1 # Inferential, non-parametric statistics to assess quality

2 # of probabilistic forecast systems

3

4 Aline de H. N. Maia[1], Holger Meinke[2], Sarah Lennox[2] and Roger Stone[2,3]

5

6 *[1] Embrapa Meio Ambiente, PO Box 69, Jaguariúna, SP, Brazil;*

7 *ahmaia@cnpma.embrapa.br*

8 *[2] Queensland Department of Primary Industries and Fisheries, Emerging*

9 *Technologies, APSRU, PO Box 102, Toowoomba, Qld 4350, Australia;*

10 *holger.meinke@dpi.qld.gov.au, sarah.lennox@dpi.qld.gov.au,*

11 *roger.stone@dpi.qld.gov.au*

12 *[3] University of Southern Queensland, Faculty of Sciences, Toowoomba, Qld 4350,*

13 *Australia*

14

1   **Abstract**

2   Many statistical forecast systems are available to interested users. In order to

3   be useful for decision-making, these systems must be based on evidence of

4   underlying mechanisms. Once causal connections between the mechanism and their

5   statistical manifestation have been firmly established, the forecasts must also provide

6   some quantitative evidence of `quality'. However, the quality of statistical climate

7   forecast systems (forecast quality) is an ill-defined and frequently misunderstood

8   property. Often, providers and users of such forecast systems are unclear about

9   what 'quality' entails and how to measure it, leading to confusion and

10  misinformation. Here we present a generic framework to quantify aspects of forecast

11  quality using an inferential approach to calculate nominal significance levels (p-

12  values) that can be obtained either by directly applying non-parametric statistical

13  tests such as Kruskal-Wallis (KW) or Kolmogorov-Smirnov (KS) or by using Monte-

14  Carlo methods (in the case of forecast skill scores). Once converted to p-values,

15  these forecast quality measures provide a means to objectively evaluate and

16  compare temporal and spatial patterns of forecast quality across datasets and

17  forecast systems. Our analysis demonstrates the importance of providing p-values

18  rather than adopting some arbitrarily chosen significance levels such as $p < 0.05$ or $p$

19  $< 0.01$, which is still common practice. This is illustrated by applying non-parametric

20  tests (such as KW and KS) and skill scoring methods (LEPS and RPSS) to the 5-phase

21  Southern Oscillation Index classification system using historical rainfall data from

22  Australia, The Republic of South Africa and India. The selection of quality measures

23  is solely based on their common use and does not constitute endorsement. We found

24  that non-parametric statistical tests can be adequate proxies for skill measures such

25  as LEPS or RPSS. The framework can be implemented anywhere, regardless of

26  dataset, forecast system or quality measure. Eventually such inferential evidence

27  should be complimented by descriptive statistical methods in order to fully assist in

28  operational risk management.

# 1. Introduction

Climate variability affects the performance of many climate sensitive systems. Agricultural systems are particularly impacted by climate variability that often results in reduced production volume or quality. Decision makers, be they farmers, policy makers or agribusiness managers, need to devise sound, adaptive risk management strategies in order to improve overall systems performance and to avoid potentially disastrous system failures such as bankruptcy, environmental collapse or famine. Such sound agricultural risk management requires objective assessments of alternative but uncertain outcomes. In highly variable climates, seasonal climate forecasting in combination with simulation models of farming systems has therefore become an important tool for risk assessments and the evaluation of management options (Hammer *et al.* 2000; Sivakumar *et al.* 2000; Ferreyra *et al.* 2001; Meinke and Stone 2005).

Hence objective criteria regarding the performance of forecast systems are required (Hartmann *et al.* 2002). We assert that for appropriate risk management, statistical forecasts must be based on evidence of underpinning mechanisms. Without a plausible explanation for the observed variability in predictors it would be inappropriate to use such forecasts in decision-making. Once some mechanistic basis has been established, other quality attributes of the forecast need to be examined. With this paper we aim to contribute to this process. Information about quality and uncertainty is as important as the forecast itself in order to establish the necessary credibility amongst users. What exactly are these attributes and how should they should be measured? A WMO report (2005) emphasises that only probabilistic forecast systems[1] should be considered for risk management. We concur. The report

---

[1] So far, most operational, probabilistic forecast systems that connect with decision-making tools such as agricultural simulation models are based on an 'analogue year' approach, whereby climate series are segregated into classes corresponding to climate indicators such as the Southern Oscillation Index (SOI), El Niño/ Southern Oscillation (ENSO) phases, sea surface temperature (SST) phases or combination of such indicators. These classes constitute 'conditional climatologies' that need to be compared to the unconditional climatology or reference distribution (Meinke and Stone 2005).

1    lists four key forecast system attributes, namely (a) consistency (whether the

2    forecasts correspond with the forecaster's judgment); (b) quality (whether the

3    forecast correspond with the observations); (c) relevancy (whether what is

4    forecasted is of concern to the user) and (d) value (whether the forecasts are/can be

5    beneficial when used). Each of these attributes deserves further attention. In

6    particular, methods that quantify the quality of probabilistic forecast systems are

7    poorly understood and often misused (Potgieter *et al.* 2003). Therefore, we focus

8    here solely on forecast quality by considering explicitly two aspects that are closely

9    related but often differentiated in the literature: discriminatory ability (DA) and skill.

10    Here we demonstrate the application of an inferential framework for the

11    evaluation of probabilistic, class-based forecast systems. The analogue-years

12    approach is a frequently used example for such class-based systems and has

13    provided valuable information for decision-makers in many world regions (e.g.

14    Messina *et al.* 1999; Singels *et al.* 1997; De Jager *et al.* 1998; Meinke and Hochman

15    2000; Nelson *et al.* 2002; Podestá *et al.* 2002; Selvaraju *et al.* 2004).

16    According to Stone *et al.* (2000), DA is the ability of the forecast system to

17    partition the unconditional probability distribution (also referred to as 'climatology') of

18    the variable of interest (e.g. rainfall, temperature, yield, drainage, runoff) into

19    conditional distributions corresponding to each class or phase within the forecast

20    system (such as, the consistently negative, consistently positive, falling, rising and

21    neutral phases of the SOI-phase system). They emphasise that DA '*does not*

22    *necessarily imply the level of forecasting skill that would be determined from a test*

23    *on independent data of forecast model performance*'. DA represents the additional

24    knowledge about future states arising from some forecast system over and above

25    the total variability of the prognostic variable (climatology in the case of our study

26    here). Note that *discriminatory ability* as defined by Stone *et al.* (2000) is different to

27    forecast discrimination (Wilks 1995; Murphy 1993). DA is concerned with

1   distributions of observations only and does not attempt to make any comparison

2   between forecasts and observations (in contrast to *forecast discrimination*).

3   Most skill measures were originally designed to quantify changes in the

4   agreement between observed and predicted **values** (accuracy) of deterministic

5   forecasts with some attempts to incorporate probabilistic properties (Mason 2004).

6   Skill measures are supposed to account for changes in accuracy, relative to using the

7   reference system as a framework (Murphy 1993; Potgieter *et al.* 2003). However, in

8   order to appropriately evaluate probabilistic forecast systems based on analogue

9   years, better skill measures are required to appropriately account for the probabilistic

10   nature of these systems (Potgieter *et al.* 2003).

11   DA is associated with variability of the observations among classes. For such

12   forecast systems, there is no single, predicted value corresponding to each

13   observation. Instead, the forecast consists of a set of possible values represented by

14   empirical distributions functions (CDF) derived from previously observed values. The

15   lack of clear distinction between sets of predicted and observed values and the

16   probabilistic nature of those predictions needs to be taken into account when

17   developing and applying skill measures.

