# *Time-window approaches to space-weather forecast metrics: a solar wind case study*

Article

Published Version

**Correspondence to:**
M. J. Owens,
m.j.owens@reading.ac.uk

# Time-Window Approaches to Space-Weather Forecast Metrics: A Solar Wind Case Study

**Mathew J. Owens[1]** ID

[1]Space and Atmospheric Electricity Group, Department of Meteorology, University of Reading, Reading, UK

**Abstract** Metrics are an objective, quantitative assessment of forecast (or model) agreement with observations. They are essential for assessing forecast accuracy and reliability and consequently act as a diagnostic for forecast development. Partly as a result of limited spatial sampling of observations, much of space-weather forecasting is focused on the time domain rather than inherent spatial variability. Thus, metrics are primarily *point-by-point* approaches, in which observed conditions at time *t* are compared directly (and only) with the forecast conditions at time *t*. Such metrics are undoubtedly useful. But in lacking an explicit consideration of timing uncertainties, they have limitations as diagnostic tools and can, under certain conditions, be misleading. Using a near-Earth solar wind speed forecast as an illustrative example, this study briefly reviews the most commonly used point-by-point metrics and advocates for complementary *time window* approaches. In particular, a scale-selective approach, originally developed in numerical weather prediction for validation of spatially patchy rainfall forecasts, is adapted to the time domain for space-weather purposes. This simple approach readily determines the time scales over which a forecast is and is not valuable, allowing the results of point-by-point metrics to be put in greater context.

## 1. Introduction

When determining how well a space-weather forecast performs, human assessment can rapidly scrutinize a large number of facets: simply looking over the observations and forecast gives an immediate *feel* for what features are reproduced and missed, how the general structure differs, over what temporal/spatial scales the forecast is applicable, whether the forecast exhibits any obvious bias, performs better within certain parameter regimes, etc. (Throughout this study, metrics are discussed with regard to *forecasting*, though the same issues and principles apply for general model diagnostics. Consequently, anywhere the term *forecast* appears, the term *model* could be directly substituted.) But this is inherently subjective, qualitative, lacking in repeatability, and simply infeasible for large volumes of data. Metrics are an automated, objective quantification of forecast performance relative to observations. As such, metrics are vitally important not just for validation of space-weather forecasts (e.g., Spence et al., 2004) but also as a diagnostic tool to inform future forecast development. (As in the majority of the space-weather literature, the term *validation* is here used to refer to the process of comparing forecasts and observations to establish accuracy and truth of the forecast. This is often referred to as *verification* in meteorology.) Different metrics quantify different, specific qualities of a forecast. Thus, while there are no right or wrong metrics per se, it is nevertheless essential to select a metric which actually measures the features of interest. This, as will be seen in the subsequent examples, is not always as straightforward as it seems. Changes to a forecast scheme made on the basis of a poorly chosen metric can potentially reduce its usefulness for an end-user, though of course the chosen metric will measure an improvement.

The space-weather community is in the process of adopting both more sophisticated forecast approaches and metrics with enhanced diagnostic capability (e.g., Jian et al., 2016; Murray, 2018; Murray et al., 2017). Many of these approaches have been adapted from numerical weather prediction (NWP; Siscoe, 2007). In NWP, there is extensive coverage by the observation network, allowing both spatial and temporal agreement to be explicitly treated. Extremely sparse observational sampling of the Sun-Earth system, however, means that space-weather forecast validation is often primarily concerned with the time domain (though errors in the time domain may well result from spatial variations). For example, while forecasts of the solar wind (such as the example of near-Earth solar wind shown in section 2.1) cover the largest spatial domain within the Sun-Earth system, they are typically validated solely against single-point in situ observations made in near-Earth space (e.g., MacNeice, 2009; MacNeice et al., 2018; Owens et al., 2008). Consequently, validation

is primarily focused on a point-by-point analysis: The observed conditions at time $t$ are compared directly (and only) with the forecast conditions at time $t$. As is illustrated in sections 2.2 and 2.3, such approaches inflict a *double penalty* for timing offsets in forecast events, due to both missing the event and generating a false alarm. On the one hand, this is a legitimate assessment of the forecast. On the other hand, it does not always provide a useful diagnostic of the forecast, and many operators will tolerate relatively small errors in event timing if the general outlook is correct. One solution is for forecasts to include a measure of their own uncertainty, as illustrated in section 2.5. However, this is not always practical. Thus, in addition to point-by-point metrics, it may be advantageous to also employ *time-window* metrics. One useful approach, outlined in section 3.1, is to specify criteria for discrete features within forecast and observation time series and to com-pare feature correspondence, including the timing. However, such feature specification requires a priori knowledge of the properties of interest, as well as repeatable signatures in said features, both from event to event and across forecast and observation data. Thus, in section 3.2 a more feature-agnostic approach is proposed, based upon NWP validation of rain forecasts. It compares forecasts and observations at a range of different spatial scales and is here adapted to the time domain as a space-weather forecast metric. It is shown that this analysis provides a useful assessment of the time scales over which a forecast is and is not valuable.

## 2. Point-By-Point Metrics

### 2.1. Example Forecast

In order to illustrate the strengths and limitations of different metrics, an example forecast is considered. The black line in Figure 1a shows hourly near-Earth solar wind speed ($V$) for Carrington rotation (CR) 2049, span-ning mid-October to mid-November 2006. Data are from the *Omni* data set of near-Earth spacecraft measure-ments (King & Papitashvili, 2005). CR 2049 was chosen as there are three distinct high-speed enhancements (HSEs) on 20 October, 28 October, and 9 November.

An illustrative forecast was produced using the *Magnetohydrodynamics Around a Sphere* (MAS; Linker et al., 1999; Riley et al., 2012) global coronal model. The inner boundary conditions are set by the observed photo-spheric magnetic field for CR 2049. Model output is available from http://www.predsci.com/mhdweb/. Typically, the MAS solution would be propagated to near-Earth space with a numerical magnetohydrody-namic solar wind model and the forecast $V$ extracted from the model grid point closest to Earth. Here, how-ever, for the purposes of demonstration, the solution was perturbed to (retrospectively) produce a closer match to the observations. Specifically, the model solar wind at 30 solar radii was sampled 5° above the sub-Earth point, as this was found to improve the representation of the HSE on 28 October. The solar wind speed was then propagated from 30 solar radii to Earth using a simple *upwind* technique (Owens & Riley, 2017) to produce the time series shown in red in Figure 1a.

### 2.2. Error Functions

Forecasts are commonly assessed using simple error functions (otherwise called cost or loss functions). The results for CR2049 are summarized in Table 1. For solar wind speed, the mean-square error (MSE) is given by

$$\text{MSE} = \frac{1}{T} \sum_{t=1}^{T} \left[ V_F(t) - V(t) \right]^2$$

where $V_F(t)$ and $V(t)$ are the forecast and observed solar wind speeds at time $t$, respectively, and $T$ is the total number of time points considered. Smaller MSE values indicate better agreement, with 0 being a perfect fore-cast. For the forecast shown in Figure 1a, the MSE is $1.30 \times 10^4$ km$^2$/s$^2$. This is usually converted to root-mean-square (RMS) error:

$$\text{RMS} = \sqrt{\frac{1}{T} \sum_{t=1}^{T} \left[ V_F(t) - V(t) \right]^2}$$

