

Probabilistic skill in ensemble seasonal forecasts

Leonard A. Smith,^{a,b} Hailiang Du,^{a*} Emma B. Suckling^a and Falk Niehörster^a

^aCentre for the Analysis of Time Series, London School of Economics, UK

^bPembroke College, Oxford, UK

*Correspondence to: H. Du, Centre for the Analysis of Time Series, London School of Economics, London, WC2A 2AE, UK.
E-mail: h.l.du@lse.ac.uk

Simulation models are widely employed to make probability forecasts of future conditions on seasonal to annual lead times. Added value in such forecasts is reflected in the information they add, either to purely empirical statistical models or to simpler simulation models. An evaluation of seasonal probability forecasts from the Development of a European Multimodel Ensemble system for seasonal to inTERannual prediction (DEMETER) and ENSEMBLES multi-model ensemble experiments is presented. Two particular regions are considered: Nino3.4 in the Pacific and the Main Development Region in the Atlantic; these regions were chosen before any spatial distribution of skill was examined. The ENSEMBLES models are found to have skill against the climatological distribution on seasonal time-scales. For models in ENSEMBLES that have a clearly defined predecessor model in DEMETER, the improvement from DEMETER to ENSEMBLES is discussed. Due to the long lead times of the forecasts and the evolution of observation technology, the forecast-outcome archive for seasonal forecast evaluation is small; arguably, evaluation data for seasonal forecasting will always be precious. Issues of information contamination from in-sample evaluation are discussed and impacts (both positive and negative) of variations in cross-validation protocol are demonstrated. Other difficulties due to the small forecast-outcome archive are identified. The claim that the multi-model ensemble provides a ‘better’ probability forecast than the best single model is examined and challenged. Significant forecast information beyond the climatological distribution is also demonstrated in a persistence probability forecast. The ENSEMBLES probability forecasts add significantly more information to empirical probability forecasts on seasonal time-scales than on decadal scales. Current operational forecasts might be enhanced by melding information from both simulation models and empirical models. Simulation models based on physical principles are sometimes expected, in principle, to outperform empirical models; direct comparison of their forecast skill provides information on progress toward that goal.

Key Words: seasonal forecasts; ensemble forecasts; forecast skill; ENSEMBLES; DEMETER

Received 20 June 2013; Revised 13 May 2014; Accepted 19 May 2014; Published online in Wiley Online Library 16 Jul 2014

1. Introduction

Skillful probabilistic forecasting of seasonal weather and climate statistics would be of value in many fields, including agriculture, health and insurance. Since the late 1990s, seasonal forecasting using dynamical models that simulate the coupled atmosphere, ocean and land surface system has become common in operational weather forecasting centres around the world. In recent years, multi-model ensembles have become popular tools to investigate and account for shortcomings due to structural model error in these simulation-model-based predictions on time-scales from days to seasons and centuries (Palmer *et al.*, 2004; Wang *et al.*,

2009; Weisheimer *et al.*, 2009). The potential for using large single-simulation model ensembles or multi-model ensembles depends critically on the forecast information that simulation models add beyond empirically based statistical approaches. Van Den Dool (2007) provides a summary of these empirical models, sometimes referred to as surrogate prediction generators (Smith, 1992) or empirical benchmarks (Suckling and Smith, 2013). Contrasting the skill of empirical models with simulation models can also be informative regarding structure model error in the simulation models.

The need for a consistent experimental design for the assessment of skill in multi-model seasonal forecasting has been embraced by two large European projects in the last decade. These

projects provided the basis for subsequent multi-model designs for operational seasonal-to-decadal forecasting (Vitart *et al.*, 2007; Kirtman *et al.*, 2013). The earlier European project, initiated in 2000, was Development of a European Multimodel Ensemble system for seasonal to interannual prediction (DEMETER: Palmer *et al.*, 2004; Doblas-Reyes *et al.*, 2005; Hagedorn *et al.*, 2005), in which a consistent framework was developed to conduct multi-model seasonal forecasting with a set of general circulation models (GCMs). A similar framework was adopted in ENSEMBLES (Hewitt and Griggs, 2004; Weisheimer *et al.*, 2009; Doblas-Reyes *et al.*, 2010), which produced the next generation of seasonal hindcast (or retrospective forecast) simulations, using updated model versions. Further details of the ENSEMBLES and DEMETER experiments can be found in Tables S1 and S2 in File S1.

The multi-model ensemble simulations from these projects provide a basis for the quantification of skill in GCM forecasts and an opportunity to assess the benefit of using multi-model ensembles (Weisheimer *et al.*, 2009; Alessandri *et al.*, 2011) over other approaches, such as forecasts based on statistical models (Smith, 1992; van Oldenborgh, 2005; Coelho *et al.*, 2006; Van Den Dool, 2007; Suckling and Smith, 2013). Furthermore, the consistency between the experimental design of the DEMETER and ENSEMBLES seasonal forecasts makes it possible to quantify the improvement of skill or, in other words, the additional information gained from the forecasts due to model development in the intervening period between the two projects. While evaluations of skill between individual model versions may exist in-house at forecast centres, the authors are unaware of any systematic comparison across centres and model versions. The analysis presented below allows direct comparisons between the relative performance and improvement in different models.

Two particular regions are considered. As a coupled atmospheric and oceanic phenomenon, the El Niño/Southern Oscillation (ENSO) in the tropical Pacific is the dominant mode of seasonal and interannual climate variability. Sea-surface temperatures (SSTs) in the Niño3.4 region at seasonal time-scales provides an indicator for the ENSO phenomenon. SSTs in the Main Development Region (MDR) over the North Atlantic provide an indicator for hurricane activity over the coming season. This article focuses on probability forecast skill in these two regions.* Probabilistic skill of seasonal forecasts from both DEMETER and ENSEMBLES are evaluated and contrasted. In each case, ensembles of GCM simulations are transformed into probabilistic distributions via kernel-dressing (see Bröcker and Smith, 2007) and blended with the climatological distribution to provide calibrated seasonal forecasts; this approach has influenced operational forecasting (Hagedorn and Smith, 2009; Met Office, 2013). Evaluating probability forecasts as probability forecasts, rather than computing summary statistics of the ensemble mean, allows clearer consideration of the uncertainties sampled by the multi-model ensemble. It is also more easily interpreted in terms of the value, or information content, of the forecast from a decision-maker's perspective.

An overview of the DEMETER and ENSEMBLES multi-model experiments used to evaluate seasonal forecast skill over the Niño3.4 and MDR regions is given in section 2 and the approach to generating probabilistic forecasts and evaluating them is described in section 3. In section 4, probabilistic skill above that of the climatological distribution is demonstrated up to a lead time of 7 months for SSTs over the Niño3.4 region and up to a lead time of 2 months for SSTs over the MDR. In section 5, forecasts from the ENSEMBLES models show improvements in skill compared with those from DEMETER for each of the models common to both projects. Broadly speaking, these results are

consistent with previous evaluations of skill from the DEMETER and ENSEMBLES projects (Weisheimer *et al.*, 2009; Alessandri *et al.*, 2011), in which improvements in the anomaly correlation, root-mean-square error (RMS) and Brier scores from DEMETER to ENSEMBLES were reported for SSTs over the tropical Pacific and some other regions up to 6 months ahead. Section 6 shows that, somewhat surprisingly, competitive results can be formed from purely empirical probability forecasts based on persistence. A similar result has been found for decadal forecasts (specifically, probability forecasts of annual mean values on lead times of 1–10 years) by Suckling and Smith (2013), who demonstrate that some empirical models often outperform the ENSEMBLES models on these decadal scales. The illustrations presented in section 7 suggest that increasing the ensemble size of future multi-model experiments could provide an efficient way of improving forecast skill, while sections 8 and 9 highlight the motivation for using proper scoring rules and the challenges involved in model combination to produce multi-model ensemble forecasts, respectively. Section 10 discusses the issues of information contamination when data are precious. The key results and conclusions are summarized in section 11.

2. The seasonal multi-model ENSEMBLES forecasts

The ENSEMBLES multi-model ensemble experiment for seasonal-to-annual forecasting comprises global coupled atmosphere–ocean climate models from the UK Met Office (UKMO), Météo-France (MF), the European Centre for Medium-Range Weather Forecasts (ECMWF), the Leibniz Institute of Marine Sciences at Kiel University (IFM-GEOMAR) and the Euro-Mediterranean Centre for Climate Change (CMCC-INGV) in Bologna (Doblas-Reyes *et al.*, 2010). In each case, the ensemble simulations include all the major radiative forcings; none of the coupled models has flux adjustments (Hewitt and Griggs, 2004; Weisheimer *et al.*, 2009; Doblas-Reyes *et al.*, 2010). A set of seasonal hindcast simulations cover the 46 year period from 1960–2005. For each launch date, the atmosphere and ocean for each model were initialized using realistic estimates of their observed states, providing an ensemble consisting of nine initial condition ensemble members for each model. Hindcast simulations were launched on the first days of February, May, August and November each year over the hindcast period and run for 7 months. This set of 46 seasonal forecasts for each launch date is analyzed below. Additionally, each model, with the exception of CMCC-INGV, was run for an extended period up to a lead time of 14 months from the November launch.

Improvements made in the ENSEMBLES multi-model forecasting system include a better representation of subgrid-scale physical processes in the simulation models, the inclusion of interannual variability in the greenhouse gas forcing and the use of improved ocean data assimilation, based on quality-controlled *in situ* ocean temperature and salinity profiles for the construction of the initial conditions (Ingleby and Huddleston, 2007; Weisheimer *et al.*, 2009). Given two simulation models from the same modelling centre, the experimental designs are sufficiently consistent to allow a direct comparison between the skill of seasonal forecasts from each version of the system. Further details of the models used for the DEMETER and ENSEMBLES projects are provided in Tables S1 and S2 of File S1.

3. Defining probabilistic forecast skill

Simulations from dynamical models are often used to make probabilistic predictions with the aim of providing useful information for decision support. Evaluating the performance of these predictions, as well as understanding the sources of skill, is crucial for guiding decision-makers in understanding in which regions and on which time-scales of interest the models are likely to be informative and, perhaps more importantly, clarifying when they are likely to be misinformative. Only proper scoring rules

*Attention was restricted to these two regions, prior to examination of forecast skill in any other regions. This approach eases interpretation of the statistical significance of the results obtained over studies that examine the entire globe and then focus analysis on areas with 'significant' skill.

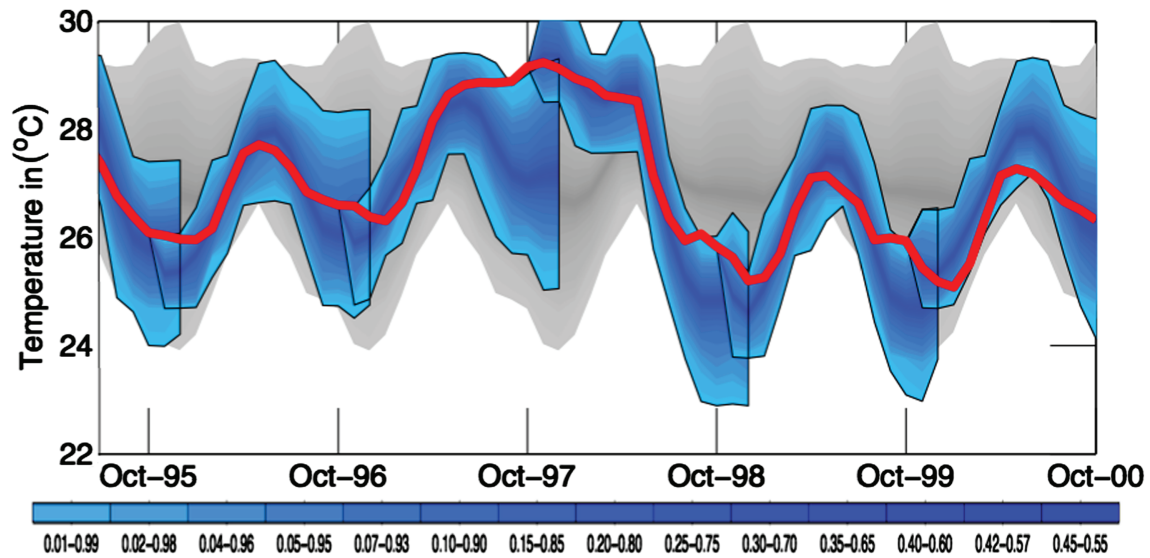


Figure 1. Probabilistic forecast distributions for the IFS (ECMWF) hindcast simulations from ENSEMBLES for the Nino3.4 index, launched in November over the period 1995–2000. The blue shaded regions indicate the forecast percentiles between 1 and 99% and the red line shows the observed outcome from the ERA40 reanalysis. The grey shaded intervals show the percentiles for the climatological distribution.

offer appropriate, clear measures of probabilistic forecast skill (Bröcker and Smith, 2006; Wilks, 2005).

