

1 **Can bias correction and statistical downscaling**
2 **methods improve the skill of seasonal precipitation**
3 **forecasts?**

4 **R. Manzanas · A. Lucero · A.**
5 **Weisheimer · J. M. Gutiérrez**

6
7 Received: date / Accepted: date

8 **Abstract** Statistical downscaling methods are popular post-processing tools which
9 are widely used in many sectors to adapt the coarse-resolution biased outputs from
10 global climate simulations to the regional-to-local scale typically required by users.
11 They range from simple and pragmatic Bias Correction (BC) methods, which di-
12 rectly adjust the model outputs of interest (e.g. precipitation) according to the
13 available local observations, to more complex Perfect Prognosis (PP) ones, which
14 indirectly derive local predictions (e.g. precipitation) from appropriate upper-air
15 large-scale model variables (predictors). Statistical downscaling methods have been
16 extensively used and critically assessed in climate change applications; however,
17 their advantages and limitations in seasonal forecasting are not well understood
18 yet. In particular, a key problem in this context is whether they serve to improve
19 the forecast quality/skill of raw model outputs beyond the adjustment of their
20 systematic biases.

21 In this paper we analyze this issue by applying two state-of-the-art BC and
22 two PP methods to downscale precipitation from a multimodel seasonal hindcast
23 in a challenging tropical region, the Philippines. To properly assess the potential
24 added value beyond the reduction of model biases, we consider two validation
25 scores which are not sensitive to changes in the mean/variance (correlation and
26 reliability categories). Our results show that, whereas BC methods maintain or
27 worsen the skill of the raw model forecasts, PP methods can yield significant skill
28 improvement (worsening) in cases for which the large-scale predictor variables con-
29 sidered are better (worse) predicted by the model than precipitation. For instance,

R. Manzanas (✉) · J. M. Gutiérrez
Meteorology Group, Institute of Physics of Cantabria (IFCA), CSIC-University of Cantabria,
Santander, 39005, Spain. E-mail: rmanzanas@ifca.unican.es

A. Lucero
Philippine Atmospheric, Geophysical and Astronomical Services Administration (PAGASA).
Quezon City, Philippines.

A. Weisheimer
Department of Physics, National Centre for Atmospheric Science (NCAS), University of Ox-
ford, Oxford OX1 3PU, UK.
European Centre for Medium-Range Weather Forecasts (ECMWF). Reading RG2 9AX, UK.

30 PP methods are found to increase (decrease) model reliability in nearly 40% of
31 the stations considered in summer (autumn). Therefore, the choice of a convenient
32 downscaling approach (either BC or PP) depends on the region and the season.

33 **Keywords** Statistical downscaling, perfect prognosis, bias correction, seasonal
34 forecasting, precipitation, skill, correlation, reliability categories

35 1 Introduction

36 Different Statistical Downscaling (SD) methods have been developed since the
37 early 1990s (see, e.g., von Storch et al, 1993) to bridge the gap between the
38 coarse-resolution biased climate information provided by Global Circulation Mod-
39 els (GCMs) and the regional-to-local scale required in different socio-economic
40 sectors such as hydrology, agriculture, energy, etc. These methods rely on em-
41 pirical/statistical models which link the local observed predictands of interest,
42 here precipitation, with explicative large-scale GCM predictors over the area of
43 interest. These models are first calibrated and tested (i.e., cross-validated) us-
44 ing data from a historical representative period (*training phase*) and subsequently
45 applied to obtain the downscaled local predictions from new GCM predictors (*pre-
46 diction/downscaling phase*). According to the nature of predictors in the training
47 phase, two different approaches for SD exist (see, e.g. Maraun et al, 2010; Gutiérrez
48 et al, 2013a): Perfect Prognosis (PP) and Model Output Statistics (MOS), the lat-
49 ter including the increasingly popular Bias Correction (BC) methods.

50 Under the PP approach, quasi-observed predictors from reanalysis are used
51 to train the statistical models (e.g. regression or analog methods), based on their
52 temporal correspondence with the observed precipitation. Afterwards, the result-
53 ing models are applied to GCM predictor data in the prediction phase. There-
54 fore, variables well represented by both reanalyses and GCMs (Wilby et al, 2004;
55 Hanssen-Bauer et al, 2005; Brands et al, 2013) accounting for a major part of
56 the variability in the predictands are typically chosen as predictors in this ap-
57 proach (usually large-scale variables at different vertical levels), whereas variables
58 directly influenced by model parameterizations and/or orography, such as precip-
59 itation, are usually discarded. As a result, one of the most time-consuming tasks
60 in PP methods is the selection of a suitable combination of predictors, which must
61 be defined over an appropriate geographical domain which encompasses the main
62 synoptic phenomena influencing the climate of the region of interest.

63 Differently, under the MOS approach, predictors are taken from the same GCM
64 for both the training and the prediction phases. In the context of seasonal forecast-
65 ing, MOS methods have been traditionally applied establishing an empirical link
66 (e.g. regression or canonical correlation analysis) between large-scale circulation
67 predictors and pairwise observations at a monthly/seasonal time-scale. However,
68 simpler MOS alternatives based on BC methods are becoming increasingly popu-
69 lar (see, e.g., Themeßl et al, 2012a). BC methods directly adjust the distribu-
70 tion of GCM predicted precipitation against local observations (e.g. local scaling
71 or quantile mapping), to ensure that their statistical properties are similar. The
72 main advantage of these methods is their simplicity, since no predictor/domain
73 screening is required (typically, GCM output from the closest model gridbox is
74 considered as unique predictor). For instance, in local scaling methods (the sim-

75 plest case of BC), a linear transformation is applied to the model output to adjust
76 the first and/or second order moments of the predicted distribution.

77 A considerable body of research on the application of SD methods to climate
78 change simulations already exists (see, e.g., Gutiérrez et al, 2013b; Vaittinada
79 et al, 2016; Maraun, 2016; San-Martín et al, 2017). Beyond the adjustment of
80 systematic biases (Maraun et al, 2015), however, the advantages and limitations
81 of these methods in seasonal forecasting are not well understood yet, in particular
82 in what refers to their effect on forecast quality/skill. To measure this skill (which
83 is understood as forecast association and reliability here), we focus on correlation
84 and reliability categories. Note that, differently to other scores such as the mean
85 absolute error and the continuous ranked probability score, these two metrics are
86 not sensitive to changes in the mean. Therefore, they allow to properly assess the
87 added value of the SD methods applied beyond the effect of bias reduction.

88 Some prospects on the potential added value of BC methods can be envis-
89 aged for the most simple ones. For instance, local scaling preserve the temporal
90 structure of the original model predictions and do not affect neither correlation
91 nor reliability. However, more sophisticated distributional BC methods such as
92 quantile mapping can introduce arbitrary temporal changes (Maraun, 2013) and
93 thus, their effect on correlation and reliability is difficult to estimate in advance.
94 Differently, PP methods do rely on the temporal correspondence between the pre-
95 dictand and the predictors considered, so there might be windows of opportunity
96 for improving correlation and/or reliability in cases where large-scale variables are
97 better predicted by the model than local precipitation.

98 In this paper we analyze this problem focusing on a challenging tropical region,
99 the Philippines, which has been identified as an ideal test-bed for SD studies due
100 to the complex topography and land-sea contrasts which determine local rainfall
101 (Moron et al, 2009; Robertson et al, 2012; Manzanas et al, 2015). Moreover, its
102 climate is largely influenced by ENSO (see, e.g., Lyon et al, 2006; Manzanas et al,
103 2014) and it is located in a region of the world where seasonal forecasts are partic-
104 ularly skillful (Manzanas et al, 2014). As a result, there may be special potential
105 for the application of SD methods to seasonal forecasts in this area. We focus on
106 downscaling methods providing daily data and refer the interested reader to the
107 existing literature (Kang et al, 2007; Robertson et al, 2012) for details on the appli-
108 cation of seasonal MOS methods in the Philippines. In particular, we analyze and
109 intercompare the results from two state-of-the-art BC (parametric and empirical
110 quantile mapping) and two PP (analogs and Generalized Linear Models, GLMs)
111 methods when applied to the seasonal hindcast provided by the ENSEMBLES
112 project (Weisheimer et al, 2009) for the period 1981-2005. To our knowledge, this
113 work provides the most comprehensive study on the added value of the BC and
114 PP approaches for downscaling of seasonal forecasts to-date.

115 The paper is organized as follows. In Section 2 we introduce the data used (both
116 predictand and predictors). Sections 3 and 4 describe the statistical downscaling
117 methods that are applied and the verification metrics which are considered to
118 assess their performance, respectively. The results obtained are presented through
119 Section 5. Finally, the most important conclusions are given in Section 6.

2 Data

2.1 Precipitation in the Philippines: Predictands and Verifying Observations

The Philippines is an archipelago of 7107 islands with complex topography (see Figure 1a) located between the monsoonal and inner tropics (4°N and 20°N). Apart from ENSO (Lyon et al, 2006; Manzanas et al, 2014), the climate of this region is affected by important large-scale processes such as the southwest summer and northeast winter monsoons of the western North Pacific Ocean (Wang, 2002), but also by local forcing related to the presence of mountains and the complex land-sea contrast (Robertson et al, 2012). As a result, the country exhibits a rich regional climate composition which has been commonly classified into four different Climatic Types (CTs) in previous studies (Coronas, 1920; Manzanas et al, 2015).

