Bai et al. 2019).

Q: The main message of the manuscript is currently, that for this specific setup (this specific RCM, this specific bias-adjustment method, this specific reference snowfall),

C8

the two identified methods if applied to bias-adjusted IPSL-WRF temperature and precipitation output can yield a representation of snowfall that well reproduces the ERA5 reference snowfall. These results are in my opinion not per se transferable to different models or to a future climate scenario or to a different reference snowfall (especially not to a true observation-based snowfall estimate). As such, the value of the work is limited for the time being in my opinion and not too informative for a broader readership. I would hence recommend to return the manuscript to the authors for major revisions. During these revisions, the mentioned points should be picked up in order to increase the relevance of the work. A couple of further issues are mentioned below. With kind regards.

A: As already mentioned in response to Reviewer's point 2, these issues are now explicitly mentioned in the Limitations paragraph in the Discussion section. We stress that we agree that all these limitations exist, but we find that it would be hardly feasible to test one statistical method by varying: Time resolution Space resolution Reference dataset Bias correction technique Physics and numerical climate model schemes

To our knowledge, such a broad validation does not exist even for simple and well established models (e.g. single threshold binary separation). As far as we can tell, most of the literature dealing with this type of statistical model for precipitation phase apportionment are validated on datasets similar to the ones we considered in this paper. We think that the Limitations section should make it clear that the results shown in the paper should be considered specific to our setting, and that a validation of the model should be performed, if this is applied to different datasets.

FURTHER ISSUES:

Q: Line 24: Very unclear what is meant.

A: We agree that this sentence was not clear. We replaced it with a simple example of a case where the binary apportionment could produce a severe bias:

C9

"Even though such binary separation of snowfall using a temperature threshold seemed a good option to retrieve snowfall data from E-OBSv20.0e, it has obvious limitations: for example, in an event characterized by abundant precipitation but a temperature associated to a roughly 50% snow fraction, snowfall would be either strongly under- or overestimated."

Q: Lines 35-36. Also rather unclear.

A: We agree about lack of clarity and we also realized this sentence was quite redundant. We will remove lines 35-38 and change the next sentence to better match the previous paragraph to:

"In order to mitigate the aforementioned biases, a BC step is usually performed. This step usually consists of a methodology designed to adjust specific statistical properties of the simulated climate variables towards a validated reference dataset in the historical period. [...]"

Q: Lines 38-40: This is actually not true, the entire set of so-called "perfect prognosis" downscaling methods is ignored here. These do not adjust the simulated variables towards observations but exploit calibrated relationships between observed (or reanalysis-simulated) large scales and observed local scales.

A: In agreement with the existing literature, we consider perfect prognosis downscaling as a part of statistical downscaling (see, e.g. Soares et al., 2019)

Q: Lines 141-142: Very unclear.

A: We will change this sentence to:

"This quantity is relevant for hydrologists, being closely related to runoff and river discharge, but also for climatologists, since it well represents the intensity of the phenomenon while, however, we remind that snowfall is not a measure of accumulation of snow on the ground."

C10

Q: Line 145: Above (line 132) you mention that only daily data are used, here you obviously employ hourly data. Please clarify.

A: We better clarify changing the sentence at line 145 to:

The initial ERA5 dataset is available at hourly frequency, while the IPSL\$\_\$WRF is available at daily frequency. Since the two time steps are different, and we have no way to disentangle a daily quantity into the sub-daily cycle, we aggregate the hourly ERA5 data into daily.

Q: Line 151: "grid step" unclear

A: We will change the sentence to: In particular, we consider DJF data from climate simulations of the historical period 1979-2005 over the same domain and at the same spatial resolution as the reanalysis dataset described in Section \ref{era5mod}.

Q: Line 181: Rather unclear what is meant by "standardized temperature anomalies" and why these are used.

A: We now specify the standardization procedure as follows: "In all the regression models discussed in the following, but not in the STM, we use as independent variable the standardized temperature anomalies, obtained by subtracting the historical mean and dividing by the historical standard deviation."