I. J. Good's logarithmic score (Ignorance: see Good, 1952; Roulston and Smith, 2002; Bröcker and Smith, 2006) is unique among several scoring rules (Wilks, 2005) designed for evaluating the skill of probabilistic forecasts. It is the only proper and local score[†] for continuous variables (see Bernardo, 1979; Raftery *et al.*, 2005; Bröcker and Smith, 2006). The Ignorance Score is defined by

$$S(p(y), Y) = -\log_2(p(Y)), \quad (1)$$

where Y is the observed outcome and $p(y)$ is the density function of the forecast distribution. Ignorance has a clear interpretation in terms of gambling returns (see Good, 1952; Kelly, 1956; Roulston and Smith, 2002): under a certain betting scenario, 'Kelly Betting' (Kelly, 1956), the Ignorance describes the expected rate at which the forecaster's wealth changes with time. Through its close relation to Shannon's information entropy, Ignorance can also be related to the amount of information expected from a forecast (see Roulston and Smith, 2002). It is easily communicated as an effective interest rate (see Hagedorn and Smith, 2009).

In practice, given K forecast-outcome pairs, $(p_i, Y_i, i = 1, \dots, K)$, the empirical Ignorance score is

$$S_E(p(y), Y) = \frac{1}{K} \sum_{i=1}^K -\log_2(p_i(Y_i)). \quad (2)$$

Relative Ignorance reflects the performance of (a set of) forecasts p from one model relative to those of a reference forecast p_{ref} :

$$S_{\text{rel}}(p(y), Y) = \frac{1}{K} \sum_{i=1}^K -\log_2[(p_i(Y_i))/p_{\text{ref}}(Y_i)]. \quad (3)$$

The Relative Ignorance of two forecast systems quantifies the information gain (in terms of bits) that the model forecast system provides over the reference system. In other words, Ignorance reflects the (average) increase in probability density that the model forecast placed on the outcome relative to that of the reference forecast. By convention, Ignorance is a negatively oriented score,

[†]'Proper' meaning that it cannot be optimized by hedging the probabilistic forecasts toward other values against the forecaster's true belief (Bröcker and Smith, 2006; Weigel *et al.*, 2008). 'Local' meaning that the score depends solely on the probability assigned to the outcome, rather than being rewarded for other features of the forecast distribution, such as its shape.

which means that the smaller the score, the more skilful the forecasts. An Ignorance score of $S_{\text{rel}} = -1$ means that, on average, forecasts from the model assign twice the probability density to the outcome compared with the reference forecast, while $S_{\text{rel}} = -2$ indicates a fourfold (2^2) increase. Suitable references could include the climatological distribution, a probability forecast from a statistical model or forecasts from another GCM. The climatological distribution provides the primary benchmark for seasonal forecast skill in this article; see however section 6.

Probability forecasts are generated from the DEMETER and ENSEMBLES simulations via kernel-dressing and are blended with climatology to produce seasonal probability forecasts (for a full description, see Bröcker and Smith, 2007 and Appendix A). The climatological distribution is estimated by kernel-dressing all available historical observations under cross-validation (see Appendix B). Figure 1 shows an example of the kernel-dressed and blended probabilistic forecast distributions for a subset (over the period 1995–2000) of the IFS (ECMWF) hindcast simulations from ENSEMBLES for the Nino3.4 index, launched in November. The blue shaded regions indicate the forecast percentiles between 1 and 99% and the red line shows the observed outcome (from the ERA40 reanalysis) for comparison. The grey shaded bands show the percentiles between 1 and 99% for the climatological distribution.

The empirical Ignorance score of the dressed and blended GCM forecasts is then computed as a function of lead time (in months) for SSTs over the MDR and Nino3.4 regions, relative to the climatology in section 4. Forecasts from each of the ENSEMBLES models are contrasted with those of DEMETER in section 5.

4. ENSEMBLES seasonal forecast skill

Figures 2 and 3 show the skill of probability forecasts from each of the models and launch dates available in the ENSEMBLES seasonal forecast project. Figure 2 shows empirical Ignorance scores for forecasts of the Nino3.4 index as a function of lead time in months, relative to climatology. Each of the four panels corresponds to a different forecast launch month (as indicated). In general, at short lead times all the models are substantially more skilful than climatology (i.e. a negative relative Ignorance) for all four initialization dates. This result is generally consistent with Weisheimer *et al.* (2009), who reported that anomaly correlation skill for the multi-model ensemble mean was found to decay with lead time over the Nino3 region, to ~ 0.5 up to 14 months ahead. At longer lead times, ENSEMBLES models show systematically less skill than at early lead times, as expected. In each case, however, the

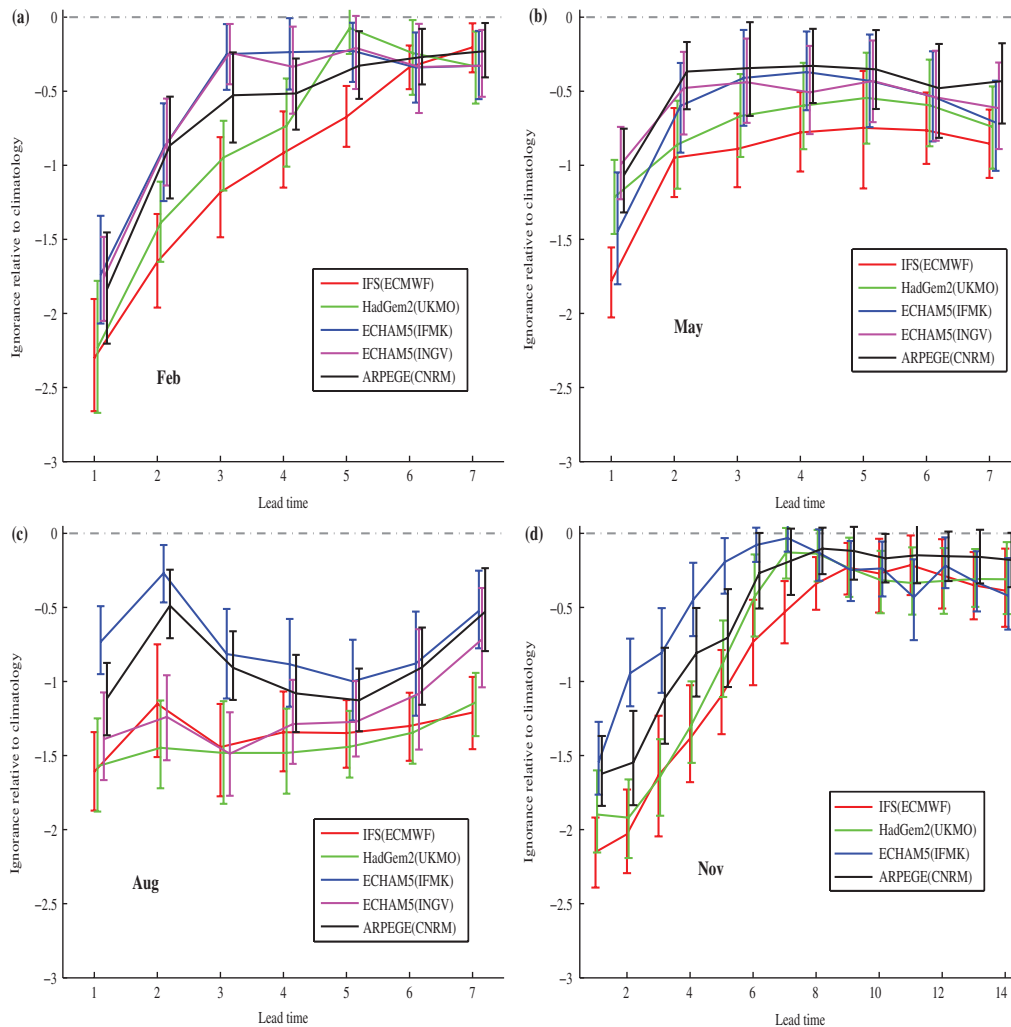


Figure 2. Ignorance score of each model from ENSEMBLES for the Nino3.4 index relative to climatology as a function of lead time in months. The four different panels show the hindcasts initialized in (a) February, (b) May, (c) August and (d) November. Zero Ignorance indicates that a model has no skill relative to climatology and negative relative Ignorance scores suggest that a model is more skilful than climatology. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. All models show significantly more skill than climatology up to a lead time of 5 months, regardless of when the forecasts are launched. For (d) the November launch, the bootstrap resampling intervals often cross the zero skill line beyond a lead time of 6 months.

simulation models demonstrate skill above the climatology up to a lead time of 7 months. For the hindcasts launched in November, some skill appears up to a lead time of 14 months (although an alternative cross-validation protocol casts some doubt on this result: see section 10). At longer lead times, relative Ignorance scores of approximately -0.25 are found for most models, which translates into the simulation models placing, on average, $\sim 19\%$ more probability density on the outcome compared with the climatological distribution. The IFS (ECMWF) and HadGEM2 (UKMO) models often score slightly lower (are more skilful) than the other three models. The sampling uncertainty across forecast launches is represented by a bootstrap resampling procedure, which resamples the set of forecast Ignorance scores for each model, with replacement. The bootstrap resampling intervals are shown as vertical bars in each of the figures as a 5–95% interval.

Figure 3 shows the Ignorance score as a function of lead time for SSTs over the MDR relative to climatology. Compared with the Nino3.4 index, hindcasts of SSTs in the MDR are less informative at all lead times, particularly for the forecasts launched in November, the performance of which decreases significantly within the first 2 months. Despite the higher Ignorance scores (lower skill), the GCM hindcasts for the MDR demonstrate significant skill relative to climatology up to 7 months ahead for most models and launch dates, with the exception of the November launch. Comparison with alternative benchmarks, like the persistence forecast, shows much larger variation than altering the cross-validation scheme.

In Figures 2 and 3, two models with similar bootstrap resampling intervals might be misinterpreted to suggest that neither model is significantly better than the other. Bootstrap resampling skill against climatology is misleading if interpreted incorrectly. One model can systematically outperform a second model on every forecast, yet the resample ranges in the skill relative to climatology may overlap. The relative Ignorance between two models, on the other hand, provides a clear result reflected in bootstrap resampling from the model–model relative scores.

Figure 4 shows the Ignorance of each of the ENSEMBLES models for the Nino3.4 index relative to the IFS (ECMWF) model. There are indeed some cases where the IFS (ECMWF) model outperforms all other models despite the overlapping bootstrap resampling intervals in Figure 2. For example, the IFS (ECMWF) model systematically outperforms the ARPEGE (CNRM), ECHAM5 (INGV) and ECHAM5 (IFMK) models, particularly at early lead times for most launch dates. In the case analyzed above, there is substantial information in the forecasts from the ENSEMBLES models for the Nino3.4 index, even at longer lead times; the IFS (ECMWF) model shows higher skill (often exceeding 0.5 bits in the first 6 months) relative to the other seasonal forecast models used in ENSEMBLES.