For a good characterization of this variability, daily precipitation from 42 gauges maintained by the Philippine Atmospheric, Geophysical and Astronomical Services Administration (PAGASA: <http://www.pagasa.dost.gov.ph>), which are uniformly distributed across the country (see Figure 1b), was considered for this work for the period 1981-2005. The percentage of missing data within this period was less than 5% in all cases (less than 1% in most of the stations) so missing values were ignored in the calibration/training and verification processes. Panels *c-f* in Figure 1 show the interannual variability of spatial average precipitation totals for each CT (see colors in the legend) for the four standard boreal seasons: winter (DJF), spring (MAM), summer (JJA) and autumn (SON). Note that precipitation along the coastlines of the northern part of the archipelago (CT1 and CT2) exhibits a strong seasonal cycle, which is driven by alternating monsoonal winds. In particular, during the southwest monsoon (summer), precipitation peaks at the stations pertaining to CT1 while CT2 is affected by relative dryness. The opposite situation occurs during the northeast monsoon (winter). During the dry months (spring), easterly winds prevail, leading to orographic precipitation along the mountain ranges in the east of the archipelago and to relatively high precipitation amounts for the stations pertaining to CT2. At the stations belonging to CT3 and CT4 (mainly situated in the center and south of the archipelago), precipitation is mainly driven by meso-scale dynamics rather than by large-scale phenomena such as the monsoon circulation, leading to a weak seasonal cycle (rains uniformly distributed along the year). For a more comprehensive description of the climate of the Philippines, the interested reader is referred to Coronas (1920); Flores and Balagot (1969); Kintanar (1984) as well as to the PAGASA website.

2.2 Model Data: Predictors

In this work we consider both reanalysis and seasonal forecast data for the upper-air variables used as predictors (zonal wind component U at 850 and 200 hPa, specific humidity Q and temperature T at 850 hPa; see Section 3) as well as for surface precipitation, the target variable.

On the one hand, and following the recommendation by Manzanas et al (2015) —who carried out an assessment of reanalysis uncertainty over the region of study,— the ERA-Interim reanalysis (Dee et al, 2011) was chosen for the training phase of the PP methods. On the other hand, seasonal forecasts were obtained

164 from four of the GCMs contributing to the ENSEMBLES multimodel seasonal
 165 hindcast (Weisheimer et al, 2009), which were produced at the following centres:
 166 The European Centre for Medium-Range Weather Forecasts (ECMWF), the Leib-
 167 niz Institute of Marine Sciences (IFM-GEOMAR), the Euro-Mediterranean Centre
 168 for Climate Change (CMCC-INGV) and Météo France (MF). Each of these models
 169 —whose main components are summarized in Table 1— ran an ensemble of nine
 170 initial conditions (nine equiprobable members), produced by perturbing the real-
 171 istic estimates of the observed initial state four times a year (the first of February,
 172 May, August and November) within the period 1960-2005, providing seven month-
 173 long retrospective forecasts. For this work, one-month lead seasonal forecasts were
 174 considered. Note that, although the ENSEMBLES models are several years older
 175 than state-of-the-art seasonal forecasting systems, they form the most homoge-
 176 neous and comprehensive multimodel ensemble publicly available to-date.

Centre	Atmospheric model and resolution	Ocean model and resolution
ECMWF	IFS CY31R1 (T159 \approx 80km/L62)	HOPE (0.3° – 1.4°/L29)
IFM-GEOMAR	ECHAM5 (T63 \approx 180km/L31)	MPI-OM1 (1.5°/L40)
CMCC-INGV	ECHAM5 (T63 \approx 180km/L19)	OPA8.2 (2.0°/L31)
MF	ARPEGE4.6 (T63 \approx 180km)	OPA8.2 (2.0°/L31)

Table 1 Main components of the four global models used in this work, which contributed to the ENSEMBLES multimodel seasonal hindcast.

177 To keep consistency among reanalysis and the ENSEMBLES models, all predic-
 178 tor data were re-gridded to a common regular 2° grid applying a nearest neighbour
 179 interpolation scheme. Moreover, daily instantaneous values at 00 UTC were chosen
 180 in all cases. The common period for the available predictands and predictors, 1981-
 181 2005, was considered for this work. Note that, according to the WMO Lead Centre
 182 for the Long Range Forecast Verification (<http://www.bom.gov.au/wmo/lrfvs>), a
 183 25-years long period is suitable for the proper verification of seasonal forecasts.

184 Finally, in order to properly harmonize the reanalysis and the ENSEMBLES
 185 model data used respectively in the training and prediction phases of the PP
 186 methods, a simple local scaling correction was applied to the latter. In particular,
 187 for every large-scale model predictor, monthly mean values were adjusted towards
 188 the corresponding reanalysis values, gridbox by gridbox, avoiding thus problems
 189 that may arise due to the models mean biases.

190 3 Downscaling Methods

191 As representative of the PP approach we considered Generalized Linear Models
 192 (GLMs) and the analog technique, and relied on the optimum downscaling con-
 193 figuration found for the region of study in Manzananas et al (2015). In particular,
 194 they used as predictors a combination of two circulation (U at 850 and 300 hPa)
 195 and two thermodynamic (Q and T at 850 hPa) variables over a domain spanning
 196 from 114°E to 132°E and from 2°N to 22°N. Here, U_{300} has been replaced by the
 197 closest available variable in the ENSEMBLES models, U_{200} .

198 GLMs were formulated by Nelder and Wedderburn (1972) in the 1970's and
 199 are an extension of the classical linear regression which allows to model the ex-

pected value for non-normally distributed variables. While GLMs have been widely used for statistical downscaling of climate change scenarios (e.g., Brandsma and Buishand, 1997; Chandler and Wheeler, 2002; Abaurrea and Asín, 2005; Fealy and Sweeney, 2007; Hertig et al, 2013), they have been rarely applied to seasonal forecasts. Given the dual (occurrence and amount) character of precipitation, we followed in this work the common two-stage implementation (see, e.g., Chandler and Wheeler, 2002; Manzanas et al, 2015) in which a GLM with Bernoulli error distribution and *logit* canonical link-function (also known as logistic regression) is used to downscale daily precipitation occurrence (as characterized by a threshold of 0.1mm) and a GLM with gamma error distribution and *log* canonical link-function is applied to downscale daily precipitation amount. A stochastic component could be introduced in both GLMs to increase the predicted variance, which is usually underestimated in deterministic ones (Enke, 1997). However, in order to keep this stochastic effect away from the validation results, the two GLMs considered in this work were deterministic, i.e., predictions were based on the expected values. For this method (denoted as PP1 hereafter), we considered as predictors the 15 leading principal components (PCs, see Preisendorfer, 1988) over the above mentioned domain. PCs were obtained, both for the reanalysis and for the seasonal forecasts, by projecting the corresponding standardized fields onto the Empirical Orthogonal Functions obtained from the reanalysis, which were computed simultaneously on all predictor variables, considering the joined vector of standardized fields. The number of PCs retained, which explain over 80% of the predictor variance, was selected as a trade-off between model parsimony and goodness-of-fit (after a sensitivity study testing models with an increasing number of PCs).

The popular analogue technique (Lorenz, 1963, 1969) estimates the local down-scaled values corresponding to a particular atmospheric configuration (as represented by a number of model predictors defined over a certain geographical domain) from the local observations corresponding to a set of similar (or analog) atmospheric configurations within a historical catalog formed by a reanalysis. Here, similarity was measured in terms of the Euclidean distance (Matulla et al, 2008), which was computed over the complete predictor fields. Analog-based methods have been applied in several previous studies to downscale precipitation in the context of seasonal forecasting (see, e.g., Frías et al, 2010; Wu et al, 2012; Shao and Li, 2013). In spite of its simplicity, the analog technique performs as well as other more sophisticated ones (Zorita and von Storch, 1999) and it is one of the most widely used. Here, a deterministic version of the technique (Zorita et al, 1995; Cubasch et al, 1996) which considers the closest analog is used. This will be referred to as PP2 hereafter.

As representative of the BC approach we used two quantile mapping methods, one parametric and one empirical. In the parametric case (referred to as BC1 henceforth) daily predicted and observed rainfall intensities are fitted to gamma distributions and then daily predicted values are corrected according to the differences of the corresponding quantiles from the fitted distributions (Piani et al, 2010; Themeßl et al, 2012a). Note that the parameters of the gamma distribution can be estimated from the first two moments and, therefore, in practice, this method is similar to a local scaling. The empirical method (denoted as BC2 hereafter) consists of calibrating the predicted empirical probability density function (PDF) by adjusting a number of quantiles based on the empirical observed PDF (see, e.g., Déqué, 2007). In particular, we proceed by adjusting percentiles 1 to 99 and

249 linearly interpolating inside this range every two consecutive percentiles. Outside
250 this range a constant extrapolation (using the correction obtained for the 1st or
251 99th percentile) is applied. Moreover, in cases when the predicted frequency of
252 dry days is larger than the observed one, the frequency adaptation proposed by
253 Themeßl et al (2012b) is applied.

254 The two BC and the two PP methods described above were separately cal-
255 ibrated/trained and applied for each of the four seasons. We followed a k -fold
256 cross-validation approach (Gutiérrez et al, 2013b) for the period 1981-2005, split-
257 ting the whole 25-year period into $k = 5$ random test sets (folds) of 5 years each.
258 Each of these sets was independently used for the prediction phase, using the re-
259 maining 20 years for training. For each model, the two BC methods were separately
260 calibrated and applied for each of the nine available ensemble members. However,
261 it is worth to notice here that other configurations were also analyzed for these
262 methods. For instance, we tested cross-validated versus not cross-validated meth-
263 ods and member- versus ensemble-wise calibrated ones (the latter considering the
264 joined nine members series), obtaining very similar results in all cases (not shown).
265 Thus, the conclusions obtained in this work for the BC methods do not depend on
266 the particular experimental configuration followed. Differently, note that the two
267 PP methods were trained just once (based on reanalysis predictor data and local
268 observed precipitation). Afterwards, the (unique) resulting statistical model was
269 separately applied to each of the nine members.

270 4 Verification Metrics

271 In order to validate the forecast quality of the raw seasonal precipitation outputs
272 from the ENSEMBLES models and the possible added value of the corresponding
273 downscaled results (beyond the adjustment of systematic biases) we considered
274 two scores recommended by the WMO Lead Centre for the Long Range Fore-
275 cast Verification (<http://www.bom.gov.au/wmo/lrfvs>): The interannual Anomaly
276 Correlation Coefficient (ACC) and a measure of reliability based on the different
277 categories introduced by Weisheimer and Palmer (2014).