Using standardized anomalies is quite a common practice in climate studies, where different variables span over very different scales (e.g. precipitation is in average about 0.1 m/day, absolute temperature is of order 250-300 K, geopotential height 5000 m). For transparency, we report to the reviewer that, in particular, we had standardized the variables as we tried to add total precipitation as a covariate, to consider the possible influence of intense precipitation on the snow fraction. Given the very different scale between the two covariates (temperature and precipitation) we standardized the variables in the various model specifications. Improvement obtained by adding precipitation were so unremarkable (practically non existent) that we did not even mention

C11

them in the manuscript, and kept the results of the models based on temperature with standardized variables.

We do not foresee affecting our results and it is quite a common practice in regression modelling even outside climatology. As we already stressed at line 182, it is important to do the standardization in the same way to reproduce the result.

Q: Chapter 3.1.2: This sub-chapter contains a large amount of rather technical information, which is appreciated, but which should be moved to some technical appendix I believe.

A: We agree and we will move this subsection to a technical appendix.

Q: Line 517: Do you have any explanation for these rather low calibrated thresholds? Is there a relation to orographic height, for instance?

A: We do not have an explanation for this. We noticed that such low thresholds appear in areas where we would expect winter precipitation to be mainly snow regardless of the specific daily temperature, given the cold continental or subarctic climate. The search algorithm is not meant to necessarily find physically meaningful values, so it is possible that it finds thresholds that improve the performance only very slightly, while even a regression without any threshold or even a binary apportionment could perform relatively well. In this sense, the threshold recovered after the recomputation via the segmented regression algorithm make more sense (temperature over continental areas seem positive, but here please consider that we are looking at anomalies).

Q: Line 526: Should be "Fig. 2" instead of "Fig. 1".

A: Thank you, we will correct accordingly

Q: Lines 688-689: Very unclear.

A: We will change this sentence to:

Since comparing IPSL\$\_\$WRF and its adjusted versions to ERA5 does not provide

C12

a one-on-one correspondence between snowfall events, it is not possible to compute correlation coefficients between reanalysis and model snowfall time series at each grid point as in Fig. \ref{boxcor}. Instead, we can study the correlation between the total 1979-2005 ERA5 snowfall and the total 1979-2005 snowfall simulated by IPSL\$\_{WRF}\$ and approximated with segmented logit and cubic spline regression at each grid point. Fig. \ref{alps\_dist} (a) shows the scatterplots of total IPSL\$\_{WRF}\$, logit segmented regression and cubic spline regression snowfall against total ERA5 snowfall for the grid points in the Alps region.

Q: Lines 707-708: Better representation of the tails is not really apparent from the figure I'd say.

A: Even though we are aware that the definition of "tail" is somehow arbitrary, we remark that when the ERA5 distribution hits the 0.95 mark, the IPSL-WRF distribution function is barely above 0.85, and when ERA5 reaches 0.99 IPSL-WRF is around the 0.95. On the other hand, the two statistical models are practically non distinguishable from ERA5, considering that this holds true for values above the 95th percentile, we think that the improvement in the tail is rather solid.

Q: Figure 1: Color scale is not very intuitive.

A: We changed the palettes to more traditional ones. The initial choice of a palette alternating different color was due to the fact that the total snowfall spans several orders of magnitude between the most and least snowy locations on the map, so that less contrasting color scales tend to flatten the variability.

Q: Figures 2 and 3: Bad color scale: White color means threshold temperatures around 0°C but also "not applicable". I'd suggest to modify the color scale.

A: We agree. We will change the color associated to 0°C to gray, to make the figures clearer.

Q: Figure 4: Legend too large. Also, the methods are named differently compared to

C13

Table1 and are sorted in a different order. Please harmonize. Also, it would be good to use the same unit in the lower panel as in Table 1 (10<sup>3</sup>mm<sup>3</sup>)

A: We proceeded to make these modifications to the figure.

Q: Figure 5: Upper panel: Please use the same sorting of methods as in Table 1.

A: We proceeded to make this modification to the figure.

Q: Figure 6: Very bad color scale, not at all intuitive. Also, the color scale should be identical for all panels to enable a comparison (same color should mean the same value in all panels). Is the unit actually m/27 years (1979-2005) or m/year? Please clarify.