5. Contrasting skill of ENSEMBLES and DEMETER

The methods and models used for the seasonal hindcast experiments in the ENSEMBLES project were developed in light of the experience gained and models available from the DEMETER

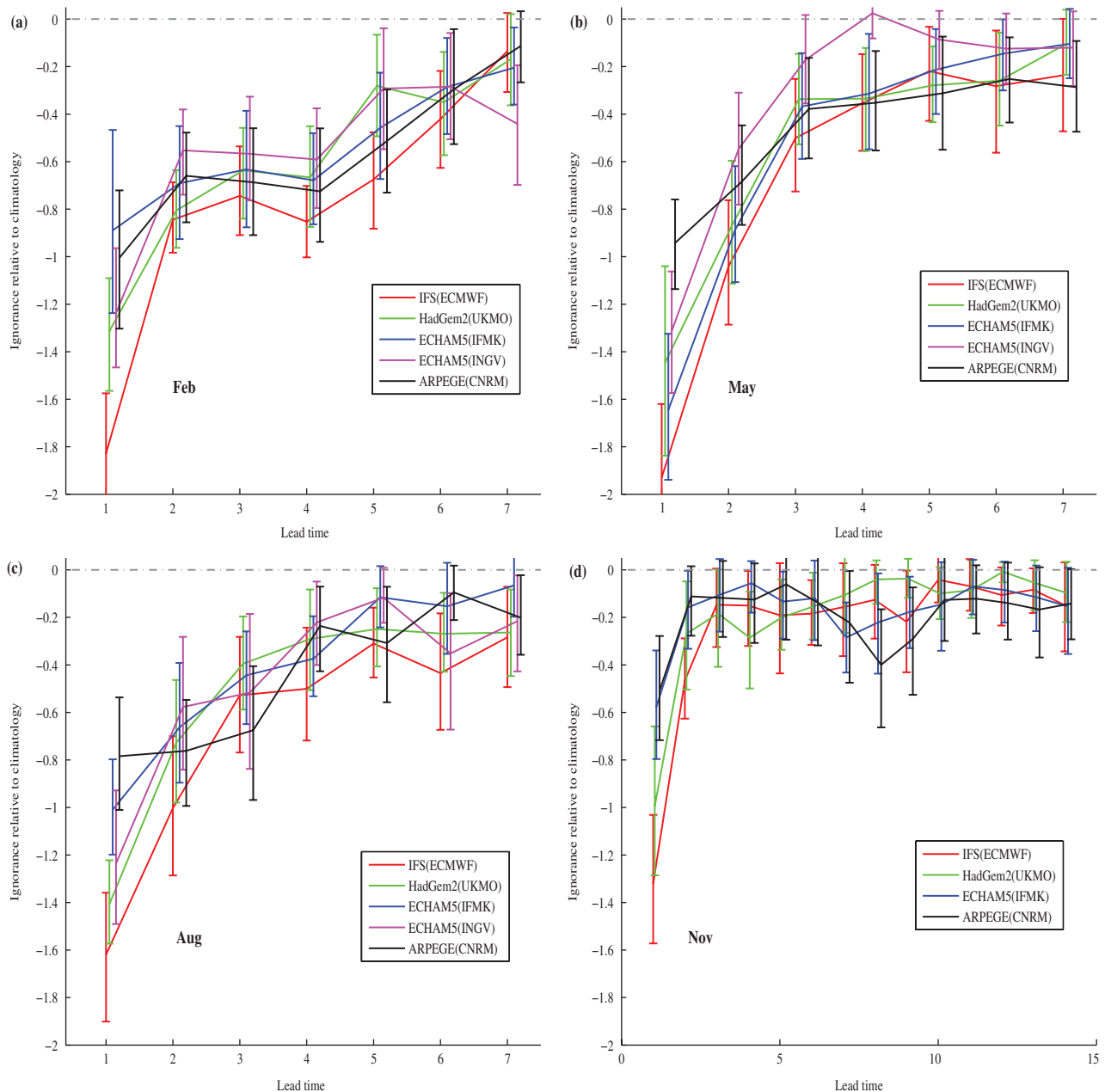


Figure 3. Ignorance score of each model from ENSEMBLES for the MDR index relative to climatology as a function of lead time in months. The four different panels show the hindcasts initialized in (a) February, (b) May, (c) August and (d) November. Zero Ignorance indicates that a model has no skill relative to climatology and negative relative Ignorance scores suggest that a model is more skilful than climatology. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. Significant skill above climatology is demonstrated for most models and launch dates at early lead times (up to 6 months for the February launches, for example), with the exception of the November forecast launches, where the bootstrap intervals overlap the zero-skill climatology beyond a lead time of 2 months.

project. The DEMETER seasonal hindcasts and ENSEMBLES hindcasts for the same verification period provide an opportunity to measure the improvement of forecast skill after 4 years of model development. Such an evaluation is aided by the similarities in the experimental design between the two projects.

Figure 5 shows the Ignorance score of each of the DEMETER model forecasts for the Nino3.4 index relative to climatology. With the exception of ECHAM5 (MPI), each model appears substantially more skilful than climatology at all lead times and for all four initialization dates. The lack of skill demonstrated by the ECHAM5 (MPI) model reflects the fact that when its ensemble members are dressed and blended with climatology (see Appendix A), they are assigned relatively little weight (i.e. the forecast is virtually the climatological distribution). There is little or no contribution from the ECHAM5 (MPI) model ensemble to the calibrated forecast beyond a lead time of 3 months. This is particularly true for the November launch, in which the forecast blending parameter as a function of lead time, α , takes values [$\alpha = 0.90, 0.81, 0.02, 0.00, 0.00, 0.00$], respectively.

In order to measure the improvement of forecast performance due to model development from the DEMETER to the ENSEMBLES project, the Ignorance of the forecast distributions derived from pairs of model simulations from each project is compared. Although seven European simulation models were used in the DEMETER project, only those models that correspond to earlier ‘versions’ of those used in ENSEMBLES are considered.

Figure 6 shows the Ignorance for seasonal forecasts of the Nino3.4 index forecasts from the ENSEMBLES models relative to those of the corresponding DEMETER models. In general, the relative Ignorance scores in Figure 6 demonstrate improvements for ENSEMBLES (negative relative Ignorance scores) for most lead times and for most models. The ECHAM5 (INGV) model is an exception to this finding; the reduction in skill for this model is consistent with Alessandri *et al.* (2010), who showed that subsurface data assimilation for ocean initialization degraded prediction skill over the tropical Atlantic. The ECHAM5 (IFMK) model shows substantial improvements, up to one bit, at early lead times, particularly for forecast launches in February and May (the

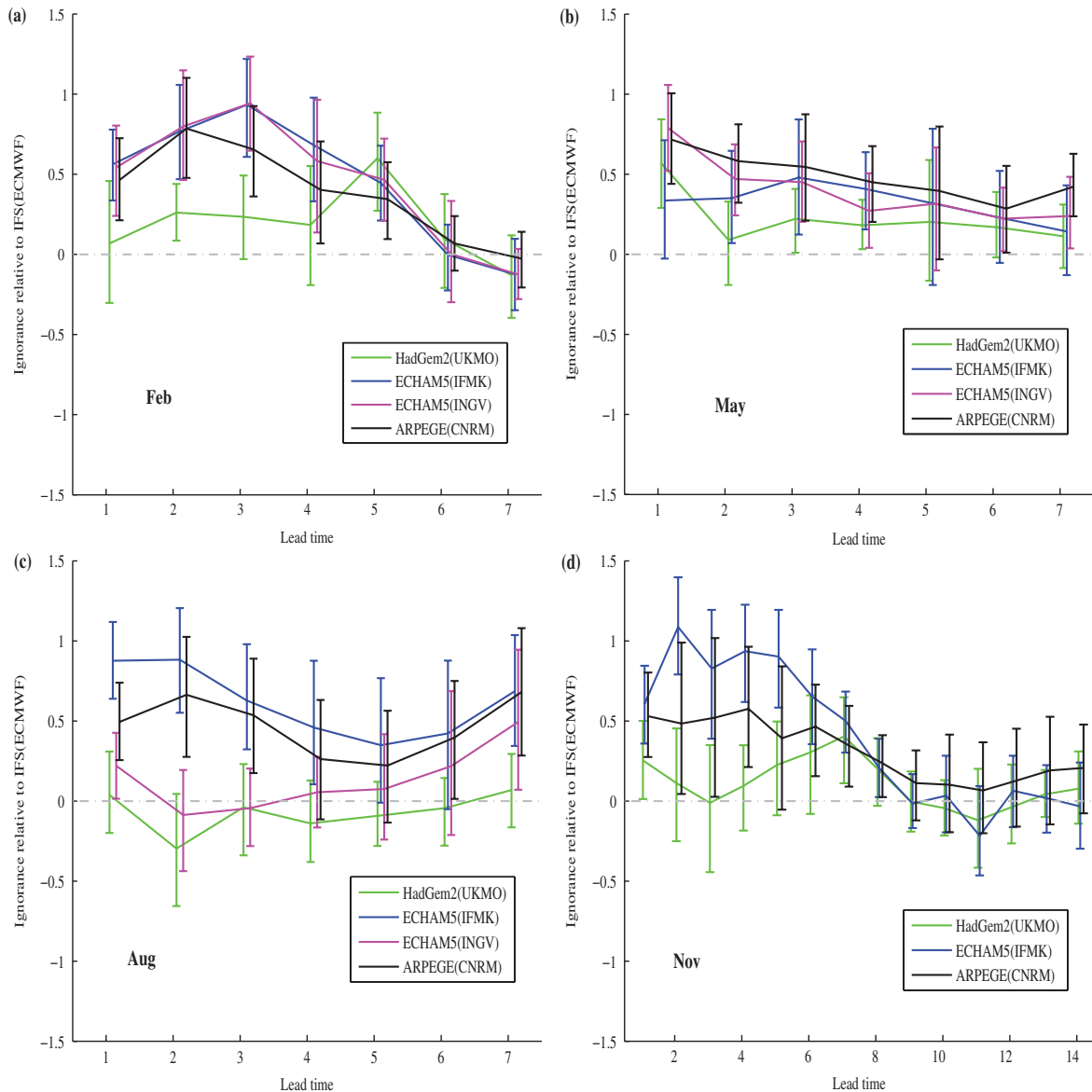


Figure 4. Ignorance score of the ENSEMBLES model forecasts for the Nino3.4 index relative to the IFS (ECMWF) model as a function of lead time in months. Zero Ignorance indicates that a model has no skill relative to the IFS (ECMWF) model and negative relative Ignorance scores suggest that a model is more skilful than the IFS (ECMWF) model. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. All models shown are typically less skilful than IFS (ECMWF) at all lead times and for most forecast launch dates. For launch dates in August, however, the IFS (ECMWF) model is shown to perform neither significantly better nor significantly worse than HadGEM2 (UKMO) and ECHAM5 (INGV).

ENSEMBLES model placing twice the probability density on the outcome compared with the DEMETER model). Improvements are also demonstrated at lead times beyond 3 months for forecasts launched in August, particularly for the ECHAM5 (IFMK) and HadGEM2 (UKMO) models.

6. Contrasting ENSEMBLES seasonal skill with persistence forecasts

In the previous sections, the climatological distribution was used as a benchmark against the performance of the ENSEMBLES and DEMETER seasonal hindcasts. Whilst comparing skill between simulations from dynamical models and climatology provides insight into the information gained from forecasting with those dynamical models, other simple empirical models can also serve as appropriate benchmarks to model performance (Smith, 1992; Suckling and Smith, 2013). A probabilistic persistence forecast provides an interesting benchmark accounting for the effects of both physical persistence and any long-term drift in the temperature of the target region. Whether the additional skill in the ENSEMBLES models over the Nino3.4 region compared with the MDR is related to the strong persistence of ENSO can be investigated by looking at the performance of forecasts

over these two regions relative to a persistence model.[‡] The persistence forecasts generated here use the observed SST value over the chosen region in the month prior to the forecast launch, persisted forward in time and transformed into a probabilistic distribution using kernel-dressing parameters that vary with lead time (as described in Suckling and Smith, 2013). While more complex persistence models could be constructed easily, this simple version is sufficient for our purpose here.

Figure 7 shows the Ignorance score of each of the ENSEMBLES models for the Nino3.4 index relative to persistence. For forecasts launched in February, most of the ENSEMBLES models are significantly more skilful than persistence at all lead times. For launch dates in August and November, little if any information is added compared with the persistence forecasts for most models at any lead time. In fact, at early lead times (up to 3 months ahead), persistence outperforms the ECHAM5 (IFMK) and ARPEGE (CNRM) models. At moderate lead times for the August launch and most lead times in the May launch, on the other hand, the IFS (ECMWF) and HadGEM2 (UKMO) models outperform persistence.

[‡]We are very grateful to an anonymous reviewer for suggesting this comparison.