278 ACC is a simple metric of forecast association which allows to assess the ability
279 of raw/downscaled precipitation to reproduce the observed interannual seasonal
280 anomalies. For each particular model, it is applied here to the deterministic forecast
281 resulting from averaging the nine (either raw or downscaled) available members.
282 In addition, a multimodel (MM) was also constructed by considering the 36 (4
283 models x 9 members) available predictions (either raw or downscaled), thus giving
284 equal weights to all models and members.

285 Reliability measures how closely the forecast probabilities of a certain event
286 correspond to the actual chance of observing that event. It is applied here for
287 probabilistic forecasts of each of the three precipitation terciles: dry (T1), normal
288 (T2) and wet (T3). For each model (the MM), probabilities are computed based on
289 the nine (36), either raw or downscaled, available members. Reliability diagrams
290 (see the illustrative examples shown in Figure 2) plot the observed frequencies of
291 the event considered (e.g. T1, T2 or T3) as a function of its forecast probabili-
292 ty, as represented by a determined number of bins (see Doblas-Reyes et al, 2008,
293 for details). For a perfectly reliable forecasting system, the curve obtained would
294 match the diagonal (perfect reliability line). Points falling within the so-called skill

295 region (in gray), i.e., the region contained between the no-resolution line (which
296 indicates the expected frequency of the event: 1/3 for terciles) and the no-skill line
297 (halfway between the no-resolution line and the diagonal) positively contribute
298 to the forecast skill (Brier Skill Score > 0). Weisheimer and Palmer (2014) pro-
299 posed a methodology to translate the information provided by these diagrams
300 to an easy-to-interpret scale with five reliability categories: *perfect* (green), *still*
301 *very useful* (blue), *marginally useful* (yellow), *not useful* (orange) and *dangerously*
302 *useless* (red). In particular, they performed a weighted linear regression as a best-
303 guess estimate on all data points in the diagram (using the number of forecasts
304 in each probability bin as weights) and defined the different reliability categories
305 based on the relative position of the so derived reliability line with respect to the
306 perfect reliability (diagonal), no-skill and no-resolution lines, as well as on the un-
307 certainty range around it (as obtained by bootstrapping with 1000 samples). Here,
308 we slightly modified this original classification by Weisheimer and Palmer (2014)
309 for a better adaptation to our particular regional study (see Section 5.3).

310 Note that the two validation metrics considered for this work are insensitive
311 to data scaling and, therefore, are suitable to assess the added value of the down-
312 scaling methods beyond the improvement of systematic biases in the mean and
313 variance. Thus, we assess here the relevant aspects which can provide added value
314 for seasonal forecasting.

315 5 Results

316 5.1 Performance of Raw Models

317 In order to obtain an estimation of the performance of the ENSEMBLES models
318 over the region of study, we carried out a regional validation considering as refer-
319 ence the observed precipitation at the 42 PAGASA stations (model precipitation
320 was bi-linearly interpolated to these gauges). Figure 3 shows the results obtained
321 in terms of local biases, which are in general strong (as compared with the observed
322 climatologies, shown in the first row). Note that in spite of local differences, all
323 models (and as a result the MM) exhibit similar spatial patterns for the different
324 seasons, which reflect their inability to properly represent the local features in this
325 region of complex orography and land-sea contrast. Notice that, by construction,
326 all the statistical downscaling methods here considered reduce the mean biases,
327 yielding absolute biases smaller than 10 mm/year in all cases (not shown). Al-
328 though this is a clear advantage for end users, here we focus on the added value in
329 terms of skill (as characterized by forecast association and reliability). The reader
330 is referred to (Maraun et al, 2015) for further information on the performance of
331 the different downscaling methods from the point of view of biases and marginal
332 statistics.

333 Figure 4 shows the local interannual ACC values obtained. In general, signifi-
334 cant correlations are found for all models throughout the year (especially in DJF
335 and MAM) except for JJA. This marked seasonality in forecast skill is a conse-
336 quence of the large influence exerted by the ENSO interannual oscillations in this
337 region (Manzanas et al, 2014). However, important local-to-regional differences
338 can be found for different models in some seasons. For instance, the ECMWF
339 model exhibits a superior performance for the CT1 region in JJA. This could be a

340 consequence of the higher resolution of this model, as compared to the other three
341 (see Table 1).

342 5.2 Correlation of Downscaled Results

343 For the different seasons (in rows) and CTs (in columns), panels in Figure 5 show
344 the interannual ACC values obtained for each of the ENSEMBLES models (see
345 the colors in the legend). Boxplots display the results along the different stations
346 for the raw/direct model output (DMO henceforward), which is indicated by a
347 light gray shadow, and for all the downscaling methods considered (right after the
348 DMO). Overall, results vary mainly among seasons, but also among CTs, models
349 and downscaling methods. For the latter, results are in general more sensitive to
350 the approach considered (BC or PP) than to the particular technique used within
351 each approach. As already explained in Section 5.1, the highest scores for the
352 DMO are obtained for DJF and MAM, whereas the worst results are found for
353 JJA, with no significant correlations for any model except for the ECMWF in the
354 CT1 region. In general, the DMO outperforms the BC methods (note that the
355 correlation gain found for the latter in some cases is limited to a few stations and
356 is counteracted by the loss found in others, so no robust signal of added value is
357 obtained for the BC approach). Nonetheless, PP methods can either improve or
358 spoil the correlations attained by the DMO, depending on the case.

359 More in detail, whereas the BC methods do not improve (or even worsen) the
360 correlations reached by the DMO in general for DJF and MAM, there are a few
361 cases in which PP methods can add important value (indicated by black dotted
362 boxes). In particular, PP methods are shown to improve raw precipitation from the
363 relatively bad performing models (those exhibiting small ACC values, as compared
364 to the rest of models), as occurs for the MF model in DJF (CT4) and the IFM-
365 GEOMAR model in MAM (CT1). Moreover, as marked with red dotted boxes, PP
366 methods can also add important local value for some particular outlier stations
367 (those in which the correlation for the raw model precipitation drops, as compared
368 with the rest of locations). See, for instance, the case of the CMCC-INGV model in
369 MAM (CT2 and CT3). Notice that, as opposite to the DMO and the BC methods
370—which depend on model precipitation at the nearest gridbox and can be affected
371 by local features such as wrong orographical gradients, land-sea interfaces, etc.,—
372 PP methods rely on large-scale predictors to infer local precipitation, which might
373 allow in turn to properly reproduce the observed interannual variability in these
374 cases.

375 With respect to JJA and SON, whereas BC methods do not clearly improve
376 (or even worsen) the correlations attained by the DMO, PP methods provide
377 in general better (worse) results than the DMO in the former (latter) season. In
378 particular, notice that PP methods yield large correlation improvements in JJA for
379 the stations pertaining to CT1 for all models (with the exception of the ECMWF),
380 which exhibit nearly-zero ACC values in this season.

381 In order to summarize the results from Figure 5 and to better quantify the
382 added value of BC and PP methods, Figure 6 shows in bar charts the percentage
383 of stations with significant ACC values for the DMO and for the different down-
384 scaling approaches (BC and PP), for the different seasons. Within each approach,
385 the two methods applied are jointly considered. Moreover, all models except the

MM (which is excluded for clarity) and all CTs are also jointly considered. This figure shows that BC methods do not outperform (or slightly reduce) the correlations attained by the DMO for any season. However, PP methods yield higher (lower) correlations than the DMO does for JJA (SON). In particular, whereas the percentage augments from 10% to 30% in JJA, it drops from more than 60% to less than 30% in SON.

5.3 Reliability of Downscaled Results

In Weisheimer and Palmer (2014), the confidence interval around the best-guess reliability line was estimated by randomly resampling members, gridboxes and years, and the 75% of the total range was considered. Here, we analyzed the sensitivity of their classification to different confidence intervals (the same bootstrapping procedure was used) and found that the ensemble size had a large influence, as higher uncertainty around the best-guess reliability line was obtained for smaller ensembles. As a result, *still very useful* (blue) categories may pass to *marginally useful* (yellow) ones due to an enlargement of the confidence region (see Weisheimer and Palmer, 2014, for details on the definition of the different categories). Therefore, in this work we considered a smaller confidence interval given by the central 50% of the total range, which is more suitable for the nine members of the ENSEMBLES models used —note that the original classification was developed for the 51 members version of the ECMWF System 4 model (Molteni et al, 2011).— Moreover, in order to introduce further discrimination power, within the original *marginally useful* (yellow) category, we differentiate those cases in which the best-guess reliability line is above the no skill line, assigning to this new category (denoted as *marginally useful +*) the dark yellow color. See, for instance, panels *g* and *h* in Figure 2 —note that both cases would correspond to the same category in the original definition.—

Figure 7 shows the reliability categories (in colors) obtained after applying the methodology described above for the different models (in columns) and seasons (in rows), by CT (note that the joined series of the different stations falling within each CT are considered). From left to right, each block shows the results for the DMO, the two BC and the two PP methods considered, for the three terciles. Overall, this figure is in good correspondence with the results found for correlation (Figures 5 and 6), with the best reliability obtained in DJF and MAM and the worst in JJA. Moreover, the results for the two BC methods are very similar to those obtained for the DMO, with slight differences due to spurious changes of category (as illustrated in the top row of Figure 2). However, the two PP methods exhibit major reliability differences with respect to the DMO, especially for JJA and SON. In particular, both PP1 and PP2 improve the results of the DMO in the former season, especially for the CT1, where *marginally useful* or *marginally useful +* categories are obtained instead of *not useful* and *dangerously useless* ones. Yet, the opposite situation is found for SON. Additionally, this figure also shows some well-known results (see, e.g., Manzananas et al, 2014), such as the higher performance attained for the extreme terciles (as compared to the normal one) and the superiority of the MM, which in general outperforms any single model.

