

Detecting failure of climate predictions

**Michael C. Runge*, USGS Patuxent Wildlife Research Center, Laurel, MD 20708 USA

Julienne C. Stroeve, National Snow and Ice Data Center, University of Colorado, Boulder, CO
80309 USA, University of College London, London, UK

Andrew P. Barrett, National Snow and Ice Data Center, University of Colorado, Boulder, CO
80309 USA

Eve McDonald-Madden, School of Geography, Planning, and Environmental Management,
University of Queensland, St. Lucia, QLD 4072 Australia

The practical consequences of climate change challenge society to formulate responses that help better achieve long-term objectives, even if those responses have to be made in the face of uncertainty^{1,2}. Such a decision-analytic focus uses the products of climate science as probabilistic predictions about the effects of management policies³. Here we present methods to detect when climate predictions are failing to capture the system dynamics. For a single model, we measure goodness of fit based on the empirical distribution function, and define failure when the distribution of observed values significantly diverges from the modelled distribution. For a set of models, the same statistic can be used to provide relative weights for the individual models, and we define failure when there is no linear weighting of the ensemble models that produces a satisfactory match to the observations. Early detection of failure of a set of predictions is important for improving model predictions and the decisions based upon them. We show that these methods (i) would have detected a range shift in northern pintail 20 years before it was actually discovered, and (ii) are

increasingly giving more weight to those climate models that forecast a September ice-free Arctic by 2055.

Recognizing the decision context of climate change issues identifies a meaningful role for empirical science, and shifts the debate to pragmatic solutions^{1,2,4}. The central role of traditional climate science in decision-making processes is to provide probabilistic predictions about outcomes of interest under various management strategies³. These predictions are, of course, made with uncertainty. The explicit articulation of this uncertainty is healthy, because it allows both risk analysis and adaptive management⁴. With risk analysis, individuals and society can examine the consequences of taking (or not taking) any action and being wrong, and so search for solutions that appropriately weigh the various risks. With adaptive management, management actions can be adjusted in response to new information that reduces uncertainty; indeed, the anticipation of this learning may influence initial actions. Both risk analysis and adaptive management require the articulation of uncertainty as a set of alternative predictions about the future. For climate forecasting, the set of coupled general circulation models (GCMs) and the various forcing scenarios provide the basis for alternative predictions about the outcomes of many potential management actions⁵.

The ability to learn and make good management decisions within an adaptive framework will depend on whether the true system dynamics are contained within, bounded by, or close to the set of models that capture current uncertainty. Two types of surprise could undermine this ability: first, the truth might not be bounded by the model set, because of a failure to anticipate some important elements of the system; or second, the system might change in unanticipated ways that lead the true dynamics outside behavior predicted by the model set. Both of these unanticipated outcomes can be considered “unknown unknowns” or “black Swans”⁶. Adaptive

management includes an internal layer of learning (“single-loop learning”⁷) that allows discernment among the existing predictions as information accrues, and adaptation of future management actions to that new understanding. In addition, a second layer of learning is needed, which examines if the system is responding as might be expected given the model set available, or if, instead, unpredicted responses are occurring. In the latter case, “double-loop learning”⁸ is triggered, in which the model set itself is reexamined, in an effort to develop new hypotheses that explain the surprising results. For example, satellite observations of arctic sea-ice extent declined faster than forecast by the World Climate Research Programme Coupled Model Intercomparison Project Phase 3 (CMIP3) models, leading to hypotheses for the discrepancy and efforts to improve subsequent models⁹. The first step in double-loop learning is the detection of the failure of the model set. Early detection of failure of a set of predictions can trigger the process of diagnosis and the process of generating new predictions, quickly turning “unknown unknowns” into “known unknowns” and leading to better ongoing management and policy interventions through adaptive management.

In this paper, we develop methods for detecting the failure of a single model and the failure of a model set. We illustrate these methods in two contexts: detecting a shift in breeding distribution for northern pintails (*Anas acuta*)¹⁰; and detecting a failure of climate models to predict the loss of Arctic sea ice⁹.