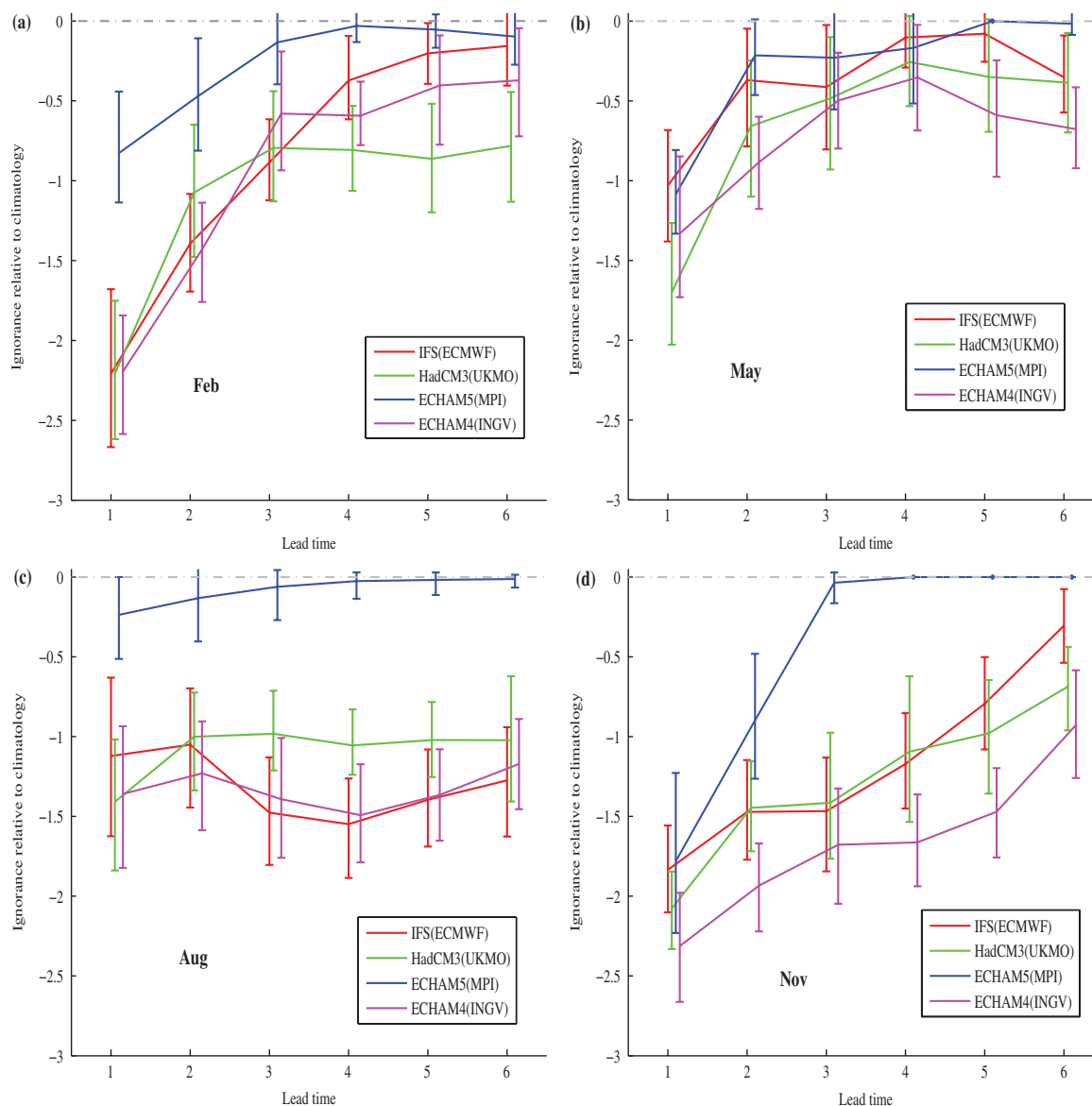


Figure 5. Ignorance score of each model from DEMETER for the Nino3.4 index relative to climatology as a function of lead time in months. Zero Ignorance indicates that a model has no skill relative to climatology and negative relative Ignorance scores suggest that a model is more skilful than climatology. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. All models, with the exception of ECHAM5 (MPI), are significantly more skilful than climatology at most lead times, particularly for forecasts launched in August and November. At lead times beyond 4 months, for forecasts launched in November, the ECHAM5 (MPI) model is given zero weight when blended with the climatological distribution.

Figure 8 shows the corresponding results for the MDR index relative to a probabilistic persistence forecast. In this case, the ENSEMBLES models and persistence have similar skill, with no one model emerging as significantly better than another. These comparable levels of skill suggest that blending statistical model output with simulation model output is likely to add value to seasonal forecasts.

7. More models or more members?

Knowledge of the relationship between ensemble size and forecast quality aids forecast system design. The cost of increasing the number of ensemble members is typically small relative to the cost of model development. The cost of increasing the ensemble size increases only (nearly) linearly. It is often true that the quality of the forecast increases with the number of ensemble members as well; however, this improvement in forecast skill depends on both the current ensemble size and the quality of that model's ultimate distribution. The seasonal forecasts from the ENSEMBLES project provide an opportunity to investigate the relationship between ensemble size and forecast quality. This analysis would be eased, for example, had one launch date included an increased number

of members so that the value of additional members could be tested more directly.[§]

Figure 9 shows the effect of decreasing the number of ensemble members on the forecast skill for the Nino3.4 index from the IFS (ECMWF) model launched in November. The skill of two-member ensembles (red) and four-member ensembles (green) is shown relative to the full nine-member ensemble (the zero line), both as a set of random draws from the nine original members without replacement (Figure 9(a)) and as the average Ignorance of all two- or four-ensemble member combinations (Figure 9(b)). In Figure 9(a), most two- and four-member combinations show less skill than the full nine-member ensemble, with only a few ensemble member combinations scoring better than the original ensemble now and then. Figure 9(b) shows that decreasing the number of ensemble members systematically decreases the average skill (i.e. increases the Ignorance score) across all lead times. This result holds both when decreasing from

[§]Quantifying the value added by including an $N + 1$ st additional ensemble member requires some subset of forecasts running more than N members. Although discussed early in the ENSEMBLES project, the decision was taken to use nine member ensembles throughout.

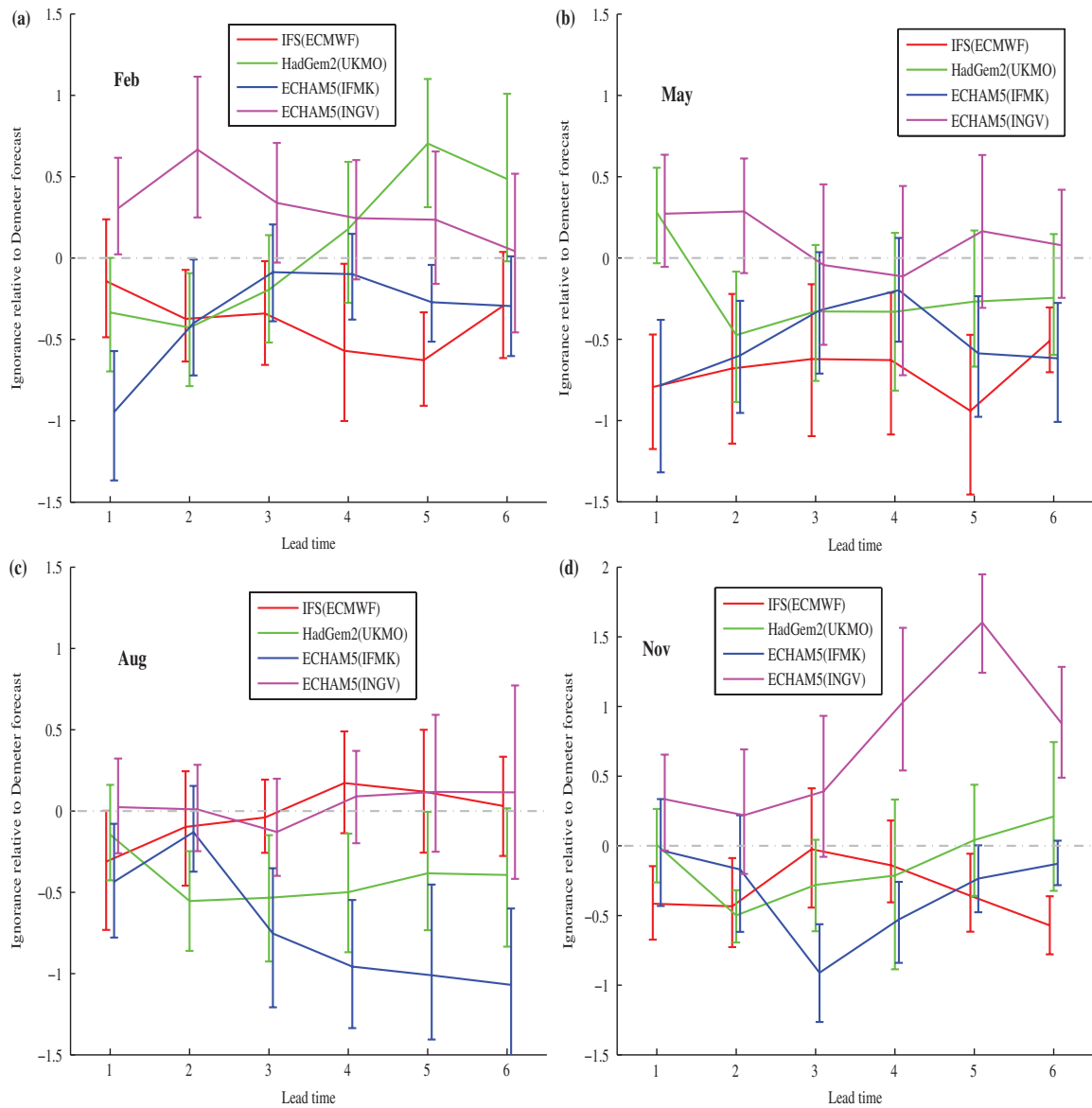


Figure 6. Ignorance score of each model from ENSEMBLES for the Nino3.4 index relative to the corresponding DEMETER forecasts as a function of lead time in months. Zero Ignorance indicates that an ENSEMBLES model has no added skill relative to the corresponding DEMETER model and negative relative Ignorance scores suggest that the ENSEMBLES model is more skilful than the corresponding DEMETER model. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. The ENSEMBLES models typically demonstrate improvements, of up to one bit in some cases, over their corresponding DEMETER models. ECHAM5 (INGV) is an exception to this improvement and is shown to perform worse in ENSEMBLES than its DEMETER model version.

nine members to four members and when decreasing from four to two ensemble members. At a lead time of 6 months, where the IFS (ECMWF) model still has non-trivial skill relative to climatology (Figure 2), for example, the two-member forecast places $\sim 7\%$ and the four-member ensemble places $\sim 3\%$ less probability density on average on the outcome[†] relative to the nine-member ensemble (Figure 9(b)). This result suggests that increasing the current ensemble size of nine would further improve the forecast performance.[‡]

A larger ensemble could be obtained by either increasing the number of ensemble members from one particular model or, alternatively, combining simulations from different models to form a multi-model ensemble (see Palmer *et al.*, 2004; Weigel *et al.*, 2008). Of course, developing a new, ideally independent model is more costly than increasing the number of ensemble members from an existing model. Combining the output of different (independent) models might, however, have the added

advantage of reducing the impact of systematic bias of any single model.^{**} One might therefore reasonably expect to obtain significantly more information by using multi-model outputs than by increasing the number of ensemble members from a single model.

Figure 10 shows the Ignorance score for a set of multi-model forecasts, in which ensemble members from each of the different ENSEMBLES models are treated equally (i.e. each ensemble member is assigned equal weight). Here, the nine-member IFS (ECMWF) forecasts define the zero line. Figure 10(a) shows the Ignorance score for forecasts built from multi-model ensembles containing four members randomly drawn from the 36 available ensemble members (nine members from each of four models) without replacement. Similarly, Figure 10(b) shows the skill of multi-model ensembles containing nine randomly drawn members. The blue line in each case shows the skill of the full

[†]Under true cross-validation (see section 10), the effect increases: a two-member forecast places $\sim 15\%$ less probability on the observed outcome.

[‡]Operational systems may typically consist of 40–50 ensemble members. Without hindcast sets, representative of operational systems, however, it is impossible to test this hypothesis fully.

^{**}In practice, numerical models developed for weather and climate simulations are far from independent, because they share common parametrizations and numerical schemes and are typically tuned towards the same training dataset. Also, they face the same technological (computational) limitation. This leads to structural similarities in the models and consequently to common shortcomings (e.g. in ‘blocking’).

