

## STANDARDIZED VERIFICATION SYSTEM (SVS) FOR LONG-RANGE FORECASTS (LRF)

### EXECUTIVE SUMMARY

#### 1. FORMULATION

The SVS is formulated in four parts:

1.1 **Diagnostics.** Information required incorporates derived diagnostic measures and contingency tables. Estimates of the statistical significance of the scores achieved are also required. Additional diagnostic measures are suggested but are not incorporated into the core SVS as yet. Use of the additional diagnostics is optional.

1.2 **Parameters.** Key variables and regions are proposed. However producers are not limited to these key parameters, and can thus all contribute regardless of the structure of individual forecast systems. The parameters to be verified are defined on three levels:

Level 1: Diagnostic measures aggregated over regions and for indices,

Level 2: Diagnostic measures evaluated at individual grid points,

Level 3: Contingency tables provided for individual grid points.

The SVS makes provision for staggered implementation of the three levels of information and the inclusion of estimates of skill significance over a two-year period.

1.3 **Verification data sets.** Key data sets of observations against which forecasts may be verified are proposed.

1.4 **System details.** Details of forecast systems employed.

1.5 **Exchange of verification information.** The SVSLRF verification results are made available through a website maintained by the Lead Centre. The functions of the Lead Centre for SVSLRF include creating and maintaining coordinated websites for the LRF verification information so that potential users would benefit from a consistent presentation of the results. The website address is <http://www.bom.gov.au/wmo/lrfvs/>.

#### 2. DIAGNOSTICS

Three diagnostic measures are incorporated in the core SVS: relative operating characteristics (ROC), reliability diagrams and accompanying measure of sharpness, and mean square skill scores (MSSS) with associated decomposition. Estimates of statistical significance in the diagnostic scores are also included in the core SVS. The three diagnostics permit direct intercomparison of results across different predicted variables, geographical regions, forecast ranges, etc. They may be applied in verification of most forecasts and it is proposed that, except where inappropriate, all three diagnostics be used on all occasions. Tabulated information at grid-point resolution is also part of the core SVS. The tabulated information will allow reconstruction of scores for user-defined areas and calculation of other diagnostic measures such as economic value.

2.1 **ROC.** To be used for verification of probability forecasts. For Level 1 information (measures aggregated over regions), the ROC curve and the standardized area under the curve (such that perfect forecasts give an area of 1 and a curve lying along the diagonal gives 0.5) should be provided. For Level 2 information (gridded values), the standardized area under the ROC curve should be provided.

2.2 **Reliability diagrams and frequency histograms.** To be used in the assessment of probability forecasts. They are required as part of Level 1 information only.

2.3 **MSSS and decomposition.** To be used in verification of deterministic forecasts. For Level 1, an overall bulk MSSS value is required and will provide a comparison of forecast performance relative to "forecasts" of climatology. The three terms of the MSSS decomposition provide valuable information on phase errors (through forecast/observation correlation), amplitude errors (through the ratio of the forecast to observed variances) and overall bias. For Level 2, quantities pertaining to the three decomposition terms should be provided. Additional terms relating to MSSS are required as part of Level 3 information.

2.4 **Contingency tables.** In addition to the derived diagnostic measures, contingency table information provided at grid points for both probability and categorical deterministic forecasts form part of the core SVS. This information constitutes Level 3 of the exchange and will allow RCCs and NMHSs (and in some cases end-users) to derive ROC, reliability, other probability-based diagnostics and scores for categorical deterministic forecasts for user-defined geographical areas.

A number of recommended contingency table-based diagnostics are listed. The Hanssen-Kuipers score is the deterministic equivalent to the area under the ROC curve, and thus provides a useful measure for comparing probabilistic and deterministic skill. The Gerrity score is a recommended score for overall assessment of forecasts using two or more categories.

### 3. PARAMETERS

The list of key parameters in the core SVS is provided below. Any verification for these key parameters should be assessed using the core SVS techniques wherever possible. Many long-range forecasts are produced which do not include parameters in the key list (for example, there are numerous empirical systems that predict seasonal rainfall over part or all of a country). The core SVS diagnostics should be used to assess these forecasts also, but full details of the predictions will need to be provided.

Forecasts can be made using different levels of post-processing, typically no-post-processing (raw or uncalibrated), simple correction of systematic errors (calibrated, i.e. calibration of mean and of variance) and more complex correction using hindcast skill (recalibrated, e.g. Model Output Statistics or perfect programme approaches). Most centres are currently issuing forecasts resulting from a simple calibration. Therefore, for the sake of comparison on the Lead Centre website scores for forecasts that were raw or calibrated (as specified in the respective skill score section) are to be submitted. It is preferable to exclude forecasts that were recalibrated, but GPCs are encouraged to apply the SVSLRF methodology and to display the results on their recalibrated forecasts on their website.

#### 3.1 Level 1: Diagrams and scores to be produced for regions

Diagrams (e.g. ROC and reliability curves) are to be supplied in digital format as specified on the Lead Centre website.

##### 3.1.1 Atmospheric parameters. Predictions for:

T2m screen temperature anomalies within standard regions:

- Tropics 20°N to 20°S;
- Northern extratropics  $\geq 20^\circ\text{N}$ ;
- Southern extratropics  $\leq 20^\circ\text{S}$ .

Precipitation anomalies within standard regions:

- Tropics 20°N to 20°S;
- Northern extratropics  $\geq 20^\circ\text{N}$ ;
- Southern extratropics  $\leq 20^\circ\text{S}$ .

##### 3.1.2 Scores and diagrams to be produced for probabilistic forecasts:

- Reliability diagram and frequency histograms;
- ROC curve and standardized area under the curve;
- Estimations of error (significance) in the scores;

The above scores and diagrams to be produced for equiprobable tercile categories.

##### 3.1.3 Score to be used for deterministic forecasts:

MSSS with climatology as standard reference forecast.

##### 3.1.4 Stratification by season

Four conventional seasons: March–April–May (MAM), June–July–August (JJA), September–October–November (SON), December–January–February (DJF).

##### 3.1.5 Lead time

Preferred minimum: two lead times, one preferably to be two weeks or more but neither greater than four months.

#### 3.2 Level 2: Grid-point data for mapping

##### 3.2.1 Grid-point verification data to be produced for each of the following variables (verification should be provided on a $2.5^\circ \times 2.5^\circ$ grid):

- T2m;
- Precipitation;
- Sea-surface temperature (SST).

##### 3.2.2 Verification parameters to be produced for deterministic forecasts

The necessary parameters for reconstructing the MSSS decomposition, the number of forecast/observation pairs, the mean square error (MSE) of the forecasts and of climatology and the MSSS are all part of the core SVS. Significance estimates for the correlation, variance, bias, MSE and MSSS terms should also be supplied.

##### 3.2.3 Verification to be provided for probability forecasts

ROC area for three tercile categories, as well as significance of the ROC scores.

##### 3.2.4 Stratification by season

If available, 12 rolling three-month periods (e.g. MAM, AMJ, MJJ). Otherwise, four conventional seasons (MAM, JJA, SON, DJF).

### 3.2.5 Lead time

Preferred minimum: two lead times, one preferably two weeks or more, but neither greater than four months.

### 3.2.6 Stratification according to the state of ENSO

Stratification by the state of ENSO should be provided if sufficient ENSO events are contained within the hindcast period used. Scores should be provided for each of three categories:

- (a) All hindcast seasons;
- (b) Seasons with *El Niño* active;
- (c) Seasons with *La Niña* active.

## 3.3 Level 3: Tabulated information to be exchanged

Tabular information to be provided for grid points of a  $2.5 \times 2.5$  grid.