18   Skill scores developed to quantify hindcast skill (e.g. LEPS skill score and RPSS)

19   of probabilistic forecast systems that produce categorical forecasts (probabilities of

20   belonging to predefined intervals or classes) are now in common use, in spite of their

21   limitations. Class-based forecast systems do not readily lend themselves to such

22   categorical evaluations without the loss of at least some valuable information by

23   reducing the full probabilistic nature of the forecast systems to some broad bands of

24   categories such as intervals defined by terciles (Potgieter *et al.* 2003). However,

25   changes in agreement between observations and predicted probabilities for the

26   predefined classes are directly related to DA, i.e. divergences between the empirical,

1   conditional CDFs corresponding to each forecast system class and the unconditional
2   CDF arising from 'climatology'. Hence, DA measures can be used as indirect skill
3   measures for class-based forecast systems.

4      Inferential methods proposed here only quantify the degree of evidence against
5   a null-hypothesis of either 'no DA' or 'no skill'. This information is essential but not
6   sufficient for sound risk management. Decision makers also require complementary
7   knowledge about the magnitude of expected change in the forecast variable. To
8   quantify this magnitude requires descriptive measures (e.g. distance among
9   cumulative distribution functions; magnitude of differences among conditional
10  median or mean values etc) rather than inferential statistical methods. However, an
11  in-depth evaluation and discussion of such descriptive measures is beyond the scope
12  of this paper. The objective of this paper is to provide a generic, inferential
13  framework for hypothesis testing that will add value to descriptive assessments of
14  forecast quality.

15      The proposed inferential approach is based on distribution-free statistical
16  methods that include both traditional non-parametric tests (e.g. Kolmogorov-Smirnov
17  test) and computationally intensive methods based on non-parametric Monte Carlo
18  techniques (e.g. bootstrapping and randomization tests). P-values derived from those
19  distribution-free procedures are used to quantify evidences of 'true' DA and skill. This
20  approach can be applied when the underlying probability distributions are unknown
21  (it is not necessary to specify a particular distribution such as normal, gamma, etc),
22  and it does not require any arbitrarily chosen level of significance. P-values range
23  between 0 and 1 and are inversely proportional to the degree of evidence against the
24  hypothesis of 'no class effect'. This approach takes into account the length of the
25  time series, the number of classes of the chosen classification system and the intra-
26  class variability. Further, given adequate spatial coverage, p-values can be mapped
27  using interpolation methods, providing a powerful and intuitive means of

1    communicating the spatial variability of DA and skill. We illustrate this approach by

2    quantifying DA and skill of the 5-phase Southern Oscillation Index (SOI) classification

3    applied to forecasting rainfall across Australia, Republic of South Africa and India.

## 4    2. Material and methods

5    We used 3-monthly rainfall totals from 3 sample stations (one each from  Australia,

6    The Republic of South Africa and India) and gridded (0.5 degree by 0.5 degree)

7    rainfall data for each of these countries to demonstrate the use of p-values to

8    measure aspects of discriminatory ability and skill of seasonal forecast systems. The

9    3 sample stations were Echuca (Australia), Bangalore (India) and Bloemfontein (The

10   Republic of South Africa). Rainfall data for Echuca was obtained from the SILO

11   Patched Point Dataset (PPD; Jeffrey et al, 2001), while Bangalore and Bloemfontein

12   rainfall data are generated from GHCN Data (NOAA; Version 2: Peterson and

13   Easterling 1994; Easterling and Peterson 1995; Easterling *et al.* 1996). Gridded data

14   from the Hadley/CRU 0.5x0.5 global rainfall grid (New *et al.* 2000) was used for the

15   spatial analyses of Australia, India and South Africa due to a lack of recent Indian

16   rainfall records.

17        All 3 sample stations had at least 94 years of daily rainfall records. They

18   represent vastly different climatic and agricultural regions.

19        The SOI 5-phase forecast system (SOI-5 FS) considered here is based on an

20   analogue year approach and is underpinned by a sound, physical understanding of

21   the ENSO cycle (Stone *et al.* 1996). The system has been used extensively for

22   decision-making in Australia (e.g. Hammer et al. 2000) and elsewhere (e.g. Hill et al.

23   2000, 2004; Selvaraju et al. 2004). Years were categorised into five analogue sets

24   according to their similarity regarding oceanic and/or atmospheric conditions as

25   measured by SOI phases just prior to the 3-months forecast period. Hence, the

26   rainfall time series were segregated into sub-series corresponding to each SOI class

27   (consistently negative, consistently positive, falling, rising and neutral), resulting in 5

28   sub-series with variable record lengths. These rainfall time series were represented

1    by their respective cumulative distribution functions (CDFs) or their complement,

2    probability of exceedance functions (POEs): a conditional CDF (or POE) for each class

3    and an unconditional CDF (or POE) for 'climatology'. Cumulative probabilities are a

4    simple and convenient way to represent probabilistic information arising from a time

5    series that exhibits no or only weak auto-correlation patterns. However, if the time

6    series shows moderate to strong auto-correlation patterns, a CDF/POE summary will

7    result in some loss of information. Yearly sequences of rainfall data from a specific

8    month or period exhibit only weak auto-correlation, thus allowing the CDF/POE

9    representation to convey seasonal climate forecast information (e.g. Selvaraju *et al.*

10   2004). Figure 1 provides an example of rainfall categorisation based on the SOI

11   classes for the three sample locations.


12


13   *a. Inferential statistical methods to quantify DA and skill*


14       The proposed inferential framework can be applied in conjunction with any

15   statistical test or skill measure. As an example we implemented the approach using:

16       (i) Two nonparametric statistical tests to quantify DA, namely the Kruskall-Wallis

17          (KW) test for comparing medians and the multi-sample Kolmogorov-Smirnov

18          (multi-sample KS) test for comparing CDFs (Conover 1980; Stone *et al.* 1996;

19          Stone *et al.* 2000).

20       (ii) Randomization tests for quantifying evidences of skill as measured by two

21          descriptive skill scores, LEPS (linear error in the probability space, LEPS score;

22          Potts *et al.* 1996) and RPSS (ranked probability skill score; Epstein 1969), which

23          is an operational forecast evaluation procedure used by the International

24          Research Institute (Goddard *et al.* 2003).

25       The KW test is a generalization of the Wilcoxon-Mann-Whitney test applied to

26   three or more groups (Stokes *et al.* 2000). It accounts for overall divergences among

1 medians of conditional CDFs[2] and is the non-parametric test equivalent to the F-test

2 used in analysis of variance. KW accounts for divergences among locations only,

3 while the multi-sample KS test, based on maximum vertical distances among CDFs,

4 takes a whole distribution approach to testing. The KS test is therefore able - unlike

5 the KW test - to detect differences due to both spread and/or shape of distributions.

6     We used p-values associated with the observed KW and multi-sample KS

7 statistics as evidences against the 'no discriminatory ability' null hypothesis. To

8 demonstrate the universal applicability of this inferential framework, we have

9 included the two popular skill measures LEPS and RPSS in spite of our reservations

10 regarding their ability to adequately represent the probabilistic features of class-

11 based forecast systems (see earlier). At the very least, the examples show how

12 quantitative descriptive measures can be converted into inferential measures, thus

13 providing a much sounder basis for comparative assessments.

14     The use of statistical tests to assess spatial variability of DA was proposed by

15 Stone (1992), who produced contour maps of significance levels of KW tests applied

16 to SOI grouping on rainfall medians over Australia. Our approach differs in one

17 important aspect: we do not propose the use of any predetermined cut-off levels for

18 evidence against the null hypothesis. Instead of choosing significance levels and

19 verifying if DA or skill is significant or not, we provide the nominal significance levels

20 (p-values). This avoids the loss of potentially valuable information and provides

21 informed users with the opportunity to form their own opinion whether or not the

22 evidence is sufficient to influence decision-making. Stone *et al.* (2000) has already

23 used the KW p-values to quantify discriminatory ability of GCM-derived analogue

24 forecast systems.