In order to summarize the results from Figure 7 and to better quantify the added value of the different approaches for statistical downscaling, Figure 8 shows

432 in stacked bar charts the percentage of reliability categories obtained from the
433 DMO and the different downscaling approaches (BC and PP) for the different
434 seasons. Within each approach, the two methods applied are jointly considered.
435 For clarity, the results from the MM and from the normal tercile are excluded from
436 this analysis. This figure shows that BC methods do not provide clear added value
437 (or even worsen the DMO) for any season. However, PP methods yield substantial
438 added value for JJA, leading to *marginally useful* or *marginally useful +* categories
439 in over 50% of the cases, as compared to less than 10% for the DMO (and for the
440 BC methods). In contrast, the opposite situation is found for the PP methods in
441 SON, with *not useful* or *dangerously useless* categories obtained in nearly 50% of
442 the cases (as compared with 10% for the DMO and 20% for the BC methods).

443 Remarkably, the good alignment between the results found for reliability and
444 those found for correlation points out the suitability and usefulness of the method-
445 ology proposed by Weisheimer and Palmer (2014) —which is slightly modified
446 here— for regional studies. Note that the original work was undertaken for the 21
447 global regions defined in Giorgi and Francisco (2000).

448 5.4 An Explanation for the Added Value of PP Methods

449 As already mentioned, PP methods rely on large-scale predictors to infer local
450 precipitation. As such, the above presented cases leading to a gain (loss) of skill
451 for the PP approach could be explained by situations where large-scale variables,
452 defined over a synoptic domain, are better (worse) predicted by the model than the
453 target precipitation, which is more affected by particular local features (as usually
454 represented by parametrizations). In order to check this premise, we focus here on
455 the climate region CT1, where PP methods were shown to improve (deteriorate)
456 the skill of the DMO in JJA (SON). Figure 9 displays the interannual ACC values
457 obtained between observed precipitation at the 13 stations pertaining to this CT
458 and the ERA-Interim and ENSEMBLES models outputs —the nearest gridbox
459 is considered— for precipitation (PR) and the different predictors used (U850,
460 U200, Q850 and T850) for the period 1981-2005. For benchmarking purposes,
461 ERA-Interim is indicated by a light gray shadow.

462 The gain of skill found in JJA for all models except the ECMWF (Figures
463 5 and 7) is in agreement with the results shown in the top panel. In particular,
464 whereas significant ACC values for precipitation are only found for the ECMWF
465 model, mostly significant correlations (similar to the benchmark provided by ERA-
466 Interim) are found for all models for U850 and T850, the large-scale predictors most
467 correlated with observed precipitation (as indicated by the reanalysis). This sug-
468 gests that PP methods might be able to exploit the model ability for reproducing
469 upper-air predictor variables to indirectly obtain improved precipitation forecasts
470 in cases of a poor skill for model precipitation.

471 The opposite situation is found for SON (bottom panel). In this season, the
472 ACC values found for precipitation are significant (although smaller than the
473 benchmark provided by ERA-Interim) in most cases. However, the results found
474 for the large-scale predictors are in general not significant. Moreover, opposite
475 correlations with observations (as compared to the reanalysis) are found in some
476 cases. The combined effect of these errors could result in wrong downscaled pre-
477 dictions, as occurs for the ECMWF model, which leads to negative ACC values

478 (see the corresponding boxplots in Figure 5) and *dangerously useless* reliability
479 categories (see the corresponding extreme terciles in Figure 7).

480 6 Conclusions

481 In order to assess the advantages and limitations of different approaches for statis-
482 tical downscaling in the context of seasonal forecasting, two state-of-the-art Bias
483 Correction (BC) and two Perfect Prognosis (PP) methods were applied to obtain
484 local precipitation at 42 stations in the Philippines, considering one-month lead
485 forecasts from the ENSEMBLES multimodel seasonal hindcast for the four boreal
486 seasons over the period 1981-2005.

487 As expected by construction, BC and PP methods were shown to be successful
488 in reducing the systematic model biases over the area of study, which are in general
489 strong (as compared to the local climatologies). In particular, both approaches lead
490 to very small biases after downscaling. However, and even though this is a clear
491 advantage for users, we focus here on the methods' ability to predict interannual
492 anomalies, which is the basis of seasonal forecasting. Therefore, we assess forecast
493 quality/skill in terms of interannual correlation and reliability categories. Note that
494 these two metrics are not sensitive to changes in the mean and allow therefore to
495 properly assess the added value of the downscaling methods beyond the effect of
496 bias reduction.

497 On the one hand, BC methods were shown to provide no added value in terms
498 of skill, maintaining or worsening both correlation and reliability. These meth-
499 ods directly transform model precipitation (by correcting different quantiles of the
500 distribution) without relying on any additional information about the underlying
501 physical phenomena (e.g. large-scale circulation). As a consequence, BC methods
502 can arbitrarily modify the temporal structure of the raw model output, with the
503 overall result of degrading the skill (Maraun, 2013). Noticeably, the conclusions
504 obtained here for the BC methods are quite general and do not depend on the par-
505 ticular experimental configuration followed. For instance, we tested cross-validated
506 versus not cross-validated methods and member- versus ensemble-wise calibrated
507 ones, obtaining very similar results in all cases.

508 On the other hand, we found that PP methods can either substantially improve
509 or deteriorate correlation and reliability. As opposite to BC ones, PP methods rely
510 on physically-based large-scale model predictors to infer local precipitation. Thus,
511 this provides an opportunity for improving the original model skill in those cases
512 for which orographic and land-sea contrasts limit the local representativeness of
513 model precipitation, but the model is yet skillful in reproducing the large-scale
514 predictors. In this work, we show that those conditions are met for certain regions
515 and/or seasons. For instance, reliability was increased by PP methods in nearly
516 40% of the stations considered in summer.

517 Therefore, we conclude that the choice of an appropriate statistical downscal-
518 ing method is not trivial and depends on factors such as the region, the season,
519 the strength of the connection between the large- and the local-scale climate and
520 the model skill for predicting surface/upper-air variables. Moreover, this selection
521 should be based on the requirements of the particular user and/or application. In
522 general, it is advisable to test the added value of PP methods as a first choice,

523 particularly in regions with complex orography and/or large local variability. How-
524 ever, BC methods could be a cost-effective and pragmatic choice in applications for
525 which the main concern is just reducing model biases, even at the cost of degrading
526 the skill.

527 **Acknowledgements** This study was partially supported by the SPECS and EUPORIAS
528 projects, funded by the European Commission through the Seventh Framework Programme
529 for Research under grant agreements 308378 and 308291, respectively. JMG acknowledges
530 partial support from the project MULTI-SDM (CGL2015-66583-R, MINECO/FEDER). Also,
531 the authors are grateful to PAGASA for the data provided.

532 References

- 533 Abaurrea J, Asín J (2005) Forecasting local daily precipitation patterns in a cli-
534 mate change scenario. *Climate Research* 28(3):183–197, DOI 10.3354/cr028183
- 535 Brands S, Herrera S, Fernández J, Gutiérrez JM (2013) How well do CMIP5
536 Earth System Models simulate present climate conditions in Europe and Africa?
537 *Climate Dynamics* 41(3):803–817, DOI 10.1007/s00382-013-1742-8
- 538 Brandsma T, Buishand TA (1997) Statistical linkage of daily precipitation in
539 Switzerland to atmospheric circulation and temperature. *Journal of Hydrology*
540 198(1-4):98–123, DOI 10.1016/S0022-1694(96)03326-4
- 541 Chandler RE, Wheater HS (2002) Analysis of rainfall variability using generalized
542 linear models: A case study from the west of Ireland. *Water Resources Research*
543 38(10):1–11, DOI 10.1029/2001WR000906
- 544 Coronas J (1920) *The climate and weather of the Philippines, 1903-1918*, Bureau
545 of Printing, Manila, pp 291–467
- 546 Cubasch U, von Storch H, Waszkewitz J, Zorita E (1996) Estimates of climate
547 change in Southern Europe derived from dynamical climate model output. *Cli-
548 mate Research* 7(2):129–149, DOI 10.3354/cr007129
- 549 Dee DP, Uppala SM, Simmons AJ, Berrisford P, Poli P, Kobayashi S, Andrae U,
550 Balmaseda MA, Balsamo G, Bauer P, Bechtold P, Beljaars ACM, van de Berg L,
551 Bidlot J, Bormann N, Delsol C, Dragani R, Fuentes M, Geer AJ, Haimberger L,
552 Healy SB, Hersbach H, Holm EV, Isaksen L, Kallberg P, Koehler M, Matricardi
553 M, McNally AP, Monge-Sanz BM, Morcrette JJ, Park BK, Peubey C, de Rosnay
554 P, Tavolato C, Thepaut JN, Vitart F (2011) The ERA-Interim reanalysis: Con-
555 figuration and performance of the data assimilation system. *Quarterly Journal
556 of the Royal Meteorological Society* 137(656):553–597, DOI 10.1002/qj.828
- 557 Déqué M (2007) Frequency of precipitation and temperature extremes over France
558 in an anthropogenic scenario: Model results and statistical correction according
559 to observed values. *Global and Planetary Change* 57(1-2):16–26, DOI 10.1016/
560 j.gloplacha.2006.11.030
- 561 Doblas-Reyes FJ, Coelho CAS, Stephenson DB (2008) How much does simplifica-
562 tion of probability forecasts reduce forecast quality? *Meteorological Applications*
563 15(1):155–162, DOI 10.1002/met.50
- 564 Enke SA W (1997) Downscaling climate model outputs into local and regional
565 weather elements by classification and regression. *Climate Research* 8(3):195–
566 207