Assessing the plausibility of a single model

The role of models in a decision context is to make predictions about system response through time and as a function of management actions. These predictions are usually probabilistic¹¹, to represent uncertainty arising from a number of sources, including environmental variation,

incomplete knowledge of system dynamics, sampling error, and incomplete control of management actions¹². Thus, a model can be viewed as a hypothesis about the distribution of the response variable of interest. We would like our probabilistic predictions to be well calibrated and sharp¹³: over time, the observations should be compatible with the modelled distribution¹⁴. For example, in forecasting rainfall, we would like the observed frequency of wet and dry years to match the predicted (hindcast or forecast) frequencies generated by GCM simulations. The empirical distribution function (EDF) tests, a class of goodness-of-fit tests, examine the agreement between two continuous distributions using a statistic that measures the distance (D_n) between the empirical cumulative distribution from the real system ($F_n(x)$, where n is the accumulated sample size) and the hypothesized cumulative distribution based on the prediction from the model ($F(x)$)¹⁵. One of the advantages of the EDF tests is that the prediction can take any form of distribution. The Kolmogorov-Smirnov (K-S) test is one of many EDF tests, and uses the distance metric

$$D_n = \max_x |F_n(x) - F(x)|. \quad (1)$$

The northern pintail (*Anas acuta*) is a waterfowl species that is important for recreational hunting in North America¹⁶. Pintails depend on ephemeral prairie wetlands for breeding and their dynamics are strongly influenced by climatic conditions¹⁷. The annual distribution of this species, as measured by the latitude of its centroid, is an indicator of the habitat conditions, with individuals breeding farther north in drier years. Because reproductive rate is also associated with habitat condition, the latitude of the breeding population is used as a predictor in setting hunting regulations¹⁸. Between 1961 and 1974 the mean latitude of the breeding distribution was

53.569 (SD=1.549) (Fig. 1A, red line). Data collected from the mid-eighties onwards shows a northerly shift in the pintail distribution, but given the variability in the data it is difficult to discern if or when this shift occurred and whether concerns should be raised about the harvest rates set using Model 1. The K-S statistic shows the observations were compatible with Model 1 (Fig. 1B, red line) until 1985 (red circle). After 1985, the distance between the observations and the predictions under Model 1 suggests a significant change in the pintail distribution. In this way, and EDF statistic can be used to identify when a single model is no longer plausible.

Assessing the plausibility of a model set

Often, a decision maker will entertain several different explanations of cause-and-effect in a system, that is, several alternative models. These models may represent a comprehensive set, in the sense that the truth is believed to be one of the models, but more commonly, the hope is merely that the set of models somehow bounds the truth. What would it mean for a model ensemble to bound the truth? We propose this means there is weighted combination of the models in the set that makes predictions consistent with the observations. If that is not the case, then the observations are falling outside anything predicted by the ensemble, which would indicate the need for careful evaluation of the model set. The EDF statistic for the best-weighted model, then, is a measure of the plausibility of the model set.

There are a number of ways that models could be weighted to form an intermediate model. One possibility is to form a linear-weighted average of the cumulative distribution functions (CDF) for each model. Another way is to average the moments of the individual distributions. In either case, the best-fitting weighted model minimizes the EDF statistic. In the examples that follow,

we have used the second weighting method, because we were particularly interested in bounding the first two moments, but the first weighting method may be appropriate in other contexts.

In 1985, when the observed pintail data indicated a divergence from the 1961-1974 model (Model 1), a possible response would have been to propose a second model with a fixed mean of 55.374 (the 5-yr moving average in 1985) and standard deviation of 1.549 (Fig. 1A, dashed blue line). This second model was not plausible between 1970 and 1985 (as judged against a K-S test with nominal $\alpha = 0.05$), but became plausible in 1985 and has remained so since (Fig. 1B, blue line). The weights in the best-fitting weighted model show the change in system dynamics (Fig. 1C): between 1971 and 1980, Model 1 received all of the weight; by 1988, all the weight had shifted to Model 2; and since 1998, the weights have fluctuated. In the period 1988-1998, having all of the weight on Model 2 raises the question whether the true dynamics have moved outside the model set and Model 2 is just the best approximation available. Nevertheless, the best-fitting weighted model remains plausible over the entire time series, suggesting that the two-model ensemble set currently bounds the true range dynamics of the northern pintail and would have performed well for setting harvest rates (Fig. 1B, black line). Use of the Anderson-Darling statistic (another in the class of EDF tests¹⁵) instead of the K-S statistic produces quite similar results, with two minor differences: first, the failure of Model 1 alone is detected in 1984 instead of 1985; and second, for one year in 1993, the tests warns that the model set may be failing. (See Supplementary Information for a comparison of the power of these two tests.)