### 3.3.1 Contingency tables

Contingency tables to be produced for verifying forecasts of tercile categories in each of the following variables:

- T2m;
- Precipitation;
- SST.

### 3.3.2 Tables to be produced for probabilistic forecast verification

The number of forecast hits and false alarms to be recorded against each ensemble member or probability bin for each of three equiprobable categories (terciles). It is recommended that the number of bins remain between 10 and 20. The forecast providers can bin according to percentage probability or by individual ensemble members as necessary. No latitude weighting of the numbers of hits and false alarms is to be applied in the contingency tables.

The user is encouraged to aggregate the tables over grid points for the region of interest and to apply methods of assessing statistical significance of the aggregated tables.

### 3.3.3 Tables to be produced for deterministic forecasts

$3 \times 3$  contingency tables comparing the forecast tercile with the observed tercile, over the hindcast period.

### 3.3.4 Stratification by season

If available, 12 rolling three-month periods (e.g. MAM, AMJ, MJJ). Otherwise four conventional seasons (MAM, JJA, SON and DJF).

### 3.3.5 Lead time

Preferred minimum: two lead times, one preferably two weeks or more, but neither greater than four months.

### 3.3.6 Stratification according to the state of ENSO

Stratification by the state of ENSO should be provided if sufficient ENSO events are contained within the hindcast period used. Scores should be provided for each of three categories:

- (a) All hindcast seasons;
- (b) Seasons with *El Niño* active;
- (c) Seasons with *La Niña* active.

## 3.4 Verification for indices (Level 1)

### 3.4.1 Indices to be verified

Niño3.4 region SST anomalies. Other indices may be added in due course.

### 3.4.2 Scores to be calculated for probabilistic forecasts

ROC area for 3 tercile categories. Where dynamical forecast models are used, the ROC scores should be calculated for the grid-point averaged SST anomaly over the Niño3.4 region. It is recommended that significance of the ROC scores should also be calculated.

### 3.4.3 Scores to be calculated for deterministic forecasts

The three terms of the Murphy decomposition of MSSS, produced with climatology as standard reference forecast. As a second (optional) control, it is recommended that damped persistence be used. Significance estimates should accompany each of the three terms.

Where dynamical models are used, the MSSS decomposition should be calculated for the grid-point averaged Niño3.4 anomaly.

#### 3.4.4 Stratification by month

Verification should be provided for each calendar month.

#### 3.4.5 Lead time

Verification for each month should be provided for six lead times. Namely zero-lead and 1-month, 2-month, 3-month, 4-month and 5-month leads. Additional lead times are encouraged if available.

### 4. STAGGERED IMPLEMENTATION

In order to ease implementation, producers may stagger the provision of the elements of the core SVS according to the following recommendation.

- (a) Verification at Levels 1 and 2 in the first year of implementation,
- (b) Verification at Level 3 by the middle of the year following implementation of Levels 1 and 2,
- (c) Level of significance by the end of the year following implementation of Levels 1 and 2.

\* \* \*

### 1. INTRODUCTION

The following sections present detailed specifications for the development of an SVS for LRF within the framework of a WMO exchange of verification scores. The SVS for LRF described herein constitutes the basis for LRF evaluation and validation, and for exchange of verification scores. It will grow as more requirements are adopted.

### 2. DEFINITIONS

#### 2.1 LRF

LRF extend from 30 days up to two years and are defined in Table 1.

Table 1

Monthly outlook	Description of averaged weather parameters expressed as a departure from climate values for that month
Three-month or 90-day 'rolling season' outlook	Description of averaged weather parameters expressed as a departure from climate values for that three-month or 90-day period
Seasonal outlook	Description of averaged weather parameters expressed as a departure from climate values for that season

#### Definition of LRF

Seasons have been loosely defined in the northern hemisphere as December-January-February (DJF) for winter (summer in the southern hemisphere), March-April-May (MAM) for spring (autumn in the southern hemisphere), June-July-August (JJA) for summer (winter in the southern hemisphere) and September-October-November (SON) for autumn (spring in the southern hemisphere). Twelve rolling seasons are also defined e.g. MAM, AMJ, MJJ. In tropical areas, the seasons may have different definitions. Outlooks over longer periods such as multi-seasonal outlooks or tropical rainy season outlooks may be provided.

It is recognized that in some countries LRF are considered to be climate products.

This attachment is mostly concerned with three-month or 90-day outlooks and seasonal outlooks.

#### 2.2 Deterministic LRF

Deterministic LRF provide a single expected value for the forecast variable. The forecast may be presented in terms of an expected category (referred to as categorical forecasts, e.g. equiprobable terciles) or may take predictions of the continuous variable (non-categorical forecasts). Deterministic LRF can be produced from a single run of a numerical weather prediction (NWP) model or a general circulation model (GCM), or from the grand mean of the members of an ensemble prediction system (EPS), or can be based on an empirical model.

The forecasts are either objective numerical values such as departure from normal of a given parameter or expected occurrences (or non-occurrences) of events classified into categories (above/below normal or above/near/below normal for example). Although equiprobable categories are preferred for consistency, other classifications can be used in similar fashion.

### 2.3 Probabilistic LRF

Probabilistic LRFs provide probabilities of occurrences or non-occurrences of an event or a set of fully inclusive events. Probabilistic LRFs can be generated from an empirical model, or produced from an EPS.

The events can be classified into categories (above/below normal or above/near/below normal for example). Although equiprobable categories are preferred for consistency, other classifications can be used in similar fashion.

### 2.4 Terminology

There is no universally accepted definition of forecast period and forecast lead time. However, the definition in Table 2 will be used here.

**Table 2**  
**Definitions of forecast period and lead time**

Forecast period	Forecast period is the validity period of a forecast. For example, LRF may be valid for a 90-day period or a season.
Lead time	Lead time refers to the period of time between the issue time of the forecast and the beginning of the forecast validity period. LRFs based on all data up to the beginning of the forecast validity period are said to be of lead zero. The period of time between the issue time and the beginning of the validity period will categorize the lead. For example, a winter seasonal forecast issued at the end of the preceding summer season is said to be of one-season lead. A seasonal forecast issued one month before the beginning of the validity period is said to be of one-month lead.

Figure 1 presents the definitions of Table 2 in graphical format.

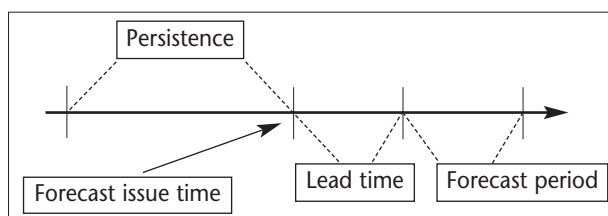


Figure 1 — Definition of forecast period, lead time and persistence as applied in a forecast verification framework

The forecast range determines how far into the future LRFs are provided; it is thus the summation of lead time and forecast period.

Persistence, for a given parameter, signifies a persisting anomaly which has been observed over the same length of time as the forecast period and immediately prior to the LRF issue time (see Figure 1). It is important to realize that only the anomaly of any given parameter can be considered in persistence. The persisting anomaly is added to the background climatology to retrieve the parameter in question. Climatology is equivalent to a uniform persisting anomaly of zero.