RMS has the advantage of being a linear measure of the magnitude of the errors with the same units as the parameter of interest. The RMS error for the forecast is 114 km/s.
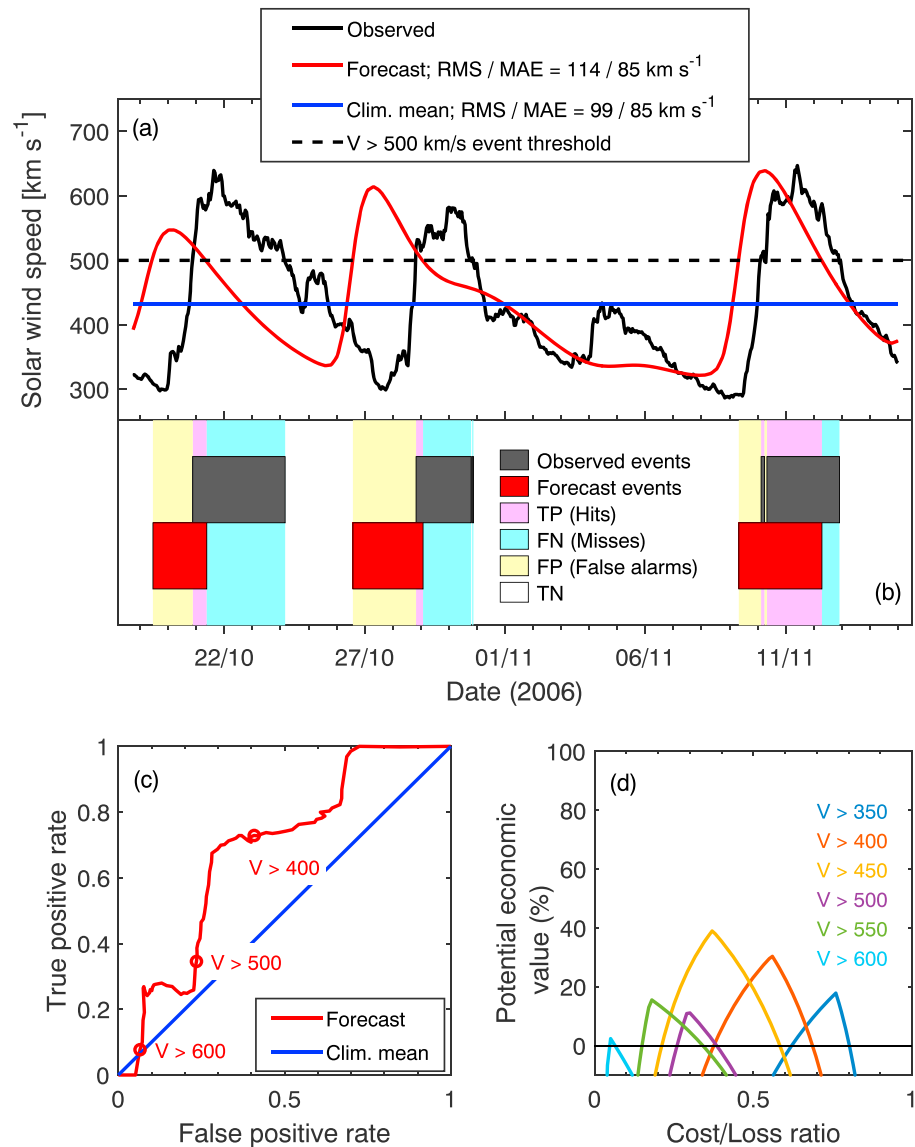
**Figure 1.** An example of a deterministic solar wind speed forecast and associated *point-by-point* metrics. (a) Time series of hourly means of near-Earth solar wind speed, *V*, for Carrington rotation 2049, spanning mid-October to mid-November 2006, as observed (black) and forecast (red). The climatological mean for this interval (blue) is also shown. (b) Solar wind speed events defined using a threshold of $V > 500$ km/s. (c) The receiver operator characteristic that plots the true positive rate against the false positive rate for a range of solar wind speed event definitions. (d) The potential economic value of the forecast at various *V* thresholds and cost/loss ratios. See text for more detail.

In isolation, these values say relatively little about the quality of the forecast. Metrics are most useful as a comparative tool. Thus, it is instructive to also consider a second solar wind speed prediction. The blue line shows the average *V* for CR 2049, 432 km/s. For validation purposes, this climatological mean would be a poor choice of comparison prediction, as it has zero variability. In practice, it would be preferable to use another simple forecast, such as 27-day recurrence (Owens et al., 2013). But for the purposes of illustrating certain issues, the climatological mean is useful here. The MSE between the observed *V* for CR 2049 and the climatological mean is $0.98 \times 10^4$ km²/s², while the RMS 98.9 km/s, both smaller than the forecast values.

An alternative measure of a similar property is the mean absolute error (MAE):

$$\text{MAE} = \frac{1}{T} \sum_{t=1}^{T} |V_F(t) - V(t)|$$

**Table 1**
*Point-By-Point Metrics for the Solar Wind Speed Forecasts Shown in Figure 1*

| | MSE[a] $(km^2/s^2)$ | RMS[b] (km/s) | MAE[c] (km/s) | $r_L$[d] | $r_S$[e] | ROC area under curve[f] |
|---|---|---|---|---|---|---|
| Forecast | $1.30 \times 10^4$ | 114.0 | 84.6 | 0.28 | 0.06 | 0.68 |
| Climatological mean | $0.98 \times 10^4$ | 98.9 | 85.0 | 0.00 | 0.00 | 0.50 |

[a]Mean-square error. [b]Root-mean-square error. [c]Mean absolute error. [d]Pearson (linear) correlation coefficient. [e]Spearman (rank-order) correlation coefficient. [f]Receiver operator characteristic area under curve.

For the V, MAE is essentially the same for the forecast (84.6 km/s) and the climatological mean (85.0 km/s).

In order to further put error functions in perspective, the *skill* of a forecast is calculated as

$$\text{Skill} = 1 - \frac{\text{MSE}}{\text{MSE}_{\text{REF}}}$$

where $\text{MSE}_{\text{REF}}$ is the MSE of a reference *baseline* model, such as the climatological mean. Skill is negative when the forecast is worse than the baseline, 0 when they are equal, and 1 for a perfect forecast. (Sometimes skill is further multiplied by 100 to express it as a percentage of a perfect forecast.) By comparing directly with a baseline model, skill potentially allows disambiguation between bad forecasts and periods/situations that are inherently difficult to forecast. For the forecast shown in Figure 1a, using the climatological mean as the reference, the forecast skill is −0.32. Thus, the forecast is deemed to be *worse* than assuming that the solar wind is always a constant 432 km/s.

The general conclusion from these error functions for this example period is that the climatological mean is at least as *good* as the forecast for CR 2049. This is, of course, an entirely correct and fair assessment. But it is obvious that it does not tell the whole story; the climatological mean lacks sharpness and discrimination, in that it does not reconstruct any of the features of the solar wind structure. It would be useless as a predictive tool for almost all applications and thus could be said to lack value. In contrast, the forecast appears to work quite well for this interval: By eye, it can be seen that the forecast produces three HSEs, as observed, and they are of comparable magnitudes and durations to the observations. By inspection of the time series, it can be seen that the error functions for the forecast are relatively high due to the approximately 1- to 2-day errors in the timings of the HSEs that result in the *double penalty* of first overpredicting V, closely followed by underpredicting V. But, depending on the application, the forecast may well still be regarded as valuable in that it enables users to make decisions that lead to beneficial outcomes (Murphy, 1993).