- 567 Fealy R, Sweeney J (2007) Statistical downscaling of precipitation for a selection of
568 sites in Ireland employing a generalised linear modelling approach. *International*
569 *Journal of Climatology* 27(15):2083–2094, DOI 10.1002/joc.1506
- 570 Flores JF, Balagot VF (1969) World Survey of Climatology, Climates of Northern
571 and Eastern Asia, vol 8, Arakawa, chap Climate of the Philippines, pp 159–213
- 572 Frías MD, Herrera S, Cofiño AS, Gutiérrez JM (2010) Assessing the skill of
573 precipitation and temperature seasonal forecasts in Spain: Windows of op-
574 portunity related to ENSO events. *Journal of Climate* 23(2):209–220, DOI
575 10.1175/2009JCLI2824.1
- 576 Giorgi F, Francisco R (2000) Uncertainties in regional climate change prediction: A
577 regional analysis of ensemble simulations with the HADCM2 coupled AOGCM.
578 *Climate Dynamics* 16(2-3):169–182, DOI 10.1007/PL00013733
- 579 Gutiérrez JM, Bedia J, Benestad R, Pagé C (2013a) Review of the different statisti-
580 cal downscaling methods for S2D prediction. Tech. rep., SPECS deliverable 5.2.1,
581 URL http://www.specs-fp7.eu/sites/default/files/u1/SPECS_D52.1.pdf
- 582 Gutiérrez JM, San-Martín D, Brands S, Manzanas R, Herrera S (2013b) Re-
583 assessing statistical downscaling techniques for their robust application under
584 climate change conditions. *Journal of Climate* 26(1):171–188, DOI 10.1175/
585 JCLI-D-11-00687.1
- 586 Hanssen-Bauer I, Achberger C, Benestad RE, Chen D, Forland EJ (2005) Sta-
587 tistical downscaling of climate scenarios over Scandinavia. *Climate Research*
588 29(3):255–268, DOI 10.3354/cr029255
- 589 Hertig E, Seubert S, Paxian A, Vogt G, Paeth H, Jacobeit J (2013) Changes of
590 total versus extreme precipitation and dry periods until the end of the twenty-
591 first century: Statistical assessments for the Mediterranean area. *Theoretical*
592 *and Applied Climatology* 111(1-2):1–20, DOI 10.1007/s00704-012-0639-5
- 593 Kang H, An KH, Park CK, Solís ALS, Stithichivapak K (2007) Multimodel
594 output statistical downscaling prediction of precipitation in the Philippines
595 and Thailand. *Geophysical Research Letters* 34(15):n/a–n/a, DOI 10.1029/
596 2007GL030730
- 597 Kintanar RL (1984) Climate of the Philippines. Tech. rep., PAGASA
- 598 Lorenz EN (1963) Deterministic nonperiodic flow. *Journal of the Atmospheric*
599 *Sciences* 20(2):130–141
- 600 Lorenz EN (1969) Atmospheric predictability as revealed by naturally occurring
601 analogues. *Journal of the Atmospheric Sciences* 26(4):636–646, DOI 10.1175/
602 1520-0469(1969)26<636:APARBN>2.0.CO;2
- 603 Lyon B, Cristi H, Verceles ER, Hilario FD, Abastillas R (2006) Seasonal reversal
604 of the ENSO rainfall signal in the Philippines. *Geophysical Research Letters*
605 33(24):n/a–n/a, DOI 10.1029/2006GL028182
- 606 Manzanas R, Frías MD, Cofiño AS, Gutiérrez JM (2014) Validation of 40 year
607 multimodel seasonal precipitation forecasts: The role of ENSO on the global
608 skill. *Journal of Geophysical Research: Atmospheres* 119(4):1708–1719, DOI
609 10.1002/2013JD020680
- 610 Manzanas R, Brands S, San-Martín D, Lucero A, Limbo C, Gutiérrez JM (2015)
611 Statistical downscaling in the tropics can be sensitive to reanalysis choice: A case
612 study for precipitation in the Philippines. *Journal of Climate* 28(10):4171–4184,
613 DOI 10.1175/JCLI-D-14-00331.1
- 614 Maraun D (2013) Bias correction, quantile mapping, and downscaling: Revis-
615 iting the inflation issue. *Journal of Climate* 26(6):2137–2143, DOI 10.1175/

- 616 JCLI-D-12-00821.1
- 617 Maraun D (2016) Bias correcting climate change simulations: A critical review.
618 *Current Climate Change Reports* 2(4):211–220, DOI 10.1007/s40641-016-0050-x
- 619 Maraun D, Wetterhall F, Ireson AM, Chandler RE, Kendon EJ, Widmann M,
620 Brienen S, Rust HW, Sauter T, Themessl M, Venema VKC, Chun KP, Good-
621 ess CM, Jones RG, Onof C, Vrac M, Thiele-Eich I (2010) Precipitation down-
622 scaling under climate change: Recent developments to bridge the gap between
623 dynamical models and the end user. *Reviews of Geophysics* 48(3):n/a–n/a, DOI
624 10.1029/2009RG000314
- 625 Maraun D, Widmann M, Gutiérrez JM, Kotlarski S, Chandler RE, Hertig E, Wibig
626 J, Huth R, Wilcke RA (2015) VALUE: A framework to validate downscaling
627 approaches for climate change studies. *Earth’s Future* 3(1):1–14, DOI 10.1002/
628 2014EF000259
- 629 Matulla C, Zhang X, Wang X, Wang J, Zorita E, Wagner S, von Storch H (2008)
630 Influence of similarity measures on the performance of the analog method for
631 downscaling daily precipitation. *Climate Dynamics* 30(2-3):133–144, DOI 10.
632 1007/s00382-007-0277-2
- 633 Molteni F, Stockdale T, Balsaseda M, Balsamo G, Buizza R, Ferranti L, Mag-
634 nusson L, Mogensen K, Palmer T, Vitart F (2011) The new ECMWF seasonal
635 forecast system (System 4). Tech. rep., ECMWF
- 636 Moron V, Lucero A, Hilario F, Lyon B, Robertson AW, DeWitt D (2009)
637 Spatio-temporal variability and predictability of summer monsoon onset
638 over the Philippines. *Climate Dynamics* 33(7-8):1159–1177, DOI 10.1007/
639 s00382-008-0520-5
- 640 Nelder JA, Wedderburn RWM (1972) Generalized linear models. *Journal of the*
641 *Royal Statistical Society Series A (Statistics in Society)* 135(3):370–384, DOI
642 10.2307/2344614
- 643 Piani C, Haerter JO, Coppola E (2010) Statistical bias correction for daily pre-
644 cipitation in regional climate models over Europe. *Theoretical and Applied Cli-*
645 *matology* 99(1-2):187–192, DOI 10.1007/s00704-009-0134-9
- 646 Preisendorfer R (1988) *Principal component analysis in meteorology and oceanog-*
647 *raphy*, 1st edn. Elsevier
- 648 Robertson AW, Qian JH, Tippett MK, Moron V, Lucero A (2012) Downscaling of
649 seasonal rainfall over the Philippines: Dynamical versus statistical approaches.
650 *Monthly Weather Review* 140(4):1204–1218, DOI 10.1175/MWR-D-11-00177.1
- 651 San-Martín D, Manzananas R, Brands S, Herrera S, Gutiérrez JM (2017) Reassessing
652 Model Uncertainty for Regional Projections of Precipitation with an Ensemble
653 of Statistical Downscaling Methods. *Journal of Climate* 30(1):203–223, DOI
654 10.1175/JCLI-D-16-0366.1
- 655 Shao Q, Li M (2013) An improved statistical analogue downscaling procedure
656 for seasonal precipitation forecast. *Stochastic Environmental Research and Risk*
657 *Assessment* 27(4):819–830, DOI 10.1007/s00477-012-0610-0
- 658 von Storch H, Zorita E, Cubasch U (1993) Downscaling of global climate change
659 estimates to regional scales: An application to Iberian rainfall in wintertime.
660 *Journal of Climate* 6(6):1161–1171, DOI 10.1175/1520-0442(1993)006<1161:
661 DOGCCE>2.0.CO;2
- 662 Themeßl MJ, Gobiet A, Heinrich G (2012a) Empirical-statistical downscaling and
663 error correction of regional climate models and its impact on the climate change
664 signal. *Climatic Change* 112(2):449–468, DOI 10.1007/s10584-011-0224-4

- 665 Themeßl MJ, Gobiet A, Heinrich G (2012b) Empirical-statistical downscaling and
666 error correction of regional climate models and its impact on the climate change
667 signal. *Climatic Change* 112(2):449–468, DOI 10.1007/s10584-011-0224-4
- 668 Vaittinada AP, Vrac M, Bastin S, Carreau J, Déqué M, Gallardo C (2016) Inter-
669 comparison of statistical and dynamical downscaling models under the EURO-
670 and MED-CORDEX initiative framework: Present climate evaluations. *Climate*
671 *Dynamics* 46(3-4):1301–1329, DOI 10.1007/s00382-015-2647-5
- 672 Wang B (2002) Rainy season of the Asian-Pacific summer monsoon. *Journal of Cli-*
673 *mate* 15(4):386–398, DOI 10.1175/1520-0442(2002)015<0386:RSOTAP>2.0.CO;2
- 674 Weisheimer A, Palmer TN (2014) On the reliability of seasonal climate forecasts.
675 *Journal of the Royal Society Interface* 11(96), DOI 10.1098/rsif.2013.1162
- 676 Weisheimer A, Doblas-Reyes FJ, Palmer TN, Alessandri A, Arribas A, Déqué
677 M, Keenlyside N, MacVean M, Navarra A, Rogel P (2009) ENSEMBLES: A
678 new multi-model ensemble for seasonal-to-annual prediction. Skill and progress
679 beyond DEMETER in forecasting tropical Pacific SSTs. *Geophysical Research*
680 *Letters* 36(21):n/a–n/a, DOI 10.1029/2009GL040896
- 681 Wilby RL, Charles S, Zorita E, Timbal B, Whetton P, Mearns L (2004) Guidelines
682 for use of climate scenarios developed from statistical downscaling methods.
683 Tech. rep., IPCC-TGCI
- 684 Wu W, Liu Y, Ge M, Rostkier-Edelstein D, Descombes G, Kunin P, Warner T,
685 Swerdlin S, Givati A, Hopson T, Yates D (2012) Statistical downscaling of cli-
686 mate forecast system seasonal predictions for the Southeastern Mediterranean.
687 *Atmospheric Research* 118:346–356, DOI 10.1016/j.atmosres.2012.07.019
- 688 Zorita E, von Storch H (1999) The analog method as a simple statistical downscal-
689 ing technique: Comparison with more complicated methods. *Journal of Climate*
690 12(8):2474–2489, DOI 10.1175/1520-0442(1999)012<2474:TAMAAS>2.0.CO;2
- 691 Zorita E, Hughes JP, Lettemaier DP, von Storch H (1995) Stochastic char-
692 acterization of regional circulation patterns for climate model diagnosis and
693 estimation of local precipitation. *Journal of Climate* 8(5):1023–1042, DOI
694 10.1175/1520-0442(1995)008<1023:SCORCP>2.0.CO;2