### 3. SVS FOR LRF

Forecasts can be made using different levels of post-processing, typically no-post-processing (raw or uncalibrated), simple correction of systematic errors (calibrated, i.e. calibration of mean and of variance) and more complex correction using hindcast skill (recalibrated, e.g. Model Output Statistics or perfect programme approaches). Most centres are currently issuing forecasts resulting from a simple calibration. Therefore, for the sake of comparison on the Lead Centre website scores for forecasts that were raw or calibrated (as specified in the respective skill score section) are to be submitted. It is preferable to exclude forecasts that were recalibrated, but GPCs are encouraged to apply the SVSLRF methodology and to display the results on their recalibrated forecasts on their website.

#### 3.1 Parameters to be verified

The following parameters are to be verified:

- (a) Surface air temperature (T2m) anomaly at screen level;
- (b) Precipitation anomaly;
- (c) SST anomaly.

In addition to these three parameters, the Niño3.4 index, defined as the mean SST anomaly over the Niño3.4 region from 170°W to 120°W and from 5°S to 5°N, inclusive, is also to be verified.

It is recommended that three levels of verification be done:

- (a) Level 1: large-scale aggregated overall measures of forecast performance (see section 3.1.1).
- (b) Level 2: verification at grid points (see section 3.1.2).
- (c) Level 3: grid point by grid point contingency tables for more extensive verification (see section 3.1.3).

Both deterministic and probabilistic forecasts are verified if available. Level 1 is applicable to the T2m anomaly, precipitation anomaly and Niño3.4 index. Levels 2 and 3 are applicable to the T2m anomaly, precipitation anomaly and SST anomaly.

### 3.1.1 Aggregated verification (Level 1)

Large-scale verification statistics are required in order to evaluate the overall skill of the models and ultimately to assess their improvements. These are bulk numbers calculated by aggregating verifications at all grid points within large regions; they will not necessarily reflect skill for any subregion. This aggregated verification is performed over three regions:

- (a) Tropics: from 20°S to 20°N, inclusive;
- (b) Northern extratropics: from 20°N to 90°N, inclusive;
- (c) Southern extratropics: from 20°S to 90°S, inclusive.

Verification of the Niño3.4 index is also part of Level 1 verification.

### 3.1.2 Grid-point verification (Level 2)

Grid-point verification is recommended for a regionalized assessment of the skill of the model, for which a  $2.5^\circ \times 2.5^\circ$  latitude/longitude grid is recommended, with origin at 0°N, 0°E. Verification should be supplied to the Lead Centre for visual rendering. The formats for supplying derived verification are specified on the Lead Centre's website.

### 3.1.3 Contingency tables (Level 3)

**Table 3**  
**Summary of the core SVS**

<i>Parameters</i>	<i>Verification regions</i>	<i>Deterministic forecasts</i>	<i>Probabilistic forecasts</i>
<b>Level 1</b>			
T2m anomaly Precipitation anomaly	Tropics Northern extratropics Southern extratropics  (section 3.1.1)	MSSS (bulk number)  (section 3.3.1)	ROC curves ROC areas Reliability diagrams Frequency histograms (sections 3.3.3 and 3.3.4)
Niño3.4 index	N/A	MSSS (bulk number)  (section 3.3.1)	ROC curves ROC areas Reliability diagrams Frequency histograms (sections 3.3.3 and 3.3.4)
<b>Level 2</b>			
T2m anomaly Precipitation anomaly SST anomaly	Grid-point verification on a $2.5^\circ \times 2.5^\circ$ grid  (section 3.1.2)	MSSS and its three-term decomposition at each grid-point  (section 3.3.1)	ROC areas at each grid point  (section 3.3.3)
<b>Level 3</b>			
T2m anomaly Precipitation anomaly SST anomaly	Grid-point verification on a $2.5^\circ \times 2.5^\circ$ grid  (section 3.1.2)	$3 \times 3$ contingency tables at each grid-point  (section 3.3.2)	ROC reliability tables at at each grid-point  (section 3.3.3)

Contingency tables allow users to perform more detailed verifications and generate statistics that are relevant for particular regions. The content and structure of the contingency tables is defined in sections 3.3.2 and 3.3.3. Data formats for supplying the contingency tables are specified on the Lead Centre's website.

### 3.1.4 Summary of the core SVS

The following gives a summary of parameters, validation regions and diagnostics that form the core SVS. The required periods, lead times and stratification against the state of ENSO are given in section 3.2.

The number of LRF realizations is far smaller than in the case of short-term numerical weather prediction forecasts. Consequently it is essential as part of the core SVS to calculate and report error bars and the level of significance (see section 3.3.5).

In order to ease implementation, participating LRF producers may stage the introduction of the core SVS by prioritizing implementation of verification at levels 1 and 2.

Other parameters and indices to be verified as well as other verification scores can be added to the core SVS in future versions.

### 3.2 Verification strategy

LRF verification should be done on a latitude/longitude grid, with areas as defined in section 3.1.1. Verification can also be done at individual stations or groups of stations. Verification on a latitude-longitude grid is performed separately from that done at stations.

A  $2.5^\circ \times 2.5^\circ$  verification latitude-longitude grid is recommended, with origin at  $0^\circ\text{N}$ ,  $0^\circ\text{E}$ . Both forecasts and the gridded verifying data sets are to be interpolated onto the same  $2.5^\circ \times 2.5^\circ$  grid.

In order to handle spatial forecasts, predictions for each point within the verification grid should be treated as individual forecasts but with all results combined into the final outcome. The same approach is applied when verification is done at stations. Categorical forecast verification can be performed for each category separately.

Similarly, all forecasts are treated as independent and combined together into the final outcome, when verification is done over a long period of time (several years for example).

Stratification of the verification data is based on the forecast period, lead time and verification area. Stratification by forecast period should, for T2m and precipitation, be by four conventional seasons for Level 1. For Levels 2 and 3 stratification should be on 12 rolling seasons (see section 2.1) if available, otherwise four conventional seasons should be used. Verification results for different seasons should not be mixed. Stratification by lead time should include a minimum of two lead times, with lead time not greater than four months. Forecasts with different lead times are similarly to be verified separately. Stratification according to the state of ENSO (where there are sufficient cases) should be as follows:

- (a) All hindcast seasons;
- (b) Seasons with El Niño active;
- (c) Seasons with La Niña active.

For Niño3.4, SST anomaly verification should be stratified according to each calendar month and lead time. Six lead times should be provided, ranging from zero to a 5-month lead.

### 3.3 Verification scores

The verification scores to be used are: MSSS and ROC.

MSSS is applicable to deterministic forecasts only, while ROC is applicable to both deterministic and probabilistic forecasts. MSSS is applicable to non-categorical forecasts (or to forecasts of continuous variables), while ROC is applicable to categorical forecasts whether deterministic or probabilistic in nature.

The verification methodology using ROC is derived from signal detection theory and is intended to provide information on the characteristics of the systems upon which management decisions can be taken. In the case of weather/climate forecasts, the decisions may relate to the most appropriate way of using a forecast system for a given purpose. ROC is applicable to both deterministic and probabilistic categorical forecasts and is useful in contrasting the characteristics of deterministic and probabilistic systems. The ROC derivation is based on contingency tables giving the hit rate and false alarm rate for deterministic or probabilistic forecasts. The events are defined as binary, which means that only two outcomes are possible, occurrence or non-occurrence. It is recognized that ROC as applied to deterministic forecasts is equivalent to the Hanssen and Kuipers score (see section 3.3.2).

The binary event can be defined as the occurrence of one of two possible categories when the outcome of the LRF system is in two categories. When the outcome of the LRF system is in three (or more) categories, the binary event is defined in terms of occurrences of one category against the remaining ones, hence the ROC has to be calculated for each possible category.