Forecasting Arctic sea ice

The rapid loss of Arctic sea ice over the past two decades has been one of the most visible and dramatic effects of global climate change¹⁹ and has led to significant concern about many aspects of the Arctic environment, including for example, the status of polar bears^{20,21}. Sea-ice extent and volume have been declining at a rate that was faster than forecast by the CMIP3 models⁹. More recent models (CMIP5) match the trends in the observed record better²² (Fig. 2A), but the question remains whether they are capturing the Arctic sea-ice dynamics well enough to support decision-making. The K-S statistics for the individual CMIP5 models (RCP8.5) are relatively stable from the early 1980s to the mid 1990s, but show substantial shifts beginning about 1995 (Fig. 2B), with one model that had previously fit the observed time series well (CESM1) falling out of favor, and several others beginning to show a better fit (HadGEM2-CC, IPSL-CM5A-MR, MRI). Throughout this time period, a linear weighting of the CMIP5 models can be found that produces a satisfactory fit to the observations, suggesting the model set is still bounding the behavior of the system (Fig. 2B, black line). Nevertheless, the sharp changes in the individual K-S statistics serve as an early indicator that the Arctic system is changing in a way that is not captured by any one of the current CMIP5 models with the RCP8.5 forcing scenario. If that trend continues, the K-S statistic for the best-fit weighted model may begin to indicate a failure of the entire model set, triggering the need for new model development. This suggests that, for the moment, the current set of models can be used by decision-makers concerned about Arctic sea ice, but a watchful eye is needed to be sure the model set still bounds the observations over the coming years.

The best fit linear weighting of the CMIP5 models changed over time (Fig. 2C), particularly after 1995, when the observed September sea ice extent began to drop relative to the multi-model

ensemble of predictions. The best fit model can be used to forecast the sea ice extent in the future (Fig. 3), with the forecast changing as the model weights are updated with each year's observation. Since 2000, the forecast September 2055 sea ice extent under the RCP8.5 emission scenario has dropped; the most recent forecast (based on data through 2015) is 0.77 million sq. km (90% prediction interval: 0.10-1.45), very close to what is considered an "ice-free" Arctic. The probability that the sea ice extent will be below 1.0 million sq. km increased from 44% using the model weights in 2000 to 71% using the model weights in 2015 (Fig. 3).

Tracking system change

The model weights, the EDF statistics for the individual models, and the EDF statistic for the best-fitting weighted model provide a way to track system change and evaluate the multi-model ensemble. A shift in model weights over time may be an indicator that the dynamics of the system are changing (or that if the system dynamics are in fact stationary, such stationarity is not captured by the models in the ensemble). If the EDF statistic for the best-fitting weighted model remains plausible, then the multi-model ensemble is bounding the behavior of the system. But if the EDF statistic for even the best-fitting weighted model is not plausible, then the ensemble is not functioning; a double-loop adaptation should be triggered, and the model set should be examined to try to explain the emerging surprises. In the case of northern pintails, this would have brought awareness to the change in system dynamics in 1985, twenty years before the effect was in fact identified and incorporated into management of hunting regulations. In the case of Arctic sea ice extent, although the model set currently bounds the observed system behavior, rapid shifts in the plausibility of individual models are an early warning that the current model set might be starting to fail.

Methods

Data and models. For northern pintails, the data are the observed latitude of the breeding population in North America, 1961-2015, taken from the Waterfowl Breeding Population and Habitat Survey. Two models, both normally distributed, were compared: Model 1 predicted a constant mean and variance (based on the mean and variance of the observed latitude, 1961-1974); Model 2 used the 5-year moving average at 1985 as the mean, and the same variance as Model 1. Both Models 1 and 2 use a fixed long-term mean, rather than a more complicated time-series model because harvest regulations for northern pintails are set assuming a fixed long-term mean for the latitude of the breeding population.

The sea ice data measure the extent of sea ice in September (million sq. km). The observational record is based on a combination of passive microwave sea ice concentrations from the NASA Team sea ice algorithm²³ and earlier satellite, aircraft, and ship observations available from the HadISST data set²⁴ that were merged to create a consistent time-series²⁵. Hindcast and forecast September sea ice extent was extracted from 11 CMIP5 models (CCSM4, 6 ensemble members; CESM1-cam5, 3; EC-EARTH, 12; GFDL-CM3, 1; HadGEM2-AO, 1; HadGEM2-CC, 1; HadGEM2-ES, 4; IPSL-CM5A-LR, 4; IPSL-CM5A-MR, 1; MIROC5, 2; MRI-CGCM3, 1), using the RCP8.5 forcing scenario. The CMIP5 models use observed greenhouse gas concentrations through 2005 and forecast concentrations thereafter. The subset of 11 was chosen from the full set of CMIP5 models based on their ability to capture basic features of the Arctic climate, as reflected in observed ice thickness distributions²⁶.