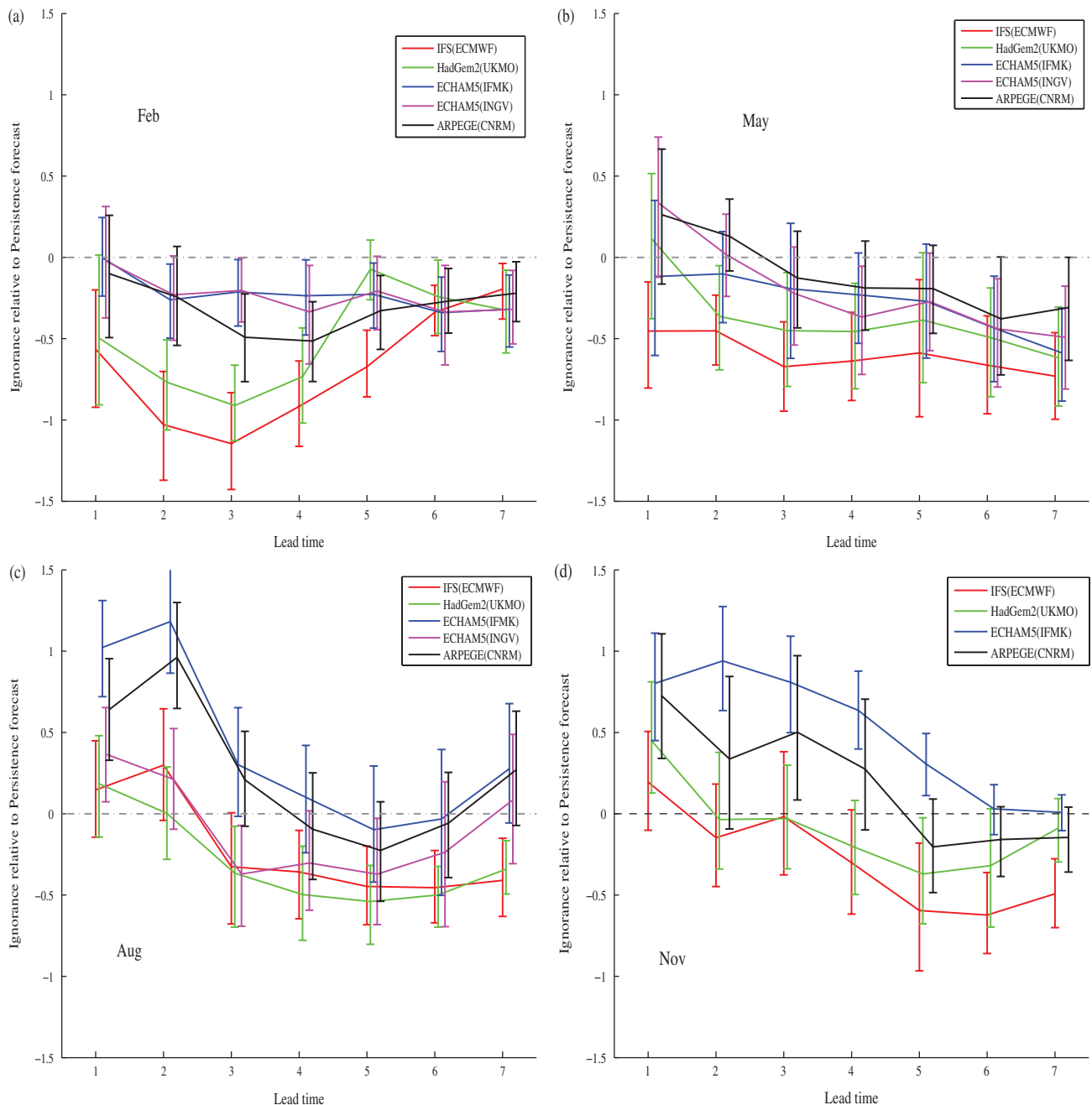


Figure 7. Ignorance score of each model from ENSEMBLES for the Nino3.4 index relative to persistence forecasts as a function of lead time in months. The four different panels show the hindcasts initialized in (a) February, (b) May, (c) August and (d) November. Scores below zero indicate that an ENSEMBLES model is more skilful than the persistence forecasts. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. ENSEMBLES model forecasts launched in February are shown to be more skilful than persistence at all lead times, whereas for forecasts launched in August the models are significantly worse than persistence at early lead times.

multi-model ensemble, containing 36 members from simulations of the IFS (ECMWF), HadGEM2 (UKMO), ECHAM5 (IFMK) and ARPEGE (CNRM) models. The four-member multi-model forecasts are shown to perform substantially worse than the nine-member IFS (ECMWF) ensemble (indicated by positive Ignorance scores), particularly over short lead times (up to 8 months). The skill of the nine-member multi-model forecasts is generally increased compared with the four-member forecasts; however, the single-model, IFS (ECMWF), forecast is still shown to be more skilful^{††} than the multi-model forecast at short lead

^{††}As noted by a referee, in this study the ‘best’ model has been identified in-sample. In this particular study, the ECMWF model is by far the highest scoring model across forecasts (see File S1) and is typically ranked first or second in over half of all skilful forecasts. Rather than resample to show that ECMWF is the best, the fraction of times it is best or second is shown in File S1. Note also Table 1 and Table 2 in this context. In practice, determining the best model *a priori*, either for a given purpose or in a multidimensional sense,

times. This is also true for the full 36 member multi-model forecast, although at longer lead times (beyond 8 months) the full multi-model ensemble is shown to outperform the IFS (ECMWF) ensemble. This result in this case suggests that increasing the ensemble size of the ‘best’ model is most likely to improve forecast skill in these regions.

8. The importance of being proper

It is sometimes said that a multi-model ensemble forecast is more skilful than any of its constituent single-model ensemble forecasts. This may be the case in terms of reducing RMS-like scores (see Palmer *et al.*, 2004; Hagedorn *et al.*, 2005; Bowler *et al.*, 2008; Weigel *et al.*, 2008; Weisheimer *et al.*, 2009; Alessandri *et al.*,

is not straightforward (if possible at all). In-sample evaluations of past model performance over relatively short hindcast periods hinder this task further.

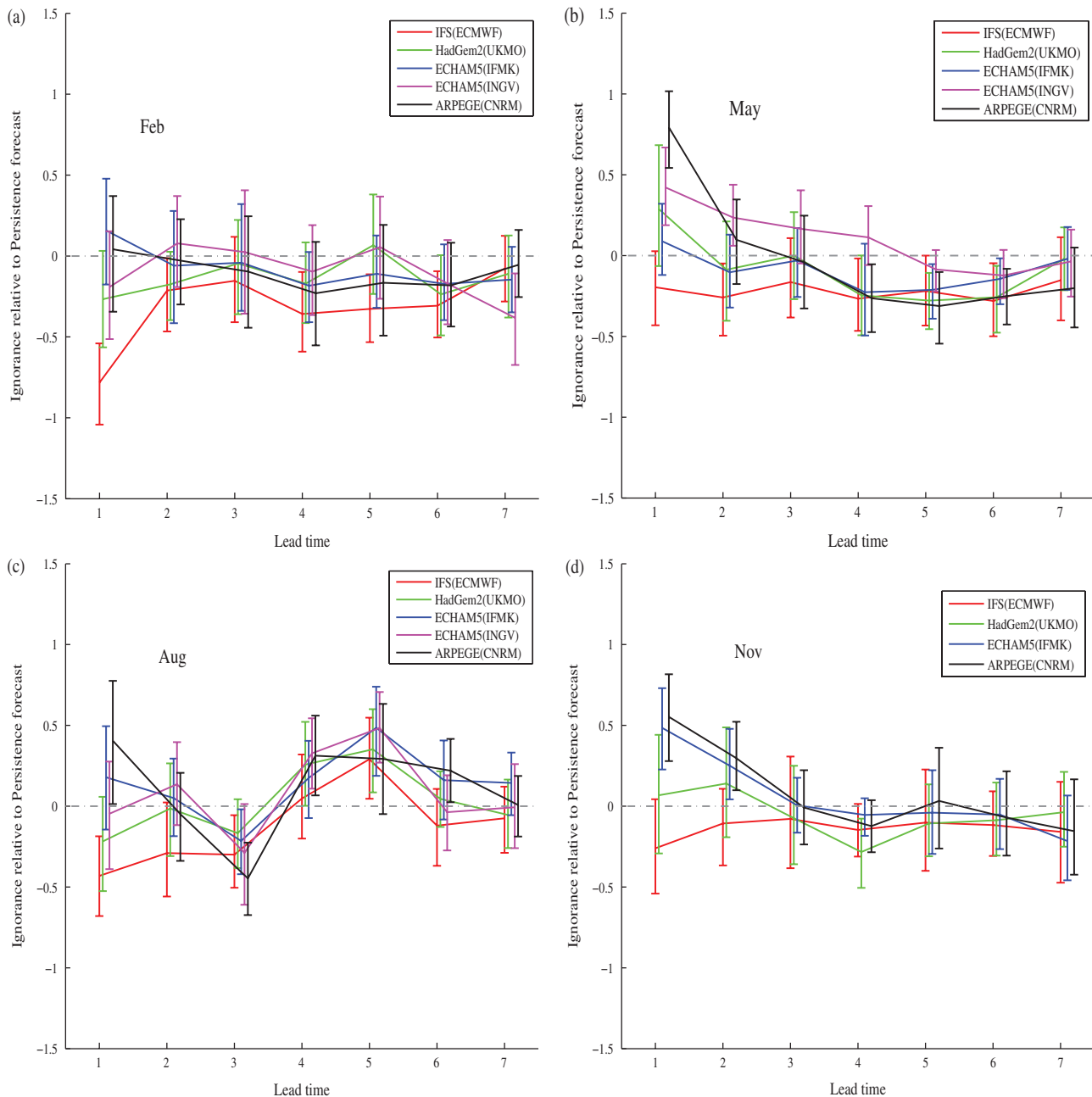


Figure 8. Ignorance score of each model from ENSEMBLES for the MDR index relative to persistence forecasts as a function of lead time in months. The four different panels show the hindcasts initialized in (a) February, (b) May, (c) August and (d) November. Scores below zero indicate that an ENSEMBLES model is more skilful than the persistence forecasts. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. While there is a tendency for Ignorance score to remain negative for several months in a row, suggesting skill, the upper (95%) resampling bound is almost always greater than zero.

2011). For probability forecasts, the definition of skill should reflect the characteristics of the forecast problem. While RMS scores are effectively optimal in linear stochastic systems, they are misleading in evaluating nonlinear forecast systems, even when the data are not precious. Indeed, RMS scores can be misleading even in the limit of an infinite forecast-verification archive (see McSharry and Smith, 1999). Improvements in RMS skill when using multi-model ensembles may be due to error cancellation from independent model contributions (see Hagedorn *et al.*, 2005; Kang and Yoo, 2006; Bowler *et al.*, 2008). For example, if some of the single-model ensembles lie below the observations and some lie above, then the ensemble mean could lie closer to the observed outcome than any single ensemble member. While such an error cancellation would reduce the RMS score, rewarding the multi-model forecast more than any single model contribution, a proper skill score (Bröcker and Smith, 2006) would not credit this ‘false’ skill. Similarly, combining ensemble members from different models may serve to reduce the variance of ensemble mean statistics, which in turn may lead to a lower RMS score.

Indeed, if the ensemble variance is large, adding ‘information-free’ ensemble members at the mean value will reduce the RMS error but need not improve a probabilistic score.

It has also been suggested that the multi-model ensemble forecast outperforms any of the single-model ensemble forecasts by reducing an apparent overconfidence in any one model (see Weigel *et al.*, 2008; Weisheimer *et al.*, 2009; Alessandri *et al.*, 2011). Such ‘improvements’ can easily be overinterpreted, however, as merely doubling the ensemble size under the same model in as much as increase the spread of the forecast distribution significantly. Another way to widen the ensemble spread is simply to blend (Bröcker and Smith, 2007) the model forecast distribution with an estimate of the climatological distribution based on the historical observations (see Appendix A for details). Two single-model forecasts may be ranked differently before and after blending with the climatological distribution. The effect of multi-model combination on seasonal forecast skill is investigated below.

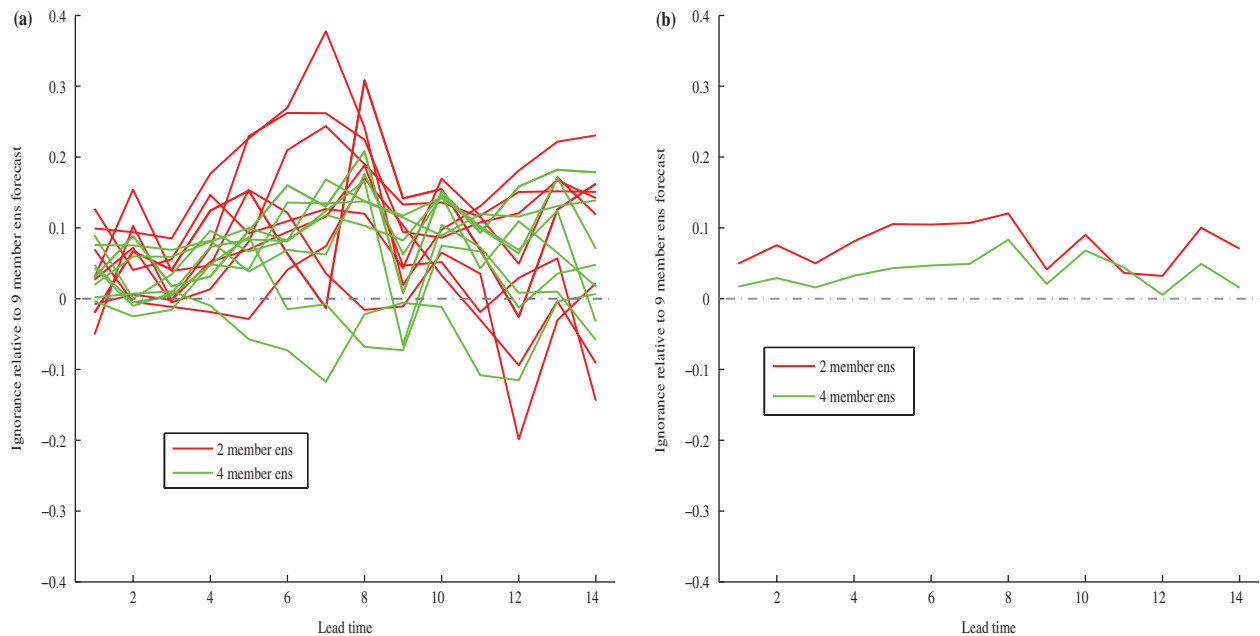


Figure 9. (a) Ignorance of the IFS (ECMWF) model as a function of lead time in months for the Nino3.4 index. The green (red) lines represent the skill of a subset of four-member (two-member) ensemble forecasts relative to the full nine-member ensemble forecast. Each four-member and two-member ensemble consists of random draws from the original nine-member ensemble. (b) Average Ignorance of all possible combinations of two-member (red) and four-member (green) ensembles. On average, the four-member ensembles are more skilful than the two-member ensemble, while both ensemble sizes are shown to perform worse on average than the full nine-member ensemble (i.e. Ignorance scores are all above zero).