A: We proceeded to change the palette with one more traditionally used for anomalies. Unfortunately, setting the colorscale in such a way that the same color has the same value over the plot would completely flatten the aspect of panels b) and c), since values in panel a) can be one order of magnitude larger.

Q: Figure 7: Legend of lower panel too small.

A: We proceeded to make this modification to the figure.

Q: Figure 8: What about the bad-performing grid cell in northern Italy in logit seq and cubic spline? What is happening here?

A: We inspected the specific grid cell in detail to assess the extreme negative value. Indeed, it seems to be one of the few points where the breakpoint search algorithm failed to converge, so that a threshold is not available and the models were then not estimated. When we took sums over time to compute snowfall totals, an 'na.rm = TRUE' option was used in the R script so that instead of an NA the sum in that grid cell resulted equal to 0, and the difference was then the negative ERA5 snowfall total in that cell. We ran the scripts without the NA removal option and produced a new figure with the concerned grid point correctly masked out as NA.

C14

Q: Figures 9 and 11, upper panels: Sorry, but even after reading the explanation several times it is not really clear to me what is displayed here. Also, I'd suggest to use a white background instead of a black background. Lower panels: Please specify the unit of the x-axis

A: We proceeded to make this modification to the figure.

#### References

Soares, P. M., Maraun, D., Brands, S., Jury, M. W., Gutiérrez, J. M., SanáŕMartín, D., ... & ObermannáŕHellhund, A. (2019). Processáŕbased evaluation of the VALUE perfect predictor experiment of statistical downscaling methods. *International Journal of Climatology*, 39(9), 3868-3893.

Bai, L., Shi, C., Shi, Q., Li, L., Wu, J., Yang, Y., ... & Meng, J. (2019). Change in the spatiotemporal pattern of snowfall during the cold season under climate change in a snowáŕdominated region of China. *International Journal of Climatology*, 39(15), 5702-5719.

Interactive comment on Nat. Hazards Earth Syst. Sci. Discuss., <https://doi.org/10.5194/nhess-2020-352>, 2020.

C15

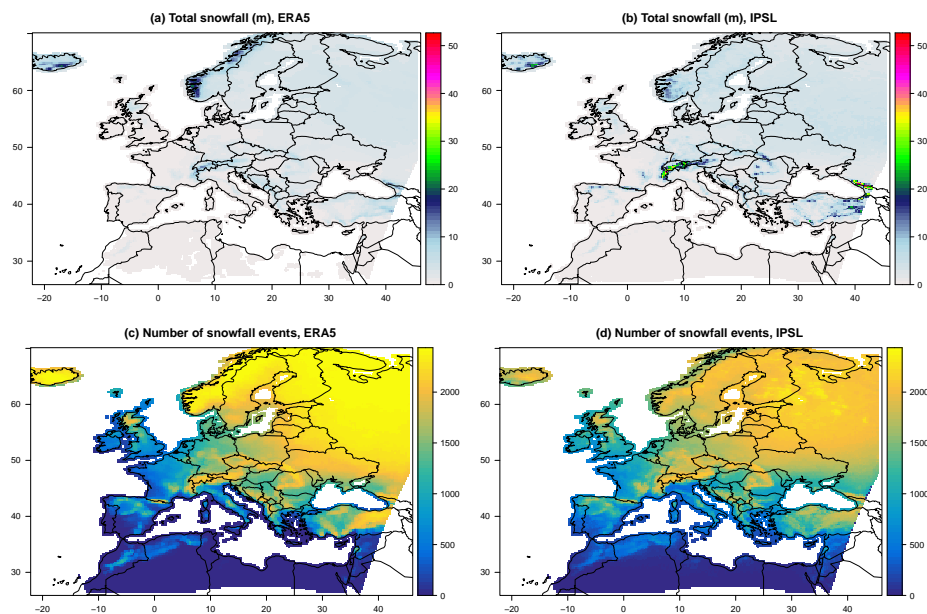


Fig. 1.