#### 3.3.1 MSSS for non-categorical deterministic forecasts

Let  $x_{ij}$  and  $f_{ij}$  ( $i=1, \dots, n$ ) denote time series of observations and continuous deterministic forecasts, respectively, for a grid point or station  $j$  over the period of verification (POV). Then, their averages for the POV,  $\bar{x}_j$  and  $\bar{f}_j$  and their sample variances  $s_{xj}^2$  and  $s_{fj}^2$  are given by:

$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}, \quad \bar{f}_j = \frac{1}{n} \sum_{i=1}^n f_{ij}$$

$$s_{xj}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2, \quad s_{fj}^2 = \frac{1}{n-1} \sum_{i=1}^n (f_{ij} - \bar{f}_j)^2$$

The MSE of the forecasts is:

$$MSE_j = \frac{1}{n} \sum_{i=1}^n (f_{ij} - x_{ij})^2$$

In the case of cross-validated (see section 3.4) POV climatology forecasts where forecast/observation pairs are reasonably temporally independent of each other (so that only one year at a time is withheld), the MSE of climatology forecasts (Murphy, 1988) is:

$$MSE_{cj} = \frac{n-1}{n} s_{xj}^2$$

The MSSS for  $j$  is defined as one minus the ratio of the MSE of the forecasts to the MSE for forecasts of climatology:

$$MSSS_j = 1 - \frac{MSE_j}{MSE_{cj}}$$

For the three domains described in section 3.1.1, it is recommended that an overall MSSS be provided, which is computed as:

$$MSSS = 1 - \frac{\sum_j w_j MSE_j}{\sum_j w_j MSE_{cj}}$$

where  $w_j$  is unity for verifications at stations and is equal to  $\cos(\theta_j)$ , where  $\theta_j$  is the latitude at grid point  $j$ , on latitude-longitude grids.

For either  $MSSS_j$  or  $MSSS$ , a corresponding root mean square skill score (RMSSS) can be obtained easily from:

$$RMSSS = 1 - (1 - MSSS)^{1/2}$$

$MSSS_j$  for fully cross-validated forecasts (with one year at a time withheld) can be expanded (Murphy, 1988) as:

$$MSSS_j = \left\{ 2 \frac{s_{fj}}{s_{xj}} r_{fxj} - \left( \frac{s_{fj}}{s_{xj}} \right)^2 - \left( \frac{[\bar{f}_j - \bar{x}_j]}{s_{xj}} \right)^2 + \frac{2n-1}{(n-1)^2} \right\} / \left\{ 1 + \frac{2n-1}{(n-1)^2} \right\}$$

where  $r_{fxj}$  is the product moment correlation of the forecasts and observations at point or station  $j$ .

$$r_{fxj} = \frac{\frac{1}{n} \sum_{i=1}^n (f_{ij} - \bar{f}_j)(x_{ij} - \bar{x}_j)}{s_{fj} s_{xj}}$$

The first three terms of the decomposition of  $MSSS_j$  are related to phase errors (through the correlation), amplitude errors (through the ratio of the forecast to observed variances) and overall bias error, respectively, of the forecasts. These terms provide the opportunity for those wishing to use the forecasts for input into regional and local forecasts to adjust or weight the forecasts as appropriate. The last term takes into account the fact that the 'climatology' forecasts are cross-validated as well.

Note that for forecasts with the same amplitude as that of observations (second term unity) and no overall bias (third term zero),  $MSSS_j$  will not exceed zero (i.e. the forecasts' MSE will not be less than for 'climatology') unless  $r_{fxj}$  exceeds approximately 0.5.

The core SVS for LRF requires grid-point values of the correlation, the ratio of the square roots of the variances, and the overall bias, i.e.

$$r_{fxj}, \frac{s_{fj}}{s_{xj}}, [\bar{f}_j - \bar{x}_j]$$

In addition, it is recommended that grid-point ( $j$ ) values of the following quantities be provided:

$$n, \bar{f}_j, \bar{x}_j, s_{fj}, s_{xj}, r_{fxj}, MSE_j, MSE_{cj}, MSSS_j$$

As an additional standard against which to measure forecast performance, cross-validated damped persistence (defined below) should be considered for certain forecast sets. A forecast of ordinary persistence, for a given parameter and target period, signifies a persisting anomaly (departure from cross-validated climatology) from a period immediately preceding the start of the lead time for the forecast period (see Figure 1). This period must have the same length as the forecast period. For example, the ordinary persistence forecast for a 90-day period made 15 days in advance would be the anomaly of the 90-day period beginning 105 days before the target forecast period and ending 16 days before. Ordinary persistence forecasts are never recommended as a standard against which to measure other forecasts if the performance or skill measures are based on MSE as here, since persistence is easily surpassed in this framework.



Damped persistence is the optimal persistence forecast in a least squares error sense. Even damped persistence should not be used in the case of extratropical seasonal forecasts, because the nature of the interannual variability of seasonal means changes considerably from one season to the next in the extratropics. In all other cases, damped persistence forecasts can be made in cross-validated mode (see section 3.4) and the skill and performance diagnostics based on MSE as described above (bulk measures, grid-point values and tables) can be computed and presented for these forecasts.

Damped persistence  $r_{\Delta,j}^m [x_{ij}(t-\Delta t) - \bar{x}_{ij}^m(t-\Delta t)]$  is the ordinary persistence anomaly  $x_{ij}(t-\Delta t) - \bar{x}_{ij}^m(t-\Delta t)$  damped (multiplied) towards climatology by the cross-validated, lagged product moment correlation between the period of persistence and the target forecast period. Thus:

$$r_{\Delta,j}^m = \frac{\frac{1}{m} \sum [x_{ij}(t-\Delta t) - \bar{x}_{ij}^m(t-\Delta t)] [x_{ij}(t) - \bar{x}_{ij}^m(t)]}{s_{xj}^m(t-\Delta t) s_{xj}^m(t)}$$

where  $t$  is the target forecast period,  $t-\Delta t$  the persistence period (preceding the lead time), and  $m$  denotes summation (for  $r_{\Delta,j}^m, \bar{x}_{ij}^m, s_{xj}^m$ ) at each stage of the cross-validation over all  $i$  except those being currently withheld (section 3.4).

⇒ MSSS, provided as a single bulk number, is mandatory for Level 1 verification in the core SVS. MSSS, together with its three-term decomposition, are also mandatory for Level 2 verification in the core SVS. For the exchange of scores via the Lead Centre website the MSSS and its decomposition term should be calculated using the raw forecasts and preferably not the calibrated ones.

### 3.3.2 Contingency tables and scores for categorical deterministic forecasts

For two- or three-category deterministic forecasts the core SVS for LRF includes full contingency tables, because it is recognized that they constitute the most informative way to evaluate the performance of the forecasts. These contingency tables then form the basis for several skill scores that are useful for comparisons between different deterministic categorical forecast sets (Gerrity, 1992) and between deterministic and probabilistic categorical forecast sets (Hanssen and Kuipers, 1965), respectively.