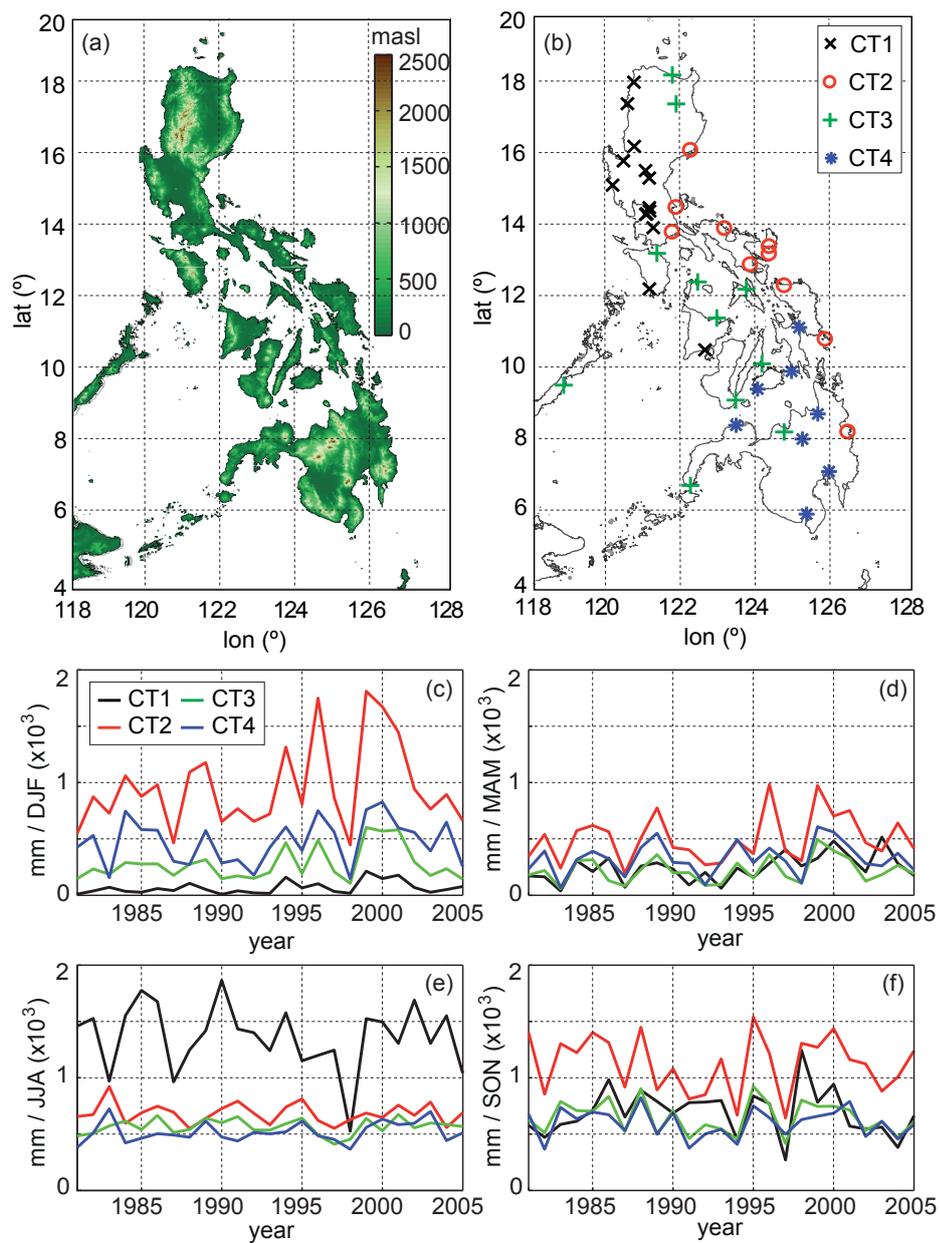


Fig. 1 (a) Topography of the Philippines. (b) Location of the 42 PAGASA gauges considered, classified into the four precipitation climatic types (CTs) defined in Coronas (1920), in colors. (c)-(f) Interannual variability of spatial average precipitation totals for each CT (see colors in the legend) for the period 1981-2005, by seasons.

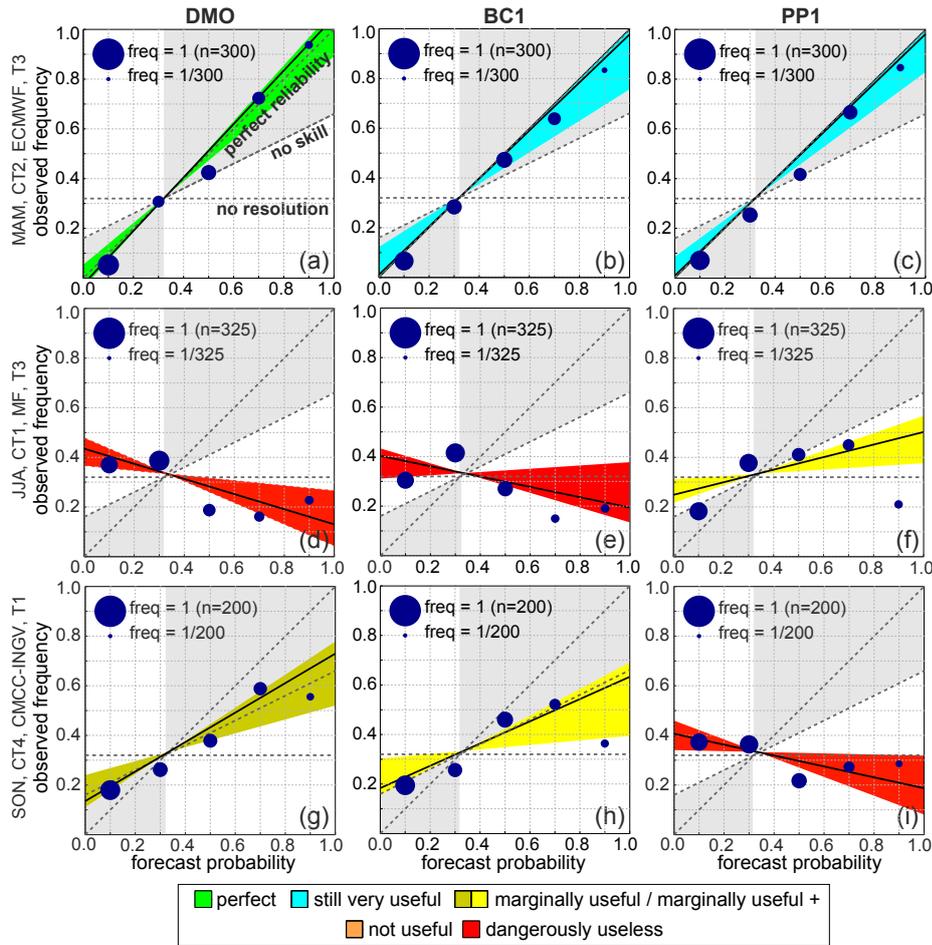


Fig. 2 Reliability diagrams for the raw/direct model output (DMO), the BC1 and the PP1 method (in columns), for three different illustrative examples of seasonal forecasts in MAM, JJA and SON (in rows), for different CTs and models (see the labels on the left-hand side). The gray area defines the region contributing positively to the forecast skill (Brier Skill Score > 0). The *perfect reliability*, *no skill* and *no resolution* lines are indicated in panel *a*. Colors correspond to the different categories used, which are based on the original scale proposed by Weisheimer and Palmer (2014) (see the text for details). Note that the joined series of the different stations falling within each CT are considered. The sample size used in each case is indicated in the upper left corner.