C16



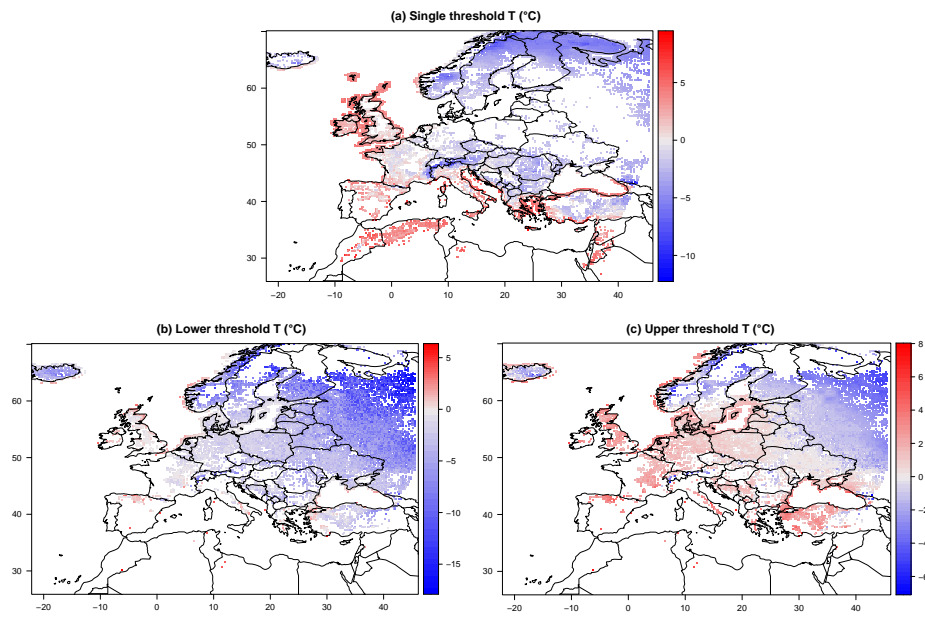


Fig. 2.

C17

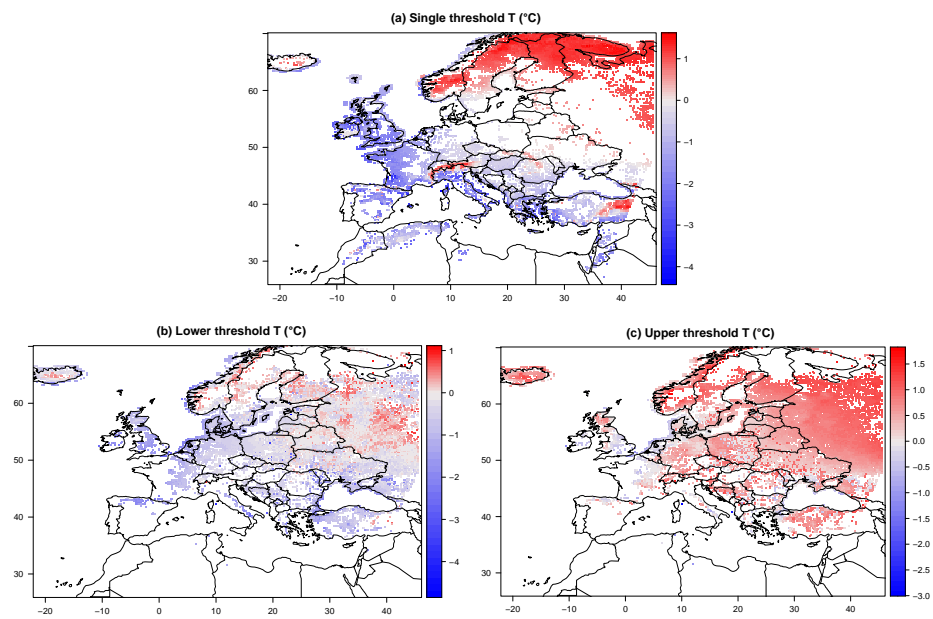


Fig. 3.