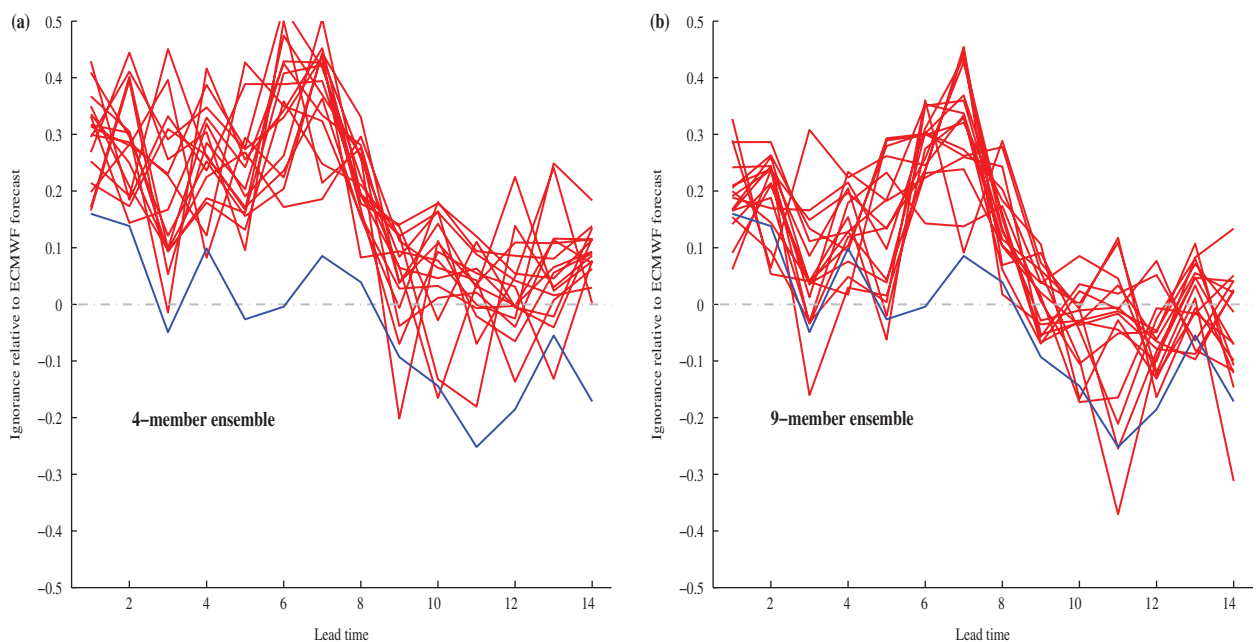


Figure 10. Ignorance of multi-model forecasts as a function of lead time in months for the Nino3.4 index, launched in November, relative to the nine-member IFS (ECMWF) forecast. The blue line represents the multi-model forecast using all 36 ensemble members from the four ENSEMBLES models, equally weighted. The red lines are multi-model forecasts using randomly drawn combinations of (a) four members and (b) nine members from the full ensemble. The four-member multi-model forecasts are shown to perform substantially worse than the nine-member IFS (ECMWF) ensemble (i.e. Ignorance scores are often above zero) and worse than the full 36 member multi-model ensemble. The nine-member multi-model forecasts perform better in general than the four-member forecasts and to a similar level of skill to the nine-member IFS (ECMWF) ensemble at lead times beyond 8 months.

9. Multiple models ensembles when data are precious

There are many ways in which forecast distributions, generated from ensembles of individual model runs, can be combined to produce a single probabilistic multi-model forecast distribution. One approach may be to assign equal weight to each model and simply sum the distributions generated from each model to obtain a single probabilistic distribution (see Hagedorn *et al.*, 2005). When different forecast models do not provide equal amounts of information, one may want to weight the models according to some measure of past performance: see, for example,

Krishnamurti *et al.* (1999), Rajagopalan *et al.* (2002) and Doblus-Reyes *et al.* (2005). The combined multi-model forecast is the weighted linear sum of the constituent distributions:

$$p_{\text{mm}} = \sum_i \omega_i p_i, \quad (4)$$

where p_i is the forecast distribution from model i and ω_i its weight, with $\sum_i \omega_i = 1$. The weighting parameters may be chosen by minimizing the Ignorance score, for example, although fitting ω_i in this way can be costly and is typically complicated by different models sharing information. Of course, the weights of individual models are expected to vary as a function of lead time.

Another, perhaps more fundamental problem of such a weighting procedure is that ω_i are likely to be over- or underfitted when the forecast-outcome archive is small (Peng *et al.*, 2002, L. A. Smith *et al.*, 2014; personal communication).

To avoid complications with fitting model weights, a simple iterative method to combine models is used below. First, a reference forecast distribution is derived from the ensemble members of one particular candidate model, in this case the IFS (ECMWF) forecasts, which were argued to provide the most skilful seasonal forecasts for the Nino3.4 index back in section 4. Each of the other candidate models, in turn, is then combined with the IFS (ECMWF) model by deriving a forecast distribution from the ensemble members of both models, equally weighted. The skill of each two-model combination is computed in terms of Ignorance relative to the IFS (ECMWF) reference forecast and is shown in Table 1 for the November launch forecasts of the Nino3.4 index. Each model combination shows the average relative Ignorance (negative scores indicate an improvement over simply using the IFS (ECMWF) forecast). The preponderance of positive values in the fifth, eighth and eleventh columns of Table 1 indicates that there is no clear improvement in skill for any two-model combination in this case. All values at lead times less than 8 months are positive. (In fact, all but two of the 42 values in these columns of the entire table are positive.) Arguably, beyond 8 months the improvements in skill are not significant; the bootstrap resampling intervals overlap with zero relative skill in each case. Table 2 shows the corresponding results when other models are combined with the UKMO model. In this case, combining with ECMWF tends to improve the average Ignorance at all lead times (negative values in the fourth and fifth columns of Table 2), but no other combination does this (all values in the eighth and eleventh columns are positive). Starting with ECMWF, combining UKMO has a much smaller effect. In cases where significant improvements are found from such a model combination, then further models could be included into the multi-model forecast by choosing those models that yield the biggest improvement in skill and adding them into the forecast one by one with equal weight until no further skill can be added. In this case, however, results suggest that the most skilful seasonal forecasts are provided by using ensemble members from a single model.

10. Establishing skill when data are precious

The DEMETER and ENSEMBLES seasonal hindcast archive contains merely 46 independent forecast-outcome pairs for each launch date. At seasonal forecast time-scales and longer, no true out-of-sample evaluation can be achieved in less than a decade, if not longer; evaluations today must necessarily be in-sample. In this case, it is desirable to strike a balance between using as many of the available data as possible to obtain the best results and holding back enough data so as to avoid information contamination (overfitting), which would lead to poor estimates of real-time operational skill.

The results shown in the previous sections used median cross-validation protocol as described in Appendix B; no additional data are held back in the evaluation of probabilistic forecast distributions beyond those excluded when determining the kernel parameters. While using median values for u , σ and α seems unlikely to allow significant information contamination, this median leave-one-out protocol is not 'true' cross-validation. In a true cross-validation protocol, more than one segment of data at a time must be removed from the fitting protocol. This reduces the chance of information contamination; it also reduces the true quality of the estimation when data are precious. Appendix B details both protocols.

Figure 11 shows the skill of forecasts from the ENSEMBLES models using true cross-validation. Figure 11(a) shows the Ignorance score for forecasts of the Nino3.4 index, launched in November. Comparing Figure 11(a) with Figure 2(d) clearly

shows a reduction in skill at longer lead times under the true cross-validation protocol, as well as a widening of the bootstrap resampling intervals in some cases. Significant skill above climatology is demonstrated only up to a lead time of 4 months. Similarly, Figure 11(b) shows the skill of the ENSEMBLES model forecasts for the MDR index. In this case, significant skill above climatology is shown to vanish beyond a lead time of 2 months.

The preferred cross-validation protocol when the data archive is small is unclear. The approach taken here is to consider more than one protocol. The true cross-validation protocol employed in this section (Figure 11) reflects the expected reduction in the skill of models simply because fewer data are used to calibrate the forecasts. The median cross-validation protocol (Figures 2 and 3) runs the risk of overfitting the dressing parameters. Only out-of-sample evaluation could establish which effect dominates in this case.

Figure 12 illustrates the effect of the different cross-validation protocols on the calculated skill of the seasonal forecasts. The figure shows Ignorance scores for the IFS (ECMWF) model from ENSEMBLES relative to climatology using the median (x -axis) and true (y -axis) cross-validation protocols for forecasts of the Nino3.4 index. Each of the four panels corresponds to a different forecast launch month (as indicated). As expected, on average the true cross-validation protocol suggests less skill (i.e. larger Ignorance scores) relative to median cross-validation. This improvement on average is not systematic across individual forecasts. The reduction of skill under true cross-validation protocol is small in most cases, giving increased confidence to results using median cross-validation. The most prominent differences are at the highest values of Ignorance, where the forecasts have little skill under either protocol. For the November launch, this typically occurs at longer lead times (beyond 7 months). The argument here is merely that it is important to consider questions of cross-validation when data are precious.

11. Conclusions

The current generation of seasonal forecasts will retire before the forecast-outcome archive grows significantly larger: seasonal verification data are precious! This complicates forecast calibration, as evaluation must be performed using cross-validation with only a small sample. Nevertheless, probabilistic seasonal forecasts based on the ENSEMBLES stream II experiment demonstrate increased skill in forecasting sea-surface temperatures in the Nino3.4 region over that of the DEMETER model simulations. Further analysis suggests that increasing the ensemble size could potentially improve forecast skill further. Such evaluations of skill, on the other hand, should be analyzed with care. RMS-based skill scores can obscure skill in nonlinear systems. The statistical characteristics reflected in RMS scores differ from those using strictly proper scoring rules, which are recommended for evaluations of such nonlinear systems as those in weather and climate dynamics. The evidence of skill presented, particularly at moderate lead times, is shown to be robust to different choices of appropriate (proper) scores (see File S1) and may prove to have non-trivial value in application. Simulation-based forecasts clearly outperform climatological probability forecasts in many cases. The fact that empirical persistence-based probability forecasts provide a significantly stronger challenge suggests that, in practice, the skill of operational forecast systems can be enhanced with information from the richer empirical models. Distinguishing the limitations of this level of skill from the limitations of our current skill scores and evaluation methodologies will also prove of great value, both in terms of informing future experimental designs for multi-model ensemble projects and for determining the value of these forecast systems to decision-makers.

Table 1. Ignorance of each two-model forecast combination, as labelled, relative to the IFS (ECMWF) forecast for each (monthly) lead time for seasonal forecasts of the Nino3.4 index, launched in November. In each case, the individual models are also blended with the climatological distribution using blending parameters that minimize the Ignorance score. Each two-model combination shows the average relative Ignorance and the 5–95% bootstrap resampling intervals, which provide an estimate of sampling uncertainty of the relative skill score. For comparison, the second column shows the skill of the (single) ECMWF model relative to climatology.