---

[2] Kruskal-Wallis is a nonparametric test for the null hypothesis that the distribution of an ordinally scaled response is the same in two or more independently sampled populations. It is sensitive to the alternative hypothesis that there is a location difference between at least a pair of populations (Stokes *et al.* 2000). It requires the assumption of same population variances when used for comparing distributions. In our case studies, we are using KW for comparing medians, thus homogeneity of variances is not required.

1    Both LEPS skill scores and RPSS quantify the departure between a categorical

2  forecast and observations in the cumulative probability space (Zhang and Casey

3  2000). Here we used calculations for LEPS skill scores and RPSS, which are based on

4  agreement between actual rainfall values and predicted probabilities for intervals

5  defined by the empirical terciles of the corresponding cross-validated forecast

6  distributions.


7    Every empirical, descriptive skill measure, including LEPS skill scores and RPSS,

8  needs to be complimented by some measure of uncertainty before the information

9  can be confidently applied in decision making (Potts et al. 1996; Zhang and Casey

10  2000; Jolliffe 2004). Beyond assessing the skill magnitude (observed skill score), it is

11  critically important for users of forecast systems to know the probability of such skill

12  arising by chance, in order to avoid making decisions based on artificial or perceived

13  skill. This probability is used to assess the true class effect, considering the time

14  series size (record length) and other sources of variability, not explained by the

15  classification system used.


16    However, appropriate null-distributions for such assessments are not readily

17  available when using the standard LEPS skill scores or RPSS. Zhang and Casey

18  (2000) therefore proposed the construction of 'statistical distributions' using quasi-

19  random experiments in order to assess the significance of forecast skill. They

20  generated 95 'quasi-random ensembles' from the rainfall time series (1900-1995) at

21  each station by shifting the observations 1 year ahead at a time and replacing, after

22  each iteration, the first observation with the last. A single 'statistical distribution' was

23  then derived for each skill measure, using the 95 skill values arising from each grid

24  location in Australia. Those 'statistical distributions' were used as 'null distributions'

25  for performing skill significance tests. However, those distributions are not

26  appropriate for assessing significance or calculating p-values because they do not

27  adequately represent the set of possible values of the skill measure under the

1   hypothesis of 'no skill'. Further, existing spatial variability of skill is misinterpreted as

2   'quasi-random variation'. Combining skill values from different locations into one

3   distribution ignores the existence of well-established, spatial correlation patterns.

4   Hence, we propose a location-by-location assessment, which allows the construction

5   of true null distributions of each skill measure at each location based on

6   randomisation techniques (Monte-Carlo analyses).


7        For each location, we calculated p-values associated with LEPS skill scores and

8   RPSS using Monte Carlo methods by randomly allocating all the observed rainfall

9   data 5000 times to the five SOI sub-classes and calculating cross-validated skill score

10  values for each random allocation in order to derive empirical null distributions for

11  both skill measures. Such distributions represent the set of possible skill score values

12  under the null hypothesis of 'no skill'. Considering the alternative hypothesis of skill

13  score for forecast system > skill score for climatology, p-values associated with

14  observed skill scores for both measures were calculated as the relative frequency of

15  skill scores that exceeded the respective observed skill scores (see example for 3

16  locations, Fig. 2).


17  *b. Spatial and temporal assessments of forecast quality*


18      To demonstrate the usefulness of p-values for spatial analyses, we mapped the

19  KW and LEPS skill score p-values[3] from all grid points for the periods analysed. This

20  allowed us to examine spatial patterns of forecast quality associated with the chosen

21  forecast system. For a temporal assessment of DA and skill we investigated month-

22  by-month changes in p-values associated with KW, multi-sample KS, LEPS skill scores

23  and RPSS at the three sample locations (Fig. 3).


24

---

[3] Mapping values arising from KS and RPSS yielded near-identical results and were therefore omitted.

# 3. Results and discussion

Here we presented a generic inferential framework for quantifying forecast quality attributes associated with probabilistic forecast systems, namely, skill and discriminatory ability. We discuss that the distinction between the two concepts becomes blurred, in the case of probabilistic forecast systems based on analogue years approach. We selected KW and KS to quantify DA because of their non-parametric nature that allows us to apply these statistical tests without the need for distributional assumptions. We selected LEPS and RPSS because they are two frequently used scoring systems (Mason 2004).

Skill scores associated with these systems are not very informative on their own, difficult to interpret and need to be accompanied by some measure of uncertainty to be useful (Hartmann *et al.* 2002; Potgieter *et al.* 2003; Jolliffe 2004). Any other statistical test or skill measure can be used with this generic inferential framework (e.g. for DA, Median or Log-Rank tests; for skill, measures such as Brier skill scores and many more, see Potgieter *et al.* 2003; Mason 2004). Here KW, multi-sample KS, LEPS and RPSS merely serve as examples to demonstrate the overall approach, the choice of the quality measure depends on the objective of the study. Our choice of LEPS or RPSS does not constitute an endorsement of these measures - they must be adequate for testing the hypothesis under investigation[4]. Using a statistical hypothesis testing approach, we converted observed skill scores into corresponding nominal significance levels (p-values). This accounts for differences in record length and number of classes and thus enables (i) investigation of temporal variability in forecast quality (e.g. length and timing of the 'autumn predictability barrier' of ENSO based forecast systems), (ii) objective comparisons of DA/skill among sites or assessments of spatial patterns of DA/skill over regions, (iii)

---

[4] When different statistical tests are available for the same hypothesis, a power analysis would provide objective criteria for choosing the most appropriate test. This could be achieved by analysing a large number of Monte Carlo samples (sets of conditional distributions) drawn from synthetically constructed distributions with known properties (divergences among conditional CDFs).

1    assessment of congruence of DA/skill magnitude as measured by different skill

2    scores and (iv) performance evaluation of different forecast systems at a location

3    and/or regionally.


4


5    *a. Converting skill scores into p-values*

6           In our examples, the relative location of the observed skill scores (Fig. 2, dark,

7    thick line) on the skill scores' empirical null distributions indicates the degree of

8    evidence against the hypothesis of 'no skill'. The higher the observed skill score value

9    relative to the null-distribution, the greater the empirical evidence of true skill of the

10   forecast system. For instance, the LEPS and RPSS skill score p-values of the SOI-5 FS

11   to predict JAS Echuca rainfall were both < 0.001. This indicates highly significant and

12   similar forecast skill, regardless of skill score used. The LEPS skill score (RPSS) p-

13   values for JAS at Bangalore, 0.009 (0.027) and for NDJ at Bloemfontein 0.002 (<

14   0.001) also indicate highly significant and similar forecast skill (Fig. 3). The results

15   show that the conversion of skill scores into p-values can overcome the issue raised

16   by Mason (2004) of some skill measures, such as Brier and RPSS, having negative

17   skill score values that can still be indicative of forecast system skill (as demonstrated

18   by low p-values shown for Bangalore, Bloemfontein and Echuca). This goes some

19   way towards overcoming the lack of equitability of some scoring systems also

20   addressed by Mason (2004).


21           Goddard and Dilley (2005) commented that Monte Carlo re-sampling would not

22   be appropriate to assess nominal significance levels (p-values) for RPSS because

23   '*forecasts drawn at random relative to the observed sequence of years typically yield*

24   *a RPSS worse than that of climatology*'. We disagree. As explained by Mason (2004),

25   such negative scores are an inherent feature of RPSS and negative scores can occur

26   even when true forecast skill exists. As the expected value of RPSS under the null

13

1 hypothesis is influenced by the forecast system employed, empirical null distributions

2 generated using Monte Carlo techniques can contain a high frequency of negative

3 values (Mason 2004) as shown in Fig 2. Such a negative bias of the RPSS expected

4 value (in contrast to LEPS) does therefore not invalidate the use of Monte Carlo

5 techniques for establishing p-values associated with skill scores.