In this particular example, other forms of point-by-point comparisons *are* able to discriminate between the predictive value of the forecast and climatological mean (see section 3.3 for an example where this is not the case). While not strictly an error function, Pearson (or linear) correlation, $r_L$, is often used in a similar manner to RMS and MAE to quantify forecast and observation agreement, where

$$r_L = \frac{\sum_{t=1}^{T} \left[ V_F(t) - \overline{V_F} \right] \left[ V(t) - \overline{V} \right]}{\sqrt{\sum_{t=1}^{T} \left[ V_F(t) - \overline{V_F} \right]^2} \sqrt{\sum_{t=1}^{T} \left[ V(t) - \overline{V} \right]^2}}$$

It is weakly positive for the forecast ($r_L = 0.28$). Spearman correlation, $r_S$, replaces the observed and forecast values at time $t$ with their ranks within their respective distributions. As a result, $r_S$ is less susceptible to outliers than $r_L$. It is effectively 0 ($r_S = 0.06$) for the forecast. The zero variance of the climatological mean results in both $r_L = 0$ and $r_S = 0$. Figure 2 summarizes these results in the form of a Taylor diagram (Riley, Linker, & Mikić, 2013; Taylor, 2001). It displays the RMS (centered by the mean values to remove forecast bias) and linear correlation between forecast and observation, along with the standard deviation of the time series under consideration. In short, the closer the forecast (red point) to the observation (black circle), the better. Thus, while the Taylor diagram does not strictly conclude that the forecast is superior to the climatological mean (blue point), the issues with the latter as predictive tool are immediately obvious. For a more realistic *baseline* forecast, this may not always be the case.
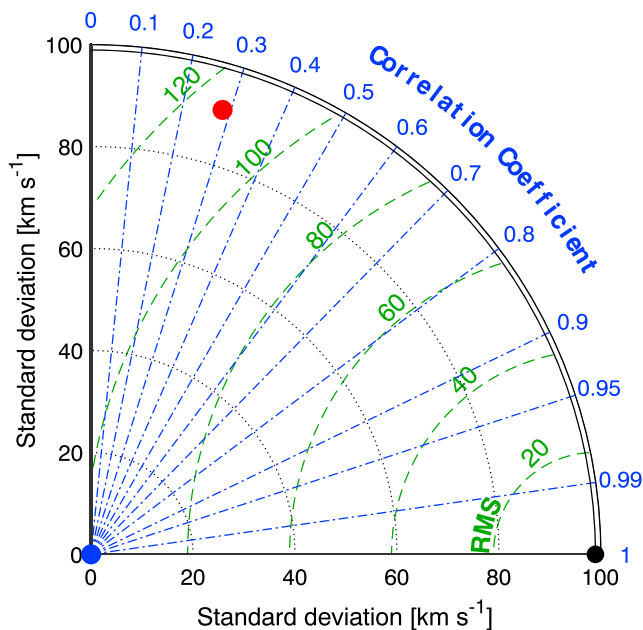


**Figure 2.** A Taylor diagram of the solar wind speed time series shown in Figure 1. The radial distance from the origin shows the standard deviation of the time series, while the azimuthal angle about the origin shows the linear correlation coefficient (note nonlinear scale) with the observed time series. The green dashed circles show contours of constant root-mean-square error (with forecast and observation mean subtracted). The black, red, and blue points show the observed, forecast, and climatological V, respectively.

**Table 2**
*A Contingency Table for the Forecast of Solar Wind Speed Events in CR 2049 Defined by a Threshold of V > 500 km/s*

| | | Event in forecast? (i.e., $V_F > 500$ km/s) | | |
|---|---|---|---|---|
| | | Yes | No | Total |
| Observed event? (i.e., $V > 500$ km/s) | Yes | $TP = 68$ | $FN = 124$ | $P = 192$ |
| | No | $FP = 109$ | $TN = 354$ | $N = 463$ |
| | Total | $P_F = 177$ | $N_F = 478$ | 655 |

*Note.* TP, FP, TN, and FN are the numbers of true positive, false positive, true negative, and false negative intervals, respectively. P and $P_F$ are the number of observed and forecast events, while N and $N_F$ are the number of observed and forecast nonevents.

In addition to potentially misleading forecast assessment, error functions can also have unintended consequences for model development. Riley, Linker, and Mikić (2013) note that changes to their coronal model that wipe out all solar wind speed variability (and thus value of the resulting forecast), are not reflected in RMS, which is essentially unchanged. Similarly, any forecast scheme trained to minimize RMS or MAE may tend preferentially toward a conservative, climatological-mean-like prediction, rather than a valuable forecast.

### 2.3. Binary Metrics

As error functions quantify the magnitude of forecast deviation from observations at every time step, they can have limitations as diagnostic tools. First, by considering every time step equally, rather than focusing on specific times or parameter ranges of interest, these metrics can be skewed toward measuring properties that are inconsequential to an operator. For example, whether the forecast correctly reproduces the details of the slow-speed wind may be unimportant, but is given equal weighting to the times of high speeds, which are important. Second, large outliers can have a relatively strong influence on error functions and especially on linear correlation. In some circumstances, this will be appropriate, as the magnitude of the extremes is of interest. In other circumstances, this may be less critical, as what matters is whether or not a given threshold is exceeded, not by how much. To address these issues, an alternative approach is to consider each time step as a binary *yes/no* state and compare observations and forecasts on this basis. For probabilistic forecasts, discussed further in section 2.5, this also involves setting a probability threshold, in addition to an event-definition threshold.

The black dashed line in Figure 1a shows a threshold of $V > 500$ km/s used to define hourly *events* in the forecast and observed time series. Figure 1b displays the timing of the subsequent forecast and observed events, sorting them into one of four categories; true positives (*TP*, or *hits*; hours for which both observed and forecast events are present), false positives (*FP* or *false alarms*; hours for which an event is forecast but not observed), false negatives (*FN*, or *misses*; hours for which an event is observed but not forecast), and true negatives (*TN*; hours for which both observation and forecast have no event). The occurrence of these classifications is summarized in a contingency table (e.g., Finley, 1884; Murphy, 1996), shown as Table 2 for the forecast and Table 3 for the climatological mean. The forecast produces approximately the correct number of events ($P_F = 177$ versus $P = 192$ observed) and nonevents ($N_F = 478$, versus $N = 463$ observed); meaning, it has little bias, whereas the climatological mean produces zero events and overestimates the nonevents ($N_F = 655$). The double penalty effect on the forecast is apparent: Because of the timing offset in the HSEs, the forecast produces both *FN* and *FP*, whereas the null prediction of the climatological mean only produces *FN*. For the forecast, the total number of false predictions, $FP + FN$, is 233, while for the climatological mean it is only 192.

From the contingency tables alone, it is not immediately clear whether the forecast is *better* than the climatological mean. It will depend on how *FP* and *FN* are weighed relative both to each other and to *TP* (and to a lesser extent, *TN*). There are a variety of ways to combine these four numbers, to emphasize different forecast aspects. The full range of combinations is not discussed here (see Thornes & Stephenson, 2001, and Reiss et al., 2016, as well as the World Meteorological Organization guide: http://www.cawcr.gov.au/projects/verification/). Two of the most useful combinations are the true positive rate ($TPR = TP/P$) and the false positive rate ($FPR = FP/N$), as together they provide a reasonable overview of a forecast. A perfect forecast would have $TPR = 1$ and $FPR = 0$. For the forecast of CR 2049, $TPR = 0.35$ and $FPR = 0.24$.

For events defined by $V > 500$ km/s, the climatological mean results in no true or false positives and so $TPR = 0$ and $FPR = 0$. If events were defined using a $V$ threshold lower than the climatological mean (e.g., $V > 400$ km/s), it would produce a prediction of events at all times, giving $TPR = 1$ and $FPR = 1$. Thus, for any event threshold, the climatological mean over the period under consideration gives $TPR = FPR$. When a forecast

**Table 3**
*The Same as Table 2 but for the Climatological Mean of Solar Wind Speed for CR 2049*

| | | Event in climatological mean? (i.e., $V_F > 500$ km/s) | | |
|---|---|---|---|---|
| | | Yes | No | Total |
| Observed event? (i.e., $V > 500$ km/s) | Yes | $TP = 0$ | $FN = 192$ | $P = 192$ |
| | No | $FP = 0$ | $TN = 463$ | $N = 463$ |
| | Total | $P_F = 0$ | $N_F = 655$ | 655 |

results in $TPR > FPR$, it is superior to the climatological mean in being able to predict the occurrence of events and nonevents.