LT	ECMWF	ECMWF and UKMO			ECMWF and CNRM			ECMWF and IFMK		
		5%	mean	95%	5%	mean	95%	5%	mean	95%
1	-2.15	-0.08	0.05	0.16	0.05	0.17	0.28	0.07	0.20	0.30
2	-2.03	-0.29	-0.07	0.10	-0.17	0.04	0.24	0.15	0.33	0.47
3	-1.63	-0.44	-0.16	0.08	-0.21	0.04	0.23	-0.09	0.18	0.37
4	-1.36	-0.17	-0.03	0.10	-0.05	0.11	0.26	0.13	0.29	0.41
5	-1.10	-0.19	0.01	0.16	-0.25	-0.04	0.16	0.09	0.28	0.42
6	-0.73	-0.16	0.01	0.17	-0.04	0.11	0.25	0.03	0.19	0.31
7	-0.53	-0.05	0.09	0.22	-0.07	0.07	0.20	0.09	0.18	0.26
8	-0.34	-0.06	0.05	0.15	-0.04	0.06	0.16	-0.04	0.06	0.15
9	-0.23	-0.14	-0.04	0.05	-0.10	0.00	0.11	-0.14	-0.04	0.04
10	-0.27	-0.16	-0.06	0.03	-0.17	-0.05	0.06	-0.14	-0.04	0.05
11	-0.22	-0.32	-0.17	-0.02	-0.22	-0.08	0.06	-0.33	-0.20	-0.08
12	-0.28	-0.20	-0.09	0.01	-0.17	-0.05	0.07	-0.13	-0.03	0.07
13	-0.35	-0.08	-0.01	0.06	-0.20	-0.03	0.11	-0.14	-0.05	0.05
14	-0.39	-0.12	-0.03	0.07	-0.12	0.00	0.13	-0.31	-0.12	0.03

Table 2. Ignorance of each two-model forecast combination, as labelled, relative to the HadGEM2 (UKMO) forecast for each (monthly) lead time for seasonal forecasts of the Nino3.4 index, launched in November. In each case, the individual models are also blended with the climatological distribution using blending parameters that minimize the Ignorance score. Each two-model combination shows the average relative Ignorance and the 5–95% bootstrap resampling intervals, which provide an estimate of sampling uncertainty of the relative skill score. For comparison, the second column shows the skill of the (single) UKMO model relative to climatology.

LT	UKMO	UKMO and ECMWF			UKMO and CNRM			UKMO and IFMK		
		5%	mean	95%	5%	mean	95%	5%	mean	95%
1	-1.90	-0.35	-0.21	-0.08	-0.02	0.08	0.17	-0.01	0.11	0.22
2	-1.92	-0.41	-0.18	0.01	0.03	0.12	0.21	0.22	0.34	0.44
3	-1.64	-0.33	-0.15	-0.01	0.00	0.13	0.26	0.14	0.28	0.40
4	-1.29	-0.24	-0.13	0.00	-0.09	0.06	0.20	0.13	0.26	0.38
5	-0.87	-0.37	-0.22	-0.09	-0.34	-0.12	0.07	0.06	0.21	0.33
6	-0.43	-0.49	-0.30	-0.11	-0.38	-0.12	0.09	-0.11	0.06	0.20
7	-0.13	-0.45	-0.31	-0.16	-0.30	-0.13	0.02	-0.09	0.00	0.08
8	-0.14	-0.26	-0.15	-0.06	-0.20	-0.05	0.06	-0.24	-0.07	0.06
9	-0.24	-0.15	-0.04	0.05	-0.21	-0.03	0.12	-0.18	-0.06	0.05
10	-0.32	-0.12	-0.02	0.08	-0.10	0.00	0.10	-0.12	-0.02	0.08
11	-0.33	-0.24	-0.05	0.12	-0.15	-0.01	0.13	-0.40	-0.16	0.03
12	-0.32	-0.22	-0.06	0.09	-0.11	0.00	0.10	-0.17	-0.03	0.11
13	-0.31	-0.13	-0.05	0.03	-0.14	-0.02	0.12	-0.17	-0.07	0.03
14	-0.31	-0.24	-0.10	0.03	-0.11	0.00	0.10	-0.39	-0.18	0.01

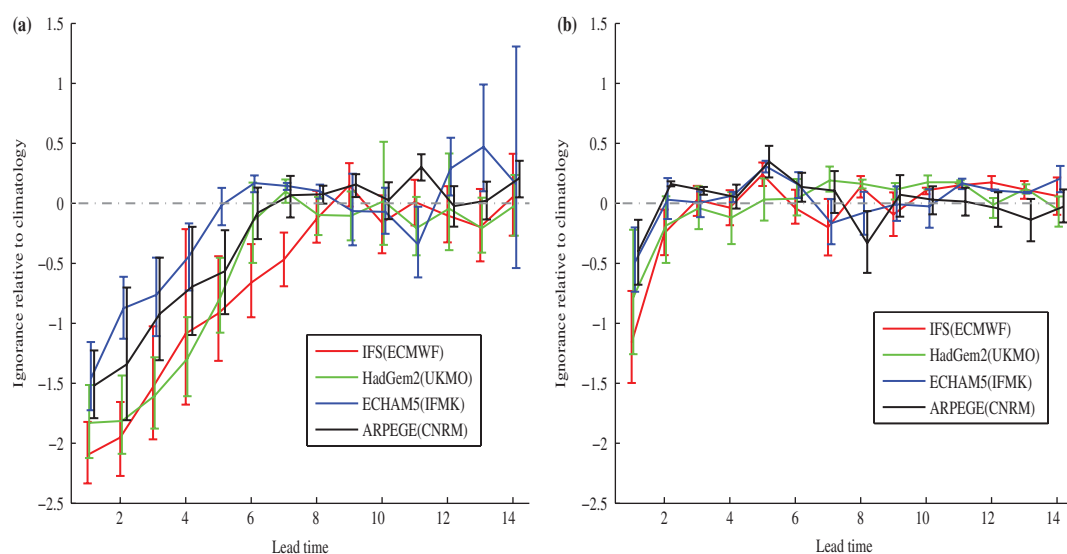


Figure 11. Ignorance score of each model from ENSEMBLES relative to climatology as a function of lead time in months using true cross-validation, for (a) forecasts of the Nino3.4 index and (b) forecasts of the MDR index launched in November. Zero Ignorance indicates that a model has no skill relative to climatology and negative relative Ignorance scores suggest that a model is more skilful than climatology. Bootstrap resampling intervals (the vertical bars) reflect the 5–95% range as estimated from 512 resamples. Skill is typically reduced compared with the median cross-validation protocol (Figures 2(d) and 3(d)), particularly at very early lead times over the MDR. The bootstrap resampling intervals are also widened in some cases.

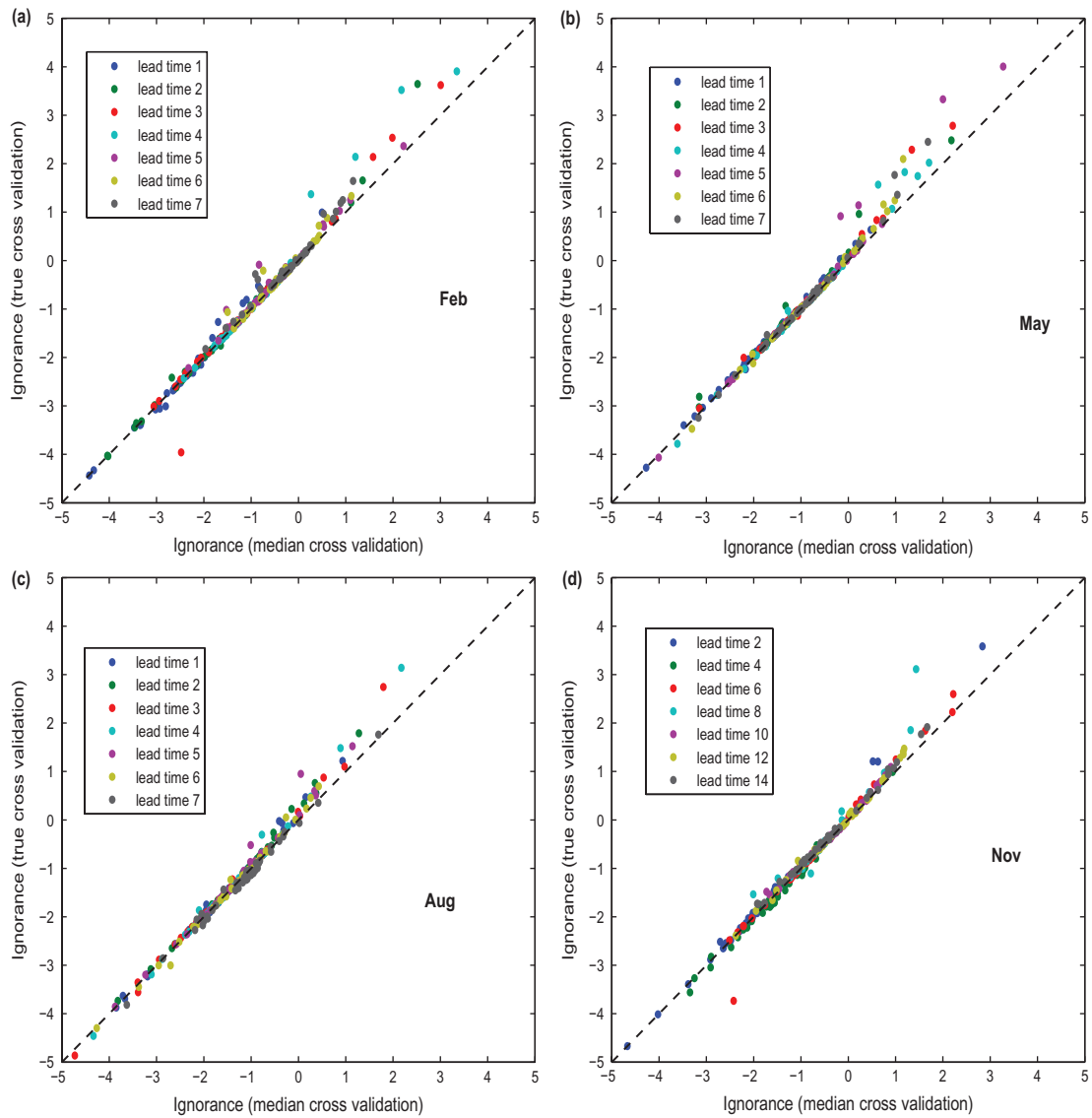


Figure 12. Comparison of Ignorance scores for the IFS (ECMWF) model from ENSEMBLES relative to climatology using the median and true cross-validation protocols for forecasts of the Nino3.4 index, launched in the months as indicated. On average, the true cross-validation protocol shows a reduction in skill (larger Ignorance scores) compared with median cross-validation, although individual forecasts can score better. The reduction of skill when using the true cross-validation protocol is most prominent at higher values of Ignorance (when the forecasts are already demonstrating poor skill under the median cross-validation protocol), which for the November launch typically occurs at longer lead times (beyond 7 months).

Appendix A: From simulation to PDF

An ensemble of simulations is transformed into a probabilistic distribution function (PDF) by a combination of kernel-dressing and blending with climatology (see Bröcker and Smith, 2007). An N -member ensemble at time t is given as $X_t = [x_t^1, \dots, x_t^N]$, where x_t^i is the value of a physical quantity (for example the SST in the MDR region) for the i th ensemble member. For simplicity, all ensemble members under a given model are treated as exchangeable. In other words, the ensemble interpretation does not depend on the ordering of the ensemble members, as long as they are generated by the same model (Bröcker and Smith, 2007). Kernel-dressing defines the model-based component of the density as

$$p(y : X, \sigma) = \frac{1}{N\sigma} \sum_i^N K\left(\frac{y - (x^i - \mu)}{\sigma}\right), \quad (\text{A1})$$

where y is a random variable corresponding to the density function p and K is the kernel, taken here to be

$$K(\zeta) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}\zeta^2\right). \quad (\text{A2})$$

Thus each ensemble member contributes a Gaussian kernel centred at $x^i - \mu$. Here μ is an offset, which accounts for any systematic ‘bias’. For a Gaussian kernel, the kernel width σ is simply the standard deviation determined empirically, as discussed below.

For any finite ensemble, there remains the chance of $\sim \frac{2}{N}$ that the outcome lies outside the range of the ensemble even when the outcome is selected from the same distribution as the ensemble itself. Given the nonlinearity of the model, such outcomes can be very far outside the range of the ensemble members. In addition to N being finite, in practice of course the simulations are not drawn from the same distribution as the outcome, as the ensemble simulation system is not perfect. To improve the skill of the probabilistic forecasts, the kernel-dressed ensemble may be blended with an estimate of the climatological distribution of the system (see Bröcker and Smith, 2007 for more details; Roulston and Smith, 2003 for an alternative kernel and Raftery *et al.*, 2005 for a Bayesian approach). The blended forecast distribution is then written as

$$p(\cdot) = \alpha p_m(\cdot) + (1 - \alpha)p_c(\cdot), \quad (\text{A3})$$

where p_m is the density function generated by dressing the model ensemble and p_c is the estimate of climatological density. The blending parameter α determines how much weight is placed

in the model. Specifying the three values (kernel width σ , kernel offset μ and weight α) at each lead time defines the forecast distribution. These parameters are fitted simultaneously by optimizing the empirical Ignorance score, using a cross-validation protocol^{††} as described in Appendix B.