The CMIP5 model results are replicate simulations taking into account temporal variation, parametric uncertainty, and uncertainty in starting conditions; each replicate is a possible future trajectory. These results, however, are not in themselves probabilistic forecasts of sea-ice extent. To develop probabilistic forecasts of sea-ice extent, we used the replicate CMIP5 results to estimate time-specific means and variances. For each of the CMIP5 models, a year-specific mean was estimated with LOESS smoothing ($\lambda = 2$, 25-yr window for α), and a corresponding year-specific variance was estimated with LOESS smoothing of the variance of the residuals. The year-specific forecast was a normal distribution with the corresponding mean and variance. This method for developing probabilistic forecasts from the CMIP5 model results, including the assumption of a normal distribution, is one possible approach and appears to work well for the sea-ice metric; other approaches and distributions have been explored³ and may be more appropriate for other metrics.

Individual model fit. To assess the fit of each model to the data, a moving window was used (10 years for the pintail data, 30 years for the sea ice data). Within the moving window, the observations were expressed as a normalized residual from the corresponding year-specific predicted distribution. An empirical cumulative distribution function was formed from the set of residuals within the window and compared against the cumulative distribution function for a standard normal distribution to calculate the Kolmogorov-Smirnov statistic (equation 1) or Anderson-Darling statistic, with an appropriate critical value¹⁵. In a decision context, the choice of the critical value is an important value judgment that reflects the relative importance of Type 1 and Type 2 errors, and the nominal critical value needs to be adjusted to account for multiple comparisons as well as the estimation of parameters²⁷. These topics are investigated in detail in

the Supplementary Information. Throughout the main body of the paper, we have used the critical value associated with a nominal Type 1 error rate (α) of 0.05.

Weighted models. Weighted models were formed from the component models with linear weighting of the first two moments. For example, to combine the 11 sea-ice models, a set of 11 weights (summing to 1) were used to weight the 11 means and the 11 standard deviations. A set of weights were evaluated by calculating the K-S statistic for the weighted model in the preceding window associated with a particular point in time. The best-fit weighted model at each point in time was found by searching for the set of weights that minimized the K-S statistic: for the pintail example using multivariate constrained optimization, specifically, sequential quadratic programming²⁸; for the sea ice example using multivariate unconstrained optimization, specifically, a gradient-based quasi-Newton method²⁹ with a cubic line search procedure²⁸.

The sea-ice forecast based on a weighted model (Fig. 3) used the year-specific means and variances from the 11 CMIP5 models, weighted by the best fit set of weights. Each of the 11 models has a forecast for the sea-ice extent in September 2055; these forecasts were weighted by the sets of weights at each point in the observational record. The quantiles were found by assuming the weighted forecast was normally distributed.