### 2.4. Forecast Summaries

Binary metrics depend on the choice of both event and probability thresholds, and thus, ways to summarize parameter space are necessary. The (often complex) relation between $FPR$ and $TPR$ for a range of event thresholds is captured by the receiver operator characteristic (ROC; Peterson et al., 1954; Mason, 1982) curve in Figure 1c. This technique is commonly used for validation of probabilistic forecasts at a range of probability thresholds (see section 2.5), including solar flare forecasts (McCloskey et al., 2018; Murray et al., 2017). However, it can also be used to summarize the deterministic $V$ forecast. In this example, all event thresholds result in $TPR > FPR$ (i.e., are above the $y = x$ line in Figure 1c) except $V > 600$ km/s, where the double penalty is strongest. The ROC can be further distilled down to the area under the curve, integrated along the horizontal axis (AUC; Mason & Graham, 2002). AUC represents a forecast's ability to correctly anticipate events and nonevents (1 being a perfect forecast, 0.5 being equal to the climatological mean). For the $V$ forecast, the AUC is 0.68.

An alternative summary can be provided by the Cost-Loss analysis (Murphy, 1977; Richardson, 2000), which determines the benefit an operator would gain from acting on a forecast. The real strength of Cost-Loss analysis is in the evaluation of probabilistic forecasts (see section 2.5), as it explicitly accounts for the fact that different operational users will act on the same forecast in a different manner. For example, if a forecast gives a low probability of a space-weather event, an operator may still choose to take mitigating action if the cost of doing so (e.g., from lost revenue), $C$, is small relative to $L$, the cost of being caught unprepared by a damaging event. In such situations, forecasts that minimize missed events, even if this means increased false alarms, are more desirable. Conversely, if $C$ is a significant fraction of $L$, an operator is unlikely to act on the basis of a low forecast probability. In such circumstances, forecasts that minimize false alarms are more desirable. This analysis has recently been applied to validation of probabilistic solar wind forecasts (Owens et al., 2014; Owens et al., 2017).

Figure 1d shows the potential economic benefit of acting on the deterministic forecast of $V$ for a range of $C/L$ values and for events defined by a range of $V$ thresholds. Potential economic benefit is measured relative to the climatological probability of an event, so that values below 0% indicate that the forecast is less useful than climatology and 100% indicates a perfect (deterministic) forecast. As shown by the ROC curve, most benefit is gained at intermediate solar wind speeds (400 to 500 km/s) and for low $C/L$ scenarios. When false alarms become costly, the forecast ceases to add value, as the double penalty effect comes into play. Despite the insight gained from binary metrics such as ROC and Cost/Loss analysis, they nevertheless operate on a strictly point-by-point comparison basis and do not account for timing errors/uncertainty. As illustrated in section 3.3, the resulting double penalty issue is even stronger for $B_Z$ forecasts, which are critical for space weather (Dungey, 1961), as large-scale $B_Z$ variations tend to be bipolar in nature.

### 2.5. Validating Probabilistic Forecasts

Ideally, a forecast would include an assessment of forecast uncertainty. Figure 3a shows an example of a probabilistic forecast of solar wind speed for CR2049. It was generated using a perturbed initial condition ensemble (Owens & Riley, 2017). The RMS and MAE of the forecast ensemble median are comparable to the deterministic $V$ forecast shown in Figure 1a. But what is of most interest here is the uncertainty estimate. Figure 3b shows the probability of $V > 500$ km/s as a function of time. For the observations, this is either 0 or 1; for the climatological mean, it is always 0; for the forecast ensemble the probability is the fraction of ensemble members for which $V > 500$ km/s at each time step (e.g., Slingo & Palmer, 2011 and references therein). For the 21 October HSE, the onset timing uncertainty is reasonable, but the forecast is too confident of no event after 22 October. For the 29 October HSE, the forecast clearly underestimates the uncertainty in the HSE arrival time and duration, as the probability peaks more than a day early and remains high ($\approx 0.75$) for around a day too long. For the 10 November HSE, there is a 3-day spread in the HSE arrival time in the probabilistic forecast, with the peak probability on the 11 November, approximately the time of the observed peak.

In order to produce the ROC curve (Figure 3c), a probability threshold is required to define events at each $V$ threshold. In general for CR 2049, higher probability thresholds produce better forecasts as given by AUC
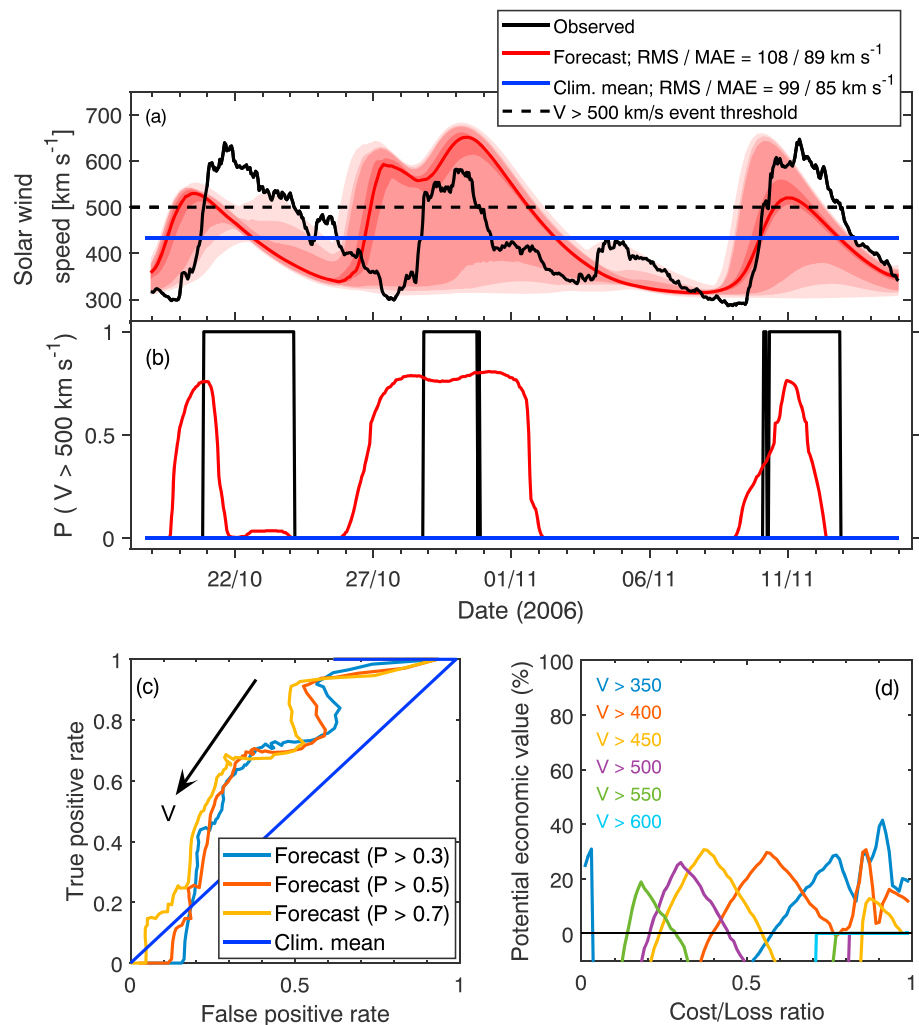
**Figure 3.** An example of a probabilistic solar wind speed forecast and associated point-by-point metrics. (a) The time series of hourly means of near-Earth solar wind speed for CR2049, mid-October to mid-November 2006, as observed (black) and forecast by the ensemble median (red), with pink-shaded areas showing 68, 90, 95, and 99.8 percentiles of the forecast ensemble. The climatological mean for this interval (blue) is also shown. (b) A threshold of $V > 500$ km/s is to define events in the time series (black dashed line), which are represented as a probability of occurrence. (c) The receiver operator characteristic for three different probability thresholds. (d) The cost-loss curves for the forecast at various action thresholds of $V$.