C18

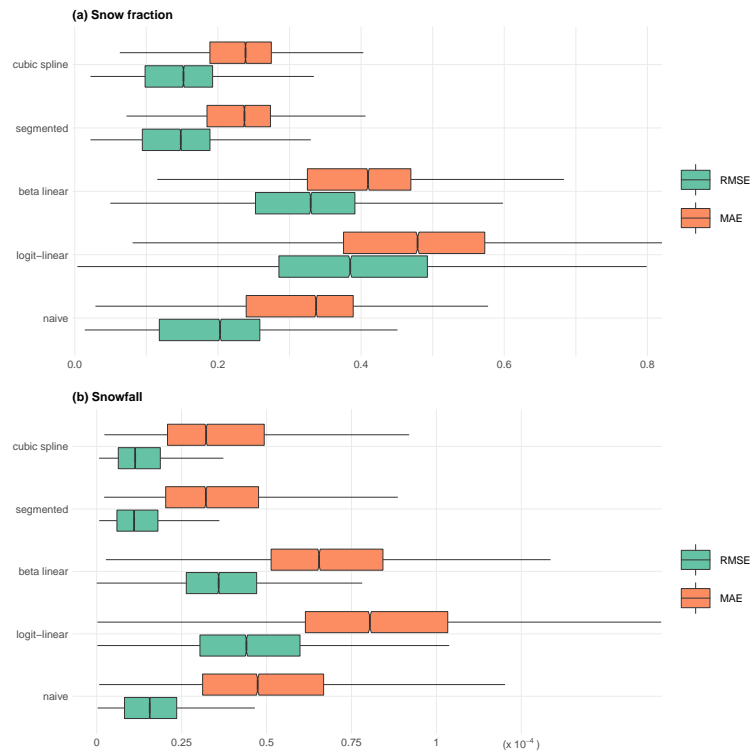


Fig. 4.

C19

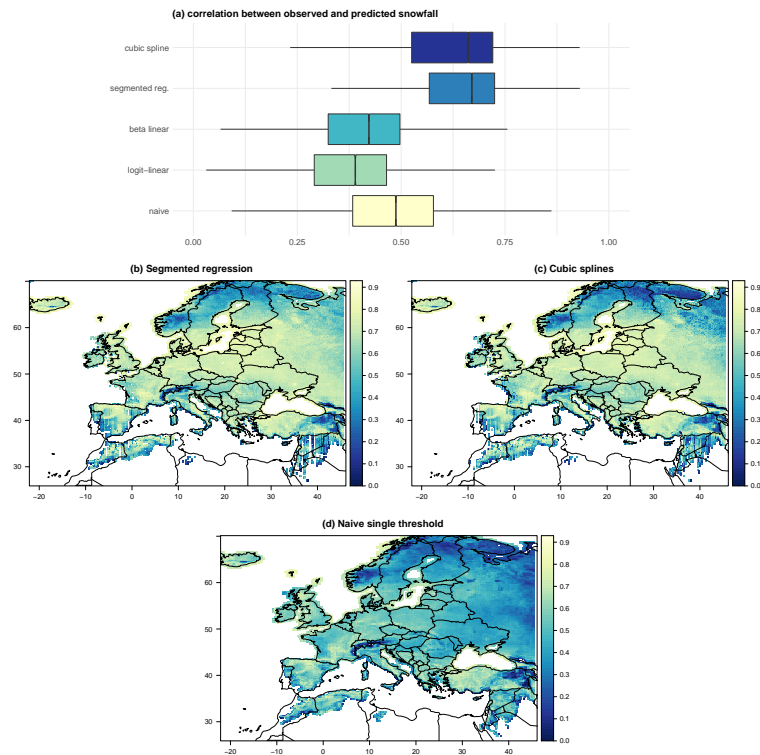


Fig. 5.