6      In this context, we need to flag the issue of how such null distributions can be

7 constructed. For class-based forecast systems using conditional distributions, random

8 reallocation of climate data to these classes of conditional probabilities is the

9 intuitively obvious method. A truly probabilistic assessment of GCM-based forecasts is

10 more difficult to obtain. Allen and Stainforth (2002) criticise the 'probabilistic' outputs

11 generated by GCMs through altering initial and boundary conditions without explicitly

12 accounting for the climate's response. They argue that climate forecasts are

13 intrinsically five-dimensional, spanning space, time and probability, a fact not

14 accounted for when compiling subjective GCM-based probability distribution of

15 forecasts. More attention to formal uncertainty analyses is required, including much

16 more rigorous sensitivity testing based on many more elaborate ensemble runs,

17 before reliable GCM-based probabilistic trajectories of future climate states can be

18 provided (Meinke *et al.* 2004). For now, the methods suggested by Stone *et al.*

19 (2000) provide a way to develop class-based forecast systems from GCM outputs,

20 allowing an immediate application of the generic inferential framework we developed

21 here.

22

23 *b. Quantifying temporal patterns of forecast quality attributes*

24      Particularly with ENSO-based forecast system, forecast quality attributes will

25 vary temporally due to the well-known, seasonal life cycle of the ENSO phenomenon.

26 Therefore, location-specific temporal analyses are necessary to evaluate when the

1    forecast system is sufficiently informative to influence decision making. Here we

2    show the usefulness of the generic inferential framework to simultaneously assess

3    temporal patterns. These tests revealed some interesting insights regarding

4    predictability and dynamics of seasonal rainfall patterns that should be investigated

5    in more detail elsewhere (Fig. 3):

6        (i) Regardless of test or skill score, p-values at all locations followed similar time

7           courses, albeit with some exceptions;

8        (ii) The autumn predictability barrier (Clarke and Shu 2000; Clarke and van

9           Gorder 2003) around MAM is clearly evident at all locations.

10      (iii) For Echuca the typical ENSO-lifecycle is evident, with a strong impact on

11          winter and spring rainfall;

12      (iv) Bloemfontein, a region that is seasonally dry in winter, shows ENSO impact

13          for the beginning of the rainy season around October, while Bangalore shows

14          a strong impact for the (northern) summer monsoon (MJJ to ASO) and a small

15          peak for the much weaker winter monsoon (OND).

16    At the 3 locations, all measures broadly identified similar trends, but differed in

17    detail. The fact that there is broad congruence between p-values regardless of test

18    or measure employed demonstrates our earlier assertion that discriminatory ability

19    can be used as a surrogate for skill for class-based, probabilistic forecasts. Further, it

20    shows the generic inferential framework's ability to reconcile vastly different

21    measures (Fig 3).

22    Although discriminatory ability and skill tests generally follow fairly consistent

23    patterns (temporally as well as spatially; Figures 3 & 4), there may be situations

24    when results may differ substantially (e.g. JJA rainfall for Echuca; Figure 3). Results

25    where the p-values from discriminatory ability tests (such as KW and the multi-

1 sample KS) are much smaller than the p-values from skill tests can be attributed to

2 procedural differences between the two types of tests. KW/KS test for differences in

3 at least 1 phase (or class) distribution. There will be instances when a single phase

4 distribution is significantly different from all the other distributions, which do not

5 differ from each other, resulting in a low p-value. For a skill test, this is unlikely to

6 yield a low p-value as such tests are designed to compare observed data with

7 forecast data. Even though these circumstances will arise occasionally, Figures 3 and

8 4 show that p-values are generally very consistent between tests, broadly indicating

9 the same temporal and spatial trends. This is particularly the case when p-values are

10 low, indicating a convergence of results when real differences between distributions

11 exist. Generally discriminatory ability tests seem to be slightly more emphatic, but

12 our results show that they serve as reliable proxies for skill measures when

13 evaluating class-based forecast systems. Here we would like to add a cautionary

14 note: there is a temptation to 'over-interpret' temporal patterns in p-values as

15 presented in Fig. 3. By definition, p-values indicate the evidence against the null-

16 hypothesis – a value of 0.2 means that in one out of five cases we would falsely

17 reject the null-hypothesis. While the broad temporal patters from all tests are similar,

18 differences are inevitable and it is probably not helpful to discuss whether or not p-

19 values of 0.3 versus 0.8 constitute evidence of 'skill' or otherwise (cf. Echuca, JJA).


20      It is up to the informed user to decide the appropriate level of 'significance' (i.e.

21 the degree of evidence against the hypothesis of 'no class effect') before using the

22 information in decision making (Nicholls 2001). We also note that Nicholls (2001)

23 argues against 'blanking out ' of areas on maps where some feature does not reach

24 statistical significance, a practice often seen in atmospheric science. This practice can

25 lead to the loss of potentially valuable information. Therefore, we argue against the

26 use of any artificial cut-off levels to determine whether or not the p-values of the

27 tests indicate sufficiently high evidence. Instead, we provide all nominal significance

28 levels (p-values) and concur with Nicholls (2001), who questions the appropriateness

1  of commonly used cut-off levels, such as p<0.05 or p<0.01. These cut-offs are no

2  more than a convention that reduce continuous probabilistic information to a

3  dichotomous response, thereby ignoring valuable information contained in the

4  nominal significance levels. Rosnell and Rosenthal (1989), cited by Nicholls (2001),

5  noted that '...*surely, God loves the .06 nearly as much as the .05*'.


6      In our examples, using traditional cut-off levels for significance testing would

7  result in conflicting conclusions for some sites and seasons. For instance, at

8  Bangalore, India, the JAS ENSO signal would be considered either as being

9  'significant' or 'not significant' if a 5% cut-off was adopted, depending on the chosen

10  test or measure. For JAS p-values ranged from 0.09 (KW) to 0.01 (LEPS). This is in

11  spite of a strong and well-established ENSO impact at this location and the fact that

12  p-values for all tests are moderate to low, but not all are below 0.05 (Fig. 3). Should

13  risk managers in the Bangalore region ignore ENSO-based forecasts during this

14  season? Alternatively, should they base their decisions on a single measure using a

15  pre-determined cut-off for significance? Risk managers must decide for themselves

16  whether or not the evidence is strong enough to influence their decisions.


17      We hypothesise that differences among p-values coming from different

18  measures (*ceteris paribus*) might be caused by differences in the ability of each test

19  or measure to detect divergences among conditional distributions regarding the

20  target attribute (e.g. median, variance, whole CDF). For example, for time series that

21  contain more than 50% of rainfall values equalling zero in all classes (all class

22  medians are zero) p-value arising from a median test would yield a value of one,

23  while a test comparing conditional CDFs could produce a low p-value, depending on

24  the differences in the right tails of conditional CDFs. Therefore, lack of agreement

25  between tests or measures can sometimes be explained by differences in the

26  underlying hypotheses tested or existing power differences between tests.

27  Understanding the differences between skill and or DA tests is important and

1    although these differences have not been addressed within this paper, it is the

2    subject of ongoing research. Based on our research here we argue that non-

3    parametric DA measures such as KW and KS are in most cases adequate surrogates

4    for skill measures and little if any additional information can be gained from using

5    skill measures originally designed for deterministic forecasts.