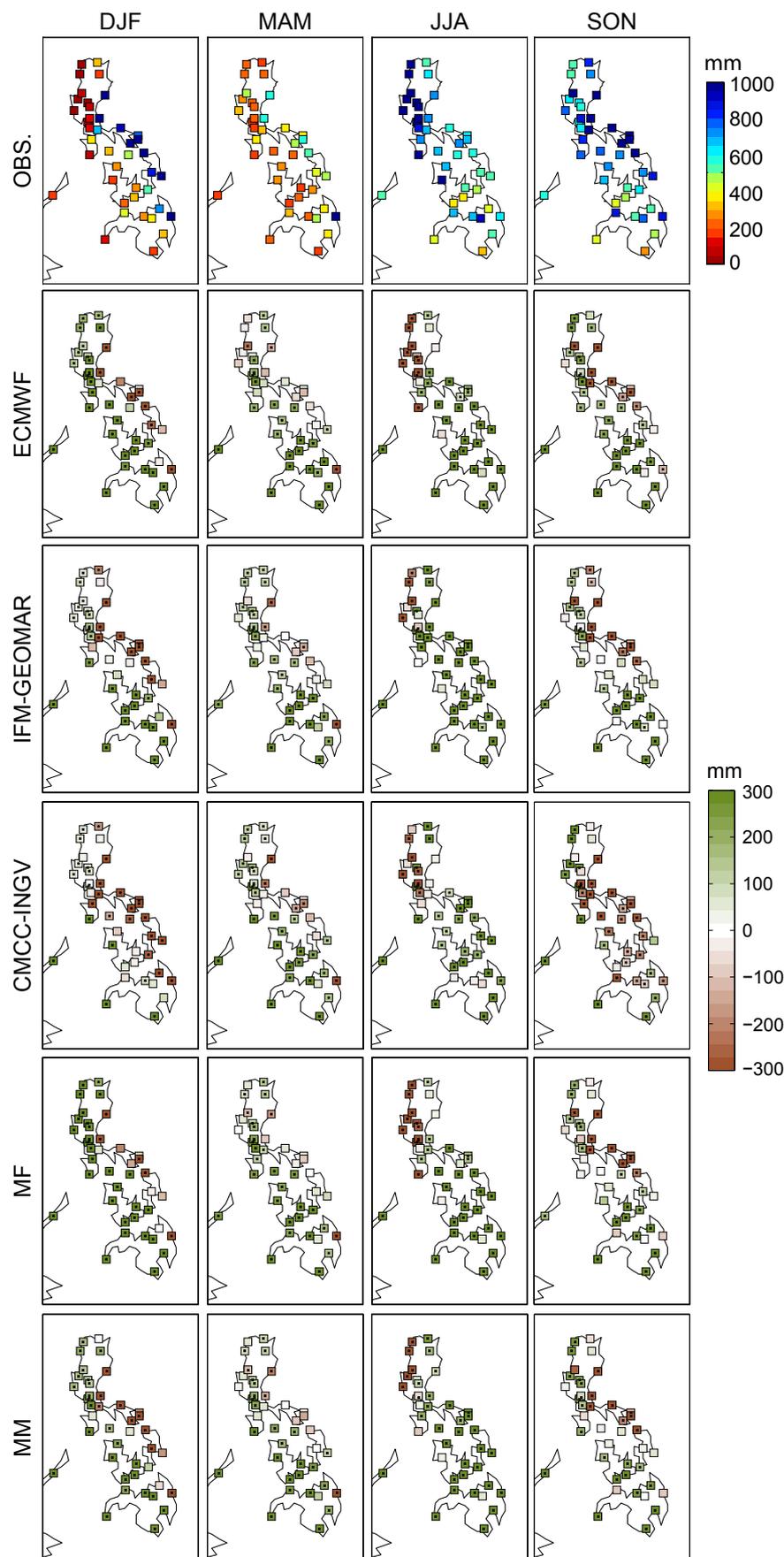


Fig. 3 First row: Observed seasonal climatologies (in mm/season) at the 42 PAGASA stations. Rest of rows: Bias (in mm/season) for the four ENSEMBLES models and the multimodel, by seasons (in columns). Significant ($\alpha = 0.05$, according to a Student's *t*-test) values are indicated with a black dot.

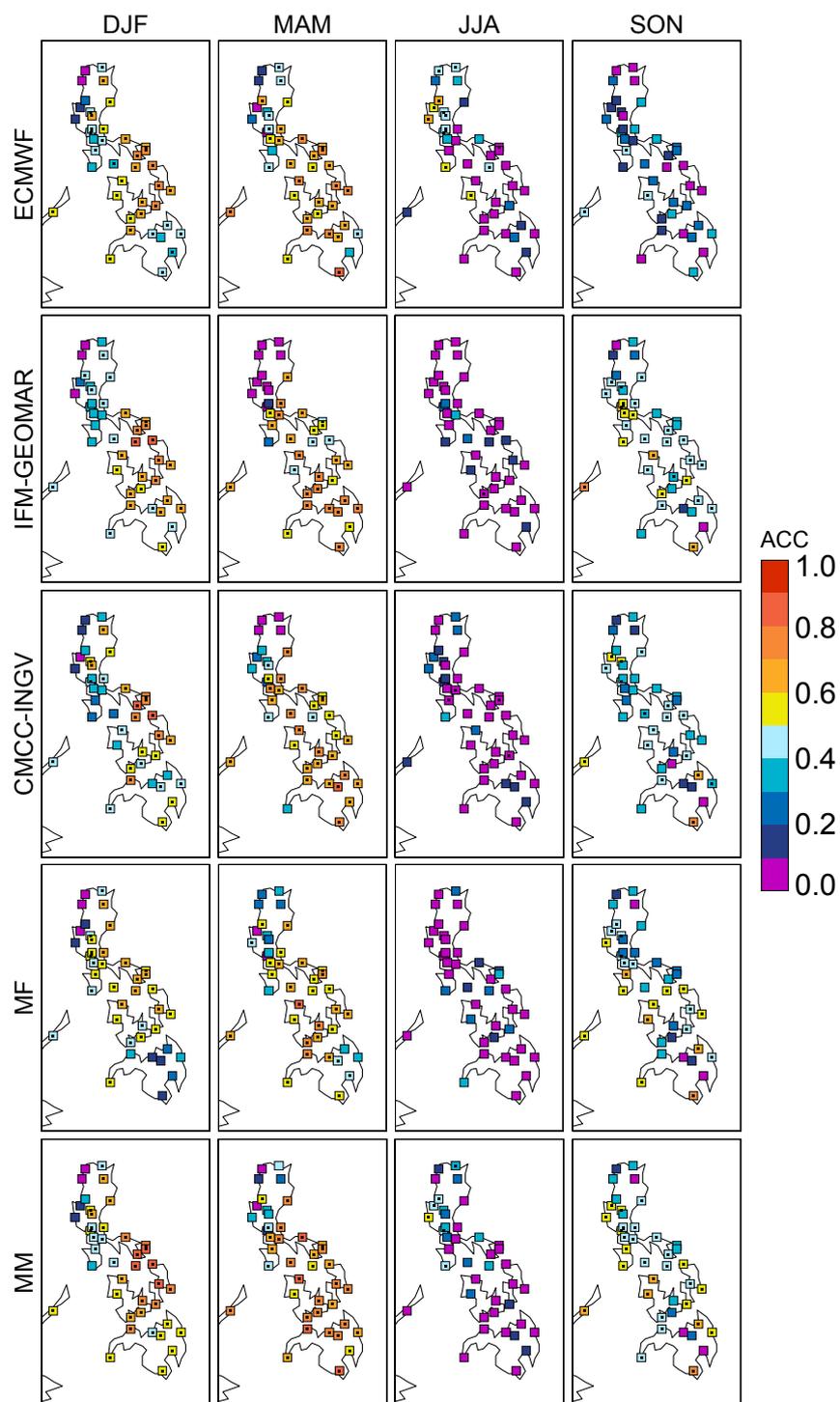


Fig. 4 Interannual ACC values obtained at the 42 PAGASA stations for the four ENSEMBLES models and the multimodel (in rows), by seasons (in columns). Significant ($\alpha = 0.05$, according to a Student's t-test) values are indicated with a black dot.

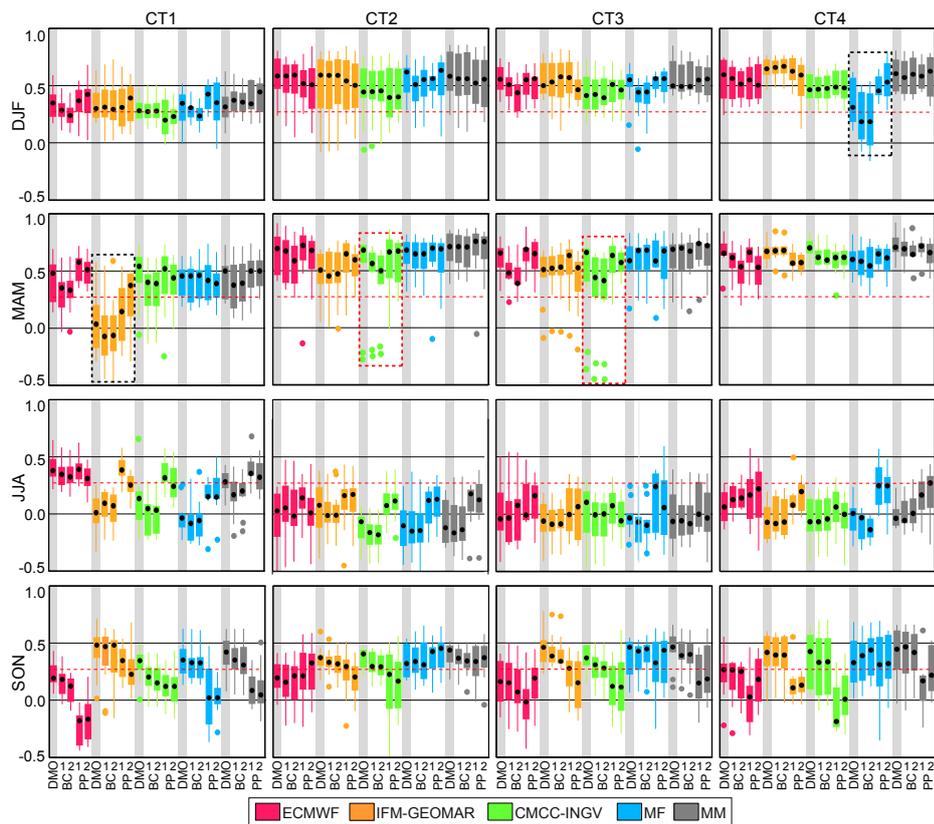


Fig. 5 Interannual ACC obtained for the different seasons (in rows) and CTs (in columns). In each panel, results for each model are shown in different colors (see the legend). From left to right, boxplots display the correlations obtained along the different stations for the DMO (indicated by a light gray shadow) and the BC1, BC2, PP1 and PP2 methods. Significant ($\alpha = 0.1$, according to a Student's t-test) values are those above the red dashed lines. Dashed boxes indicate particular situations which are described in the text.