C20

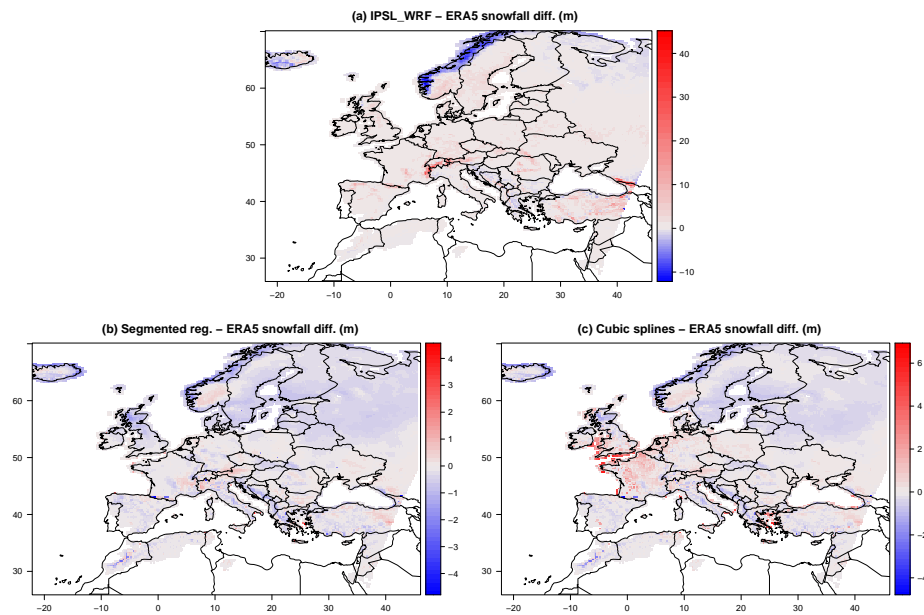


Fig. 6.

C21

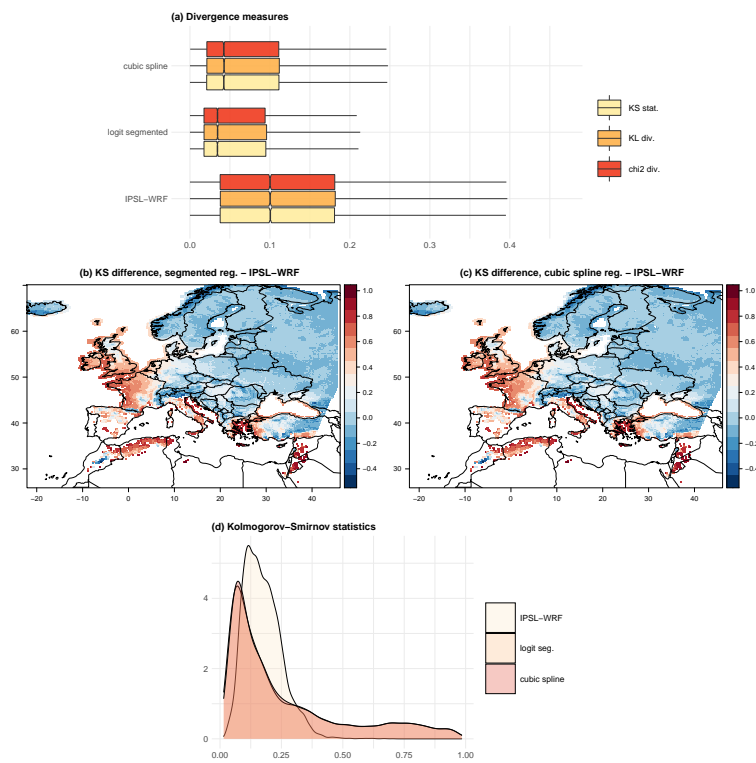


Fig. 7.