6


7    *c. Quantifying spatial patterns of DA and skill over a region*

8         As we have shown, quality measures of forecast systems vary temporally and

9    spatially. The spatial patterns of DA and skill for the SOI-5 FS based on KW and LEPS

10   p-values (Fig. 4) were consistent with known, typical ENSO impacts. Again, DA and

11   skill measures showed similar spatial patterns, regardless of season - high divergence

12   among conditional CDFs is highly likely to lead to improved agreement between

13   `predicted' and `observed' values as captured by a measure of skill. However, high

14   DA does not necessarily imply high skill (Stone *et al.* 2000) due to, for instance,

15   possible changes in skill over time, a factor not accounted for when calculating DA,

16   but considered during the cross-validation procedure when calculating skill. Hence, it

17   is not unexpected that there appears to be a general tendency for p-values

18   associated with DA to be slightly lower than those associated with skill.


19        Parametric approaches were initially developed during a time when computer

20   power was not available. Due to their reliance on distributional assumptions they are

21   a convenient way to quickly perform hypothesis tests and calculate associated

22   nominal significance levels (p-values) using known, analytically-derived null

23   distributions. Initially, most of the known parametric methods were based on the

24   assumption of normality. Nowadays there are increasing numbers of parametric

25   methods available that are based on a range of distribution types, such as Tweedie

26   family (e.g. Tweedie 1984; Jørgensen 1987), which include the Normal, gamma, and

1    Poisson distributions as special cases and more. Although this greatly broadens the

2    applicability of parametric approaches, spatial assessments of forecast quality still

3    require case-by-case evaluation before parametric methods can be applied. The

4    historical limitation imposed by the lack of computer power no longer holds and it is

5    therefore no longer necessary to make assumptions about distributions – these can

6    now be constructed via non-parametric Monte Carlo approaches, such as

7    bootstrapping and randomisation techniques. This flexible approach is of particular

8    importance for climate science where data sources are varied, underlying

9    distributions can come in many shapes and predictor/predictand relationships are

10   often non-linear (Von Storch and Zwiers 1999).

11

## 4. Conclusions

13       Using an inferential, non-parametric framework we have evaluated aspects of

14   forecast quality using 3-monthly rainfall forecasts for a range of locations in

15   Australia, The Republic of South Africa and India. The approach taken is generic and

16   independent of location, season, data source, statistical test or skill score and

17   provides intuitively simple, but powerful methods that objectively quantify

18   discriminatory ability and skill of probabilistic forecast systems. Forecast quality

19   measures, once converted to nominal significance levels (p-values), can provide the

20   means for investigating temporal and spatial patterns of discriminatory ability and

21   skill. In addition, this allows comparisons among different probabilistic forecast

22   systems according to objective quality criteria – a key issue to further improve risk

23   management in climate-sensitive agricultural systems. In a subsequent step, this

24   framework should be supplemented with descriptive statistical tools that quantify the

25   magnitude of difference between target forecast quantities derived from forecast

26   probability distributions, in addition to the evidence against the null-hypothesis

1 provided by p-values. From a risk management perspective, the ultimate
2 responsibility for the utility of these tools resides with the decision makers. They
3 must be able to properly interpret the relevance of information obtained using these
4 approaches. This requires an ability to handle intrinsically uncertain information
5 (probabilities) as well as asking the relevant questions (Hoffman and Kaplan 1999).

6 The most insightful results were obtained when two very different skill scoring
7 systems, namely LEPS and RPSS converged in terms of p-values and the resulting
8 evidence against the null-hypothesis was similar for both scoring systems and,
9 indeed, even for all four DA and skill measures. This provides strong supporting
10 evidence of the general applicability of this generic, inferential framework to explore
11 and quantify forecast quality. We concur with comments made by Zhang and Casey
12 (2000), Potgieter *et al.* (2003) and Mason (2004), who all stated the need to
13 consider a set of measures due to the multidimensional nature of forecast quality.

14 Finally, the generic inferential framework proposed here might also be useful
15 for evaluating similarities among different sets of quality measures: the approach
16 proposed by Potgieter *et al.* (2003) based on Principal Component Analysis and
17 clustering techniques might reveal more insights when applied to p-values arising
18 from those measures. Further, the value of adapting skill measures based on tercile-
19 based categories in order to provide continuous assessments of forecast quality can
20 now be objectively evaluated. All these issues are subjects of on-going research.

1 **Acknowledgements**

8

9 **References**

10 Allen, M. R., and D. A. Stainforth, 2002: Towards objective probabilistic climate

11     forecasting. *Nature*, **419**, 228.

12 Clarke, A.J. and L. Shu, 2000: Biennial winds in the far western equatorial Pacific

13     phase-locking El Niño to the seasonal cycle, *Geophys. Res. Lett.*, **27(6)**, 771-

14     774.

15 Clarke, A.J. and S. van Gorder, 2003: Improving El Niño prediction using a space-

16     time integration of Indo-Pacific winds and equatorial Pacific upper ocean heat

17     content. *Geophys. Res. Lett.*, **30(7)**, 1399-1402, doi:10-1029/2002GLO16673.

18 Conover, W. J., 1980: *Practical Nonparametric Statistics*. 3d ed. John Wiley & Sons,

19     584 pp.

20 De Jager, J. M., A. B. Potgieter and W. J. van der Berg, 1998: Framework for

21     forecasting the extent and severity of drought in maize in the Free State

22     province of South Africa. *Agricultural Systems*, **57**, 351-365.

23 Easterling, David R., Thomas C. Peterson, and Thomas R. Karl, 1996: On the

24     development and use of homogenized climate data sets. *Journal of Climate*, **9**,

25     1429-1434.

1    Easterling, David R. and Thomas C. Peterson, 1995: A new method for detecting and

2        adjusting for undocumented discontinuities in climatological time series.

3        *International Journal of Climatology,* **15**, 369-377.

4    Epstein, E.S., 1969: A scoring system for probability forecasts of ranked categories.

5        *J. Appl. Meteor.,* **8**, 985-987.

6    Ferreyra, R.A., G. P. Podestá, C. D. Messina, D. Letson, J. Dardanelli, E. Guevara,

7        and S. Meira, 2001: A linked-modeling framework to estimate maize production

8        risk associated with ENSO-related climate variability in Argentina. *Agricultural*

9        *and Forest Meteorology,* **107**, 177–192.

10   Goddard, L. and M. Dilley, 2005: El Niño: Catastrophe or opportunity. *J. Climate,* **18**,

11        651-665.

12   Goddard, L., A. G. Barnston, and S. J. Mason, 2003: Evaluation of the IRI's "Net

13        Assessment" seasonal climate forecasts 1997–2001. *Bull. Amer. Meteor. Soc.,*

14        **84**, 1761-1781.

15   Hammer, G.L., N. Nicholls, C. Mitchell, 2000 (eds): *Applications of seasonal climate*

16        *forecasting in agriculture and natural ecosystems: The Australian experience.*

17        Kluwer Academic Publishers, 469 pp.

18   Hartmann, H.C., T.C. Pagano, S. Sorooshian, and R. Bales, 2002: Confidence

19        builders: evaluating seasonal climate forecasts from a user perspective. *Bull.*

20        *Amer. Meteor. Soc.,* **83**, 683-698.

21   Hill, S. J. H., J. Park, J.W. Mjelde, W. Rosenthal, H. A. Love, and S. W. Fuller, 2000:

22        Comparing the value of Southern Oscillation Index-based climate forecast

23        methods for Canadian and US wheat producers. *Agric. Forest Met.,* **100**, 261-

24        272.