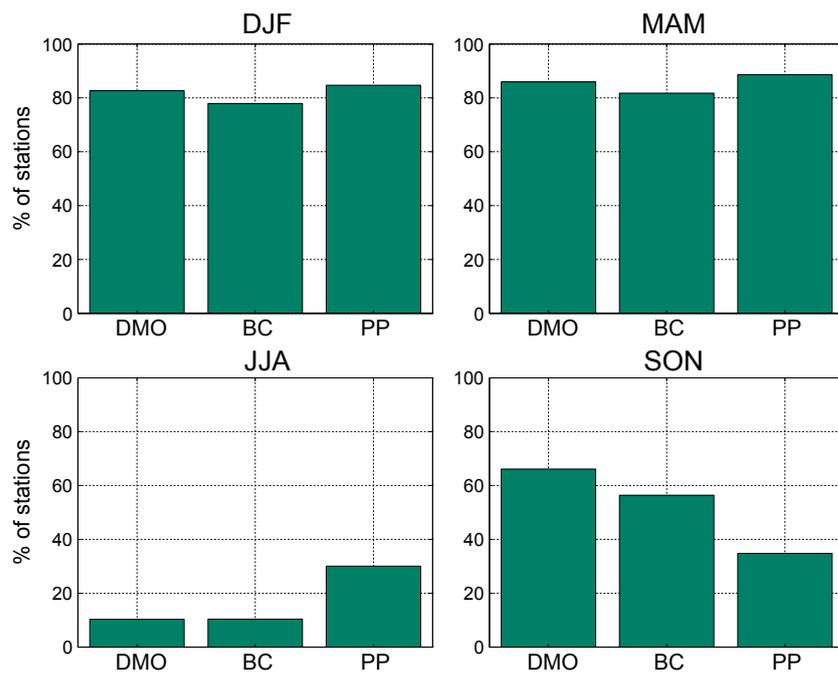


Fig. 6 Summary of Figure 5 showing in bar charts the percentage of stations with significant ($\alpha = 0.1$, according to a Student's t-test) interannual ACC for the DMO and the BC and PP downscaling approaches, for the different seasons. Within each approach, the two methods considered are jointly analyzed. Moreover, all models except the MM (which is excluded for clarity) and all CTs are also jointly considered.

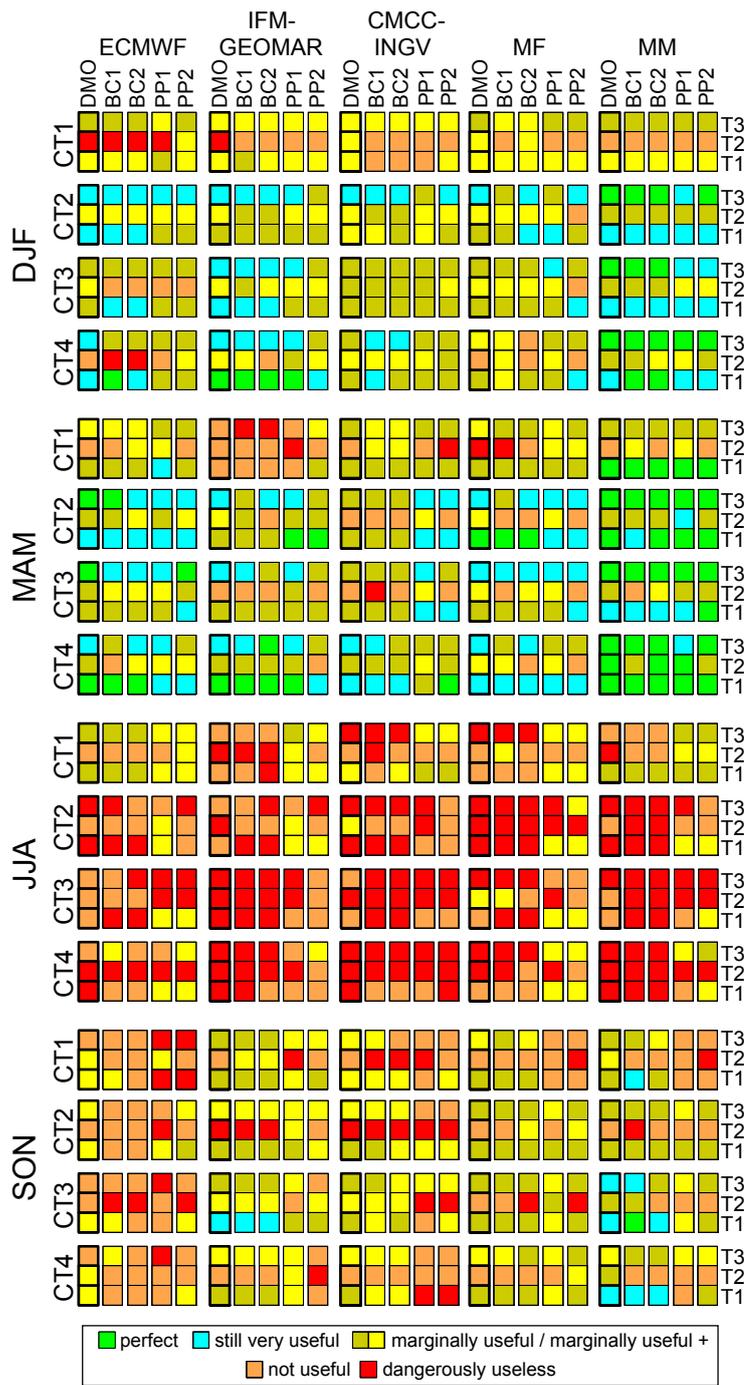


Fig. 7 Reliability categories obtained for the different ENSEMBLES models (in columns) along the different seasons and CTs (in rows). Each block shows the results obtained for the DMO, the two BC and the two PP methods considered, for the three terciles (T1, T2 and T3). Colors correspond to the different categories used, which are based on the original classification proposed by Weisheimer and Palmer (2014) (see the text for details).

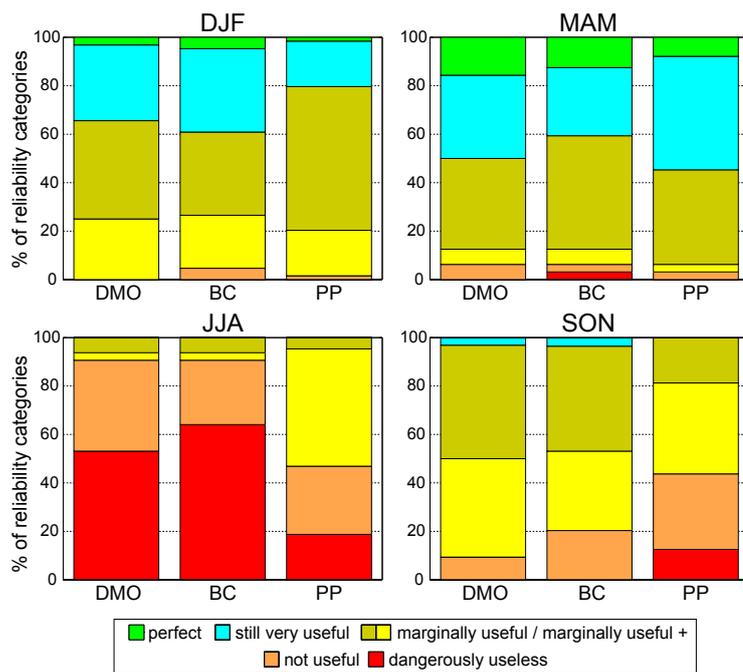


Fig. 8 Stacked bar charts with the percentage of reliability categories (in colors) for the DMO and the BC and PP approaches (within each approach, the two methods considered are jointly analyzed) for the different seasons. For clarity, results from the MM and from the normal tercile (T2) are excluded from this analysis.

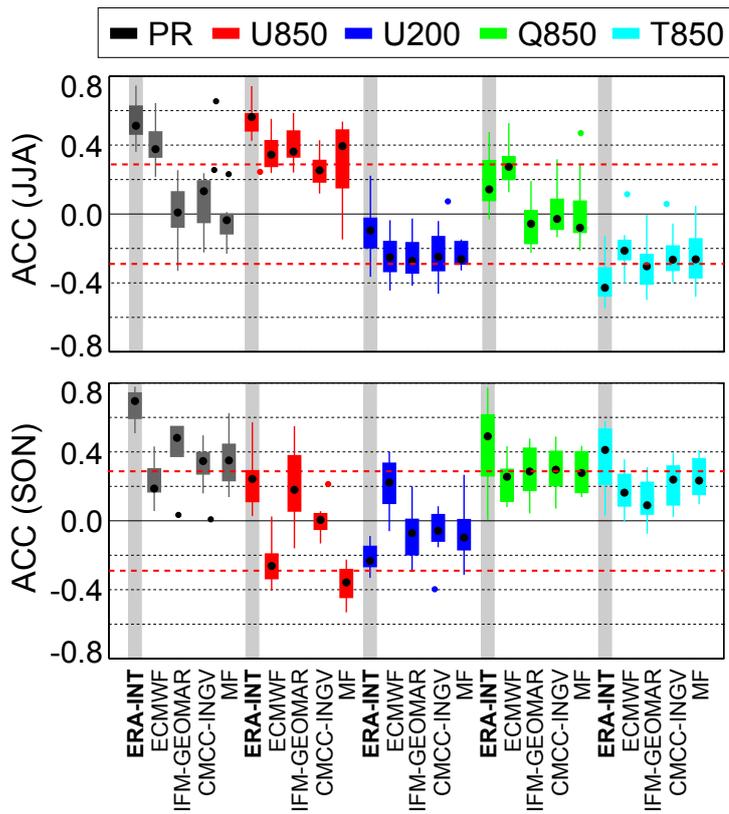


Fig. 9 Interannual ACC values between observed precipitation at the 13 stations pertaining to CT1 and the corresponding ERA-Interim and ENSEMBLES models outputs —the nearest gridbox is considered— for precipitation (PR) and the different predictors used (U850, U200, Q850 and T850) for (top) JJA and (bottom) SON. Significant ($\alpha = 0.1$) positive (negative) values are those above (below) the upper (lower) red dashed line.