The contingency tables should be provided for every combination of parameter, lead time, target month or season, and ENSO stratification (when appropriate) at every verification point for both the forecasts and (when appropriate) damped persistence. The definition of ENSO events is provided on the Lead Centre's website. If  $x_i$  and  $f_i$  now denote an observation and corresponding forecast of category  $i$  ( $i = 1, \dots, 3$ ), let  $n_{ij}$  be the count of those instances with forecast category  $i$  and observed category  $j$ . The full contingency table is defined as the nine  $n_{ij}$ . Graphically, the nine cell counts are usually arranged with the forecasts defining the table rows and the observations the table columns:

**Table 4**  
**General three-by-three contingency table**

		Observations			
		<i>Below normal</i>	<i>Near normal</i>	<i>Above normal</i>	
Forecasts	<i>Below normal</i>	$n_{11}$	$n_{12}$	$n_{13}$	$n_{1\bullet}$
	<i>Near normal</i>	$n_{21}$	$n_{22}$	$n_{23}$	$n_{2\bullet}$
	<i>Above normal</i>	$n_{31}$	$n_{32}$	$n_{33}$	$n_{3\bullet}$
		$n_{\bullet 1}$	$n_{\bullet 2}$	$n_{\bullet 3}$	$T$

In Table 4,  $n_{i\bullet}$  and  $n_{\bullet i}$  represent the sum of the rows and columns respectively;  $T$  is the total number of cases. Generally at least 90 forecast/observation pairs are required to estimate properly a three-by-three contingency table. Thus it is recommended that the tables provided be aggregated by users over windows of target periods, such as several adjacent months or overlapping three-month periods, or over verification points. In the latter case, weights  $W_i$  should be used in summing  $n_{ij}$  over different points  $i$  (see discussion on Table 5).  $W_i$  is defined as:

$W_i = 1$  when verification is done at stations or at single grid points within a limited graphical region;

$W_i = \cos(\theta_i)$  at grid point  $i$ , when verification is done on a grid,  $\theta_i$  being the latitude at grid point  $i$ .

On a 2.5 degree latitude-longitude grid, the minimum acceptable sample is easily attained, even with a record as short as  $n = 10$ , by aggregating over all grid points within a 10-degree box. Alternatively, in this case, an adequate sample can be achieved by aggregation over three adjacent months or overlapping three-month periods and within a 5-degree box. Regardless of this, scores derived from any contingency table should be accompanied by error bars, confidence intervals or level of significance.

Contingency tables such as the one in Table 4 are mandatory for Level 3 verification in the core SVS.

Relative sample frequencies  $p_{ij}$  are defined as the ratios of the cell counts to the total number of forecast/observation pairs  $N$  ( $n$  is reserved to denote the length of the POV):

$$p_{ij} = n_{ij} / N$$

The sample probability distributions of forecasts and observations, respectively, then become:

$$p(f_i) = \sum_{j=1}^3 p_{ij} = \hat{p}_i; \quad i = 1, \dots, 3$$

$$p(x_i) = \sum_{j=1}^3 p_{ji} = p_i; \quad i = 1, \dots, 3$$

A recommended skill score for the three-by-three table which has many desirable properties and is easy to compute is the Gerrity skill score (GSS). The definition of the GSS uses a scoring matrix  $s_{ij}$  ( $i = 1, \dots, 3$ ), which is a tabulation of the reward or penalty every forecast/observation outcome represented by the contingency table will be accorded:

$$GSS = \sum_{i=1}^3 \sum_{j=1}^3 p_{ij} s_{ij}$$

The scoring matrix is given by:

$$s_{ii} = \frac{1}{2} \left( \sum_{r=1}^{i-1} a_r^{-1} + \sum_{r=i}^2 a_r \right)$$

$$s_{ij} = \frac{1}{2} \left[ \sum_{r=1}^{i-1} a_r^{-1} - (j-1) + \sum_{r=j}^2 a_r \right]; \quad 1 \leq i < 3, i < j \leq 3$$

where:  $a_i = \frac{1 - \sum_{r=1}^i p_r}{\sum_{r=1}^i p_r}$

Note that the GSS is computed using the sample probabilities, not those on which the original categorizations were based (i.e. 0.33, 0.33, 0.33).

Alternatively, the GSS can be computed by the numerical average of two of the three possible two-category, unscaled Hanssen and Kuipers scores (introduced below) that can be computed from the three-by-three table. The two are computed from the two two-category contingency tables formed by combining categories on either side of the partitions between consecutive categories: (a) above normal and a combined near and below normal category; and (b) below normal and a combined near and above normal category.

The easy construction of the GSS ensures its consistency from categorization to categorization and with underlying linear correlations. The score is also equitable, does not depend on the forecast distribution, does not reward conservatism, utilizes off-diagonal information in the contingency table, and penalizes larger errors more. For a limited subset of forecast situations, it can be manipulated by the forecaster to his/her advantage (Mason and Mimmack, 2002), but this is not a problem for objective forecast models that have not been trained to take advantage of this weakness. For all these reasons it is the recommended score.

An alternative score to the GSS to be considered is LEPCAT (Potts, et al., 1996).

Table 5 shows the general form for the three possible two-by-two contingency tables referred to above (the third is the table for the near normal category and the combined above and below normal category). In Table 5,  $T$  is the grand sum of all the proper weights applied on each occurrence and non-occurrence of the events.

**Table 5**  
**General ROC contingency table for deterministic forecasts**

The two-by-two table in Table 5 may be constructed from the three-by-three table described in Table 4 by summing the appropriate rows and columns.

		Observations		
		Occurrences	Non-occurrences	
Forecasts	Occurrences	O <sub>1</sub>	NO <sub>1</sub>	O <sub>1</sub> +NO <sub>1</sub>
	Non-occurrences	O <sub>2</sub>	NO <sub>2</sub>	O <sub>2</sub> +NO <sub>2</sub>
		O <sub>1</sub> +O <sub>2</sub>	NO <sub>1</sub> +NO <sub>2</sub>	T

In Table 5:

O<sub>1</sub> represents the correct forecasts or hits:  $O_1 = \sum W_i (OF)_i$ , (OF) being 1 when the event occurrence is observed and forecast; 0 otherwise. The summation is over all grid points or stations;

NO<sub>1</sub> represents false alarms:  $NO_1 = \sum W_i (NOF)_i$ , (NOF) being 1 when the event occurrence is not observed but was forecast; 0 otherwise. The summation is over all grid points or stations;

O<sub>2</sub> represents the misses:  $O_2 = \sum W_i (ONF)_i$ , (ONF) being 1 when the event occurrence is observed but not forecast; 0 otherwise. The summation is over all grid points or stations;

$NO_2$  represents the correct rejections:  $NO_2 = \sum W_i (NONF)_i$ , (NONF) being 1 when the event occurrence is not observed and not forecast; 0 otherwise. The summation is over all grid points or stations;

$W_i = 1$  when verification is done at stations or at single grid points;  $W_i = \cos(\theta_i)$  at grid point  $i$  when verification is done on a grid,  $\theta_i$  being the latitude at grid point  $i$ .

When verification is done at stations, the weighting factor is one. Consequently, the number of occurrences and non-occurrences of the event are entered in the contingency table as per Table 5.

However, when verification is done on a grid, the weighting factor is  $\cos(\theta_i)$ , where  $\theta_i$  is the latitude at grid point  $i$ . Consequently, each number entered in the contingency table as per Table 6, is, in fact, a summation of the weights properly assigned.

Using stratification by observations (rather than by forecast), the hit rate (HR) is defined (referring to Table 5) as:

$$HR = O_1 / (O_1 + O_2)$$

The range of values for HR goes from 0 to 1, the latter value being desirable. An HR of 1 means that all occurrences of the event were correctly forecast.

The false alarm rate (FAR) is defined as:

$$FAR = NO_1 / (NO_1 + NO_2)$$

The range of values for FAR goes from 0 to 1, the former value being desirable. A FAR of zero means that in the verification sample, no non-occurrences of the event were forecast to occur.