(though it is not a simple linear relation). On this basis alone, it may be tempting to conclude that the probabilistic forecast is most beneficial in operational situations where few false alarms are present (i.e., high $C/L$ ratios). However, that is not generally the case (as shown in Figure 3c and discussed below). What the ROC is actually revealing is simply that higher probability thresholds reduce the total number of forecast events and, in the presence of timing errors, minimize the double penalties described in the previous section. Thus again, even with probabilistic forecasts, point-by-point metrics can favor overly conservative forecasts.

From Figure 3b, it can be seen that for $V > 500$ km/s, there are no periods where the forecast probability of an event exceeds around 0.75. This means that operational settings in which forecast certainty is critical (i.e., where false alarms are costly), the forecast will not be useful. This is demonstrated in the cost-loss analysis in Figure 3d, where for $V > 500$ km/s, there is no economic benefit to acting on the forecast when $C/L > 0.5$. At lower speed thresholds, for example, 400 km/s, there are times when the forecast correctly predicts 0 probability of an event (6 to 8 November) and 1 probability of an event (22 October). This results in a valuable forecast for higher $C/L$ values, unlike the similar deterministic forecast.
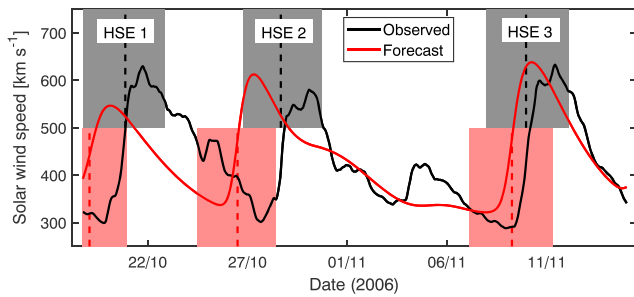
**Figure 4.** High-speed enhancement (HSE) analysis applied to the solar wind speed observed (black) and forecast (red) for CR2049, mid-October to mid-November 2006. All data have been 8-hr smoothed. The black- and red-shaded intervals show times when observed and forecast $V$ meet the criteria for a HSE, respectively. The dashed vertical lines show the times of maximum $V$ gradient.

Clearly, forecasts should intrinsically account for uncertainty, including the timing of features. However, forecasts often do a poor job in this respect (as shown in the example above), and uncertainty can be costly to estimate. For example, estimating the timing uncertainty in a coronal mass ejection (CME) forecast through a numerical model ensemble (Mays et al., 2015; Riley, Linker, & Mikić, 2013; Riley, Linker, Mikič, Zank, et al., 2013) will require a minimum of an order-of-magnitude more computing resources. Additionally, an operator may tolerate a greater timing error than the estimated forecast timing uncertainty. Thus, it is also desirable to use metrics that explicitly allow for timing uncertainty.

## 3. Time-Window Metrics
### 3.1. Feature-Based Metrics

One approach to dealing with timing uncertainty is to define discrete features (also known as objects or events) on the basis of extended spatial information or time history (rather just using a simple threshold on a point-by point basis, as in the case of binary metrics) and compare their properties, including timing (e.g., Ebert & Gallus, 2009). For example, Owens et al. (2005) defined a HSE as a net 100 km/s increase in $V$ over a 2-day interval in 8-hr smoothed data (computed as the mean in a rolling 8-hr window). The smoothing allows the analysis to be readily applied to both observations and numerical solar wind model output. The HSE lasts as long as these criteria are met, with the characteristic time of the HSE being the time of maximum $V$ gradient. Reiss et al. (2016) and MacNeice (2009) used similar definitions. Figure 4 shows the analysis applied to the CR2049 observations and forecast. In practice, when applying the analysis to years of data, observed and forecast HSEs are paired up algorithmically. In this instance, there are three observed and forecast HSEs, with forecast/observed pairs overlapping in time, so the pairing is trivial. Results are summarized in Table 4.

During this short interval of comparison, the forecast produces approximately the correct magnitude of HSE (in 8-hr smoothed data), but the timing of HSEs is systematically biased early. Clearly, this approach provides quantitative diagnostic information about *why* the RMS and MAE are high for this forecast relative to the climatological mean. The limitation in this kind of analysis is that features of interest have to be rigorously defined a priori. For solar wind speed, this is reasonable, but for $B_Z$, it may be more difficult, particularly regarding time scale and magnitude, as further discussed in section 3.3.

An alternate approach to timing uncertainties is to consider the peak value within a fixed time window (e.g., maximum $V$ in a 24-hr window of 1-hr data). This can provide useful information if, again, tailored to the specific needs of the operational setting. But there are a number of considerations with applying this approach more generally. First, different time windows will, of course, be more or less appropriate for different forecast applications. Second, for a fixed time window, the same peak value can result from a single data spike, multiple peaks, or the whole window being elevated. Third, changing the time resolution of the data can affect the peak values in different ways: The peak value of the single data spike will be dramatically reduced, whereas broader peaks will be less affected. A method to effectively summarize this parameter space for a binary forecast is described in the next section.

**Table 4**
*Results of the High-Speed Enhancement Analysis Applied to the Observed and Forecast Solar Wind Speed for CR 2049, Mid-October to mid-November 2006*

| HSE | Observed | Forecast | $\Delta T$ (days) | $|\Delta T|$ (days) | $V_{MAX}$ obs (km/s) | $V_{MAX}$ for (km/s) | $\Delta V_{MAX}$ (km/s) | $|\Delta V_{MAX}|$ (km/s) |
|---|---|---|---|---|---|---|---|---|
| 1 | 20 October 2006 T21:00 | 19 October 2006 T02:00 | 1.79 | 1.79 | 630 | 547 | 82.9 | 82.9 |
| 2 | 28 October 2006 T16:00 | 26 October 2006 T12:00 | 2.17 | 2.17 | 580 | 612 | −32.1 | 32.1 |
| 3 | 9 November 2006 T23:00 | 9 November 2006 T6:00 | 0.71 | 0.71 | 633 | 638 | −5.1 | 5.1 |
| Mean | – | – | 1.56 | 1.56 | 614 | 599 | 15.2 | 40.1 |