Appendix B: Information contamination and cross-validation

Ideally, forecast performance is evaluated ‘out-of-sample’, with new data unknown at the time the model parameters were determined (much less data seen by the analyst). Given a large forecast-outcome archive, cross-validation reduces information contamination and overfitting when working in-sample (that is, when evaluating a model on the sample used to fit the parameters of that model) by dividing the archive into two sets: a training set, used to build the forecast model and fit the parameters, and a testing set, used to obtain an estimate of the skill and likely performance of the model. The process can be repeated to examine the robustness of the results, but information from the test set(s) must not be used to ‘improve’ the forecast model. When the archive is small and will increase only slowly, one does not have the luxury of this approach. Calibration and evaluation are at best performed under more complex cross-validation; the ideal protocol is not clear and the results can be expected to change with the protocol. A median protocol and a true leave-one-out protocol are defined below.

First, define the forecast probability distribution to be $p(x, X_t, \Theta)$, $t = 1, \dots, N$, where X represents the ensemble forecast at time t , Θ represents a vector of parameters (including the kernel width σ , offset μ and blending parameter α) to be fitted and N is the number of forecasts. The corresponding outcomes are defined to be s_t . For each forecast at time $j = 1, \dots, N$, leave out one pair of forecast-outcome data (X_j, s_j) and use the remaining forecast-outcome data pairs to determine the parameter Θ_j by minimizing the empirical score (in this article Ignorance is used). The median value, $\bar{\Theta}$, of the set of N Θ_j is then used in the forecast model. This ‘median protocol’ maintains a large learning set with only slight information contamination.

The leave-one-out protocol described in the previous paragraph is not pure cross-validation, as $\bar{\Theta}$ arguably contains information from every (X_j, s_j) when the median is taken. To achieve pure cross-validation, the following protocol is adopted. For each forecast at time j , first leave out (X_j, s_j), then for the remaining set apply the median cross-validation protocol described above to obtain N parameter values $\bar{\Theta}_j$. The value $\bar{\Theta}_j$ at each time j is then independent of (X_j, s_j). The forecast empirical Ignorance is then given by $\sum_{j=1}^N -\log_2 p(s_j, X_j, \bar{\Theta}_j)$. This protocol ensures that the parameters $\bar{\Theta}_j$ have no explicit dependence on the datum used to evaluate them, at the cost of smaller learning set(s). Even in this case, the datum was known to the analyst. Indeed, use of a common archive in DEMETER and in ENSEMBLES (Stream Two) clouds the possibility of assigning clear statistical significance to estimates of expected skill.

Acknowledgements

We happily acknowledge constructive conversations with Jochen Broecker, James Hansen, Tim Palmer, Erica Thompson, Antje Weisheimer, Edward Wheatcroft and Daniel S. Wilks. This research was supported by the EU Framework 6 ENSEMBLES project; it was also supported both by the LSE’s Grantham

^{††}As only 46 years of data are used in this case, any estimation of the two parameters lacks robustness. If one had 4000 years of data, one could draw multiple 46 year data sets from them and estimate the parameters for each sample set. In experiments with simple systems, it turns out that the variation of such estimates is large (see L. A. Smith *et al.*, 2014; personal communication). Note that a 46 year hindcast archive of the full ensemble system may not be available to aid the construction of operational forecast systems.

Research Institute on Climate Change and the Environment and the ESRC Centre for Climate Change Economics and Policy, funded by the Economic and Social Research Council and Munich Re. LAS gratefully acknowledges support from Pembroke College, Oxford.

Supporting information

The following supporting information is available as part of the online article:

File S1. Probabilistic skill in ensemble seasonal forecasts.

References

- Alessandri A, Borrelli A, Masina S, Pietro PD, Carril A, Cherchi A, Gualdi S, Navarra A. 2010. The INGV-CMCC seasonal prediction system: Improved ocean initial conditions. *Mon. Weather Rev.* **138**: 2930–2952.
- Alessandri A, Borrelli A, Navarra A, Arribas A, Déqué M, Rogel P, Weisheimer A. 2011. Evaluation of probabilistic quality and value of the ensembles multimodel seasonal forecasts: Comparison with DEMETER. *Mon. Weather Rev.* **139**: 581–607.
- Bernardo JM. 1979. Expected information as expected utility. *Ann. Stat.* **7**: 686–690.
- Bowler NE, Arribas A, Mylne KR. 2008. The benefits of multi-analysis and poor-mans ensembles. *Mon. Weather Rev.* **136**: 4113–4129.
- Bröcker J, Smith LA. 2006. Scoring probabilistic forecasts: On the importance of being proper. *Weather and Forecasting* **22**: 382–388.
- Bröcker J, Smith LA. 2007. From ensemble forecasts to predictive distribution functions. *Tellus A* **60**: 663–678.
- Coelho CAS, Stephenson DB, Balmaseda M, Doblas-Reyes FJ, van Oldenborgh GJ. 2006. Towards an integrated seasonal forecasting system for South America. *J. Clim.* **19**: 3704–3721.
- Doblas-Reyes FJ, Hagedorn R, Palmer TN. 2005. The rationale behind the success of multi-model ensembles in seasonal forecasting. Part II: Calibration and combination. *Tellus A* **57**: 234–252.
- Doblas-Reyes FJ, Weisheimer A, Palmer TN, Murphy JM, Smith D. 2010. *Forecast Quality Assessment of the ENSEMBLES Seasonal-to-decadal Stream 2 Hindcasts, Technical Memorandum 621*. ECMWF: Reading, UK.
- Good IJ. 1952. Rational decisions. *J. R. Stat. Soc.* **XIV**: 107–114.
- Hagedorn R, Smith LA. 2009. Communicating the value of probabilistic forecasts with weather roulette. *Meteorol. Appl.* **16**: 143–155.
- Hagedorn R, Doblas-Reyes FJ, Palmer TN. 2005. The rationale behind the success of multi-model ensembles in seasonal forecasting. Part I: Basic concept. *Tellus A* **57**: 219–233.
- Hewitt CD, Griggs DJ. 2004. Ensembles-based predictions of climate changes and their impacts. *Eos Trans. Am. Geophys. Union* **85**: 566.
- Ingleby B, Huddleston M. 2007. Quality control of ocean temperature and salinity profiles –historical and real-time data. *J. Mar. Syst.* **65**: 158–175.
- Kang I-S, Yoo J. 2006. Examination of multi-model ensemble seasonal prediction methods using a simple climate system. *Clim. Dyn.* **26**: 285–294.
- Kelly JL Jr. 1956. A new interpretation of information rate. *Bell Syst. Tech. J.* **35**: 917–926.
- Kirtman BP, Min D, Infanti JM, Kinter JL III, Paolino DA, Zhang Q, van den Dool H, Saha S, Mendez MP, Becker E, Peng P, Tripp P, Huang J, DeWitt DG, Tippett MK, Barnston G, Li S, Rosati A, Schubert SD, Rienecker M, Suarez M, Li ZE, Marshall J, Lim Y-K, Tribbia J, Pegion K, Merryfield WJ, Denis B, Wood EF. 2013. The North American Multi-Model Ensemble (NMME): Phase-1 seasonal to interannual prediction, phase-2 toward developing intra-seasonal prediction. *Bull. Am. Meteorol. Soc.* **95**: 585–601.
- Krishnamurti TN, Kishtawal CM, LaRow TE, Bachiochi DR, Zhang Z, Williford CE, Gadgil S, Surendran S. 1999. Improved weather and seasonal climate forecasts from multimodel superensemble. *Science* **285**: 1548–1550.
- McSharry PE, Smith LA. 1999. Better nonlinear models from noisy data: Attractors with maximum likelihood. *Phys. Rev. Lett.* **83**: 4285–4288.
- Palmer TN, Alessandri A, Andersen U, Cantelaube P, Davey M, Décluse P, Déqué M, Diez E, Doblas-Reyes FJ, Feddersen H, Graham R, Gualdi S, Guérémy J-F, Hagedorn R, Hoshen M, Keenlyside N, Latif M, Lazar A, Maisonnave E, Marletto V, Morse AP, Orfila B, Rogel P. 2004. Development of a European multimodel ensemble system for seasonal-to-interannual prediction (DEMETER). *Bull. Am. Meteorol. Soc.* **85**: 853–872.
- Peng P, Kumar A, van den Dool H, Barnston AG. 2002. An analysis of multimodel ensemble predictions for seasonal climate anomalies. *J. Geophys. Res.* **107**: 4710, doi: 10.1029/2002JD002712.
- Raftery AE, Gneiting T, Balabdaoui F, Polakowski M. 2005. Using Bayesian model averaging to calibrate forecast ensembles. *Mon. Weather Rev.* **131**: 1155–1174.

- Rajagopalan B, Lall U, Zebiak SE. 2002. Categorical climate forecasts through regularization and optimal combination of multiple GCM ensembles. *Mon. Weather Rev.* **130**: 1792–1811.
- Roulston MS, Smith LA. 2002. Evaluating probabilistic forecasts using information theory. *Mon. Weather Rev.* **130**: 1653–1660.
- Roulston MS, Smith LA. 2003. Combining dynamical and statistical ensembles. *Tellus* **55A**: 16–30.
- Smith DM, Eade R, Dunstone NJ, Fereday D, Murphy JM, Pohlmann H, Scaife AA. 2010. Skilful multi-year predictions of Atlantic hurricane frequency. *Nat. Geosci.* **3**: 846–849.
- Smith LA. 1992. Identification and prediction of low-dimensional dynamics. *Physica D* **58**: 50–76.
- Smith LA. 1997. In *Proceedings International School of Physics 'Enrico Fermi', Course CXXXIII*, Cini G. (ed.): 177–246. Societa Italiana di Fisica: Bologna, Italy.
- Suckling EB, Smith LA. 2013. An evaluation of decadal probability forecasts from state-of-the-art climate models. *J. Clim.* **26**: 23.
- The Met Office. 2013. *3-month Outlook for Contingence Planning: User Guidance*, HM Government document. The Met Office: Devon, UK. http://www.metoffice.gov.uk/media/pdf/g/o/3-month_Outlook_User_Guidance-150.pdf (accessed 10 July 2014).
- Van Den Dool HM. 2007. *Empirical Methods in Short-term Climate Prediction*. Oxford University Press: Oxford, UK.
- van Oldenborgh GJ, Balmaseda MA, Ferranti L, Stockdale TN, Anderson DLT. 2005. Did the ECMWF seasonal forecast model outperform statistical ENSO forecast models over the last 15 years? *J. Clim.* **18**: 3240–3249.
- Vitart F, Huddleston MR, Déqué M, Peake D, Palmer TN, Stockdale TN, Davey MK, Ineson S, Weisheimer A. 2007. Dynamically-based seasonal forecast of Atlantic tropical storm activity issued in June by EUROSIP. *Geophys. Res. Lett.* **34**: L16815, doi: 10.1029/2007GL030740.
- Wang B, Lee J-Y, Kang I-S, Shukla J, Park CK, Kumar A, Schemm J, Cocke S, Kug JS, Luo JJ, Zhou T, Wang B, Fu X, Yun WT, Alves O, Jin E, Kinter J, Kirtman B, Krishnamurti T, Lau N, Lau W, Liu P, Pegion P, Rosati T, Schubert S. 2009. Advance and prospectus of seasonal prediction: Assessment of the APCC/CliPAS 14-model ensemble retrospective seasonal prediction (19802004). *Clim. Dyn.* **33**: 93–117.
- Weigel AP, Liniger MA, Appenzeller C. 2008. Can multi-model combination really enhance the prediction skill of probabilistic ensemble forecasts? *Q. J. R. Meteorol. Soc.* **134**: 241C260.
- Weisheimer A, Doblas-Reyes FJ, Palmer TN, Alessandri A, Arribas A, Déqué M, Keenlyside N, MacVean M, Navarra A, Rogel P. 2009. ENSEMBLES: A new multi-model ensemble for seasonal-to-annual predictions and Skill and progress beyond DEMETER in forecasting tropical Pacific SSTs. *Geophys. Res. Lett.* **36**: L21711, doi: 10.1029/2009GL040896.
- Wilks DS. 2005. *Statistical Methods in the Atmospheric Sciences, International Geophysics* **91** (2nd edn). Academic Press: Oxford, UK.