The Hanssen and Kuipers score (see Hanssen and Kuipers, 1965 and Stanski, *et al.*, 1989) is calculated for deterministic forecasts. The Hanssen and Kuipers score (KS) is defined as:

$$KS = HR - FAR = \frac{O_1 NO_2 - O_2 NO_1}{(O_1 + O_2)(NO_1 + NO_2)}$$

The range of values for KS goes from -1 to +1, the latter value corresponding to perfect forecasts (HR being 1 and FAR being 0). KS can be scaled so that the range of possible values goes from 0 to 1 (1 being for perfect forecasts):

$$KS_{scaled} = \frac{KS + 1}{2}$$

The advantage of scaling KS is that it becomes comparable to the area under the ROC curve for probabilistic forecasts (see section 3.3.3) where a perfect forecast system has an area of 1 and a forecast system with no information has an area of 0.5 (HR being equal to FAR).

⇒ Contingency tables for deterministic categorical forecasts (such as in Table 4) are mandatory for Level 3 verification in the core SVS. These contingency tables can provide the basis for the calculation of several scores and indices such as the GSS, the LEPCAT or the scaled Hanssen and Kuipers score and others.

### 3.3.3 ROC for probabilistic forecasts

Tables 6 and 7 show contingency tables (similar to Table 5) that can be built for probabilistic forecasts of binary events.

Table 6

General ROC contingency table for probabilistic forecasts of binary events with definitions of the different parameters. This contingency table applies when probability thresholds are used to define the different probability bins

Bin number	Forecast probabilities	Observed occurrences	Observed non-occurrences
1	0-P <sub>2</sub> (%)	O <sub>1</sub>	NO <sub>1</sub>
2	P <sub>2</sub> -P <sub>3</sub> (%)	O <sub>2</sub>	NO <sub>2</sub>
3	P <sub>3</sub> -P <sub>4</sub> (%)	O <sub>3</sub>	NO <sub>3</sub>
...	...	...	...
n	P <sub>n</sub> -P <sub>n+1</sub> (%)	O <sub>n</sub>	NO <sub>n</sub>
...	...	...	...
N	P <sub>N</sub> -100 (%)	O <sub>N</sub>	NO <sub>N</sub>

In Table 6:

$n$  = the number of the  $n^{\text{th}}$  probability interval or bin  $n$ ;  $n$  goes from 1 to  $N$ ;

$P_n$  = the lower probability limit for bin  $n$ ;

$P_{n+1}$  = the upper probability limit for bin  $n$ ;

$N$  = the number of probability intervals or bins;

$O_n = \sum W_i (O)_i$ , (O) being 1 when an event corresponding to a forecast in bin  $n$ , is observed as an occurrence; 0 otherwise. The summation is over all forecasts in bin  $n$ , at all grid points or stations;

$NO_n = \sum W_i (NO)_i$ , (NO) being 1 when an event corresponding to a forecast in bin  $n$ , is not observed; 0 otherwise. The summation is over all forecasts in bin  $n$ , at all grid points  $i$  or stations  $i$ ;

$W_i = 1$  when verification is done at stations or at single grid points;  $W_i = \cos(\theta_i)$  at grid point  $i$ , when verification is done on a grid;  $\theta_i$  being the latitude at grid point  $i$ .

Table 7

General ROC contingency table for probabilistic forecasts of binary events with definitions of the different parameters. This contingency table applies when the different probability bins are defined as a function of the number of members in the ensemble

Bin number	Member distribution	Observed occurrences	Observed non-occurrences
1	F = 0, NF = M	O <sub>1</sub>	NO <sub>1</sub>
2	F = 1, NF = M-1	O <sub>2</sub>	NO <sub>2</sub>
3	F = 2, NF = M-2	O <sub>3</sub>	NO <sub>3</sub>
...		...	...
n	F = n-1, NF = M-n+1	O <sub>n</sub>	NO <sub>n</sub>
...		...	...
N	F = M, NF = 0	O <sub>N</sub>	NO <sub>N</sub>

In Table 7:

$n$  = the number of the  $n^{\text{th}}$  bin;  $n$  goes from 1 to  $N = M + 1$ ;

$F$  = the number of members forecasting occurrence of the event;

$NF$  = the number of members forecasting non-occurrence of the event. The bins may be aggregated;

$O_n = \sum W_i (O)_i$ , (O) being 1 when an event corresponding to a forecast in bin  $n$  is observed as an occurrence; 0 otherwise. The summation is over all forecasts in bin  $n$ , at all grid-points  $i$  or stations  $i$ ;

$NO_n = \sum W_i (NO)_i$ , (NO) being 1 when an event corresponding to a forecast in bin  $n$  is not observed; 0 otherwise. The summation is over all forecasts in bin  $n$ , at all grid points  $i$  or stations  $i$ ;

$W_i = 1$  when verification is done at stations or at single grid points;  $W_i = \cos(\theta_i)$  at grid point  $i$ , when verification is done on a grid,  $\theta_i$  being the latitude at grid point  $i$ .

To build the contingency table in Table 6, probability forecasts of the binary event are grouped in categories or bins in ascending order, from 1 to  $N$ , with probabilities in bin  $n-1$  lower than those in bin  $n$  ( $n$  goes from 1 to  $N$ ). The lower probability limit for bin  $n$  is  $P_n$  and the upper limit is  $P_{n+1}$ . The lower probability limit for bin 1 is 0%, while the upper limit in bin  $N$  is 100%. The summation of the weights on the observed occurrences and non-occurrences of the event corresponding to each forecast in a given probability interval (bin  $n$  for example) is entered in the contingency table.

Tables 6 and 7 outline typical contingency tables. It is recommended that the number of probability bins remain between 10 and 20. The forecast providers can bin according to per cent thresholds (Table 6) or ensemble members (Table 7) as necessary. Table 7 gives an example of a table based on ensemble members.

The HR and FAR are calculated for each probability threshold  $P_n$  (see Tables 6 and 7). The HR for probability threshold  $P_n$  ( $HR_n$ ) is defined as (referring to Tables 6 and 7):

$$HR_n = \left( \sum_{i=n}^N O_i \right) / \left( \sum_{i=1}^N O_i \right)$$

and  $FAR_n$  is defined as:

$$FAR_n = \left( \sum_{i=n}^N NO_i \right) / \left( \sum_{i=1}^N NO_i \right)$$

where  $n$  goes from 1 to  $N$ . The range of values for  $HR_n$  goes from 0 to 1, the latter value being desirable. The range of values for  $FAR_n$  goes from 0 to 1, zero being desirable. Frequent practice is for probability intervals of 10% (10 bins, or  $N = 10$ ) to be used. However the number of bins ( $N$ ) should be consistent with the number of members in the EPS used to calculate the forecast probabilities. For example, intervals of 33% for a nine-member ensemble system could be more appropriate.

The HR and FAR are calculated for each probability threshold  $P_n$ , giving  $N$  points on a graph of HR (vertical axis) against FAR (horizontal axis) to form the ROC curve. This curve, by definition, must pass through the points (0,0) and (1,1) (for events being predicted only with >100% probability (never occur) and for all probabilities exceeding 0%, respectively). No-skill forecasts are indicated by a diagonal line (where  $HR = FAR$ ); the further the curve lies towards the upper left-hand corner (where  $HR = 1$  and  $FAR = 0$ ) the better.

The area under the ROC curve is a commonly used summary statistic representing the skill of the forecast system. The area is standardized against the total area of the figure such that a perfect forecast system has an area of one and a curve lying along the diagonal (no information) has an area of 0.5. The normalized ROC area has become known as the ROC score. Not only can the areas be used to contrast different curves, but they are also a basis for Monte Carlo significance tests. It is proposed that Monte Carlo testing be done within the forecast data set itself. For the core SVS for LRF, the area under the ROC curve can be calculated using the trapezium rule, although other techniques are available to calculate the ROC score (see Mason, 1982).