*Note.* In both case, three HSEs were identified. The symbol $\Delta$ indicates the (observed-forecast) value.
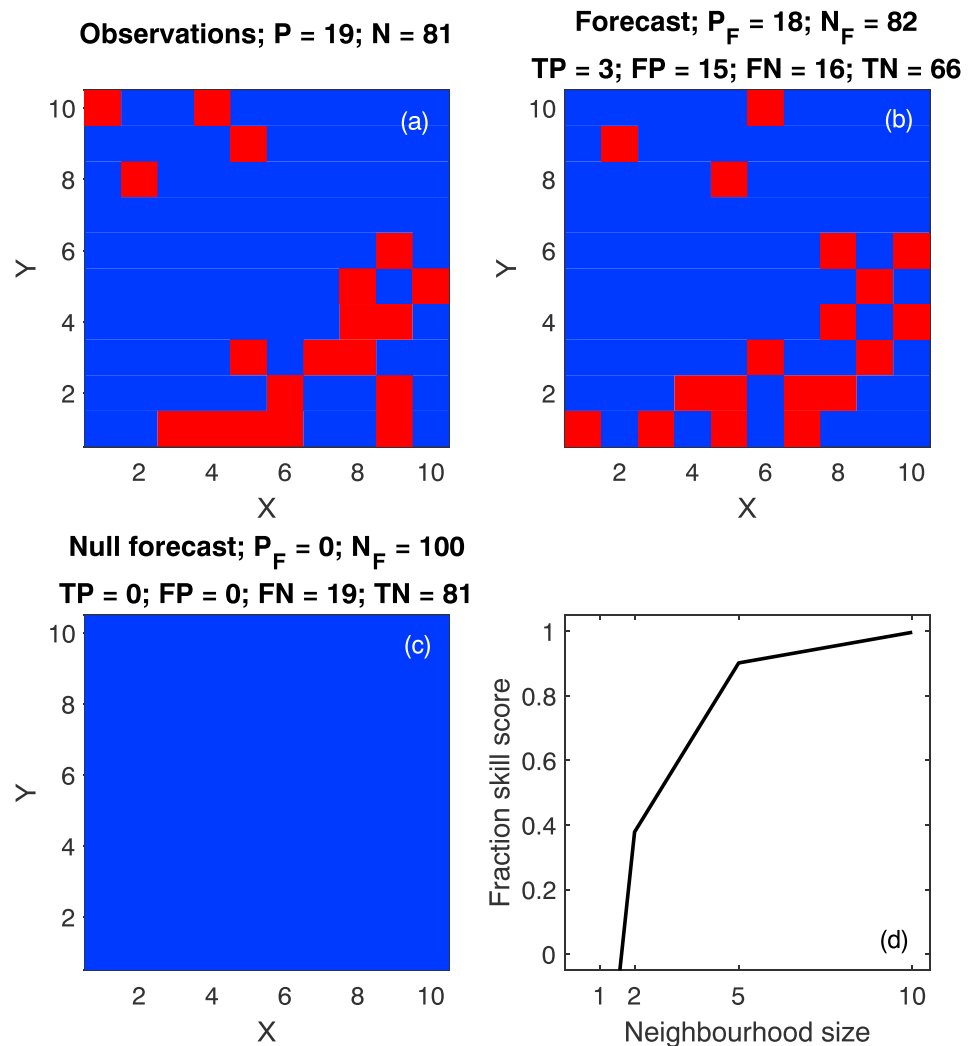
**Figure 5.** Spatial distributions of hypothetical (a) observed and (b) forecast rain. Red is a positive observation/forecast at a given position; blue is negative. (c) A null forecast predicts no rain anywhere. (d) The fraction skill score for different spatial scales (or neighborhood sizes).

### 3.2. Scale-Selective Metrics

In validation of forecasts from NWP, double penalties are also a ubiquitous issue, resulting from both spatial and temporal offsets between forecast and observation. A particularly apposite example is convective rain, which is inherently patchy on the spatial scales measureable by radar and forecastable by NWP. This can lead to misdiagnosis of forecasts if performed on a point-by-point basis at the grid-cell level. Hypothetical rain observations and forecast for a $10 \times 10$ grid are shown in Figure 5. The forecast has little bias over the whole domain (forecast and observation predict 18%, and 19% of grid points, respectively, will contain rain) and captures much of the large-scale structure, with a front of rain in the bottom-right corner of the domain. There is, however, little correspondence at the individual grid-point level. Making a simple point-by-point comparison of the forecast and observations reveals $FPR > TPR$; meaning, it performs worse than climatology. In fact, even a completely null prediction, where rain is never predicted anywhere, is found to be superior in this instance.

Roberts and Lean (2008) and Roberts (2008) suggest a scale-selective approach to address this issue. This considers how well the forecast captures the observed rain on increasing larger spatial scales, or *neighborhood* sizes, $n$ (Theis et al., 2005). In the example shown in Figure 5, the available neighborhood sizes would be $n = 1$ (where each neighborhood is one grid point, resulting in the original distribution of observed and

forecast rain), to $n = 2$ (where each neighborhood contains $2 \times 2$ grid points), $n = 5$ (25 grid points), and $n = 10$ (100 grid points, the entire domain). At each $n$, the fraction, $f$, of grid points within each neighborhood that contains rain is computed. For $n = 1$, each $f$ will be either exactly 0 or exactly 1. For higher values of $n$, $f$ will take a value between 0 and 1. For the example shown, at $n = 10$ the observed $f = 0.19$, while the forecast $f_F = 0.18$. For each $n$, the fraction MSE, fMSE, can be computed:

$$\text{fMSE(n)} = \frac{1}{N_x N_y} \sum_{x=1}^{N_x} \sum_{y=1}^{N_y} [f(x,y) - f_F(x,y)]^2$$

where $x$ and $y$ are the neighborhood number in the $x$ and $y$ directions and $N_X$ and $N_Y$ are the total number of neighborhoods in the $x$ and $y$ directions, respectively. Thus, for the example shown, $N_X = N_Y = 10/n$. The fraction skill score (FSS) is computed by comparing the forecast fMSE with the fMSE of a reference (or baseline) forecast, in this case the null rain forecast:

$$\text{FSS}(n) = 1 - \frac{\text{fMSE}(n)}{\text{fMSE}_{REF}(n)}$$

Figure 5d shows how the FSS varies with $n$. As discussed above, at $n = 1$, FSS is negative as the total number of false grid points (i.e., $FN + FP$) is higher for the forecast than for the null prediction. But as neighborhood size increases, FSS becomes increasingly positive, as the forecast captures the large-scale spatial structure of the observed rainfall. At $n = 10$, FSS approaches 1 as the forecast bias is very low, whereas the null prediction bias is high. The overall conclusion is that if an operator is interested in spatial scales greater than those represented by single grid points, the forecast is valuable (relative to a null forecast).

This same scale-selective approach can be adapted to the time domain for space-weather purposes. For the $V$ time series, the fMSE for neighborhood size $n$ becomes

$$\text{fMSE}(n) = \frac{n}{T} \sum_{t=1}^{T/n} [f(t) - f_F(t)]^2$$

where $f(t)$ and $f_F(t)$ are the fraction of observed and forecast hours in time bin $t$ for which $V > 500$ km/s. The first panel of Figure 6 shows the observed and forecast $f$ as a function of time for the CR 2049 solar wind speed, with events (red) and nonevents (blue) defined using $V > 500$ km/s. At this 1-hr neighborhood size, this is equivalent to the original point-by-point analysis (i.e., the same as Figure 1b) and $f$ is either exactly 0 or 1. The fMSE of the forecast is 0.370, whereas for the climatological mean, fMSE = 0.3048. Thus, for $n = 1$, FSS = $-0.21$.

The second panel of Figure 6 shows a neighborhood size of 45 hr. There are still neighborhoods with $f = 0$ and $f = 1$, but there are now also intermediate values. By eye, the agreement is still far from perfect, but the *smearing* of events in time means that there are fewer intervals that are so starkly wrong, that is, where $|f - f_F| = 1$. The fMSE of both the forecast and climatological mean have dropped (to 0.184 and 0.230, respectively) and the FSS is now weakly positive (0.2). The third and fourth panels show neighbor sizes of 105 and 210 hr, respectively. The agreement between forecast and observation has been greatly enhanced, though at these long temporal scales, a lot of information has also been lost.