C22

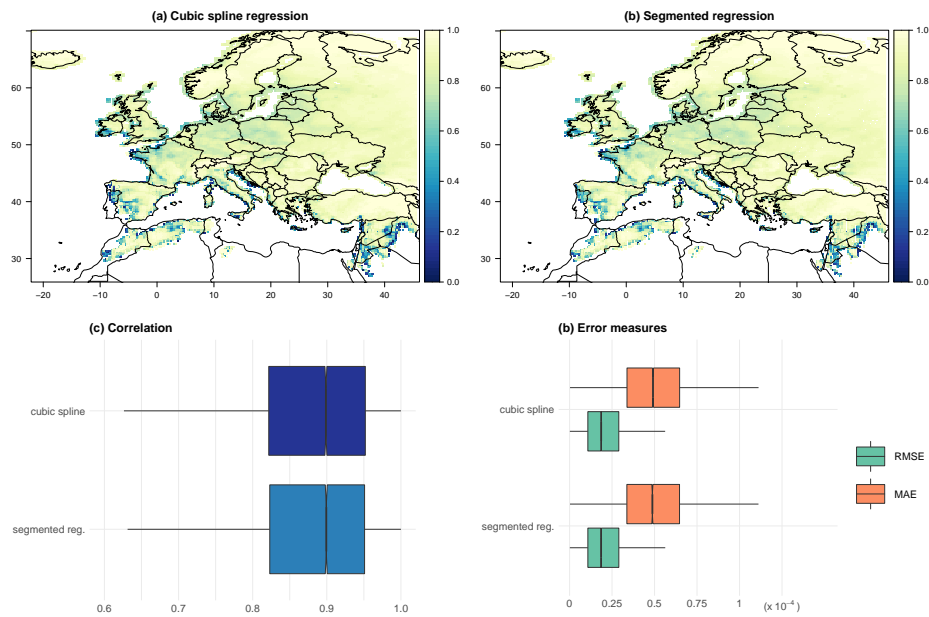


Fig. 8.

C23

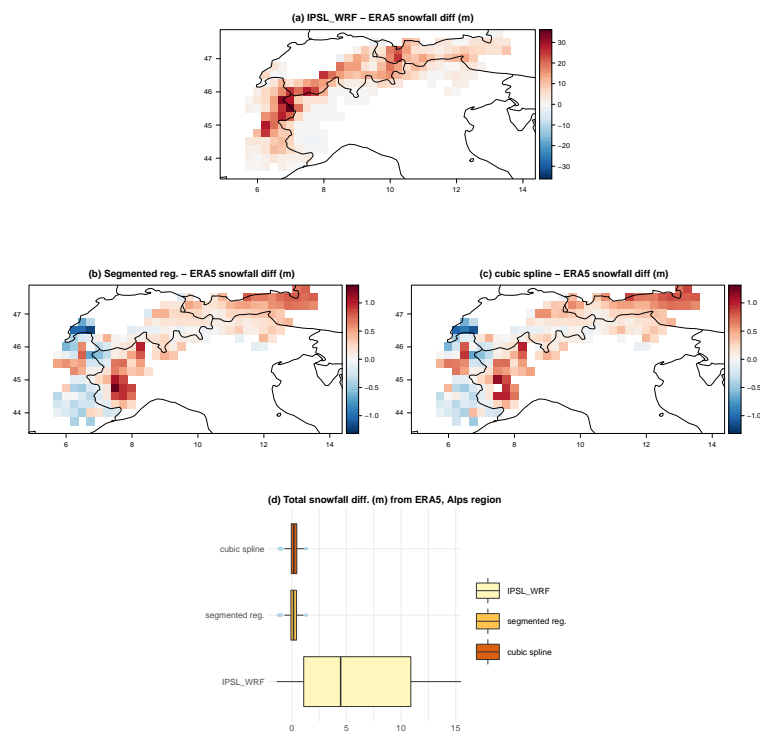


Fig. 9.

C24

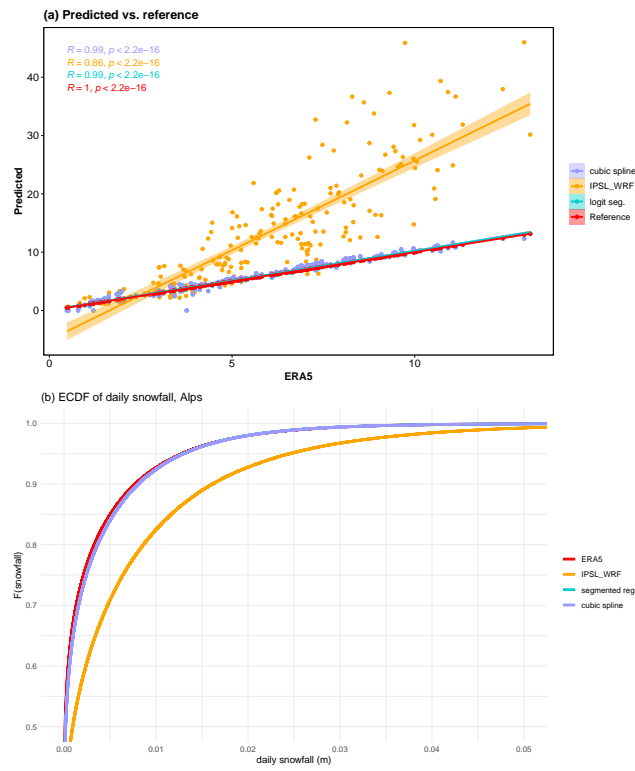


Fig. 10.