⇒ Contingency tables for probabilistic forecasts (such as in Tables 6 and 7) are mandatory for Level 3 verification in the core SVS. ROC curves and ROC areas are mandatory for Level 1 verification in the core SVS, while ROC areas are only mandatory for Level 2 verification in the core SVS.

### 3.3.4 Reliability diagrams and frequency histograms for probabilistic forecasts

It is recommended that the construction of reliability curves (including frequency histograms to provide indications of sharpness) be done for large-sample probability forecasts aggregated over the tropics and, separately, the two extratropical hemispheres. Given frequency histograms, the reliability curves are sufficient for the ROC curve, and have the advantage of indicating the reliability of the forecasts, which is a deficiency of the ROC. It is acknowledged that the ROC curve is frequently the more appropriate measure of forecast quality than the reliability diagram in the context of verification of LRF because of the sensitivity of the reliability diagram to small sample sizes. However, because measures of forecast reliability are important for modellers, forecasters and end-users, it is recommended that in exceptional cases when forecasts are spatially aggregated over the tropics and over the two extratropical hemispheres, reliability diagrams be constructed in addition to ROC curves.

The technique for constructing the reliability diagram is somewhat similar to that for the ROC. Instead of plotting the HR against the FAR for the accumulated probability bins, the HR is calculated only from the sets of forecasts for each probability bin separately, and is plotted against the corresponding forecast probabilities. The HR for each probability bin ( $HR_n$ ) is defined as:

$$HR_n = \frac{O_n}{O_n + NO_n}$$

This equation should be contrasted with the hit rate used in constructing the ROC diagram.

Frequency histograms are constructed similarly from the same contingency tables as those used to produce reliability diagrams. Frequency histograms show the frequency of forecasts as a function of the probability bin. The frequency of forecasts for probability bin  $n$  ( $F_n$ ) is defined as:

$$F_n = \frac{O_n + NO_n}{T}$$

where  $T$  is the total number of forecasts and  $T = \sum_{n=1}^N (O_n + NO_n)$

⇒ Reliability diagrams and frequency histograms are mandatory for Level 1 verification in the core SVS.

### 3.3.5 Level of significance

Because of the increasing uncertainty in verification statistics with decreasing sample size, significance levels and error bars should be calculated for all verification statistics. Recommended procedures for estimating these uncertainties are detailed below.

#### ROC AREA

In certain special cases, the statistical significance of the ROC area can be obtained from its relationship to the Mann-Whitney U statistic. The distribution properties of the U statistic can be used only if the samples are independent. This

assumption of independence will be invalid when the ROC is constructed from forecasts sampled in space because of the strong spatial (cross) correlation between forecasts (and observations) at nearby grid points or stations. However, because of the weakness of serial correlation of seasonal climate anomalies from one year to the next, an assumption of sequential independence may frequently be valid for LRF, so the Mann–Whitney U statistic may be used for calculating the significance of the ROC area for a set of forecasts from a single point in space. An additional assumption for using the Mann–Whitney U test is that the variance of the forecast probabilities (not that of the individual ensemble predictions per se) for when non-events occurred is the same as for when events occurred. The Mann–Whitney U test is, however, reasonably robust to violations of homoscedasticity, which means that the variance of the error term is constant across the range of the variable, hence significance tests in cases of unequal variance are likely to be only slightly conservative.

If the assumptions for the Mann–Whitney U test cannot be held, the significance of the ROC area should be calculated using randomization procedures. Because the assumptions of permutation procedures are the same as those of the Mann–Whitney U test, and because standard bootstrap procedures assume independence of samples, alternative procedures such as moving block bootstrap procedures (Wilks, 1997) should be conducted to ensure that the cross- and/or serial-correlation structure of the data is retained.

#### ROC CURVES

Confidence bands for the ROC curve should be indicated, and can be obtained either by appropriate bootstrap procedures, as discussed above, or, if the assumption of independent forecasts is valid, from confidence bands derived from a two-sample Kolmogorov-Smirnov test comparing the empirical ROC with the diagonal.

#### MSSS

Appropriate significance tests for the MSSS and the individual components of the decomposition again depend upon the validity of the assumption of independent forecasts. If the assumption is valid, significance tests could be conducted using standard procedures (namely the F ratio for the correlation and for the variance ratio, and the t test for the difference in means), otherwise bootstrap procedures are recommended.

⇒ Level of significance will be mandatory in the core SVS once guidelines for calculation have been established for the complete suite of scores. A phased-in introduction of level of significance in the SVS may be used (see section 3.1.4).

### 3.4 Hindcasts

In contrast to short- and medium-range dynamical numerical weather prediction (NWP) forecasts, LRF are produced relatively few times a year (for example, one forecast for each season or one forecast for the following 90-day period, issued every month). Therefore the verification sampling for LRF may be limited, possibly to the point where the validity and significance of the verification results may be questionable. Providing verification for a few seasons or even over a few years only may be misleading and may not give a fair assessment of the skill of any LRF system. LRF systems should be verified over as long a period as possible in hindcast mode. Although there are limitations on the availability of verification data sets and in spite of the fact that validating numerical forecast systems in hindcast mode requires large computer resources, the hindcast period should be as long as possible. The recommended period for the exchange of scores is advertised on the Lead Centre website (<http://www.bom.gov.au/wmo/lrfvs/>).

Verification in hindcast mode should be achieved in a form as close as possible to the real-time operating mode in terms of resolution, ensemble size and parameters. In particular, dynamical/empirical models must not make any use of future data. Validation of empirical models and dynamical models with postprocessors (including bias corrections), and calculation of period of verification means, standard deviations, class limits, etc. must be done in a cross-validation framework. Cross-validation allows the entire sample to be used for validation (assessing performance, developing confidence intervals, etc.) and almost the entire sample for model and postprocessor building and for estimation of the period of verification climatology. Cross-validation procedures are as follows:

1. Delete 1, 3, 5, or more years from the complete sample;
2. Build the statistical model or compute the climatology;
3. Apply the model (e.g. make statistical forecasts or postprocess the dynamical forecasts) or the climatology for one (usually the middle) year of those deleted and verify;
4. Replace the deleted years and repeat 1–3 for a different group of years;
5. Repeat 4 until the hindcast verification sample is exhausted.

Ground rules for cross-validation are that every detail of the statistical calculations be repeated, including redefinition of climatology and anomalies, and that the forecast year predictors and predictands are not serially correlated with their counterparts in the years reserved for model building. For example, if adjacent years are correlated but every other year is effectively not, three years must be set aside and forecasts made only on the middle year (see Livezey, 1999 for estimation of the reserved window width).

The hindcast verification statistics should be updated once a year based on accumulated forecasts.

⇒ Verification results over the hindcast period are mandatory for the exchange of LRF verification scores. Producing centres have to send new hindcast verification results as soon as their forecast system is changed.

### 3.5 Real-time monitoring of forecasts

It is recommended that there be regular monitoring of real-time LRF. It is acknowledged that this real-time monitoring is neither as rigorous nor as sophisticated as hindcast verification; nevertheless it is necessary for forecast production and dissemination. It is also acknowledged that the sample size for this real-time monitoring may be too small to assess the overall skill of the models. However, it is recommended that the forecast and the observed verification for the previous forecast period be presented in visual format to the extent possible given the restrictions on availability of verification data.

Real-time monitoring of forecast performance is an activity for the GPCs rather than the Lead Centre. GPCs are free to choose the format and content of real-time monitoring information.