Figure 7a shows how the FSS varies with $n$ and $V$ thresholds. In order to avoid aliasing between features in the $V$ time series and the neighborhood boundaries, the boundaries are slid across the time series to consider all possible neighborhood combinations for a given value of $n$. The mean FSS for a given $n$ is shown. For the CR2049, the $V$ forecast is generally most valuable for lower $V$ thresholds. However, at the very lowest threshold, $V > 350$ km/s, the forecast has little value as it fails to capture the lowest observed solar wind speeds during this interval. Across $V$ thresholds, forecast skill increases very gradually from $n = 1$ hr to $n = 20$–30 hr, when it rises more sharply. For $V > 500$ km/s, the forecast becomes more valuable than the climatological mean at neighborhood sizes of around 20 hr or longer. This time scale is roughly comparable the average timing error for HSEs (see Table 4) and indicates where the false alarm and missed events begin to cancel out, removing the double penalty effect. The fact that most $V$ thresholds converge to FSS $\approx 1$ at the maximum neighborhood size ($n = 630$ hr) shows that there is little bias in the occurrence of such events. For $V > 600$ km/s, FSS converges to values less than 1, highlighting an occurrence bias in the forecast for such an event definition (with the forecast slightly overpredicting occurrence of $V > 600$ km/s).
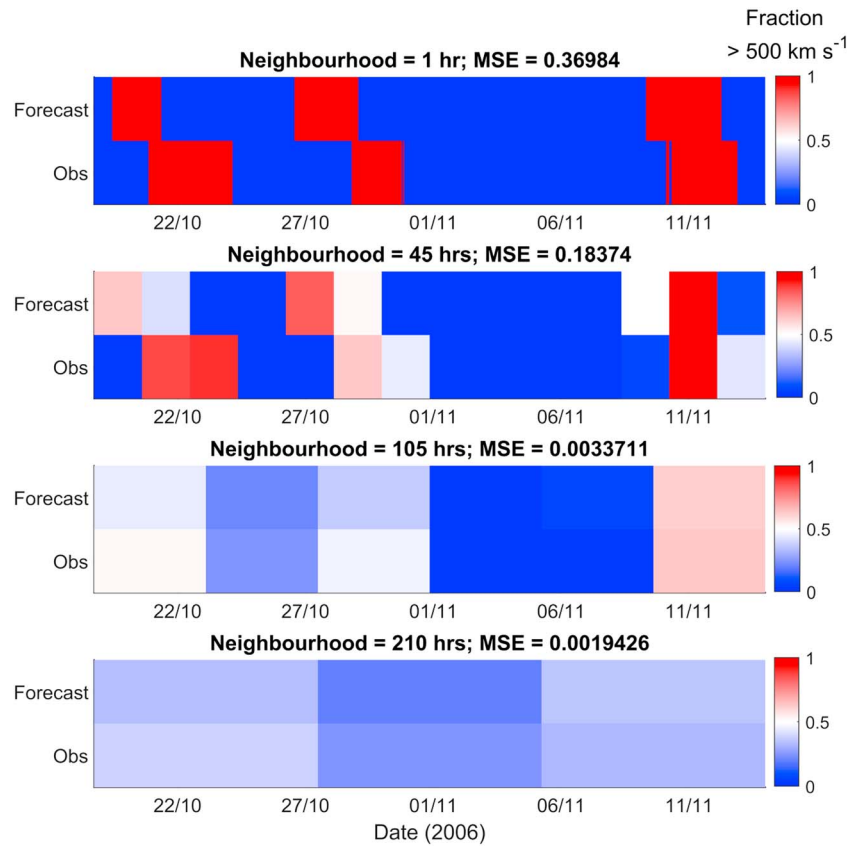
**Figure 6.** Scale-selective metrics applies the observed and forecast solar wind speed for CR2049, mid-October to mid-November 2006. The color scale shows the fraction of individual hours within a neighborhood that exceed a speed threshold of 500 km/s, from 0 (blue) to 1 (red). The first panel shows a neighbor size of 1 hr and thus is simply the threshold applied to the original observations and forecast (i.e., the same as Figure 1b). The second, third, and fourth panels show neighbor sizes of 45, 105, and 210 hr.
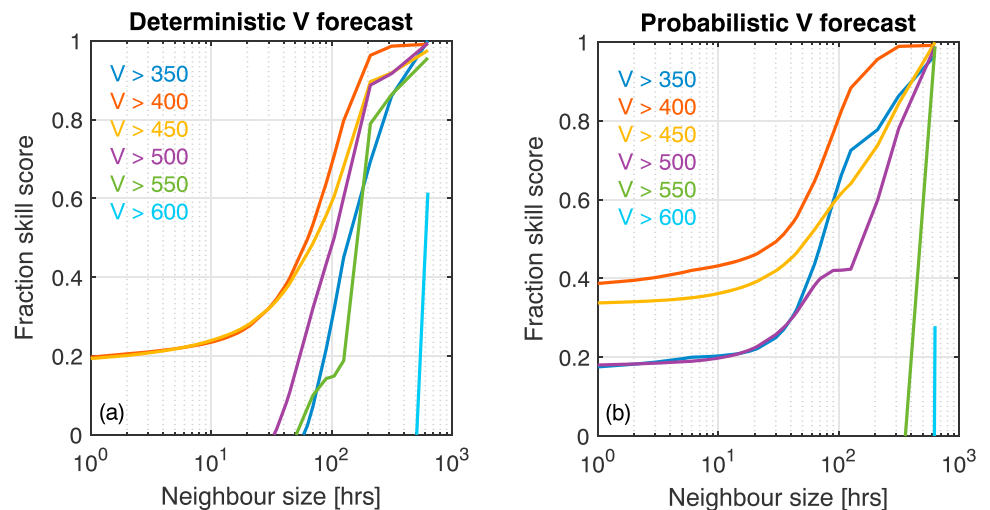


**Figure 7.** Fraction skill score for the forecast solar wind speed for mid-October to mid-November 2006 for a range of speed thresholds and neighbor sizes. The climatological mean is used as the baseline forecast. (a) The deterministic forecast of *V*; (b) the probabilistic forecast of *V*.
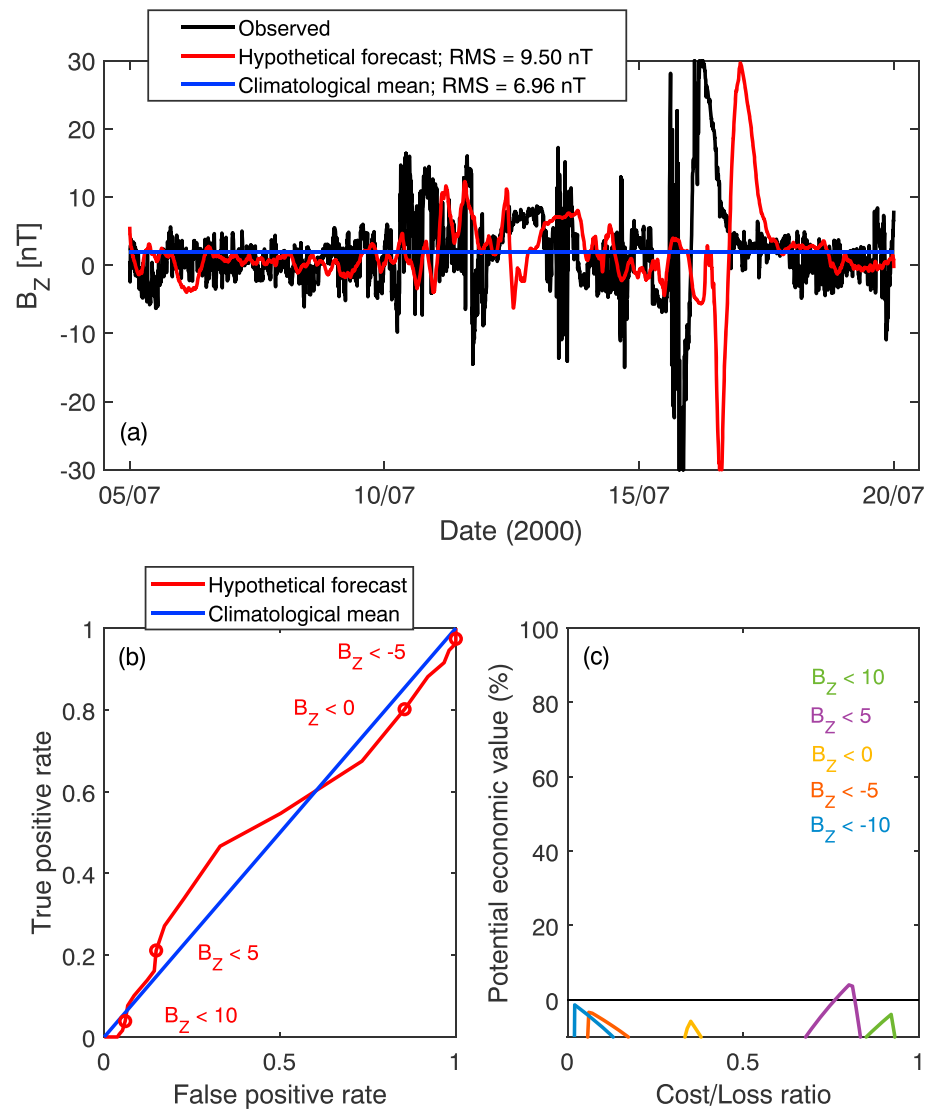
**Figure 8.** Point-by-point metrics for a hypothetical forecast of the out-of-ecliptic heliospheric magnetic field component, $B_Z$. (a) Time series of observed $B_Z$ (black) for 15 days around the *Bastille Day* coronal mass ejection of July 2000. A hypothetical forecast (black) has been produced by smoothing and shifting the observations by 18 hr. The blue line shows the mean $B_Z$ for this period. (b) The receiver operator characteristic for the forecast and climatological mean. (c) Cost-loss analysis for different $B_Z$ thresholds.