References

- 1 McDonald-Madden, E., Runge, M. C., Martin, T. G. & Possingham, H. Optimal timing for managed relocation of species faced with climate change. *Nature Climate Change* 1, 261-265 (2011).
- 2 Conroy, M. J., Runge, M. C., Nichols, J. D., Stodola, K. W. & Cooper, R. J. Conservation in the face of climate change: The roles of alternative models, monitoring, and adaptation in confronting and reducing uncertainty. *Biol. Conserv.* 144, 1204-1213, doi:10.1016/j.biocon.2010.10.019 (2011).
- 3 Terando, A., Keller, K. & Easterling, W. E. Probabilistic projections of agro-climate indices in North America. *J. Geophys. Res.* 117, D08115, doi:10.1029/2012JD017436 (2012).
- 4 Lawler, J. J. *et al.* Resource management in a changing and uncertain climate. *Frontiers in Ecology and the Environment* 8, 35-43 (2010).
- 5 IPCC. *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change.* (Cambridge University Press, 2013).
- 6 Taleb, N. N. *The Black Swan: The Impact of the Highly Improbable.* (Random House, 2007).
- 7 Pahl-Wostl, C. A conceptual framework for analysing adaptive capacity and multi-level learning processes in resource governance regimes. *Global Environmental Change* 19, 354-365 (2009).
- 8 Argyris, C. & Schön, D. A. *Organizational Learning: a Theory of Action Perspective.* (Addison-Wesley, 1978).
- 9 Stroeve, J., Holland, M. M., Meier, W., Scambos, T. & Serreze, M. Arctic sea ice decline: Faster than forecast. *Geophysical Research Letters* 34, L09501 (2007).
- 10 Miller, M. R. & Duncan, D. C. The northern pintail in North America: status and conservation needs of a struggling population. *Wildl. Soc. Bull.* 27, 788-800 (1999).
- 11 Dawid, A. P. Statistical theory: the prequential approach. *Journal of the Royal Statistical Society A* 147, 278-292 (1984).
- 12 Williams, B. K., Johnson, F. A. & Wilkins, K. Uncertainty and the adaptive management of waterfowl harvests. *The Journal of Wildlife Management* 60, 223-232 (1996).
- 13 Gneiting, T., Balabdaoui, F. & Raftery, A. E. Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69, 243-268 (2007).
- 14 Seillier-Moiseiwitsch, F. & Dawid, A. P. On testing the validity of sequential probability forecasts. *Journal of the American Statistical Association* 88, 355-359 (1993).
- 15 Stephens, M. A. in *Goodness-of-fit Techniques* Vol. 68 (eds Ralph B D'Agostino & Michael A Stephens) Ch. 4, 97-193 (Marcel Dekker, 1986).
- 16 Austin, J. E. & Miller, M. R. in *The Birds of North America* (ed A Poole) Issue 163 (Cornell Laboratory of Ornithology, 1995).

- 17 Hestbeck, J. B. Response of northern pintail breeding populations to drought, 1961-
92. *The Journal of Wildlife Management* 59, 9-15, doi:10.2307/3809109 (1995).
- 18 U.S. Fish and Wildlife Service. *Adaptive Harvest Management: 2014 Hunting Season*.
(United States Department of Interior, 2014).
- 19 Stroeve, J. *et al.* The Arctic's rapidly shrinking sea ice cover: a research synthesis.
Clim. Change 110, 1005-1027, doi:10.1007/s10584-011-0101-1 (2012).
- 20 Amstrup, S. C. *et al.* Greenhouse gas mitigation can reduce sea-ice loss and increase
polar bear persistence. *Nature* 468, 955-958 (2010).
- 21 Hunter, C. M. *et al.* Climate change threatens polar bear populations: a stochastic
demographic analysis. *Ecology* 91, 2883-2897 (2010).
- 22 Stroeve, J. C. *et al.* Trends in Arctic sea ice extent from CMIP5, CMIP3 and
observations. *Geophysical Research Letters* 39, L16502 (2012).
- 23 Cavalieri, D. J., Parkinson, C. L., Gloersen, P. & Zwally, H. J. *Sea ice concentrations*
from Nimbus-7 SMMR and DMSP SSM/I-SSMIS passive microwave data. (NASA
DAAC at the National Snow and Ice Data Center, 1996, updated yearly).
- 24 Rayner, N. A. *et al.* Global analyses of sea surface temperature, sea ice, and night
marine air temperature since the late nineteenth century. *Journal of Geophysical*
Research: Atmospheres (1984–2012) 108, 4407 (2003).
- 25 Meier, W. N., Stroeve, J. & Fetterer, F. Whither Arctic sea ice? A clear signal of
decline regionally, seasonally and extending beyond the satellite record. *Annals of*
Glaciology 46, 428-434 (2007).
- 26 Stroeve, J., Barrett, A., Serreze, M. & Schweiger, A. Using records from submarine,
aircraft and satellites to evaluate climate model simulations of Arctic sea ice
thickness. *The Cryosphere* 8, 1839-1854, doi:10.5194/tc-8-1839-2014 (2014).
- 27 Babu, G. J. & Rao, C. R. Goodness-of-fit tests when parameters are estimated.
Sankhyā: The Indian Journal of Statistics 66, 63-74 (2004).
- 28 Fletcher, R. *Practical Methods of Optimization*. (John Wiley and Sons, 1987).
- 29 Shanno, D. F. Conditioning of quasi-Newton methods for function minimization.
Mathematics of Computation 24, 647-656 (1970).

Acknowledgments

E.M. was supported by an ARC DECRA Fellowship and by the ARC Centre for Excellence in Environmental Decisions.