25   Hill, H. S. J., J. W. Mjelde, H. A. Love, D. J. Rubas, S. W. Fuller, W. Rosenthal, and G.

26        Hammer, 2004: Implications of seasonal climate forecasts on world wheat

1    trade: a stochastic, dynamic analysis. *Canadian Journal of Agricultural*

2    *Economics*, **52**(3), 289-312.

3    Hoffman, F.O. and S. Kaplan, 1999: Beyond the Domain of Direct Observation: How

4    to Specify a Probability Distribution that Represents the "State of Knowledge"

5    About Uncertain Inputs. *Risk Analysis*, **19(1)**, 131-134.

6    Jeffrey, S.J., J. O. Carter, K. M. Moodie, and A. R. Beswick, 2001: Using spatial

7    interpolation to construct a comprehensive archive of Australian climate data.

8    *Environmental Modelling and Software*, **16/4**, 309-330. [Available at

9    http://www.nrm.qld.gov.au/silo/ppd/index.html].

10    Jolliffe, I.T., 2004: Estimation of uncertainty in verification measures. *Proc.*

11    *International Verification Methods Workshop*, Montreal, Canada. [Available at

12    http://www.bom.gov.au/bmrc/wefor/staff/eee/verif/Workshop2004/home.html]

13    Jørgensen, B., 1987: Exponential dispersion models. *Journal of the Royal Statistical*

14    *Society, Series B*, **49**, 127-162.

15    Maia, A. H. N., H. Meinke and S. Lennox, 2004: Assessment of probabilistic forecast

16    'skill' using p-values. *Proc. 4th International Crop Science Congress*, Brisbane,

17    Australia.

18    Manly, B.F., 1981: *Randomization tests and Monte Carlo methods in biology*. John

19    Willey & Sons.

20    Mason, S. J., 2004: On using 'climatology' as a reference strategy in the Brier and

21    ranked probability skill scores. *Mon. Wea. Rev.*, **132**, 1891-1895.

22    Meinke, H. and Z. Hochman, 2000: Using seasonal climate forecasts to manage

23    dryland crops in northern Australia. In: G.L. Hammer, N. Nicholls, and C.

24    Mitchell (eds.), Applications of Seasonal Climate Forecasting in Agricultural and

25    Natural Ecosystems – The Australian Experience. Kluwer Academic, The

26    Netherlands, p. 149-165.

1  Meinke H. and R. C. Stone, 2005: Seasonal and inter-annual climate forecasting: the
2      new tool for increasing preparedness to climate variability and change in
3      agricultural planning and operations. *Climatic Change*, **70**, 221-253.

4  Meinke, H., L. Donald, P. deVoil, B. Power, W. Baethgen, M. Howden, R. Allan, and
5      B. Bates, 2004: How predictable is the climate and how can we use it in
6      managing cropping risks? Invited Symposium Paper, *Proc. of the 4th*
7      *International Crop Science Congress*, 26 Sep – 1 Oct 2004, Brisbane, Australia,
8      CD-ROM. [Available online at www.cropscience.org.au]

9  Messina, C.D., J. W. Hansen, and A. J. Hall, 1999: Land allocation conditioned on
10     ENSO phases in the Pampas of Argentina. *Ag. Systems*, **60**, 197-212.

11  Murphy, A.H., 1993: What is a good forecast? An essay on the nature of goodness in
12     weather forecasting. *Weather and Forecasting*, **8**, 281-293.

13  Nelson, R.A., D. P. Holzworth, G. L. Hammer, and P. T. Hayman, 2002: Infusing the
14     use of seasonal climate forecasting into crop management practice in North
15     East Australia using discussion support software. *Ag. Systems*, **74**, 393-414.

16  New, M. G., M. Hulme and P. D. Jones, 2000: Representing twentieth-century space-
17     time climate variability. Part II: Development of 1901-1996 monthly grids of
18     terrestrial surface climate. *J. Climate*, **13**, 2217-2238.

19  Nicholls, N., 2001. The insignificance of significance testing. *Bull. Amer. Meteor. Soc.*,
20     **82**, 981-986.

21  Peterson, Thomas C. and David R. Easterling, 1994: Creation of homogeneous
22     composite climatological reference series. *International Journal of Climatology*,
23     **14**, 671-679.

24  Potgieter, A.B., Y. Everingham, and G. L. Hammer, 2003: Measuring quality of a
25     commodity forecasting from a system that incorporates seasonal climate
26     forecasts. *Int. J. Remote Sensing*, **23**, 1195-1210.

Potgieter, A.B., G. L. Hammer, H. Meinke, R. C. Stone, and L. Goddard, 2005: Three putative types of El Nino revealed by spatial variability in impact on Australian wheat yield. *Journal of Climate*, **18**, 1566-1574.

Potts, J.M., C. K. Folland, I. T. Jolliffe, and D. Sexton, 1996: Revised "LEPS" scores for assessing climate model simulations and long-range forecasts. *J. Climate*, **9**, 34-53.

Podestá, G., D. Letson, C. Messina, F. Rocye, A. Ferreyra, J. Jones, I. Llovet, J. Hansen, M. Grondona, and J. O'Brien, 2002: Use of ENSO-related climate information in agricultural decision making in Argentina: a pilot experience. *Ag. Systems*, **74**, 371-392.

Rosnell, R. L., and R. Rosenthal, 1989: Statistical procedures and the justification of knowledge and psychological science. *Amer. Psychol.*, **44**, 1276-1284.

Selvaraju, R., H. Meinke, and J. Hansen, 2004: Climate information contributes to better water management of irrigated cropping systems in Southern India. *Proc 4th International Crop Science Congress*, Brisbane, Australia.

Singels, A. and A. B. Potgieter, 1997: A technique to evaluate ENSO-based maize production strategies. *South African Journal of Plant and Soil*, **14**, 93-97.

Sivakumar, M.V.K., R. Gommes, and W. Baier, 2000: Agrometeorology and sustainable agriculture, *Agric. Forest Met*, **103**, 11-26.

Stokes, M. E., C. S. David, and G. G. Koch, 2000: *Categorical data analysis using the SAS® System*. 2nd ed. SAS Institute Inc., p. 165.

Stone, R. C., 1992: SOI phase relationship with rainfall probabilities over Australia. Unpublished PhD Thesis, University of Queensland, Brisbane, Australia.

Stone, R. C., G. L. Hammer, T. Marcussen, 1996: Prediction of global rainfall probabilities using phases of the Southern Oscillation Index. *Nature*, **384**, 252-55.

1   Stone, R. C., I. Smith, and P. McIntosh, 2000: Statistical methods for deriving

2        seasonal climate forecasts from GCM´s. In: *Applications of seasonal climate*

3        *forecasting in agricultural and natural ecosystems.* Hammer, G.L., N. Nichols, C.

4        Mitchell, Eds., Kluwer Academic Publishers, 135-147.

5   Tweedie, M. C. K., 1984: An index which distinguishes between some important

6        exponential families. In 'Statistics Applications and New Directions', *Proc. of the*

7        *Indian Statistical Institute Golden Jubilee International Conference,* Indian

8        Statistical Institute, Calcutta, J. K. Ghosh and J. Roy (eds), 579-604.

9   Von Storch, H. and F.W. Zwiers, 1999: *Statistical Analysis in Climate Research,*

10       Cambridge University Press, 484 pp.

11   Wilks, Daniel S., 1995: Statistical Methods in the Atmospheric Sciences, Academic

12       Press, 467 pp.

13   WMO, 2005: Proceedings of the meeting of the CCl OPAG3 expert team on

14       verification. [Available on line at:

15       http://www.wmo.ch/web/wcp/clips2001/html/Report_ETVerification_080205_2.

16       pdf]

17   Zhang, H. and T. Casey, 2000: Verification of Categorical Forecasts. *Weather and*

18       *Forecasting,* **15**, 80-89.