C25

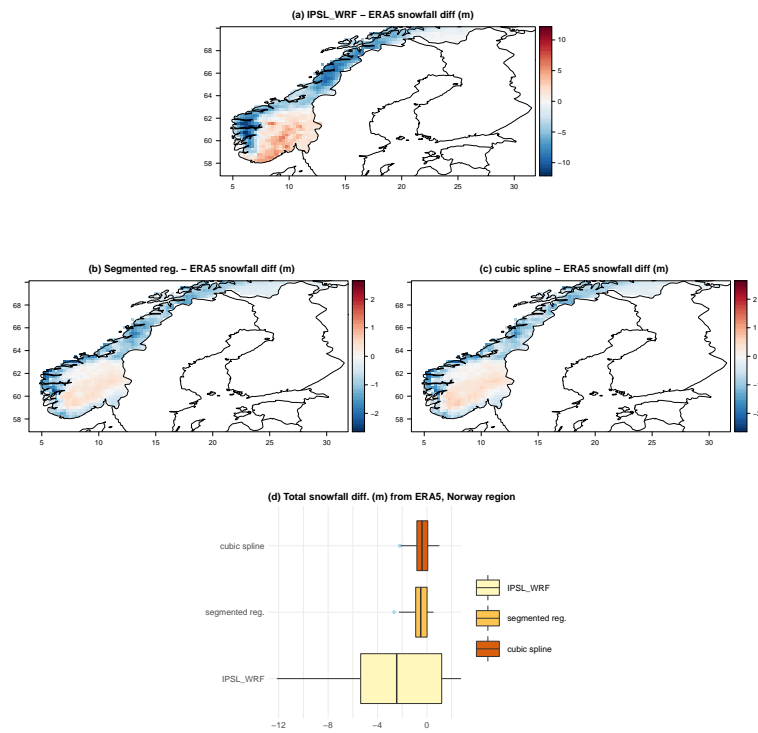


Fig. 11.

C26

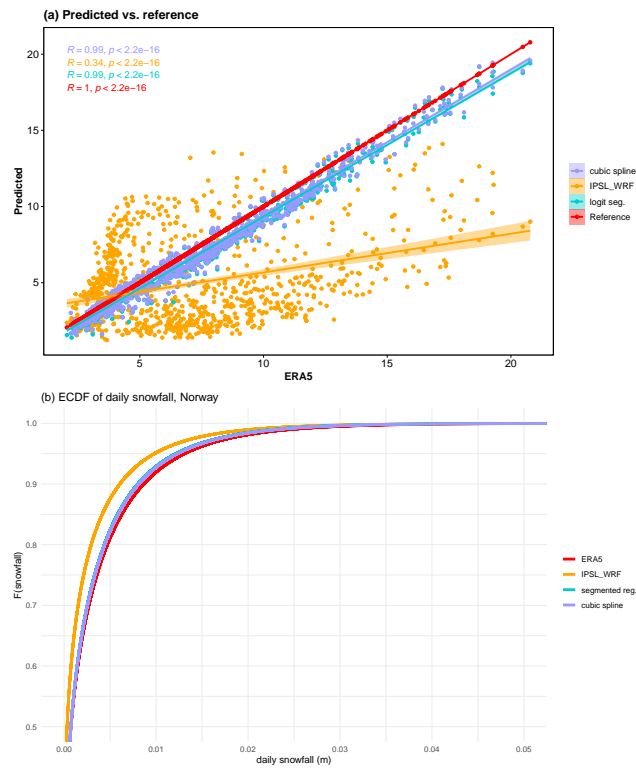


Fig. 12.