Author Contributions

M.C.R. and E.M. conceived of the methods; J.C.S. and A.B. extracted the sea-ice forecasts from the CMIP5 models; M.C.R. analyzed the pintail and sea ice data and prepared the figures; and M.C.R., E.M., and J.C.S. co-wrote the paper.

Competing Financial Interests

The authors declare no competing financial interests.

Disclaimer

This draft manuscript is distributed solely for purposes of scientific peer review. Its content is deliberative and predecisional, so it must not be disclosed or released by reviewers. Because the manuscript has not yet been approved for publication by the U.S. Geological Survey (USGS), it does not represent any official USGS finding or policy.

Figure Legends

Figure 1. Analysis of trends in distribution of northern pintails (*Anas acuta*), 1961-2015. (A) Latitude of the centroid of the pintail breeding distribution. Model 1 (red line) is the average for the period 1961-1974. Model 2 (blue line) is the 5-year moving average, 1965-1985. (B)

Kolmogorov-Smirnov fit statistics for Model 1 (red), Model 2 (blue), and the best weighted model (black), using a 10-year moving window of the data. The nominal critical value ($\alpha = 0.05$) is shown as a dashed line. (C) Weights on Models 1 (red) and 2 (blue) that provide the best fit to a 10-year moving window of observations, as measured by the Kolmogorov-Smirnov statistic.

Figure 2. Analysis of trends in the extent of sea ice in the Arctic, 1953-2015. (A) September Arctic sea ice extent (million sq. km)²². The thin lines show loess means from 11 CMIP5 models (RCP8.5). The thick red line shows the observed record. (B) Kolmogorov-Smirnov fit statistics for the individual CMIP5 models and the best weighted model (thick black line), using a 30-year moving window of the data. The critical value (nominal $\alpha = 0.05$) is shown as a dashed line. (C) Weights on the individual CMIP5 models that provide the best fit to a 30-year moving window of observations, as measured by the Kolmogorov-Smirnov statistic.

Figure 3. Forecast extent of sea ice in the Arctic in 2055 as a function of the weights on the 11 CMIP5 models over the course of the observed record, and assuming the RCP8.5 forcing scenario. The boxplots show the 5%, 25%, 50%, 75%, and 95% quantiles of the weighted model. The dashed line (at 1.0 million sq. km) is frequently cited as the threshold for an “ice-free” Arctic. The black line shows the probability that the sea-ice extent will be less than 1.0 million sq. km in 2055, based on the best-fit weighted model.