1    **Figure Captions**

2    Figure 1: Time series and probability of exceedance plots for JAS rainfall by

3    May/June SOI phase at (a) Echuca, Victoria (36.17°S, 144.76°E), and (b) Bangalore,

4    India (13.00°N,77.60°E) and for NDJ rainfall by September/October SOI phase at (c)

5    Bloemfontein, Republic of South Africa (29.10°S,26.30°E).

6    Figure 2: Empirical null distribution for the three-category LEPS skill scores (top row)

7    and RPSS (bottom row) arising from the SOI forecast system for predicting (from left

8    to right) JAS rainfall at Echuca and Bangalore, and NDJ rainfall at Bloemfontein.  The

9    LEPS skill scores are defined on the range −1 to +1 and the RPSS are −Inf to +1.

10    Dark, thick lines indicate the location of the observed skill score values. The area to

11    the right of the dark, thick lines correspond to the respective p-value.
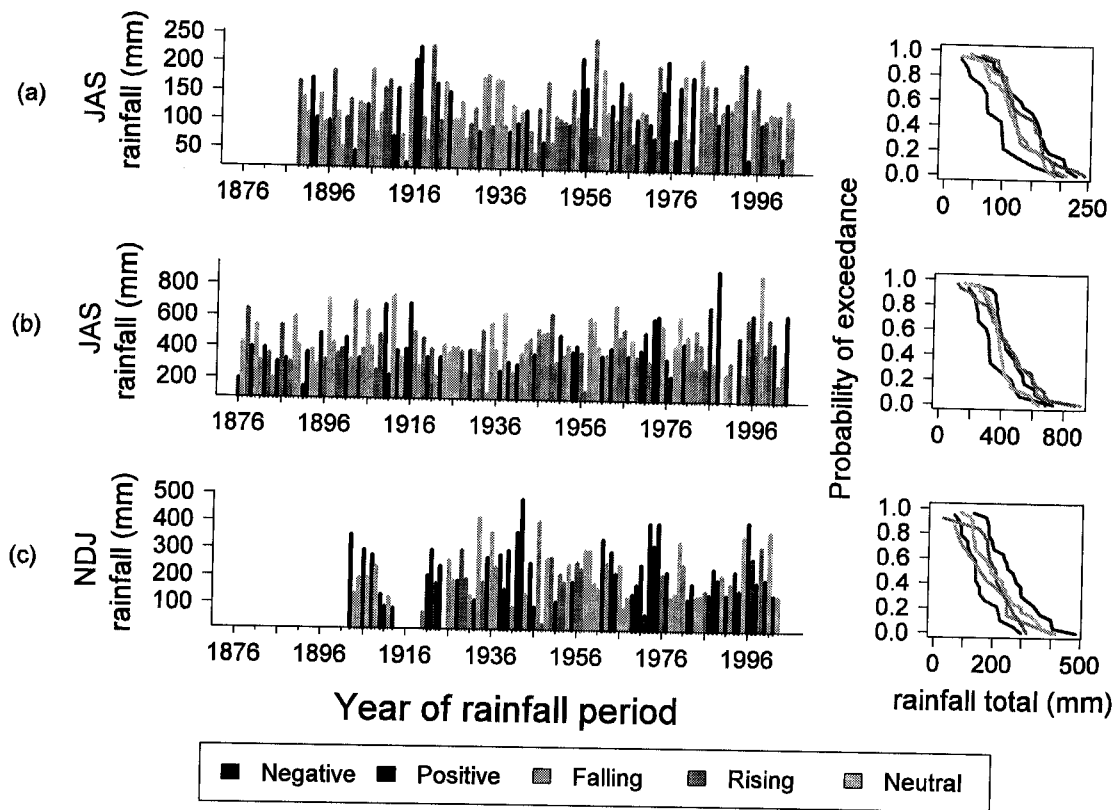
12    Figure 3: Annual patterns of discriminatory ability and skill for Echuca, Bangalore and

13    Bloemfontein based on p-values derived from the Kruskal-Wallis (KW) and

14    Kolmogorov-Smirnov (multi-sample KS), cross-validated RPSS and cross-validated

15    LEPS skill scores.

16    Figure 4: Discriminatory ability (DA, top row) and skill (bottom row) of the SOI-5 FS

17    based on p-values derived from the Kruskal-Wallis test and cross-validated LEPS skill

18    scores for July-September (JJA) rainfall in Australia and India, and November-

19    January (NDJ) rainfall in the Republic of South Africa. For computational reasons,

20    grids that have a large proportion (>33%) of dry seasons are removed from the