## 4. VERIFICATION DATA SETS

The same data should be used to generate both climatology and verification data sets, although the forecast issuing Centres/Institutes own analyses or reanalyses and subsequent operational analyses may be used when other data are not available.

Many LRF are produced that are applicable to limited or local areas. It may not be possible to use the data in either the recommended climatology or verification data sets for validation or verification purposes in these cases. Appropriate data sets should then be used with full details provided.

Verification should be done using the recommended data sets as listed on the Lead Centre website (<http://www.bom.gov.au/wmo/lrfvs/>).

## 5. SYSTEM DETAILS

Information must be provided on the system being verified. This information should include (but is not restricted to):

1. Whether the system is numerical, empirical or hybrid,
2. Whether the system is deterministic or probabilistic,
3. Model type and resolution,
4. Ensemble size,
5. Specifications of boundary conditions,
6. List of parameters being assessed,
7. List of regions for each parameter,
8. List of forecast ranges (lead times) and periods for each parameter,
9. Period of verification,
10. The number of hindcasts or predictions incorporated in the assessment and the dates of these hindcasts or predictions,
11. Details of climatological and verification data sets used (with details on quality control when these are not published),
12. If appropriate, resolution of fields used for climatologies and verification.

Verification data for the aggregated statistics and the grid-point data should be provided on the Web. The contingency tables should be made available on the Web or by anonymous FTP. Real-time monitoring should be done as soon as possible and made available on the Web.

## 6. SVS FOR LRF LEAD CENTRE

The Fourteenth WMO Congress endorsed the designation by CBS-Ext.(02) of WMC Melbourne and the Canadian Meteorological Centre in Montreal as Co-Lead Centres for verification of long-range and SI forecast activities. The Co-Lead Centre functions include creating and maintaining coordinated websites for LRF verification information, so that potential users will benefit from a consistent presentation of the results. The goal is to help the RCCs and NMHSs to have a tool for improving the long-range forecasts delivered to the public. Congress urged all Members to actively participate in that activity as either users or producers of LRF verification information to assure the use of the best available products.



## 6.1 Role of the Lead Centre

6.1.1 The purpose of the Lead Centre is to create, develop and maintain the website (the "SVSLRF website") to provide access to the LRF verification information. The address of the website is <http://www.bom.gov.au/wmo/lrfvs/>. The website will:

- (a) Provide access to standardized software for calculating scoring information (ROC curves, areas, contingency table scores, hit rates, ...).
- (b) Provide consistent graphical displays of the verification results from participating centres through processing of digital versions of the results;
- (c) Contain relevant documentation and links to the websites of global-scale producing centres;
- (d) Provide some means for the collection of feedback from NMHSs and RCCs on the usefulness of the verification information;
- (e) Contain information and, preferably, provide access to available verification data sets;

6.1.2 The Centre will also:

- (a) Produce monthly verification data sets in common format on  $2.5^\circ \times 2.5^\circ$  grids where appropriate;
- (b) Liaise with other groups involved in verification (e.g. WGSIP and CCI) on the effectiveness of the current standardized verification system (SVS) and identify areas for future development and improvement;
- (c) Provide periodic reports to CBS and other relevant commissions assessing the effectiveness of the SVS;
- (d) Facilitate the availability of information to assess the skill of long-range forecasts but not to provide a direct intercomparison between the GPCs' models.

6.1.3 Detailed tasks of the Lead Centre

6.1.3.1 The Lead Centre will provide access to verification data sets on the SVSLRF website. The verification data sets will be in GRIB Edition 1 format. They will be translated to GRIB Edition 2 format when the encoder/decoder becomes widely available. RSMC Montreal will be responsible for preparing the verification data sets. These will be updated on the SVSLRF website on a yearly basis provided that new data are available. The choice of verification data sets will be revised as new ones become available and as recommended by the appropriate CBS expert team.

6.1.3.2 The Lead Centre will develop and provide specifications defining the format of the data to be sent to the Lead Centre for graphics preparation. There is no need to specify standards for graphics to be sent to the SVSLRF website because all graphics will be generated by the Lead Centre. WMC Melbourne will develop the infrastructure to generate all graphics posted on the SVSLRF website.

6.1.3.3 The Lead Centre will be responsible for making available the digital verification information as specified at Levels 1, 2 and 3 (see section 3.1).

6.1.3.4 The Lead Centre will ensure that clear and concise information explaining the verification scores, graphics and data is available and kept up to date on the SVSLRF website. The production of this documentation will be shared between the two Co-Lead Centres. Links to the participating GPCs will be listed on the SVSLRF website. The content of the documentation and information on interpretation and use of the verification data will be determined in consultation with the appropriate CBS expert team.

6.1.3.5 The Lead Centre will consult with the GPCs to make sure that the verification data are correctly displayed before making available their verification results on the SVSLRF website.

6.1.3.6 The Lead Centre will ensure that the verification results placed on the SVSLRF website come from GPCs (officially recognized by CBS) with operational LRF commitments.

6.1.3.7 The Lead Centre will provide and maintain software to calculate the verification scores. Development of the software will be the responsibility of RSMC Montreal. The software code will be available on the SVSLRF website, and will be in FORTRAN language. However, it is recognized that the use of this software is not mandatory.

6.1.3.8 The Lead Centre will publicize the SVSLRF website to other bodies involved in verification (such as WGSIP and CCI) and establish contacts in order to receive feedback and facilitate discussion for further development and improvement.

6.1.3.9 Once the SVSLRF website is operational, the Lead Centre will provide progress reports every two years to CBS, prior to its sessions.



## 7. REFERENCES

- Gerrity, J. P. Jr., 1992: A note on Gandin and Murphy's equitable skill score. *Monthly Weather Review*, 120, pp. 2707-2712.
- Hanssen A. W. and W. J. A. Kuipers, 1965: On the relationship between the frequency of rain and various meteorological parameters. Koninklijk Nederlands Meteorologisch Instituut, *Mededelingen en Verhandelingen*, No. 81, pp. 2-15.
- Livezey, R. E., 1999: Chapter 9: Field intercomparison. *Analysis of Climate Variability: Applications of Statistical Techniques*, H. von Storch and A. Navarra, Eds, Springer, pp. 176-177.
- Mason I., 1982: A model for assessment of weather forecasts. *Australian Meteorological Magazine*, 30, pp. 291-303.
- Mason, S. J., and G. M. Mimmack, 2002: Comparison of some statistical methods of probabilistic forecasting of ENSO. *J. Climate*, 15, pp. 8-29.
- Murphy, A. H., 1988: Skill scores based on the mean square error and their relationships to the correlation coefficient. *Monthly Weather Review*, 116, pp. 2417-2424.
- New, M. G., M. Hulme and P. D. Jones, 2000: Representing twentieth-century space-time climate variability. Part II: Development of 1901-1996 monthly grids of terrestrial surface climate. *J. Climate*, 13, 2217-2238.
- 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, pp. 34-53.
- Reynolds, R. W. and T. M. Smith, 1994: Improved global sea surface temperature analyses using optimum interpolation. *J. Climate*, 7, 929-948.
- Smith T. M., R. W. Reynolds, R. E. Livezey and D. C. Stokes, 1996: Reconstruction of historical sea-surface temperatures using empirical orthogonal functions, *J. Climate*, 9, pp. 1403-1420.
- Stanski H. R., L. J. Wilson and W. R. Burrows, 1989: Survey of common verification methods in meteorology. *World Weather Watch Technical Report No. 8*, WMO/TD-No. 358, 114 pp.
- Wilks, D. S., 1997: Resampling hypothesis tests for autocorrelated fields. *J. Climate*, 10, pp. 65-82.
-