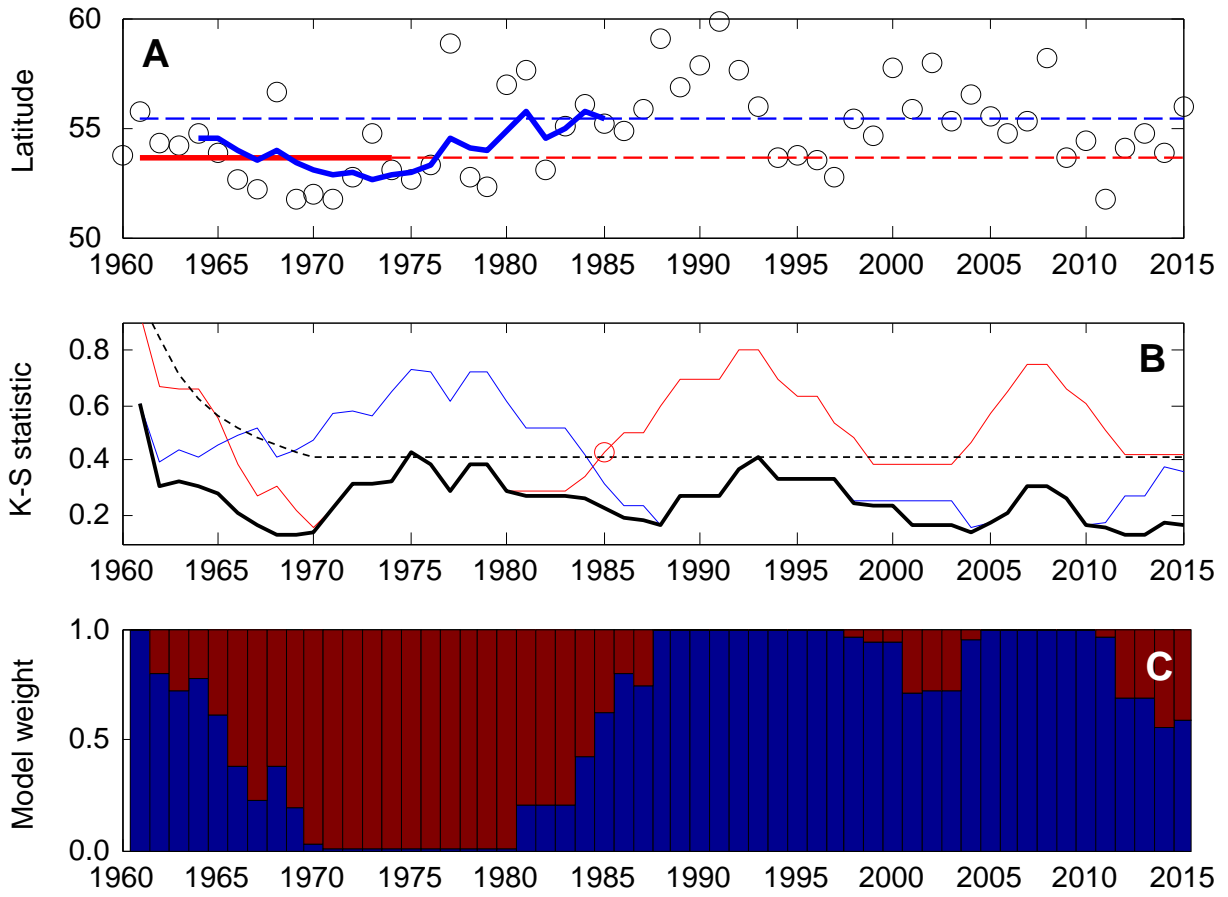


Figure 1.