21    LEPS analysis and therefore appear as white grids in the LEPS maps.

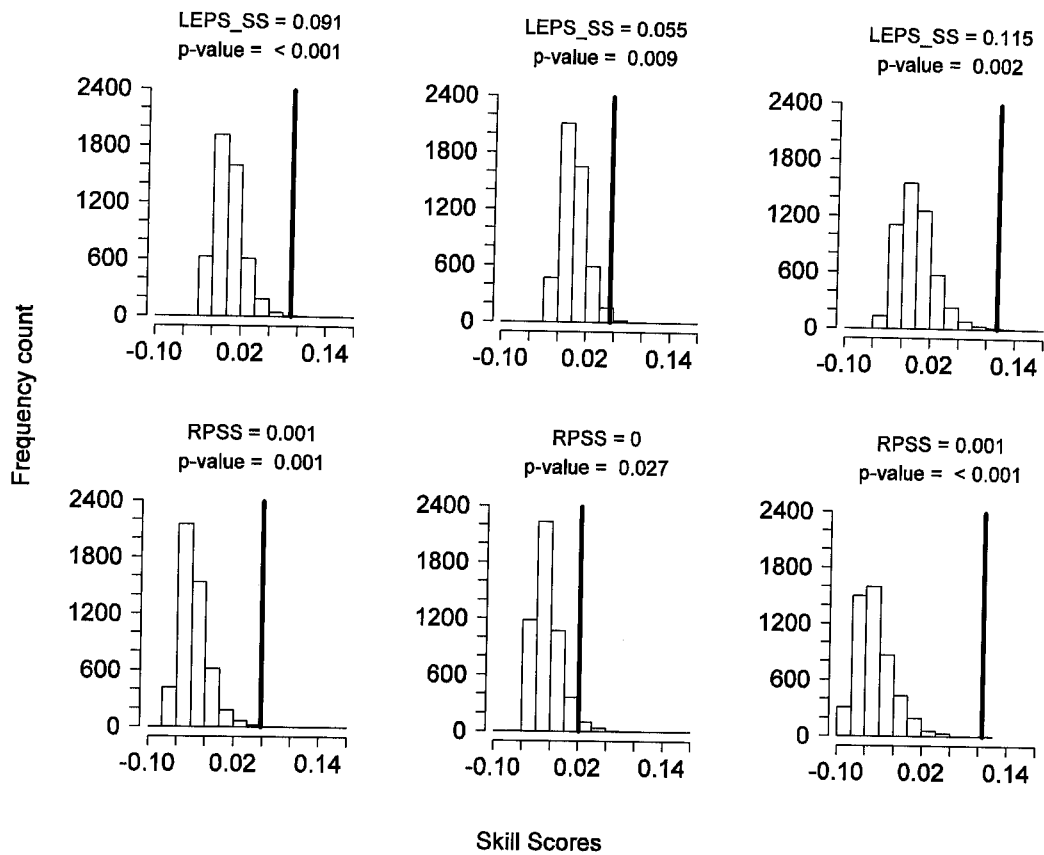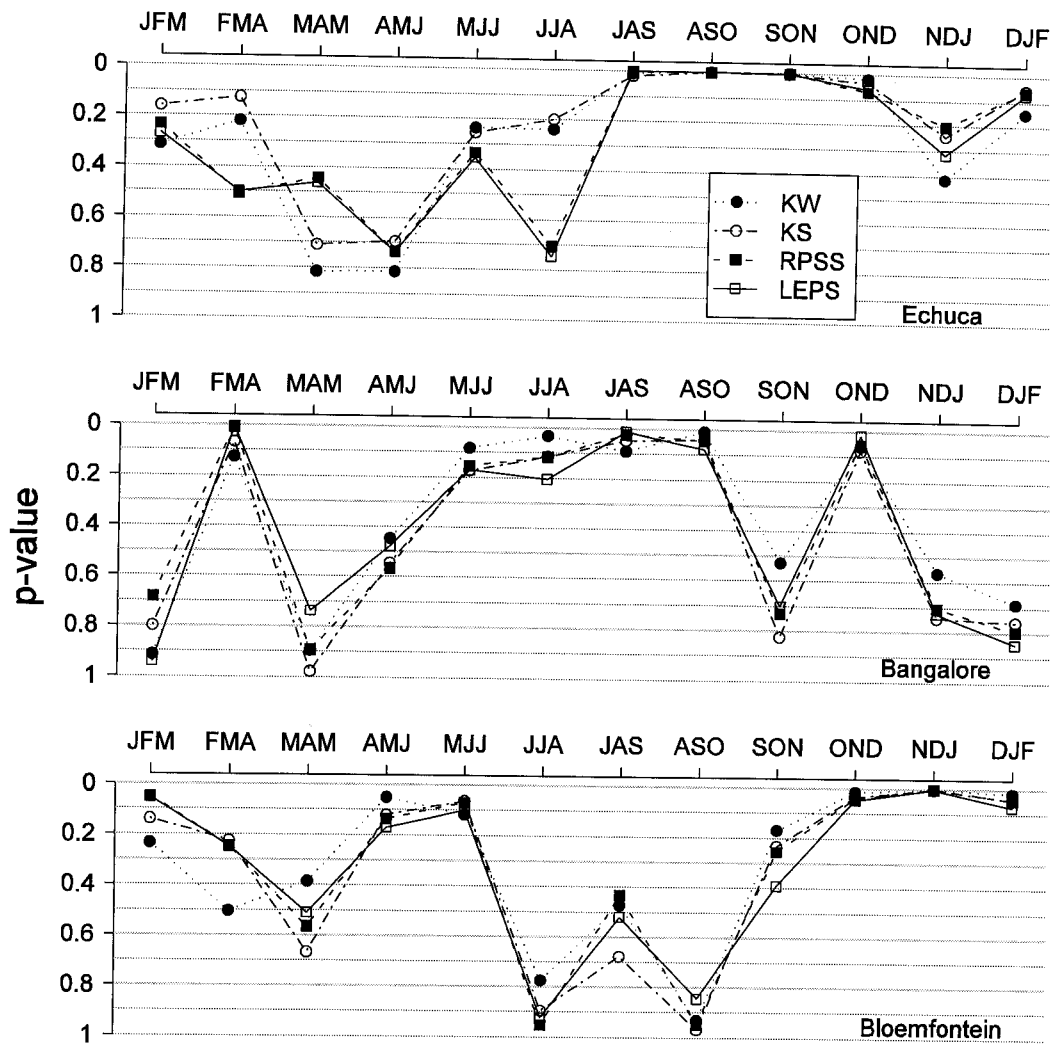22

23

24

Figure 1: Time series and probability of exceedance plots for JAS rainfall by May/June SOI phase at (a) Echuca, Victoria (36.17°S, 144.76°E), and (b) Bangalore, India (13.00°N,77.60°E) and for NDJ rainfall by September/October SOI phase at (c) Bloemfontein, Republic of South Africa (29.10°S,26.30°E).
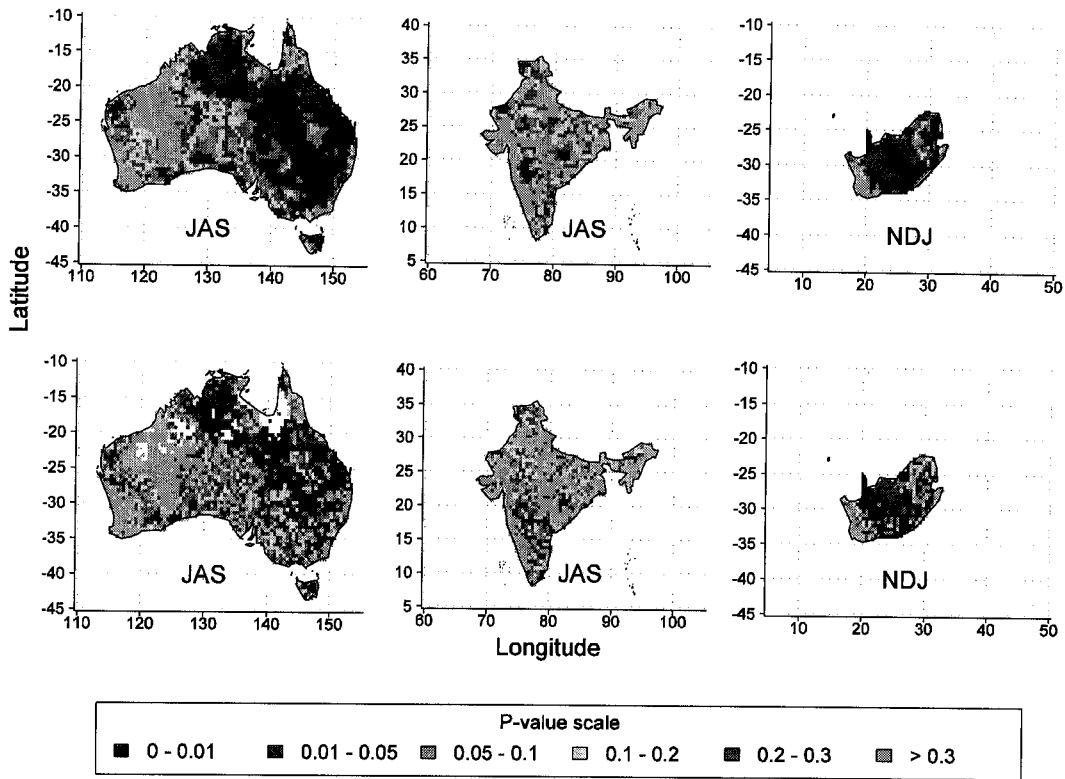
Figure 2: Empirical null distribution for the three-category LEPS skill scores (top row) and RPSS (bottom row) arising from the SOI forecast system for predicting (from left to right) JAS rainfall at Echuca and Bangalore, and NDJ rainfall at Bloemfontein. The LEPS skill scores are defined on the range −1 to +1 and the RPSS are −Inf to +1. Dark, thick lines indicate the location of the observed skill score values. The area to the right of the dark, thick lines correspond to the respective p-value.

Figure 3: Annual patterns of discriminatory ability and skill for Echuca, Bangalore and Bloemfontein based on p-values derived from the Kruskal-Wallis (KW) and Kolmogorov-Smirnov (multi-sample KS), cross-validated RPSS and cross-validated LEPS skill scores.

Figure 4: Discriminatory ability (DA, top row) and skill (bottom row) of the SOI-5 FS based on p-values derived from the Kruskal-Wallis test and cross-validated LEPS skill scores for July-September (JAS) rainfall in Australia and India, and November-January (NDJ) rainfall in the Republic of South Africa. For computational reasons, grids that have a large proportion (>33%) of dry seasons are removed from the LEPS analysis and therefore appear as white grids in the LEPS maps.