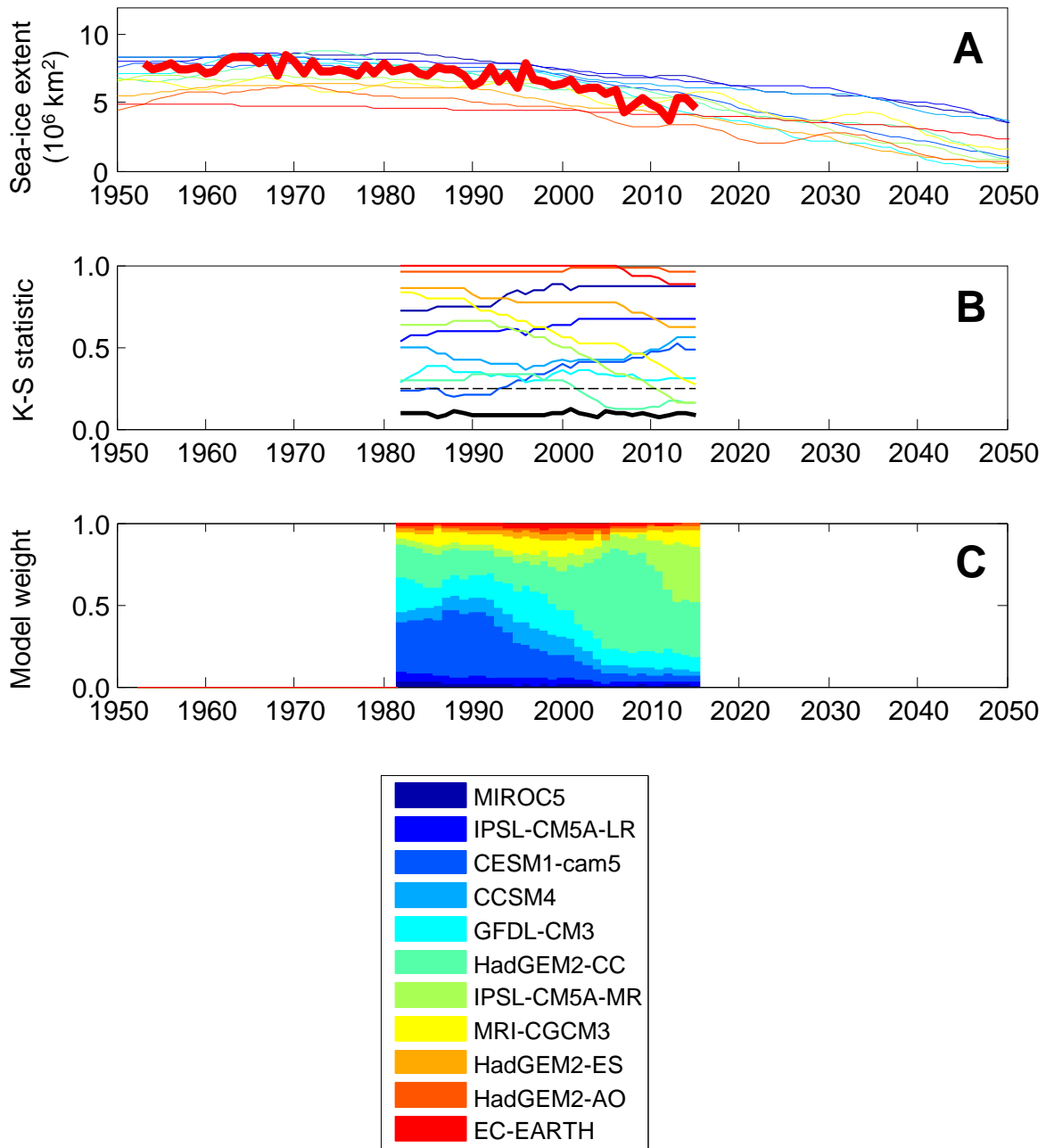


Figure 2.

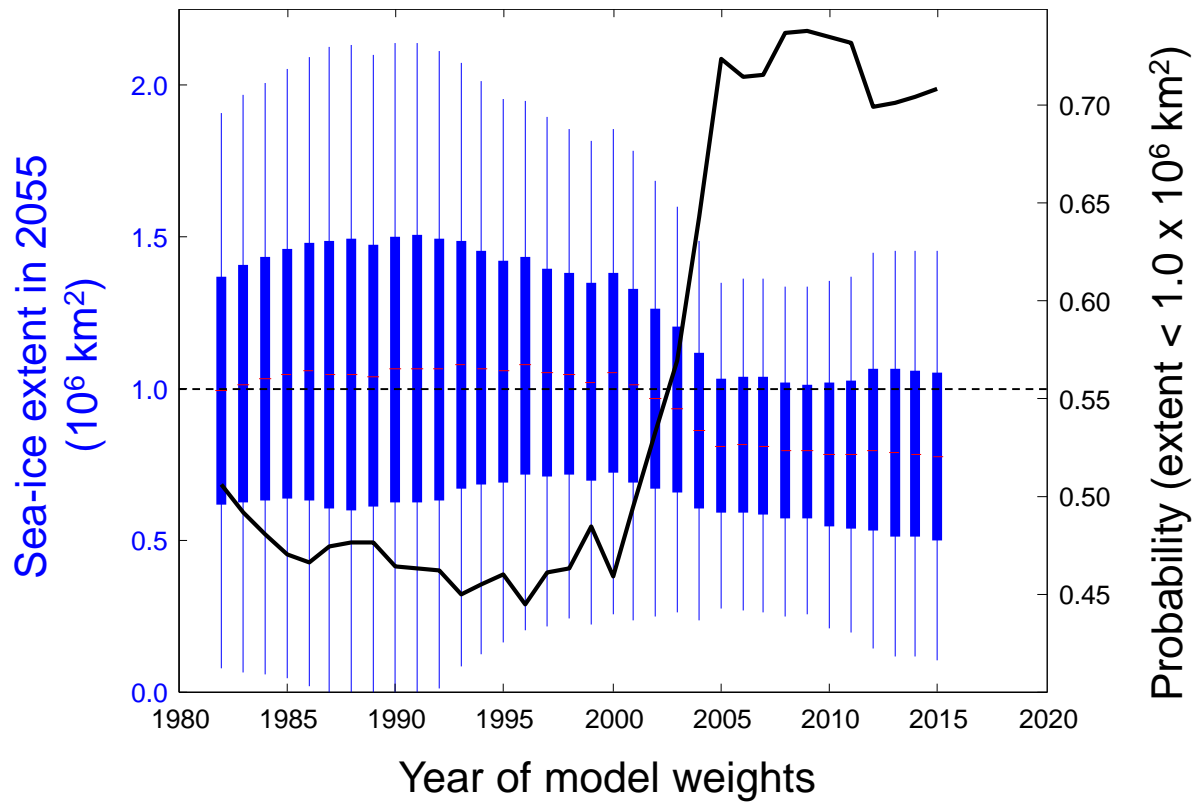


Figure 3.