The same basic approach can also be applied to a probabilistic forecast. However, in the fMSE calculation the forecast fraction of hours above the threshold $V$, $f_F$, is replaced by $p_F$, the average probability of $V$ above the threshold in a given neighborhood. Thus, the uncertainty information is preserved, without the need to investigate different probability thresholds. Figure 7b shows how the FSS varies with neighborhood size and $V$ threshold for the probabilistic forecast of CR 2049. The general trends are similar to the deterministic forecast. But it is clear that the probabilistic forecast provides significantly higher FSS at lower $n$, particularly for $V$ thresholds below 500 km/s. This is because it intrinsically involves an (imperfect) estimate of timing error and thus some reduction of the double penalty. The rapid rise in FSS with $n$ is consequently less apparent. As the probabilistic forecast includes an increased occurrence of low speed solar wind compared to the deterministic forecast, albeit at low probability, the probabilistic forecast at $V > 350$ km/s is now valuable relative to the climatological mean. At the very highest event thresholds, $V > 550$ km/s and $V > 600$ km/s, there are insufficient events of high probability in this short interval, resulting in low FSS and a high bias for
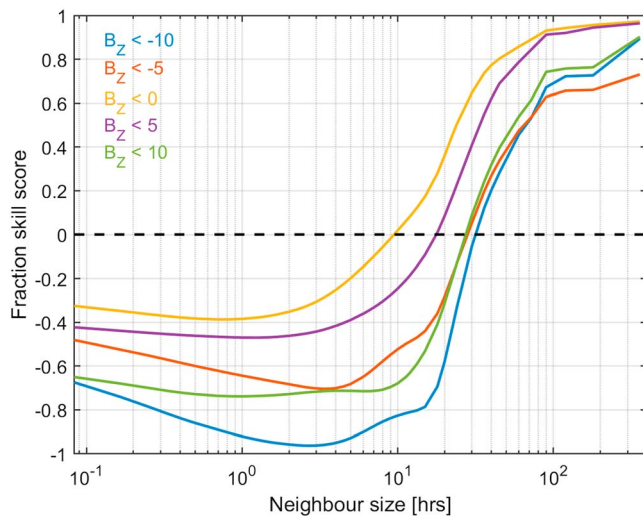
AGU
100
ADVANCING EARTH
AND SPACE SCIENCE

**Space Weather**

10.1029/2018SW002059

**Figure 9.** Fraction skill score for the $B_Z$ forecast for 15 days around the *Bastille Day* coronal mass ejection of July 2000, for a range of $B_Z$ thresholds and neighbor sizes. The climatological mean is used as the baseline forecast.

$V > 600$ km/s occurrence. Thus, a great deal of diagnostic information can be obtained from the simple FSS analysis, which is complementary to point-by-point approaches.

### 3.3. $B_Z$ Forecasts

Thus far, these issues have been illustrated exclusively with an example of a solar wind speed forecast. $V$ is one of the more accurately forecast solar wind parameters (e.g., MacNeice, 2009; Owens et al., 2008) and is always positive in value. Perhaps the solar wind parameter of greatest importance for space weather is the out-of-ecliptic component of the heliospheric magnetic field, $B_Z$, which is fundamentally less predictable than $V$ (Lockwood et al., 2016), both due to its stochastic nature and the difficulty in making remote observations of this parameter (DeForest et al., 2017). Validation of $B_Z$ forecasts is complicated by the bipolar variations associated with geoeffective coronal mass ejections, which will be particularly susceptible to double penalties. This is illustrated in Figure 8, where a hypothetical forecast of $B_Z$ for the Bastille Day interplanetary coronal mass ejection (ICME), in July 2000, has been produced by smoothing and shifting the observed time series by 18 hr, representative of current ICME forecast timing errors (Riley et al., 2018; Tucker-Hood et al., 2015). By accurately reproducing the magnitude and direction of the magnetic field within the ICME and sheath region, such a forecast is a far more accurate than any current capability (e.g., Savani et al., 2015). Yet all the point-by-point metrics, whether they be error functions (even $r_L = -0.1$) or binary metrics, show the forecast to be significantly worse than assuming that $B_Z$ is approximately 0 at all times. (The total area under the ROC curve is slightly larger than 0.5, but the sampling of $B_Z$ space is uneven. For negative $B_Z$ thresholds, the conditions of interest for space weather, the forecast lies below the $y = x$ line and hence is deemed worse than the climatological mean.)

A features-based metric, equivalent to the HSEs, would clearly work well in this instance. But the difficulty is in rigorously defining a useable definition: Time scales that would pick out a feature in the body of this ICME may exclude negative $B_Z$ intervals in other ICMEs or in the ICME sheath, which involve higher frequency variations. The more feature-agnostic approach of the scale-selective FSS is preferable. Figure 9 shows the FSS of the $B_Z$ forecast over a range of time scales (or neighborhood sizes) and for a range of $B_Z$ thresholds. For neighborhoods smaller than 10 hr, the forecast is worse than assuming $B_Z \approx 0$ at all times, as the point-by-point analyses concluded. But as the time scale is increased to around 10–30 hr, the forecast is shown to be skilful relative to the climatological mean, as one would likely conclude by eye.

## 4. Summary

This study briefly reviewed some of the commonly used metrics for space-weather forecast and model validation. Simple error functions, like RMS and MAE, are the mainstay of forecast validation. They compare forecasts and observations on a strictly point-by-point basis. They are undoubtedly a valuable tool for forecast comparison. But there are limitations in their use as forecast diagnostics and they can, in some circumstances, give misleading results about the value or usefulness of a forecast. In particular, by treating each time point entirely independently, timing uncertainties are not explicitly accounted for. Thus, when timing errors are present in the forecasts, they can be hit with *double penalties*, for both missing the observed event and issuing a false alarm. While there is nothing inherently wrong with this form of assessment, it can systematically favor overly conservative forecasts, which may not be beneficial. Binary metrics, in which a forecast is converted to series of yes/no predictions, reduce the emphasis on event magnitude and hence somewhat reduce the effect of double penalties for timing errors. These kinds of approaches are summarized by the ROC and the Cost-Loss analysis. These can provide useful insight into the operational circumstances in which a particular forecast is most useful (e.g., in settings where false alarms are not a major issue).

A neat, simple, solution to the double penalty problem is for all forecasts to include an accurate assessment of uncertainty. As shown here, even relatively coarse estimates of uncertainty can add value to existing

forecasts. But there are a number of reasons why this is not always practical. Instead, this study has advocated a more pragmatic solution of time-window metrics alongside the more traditional point-by-point approaches. Defining discrete, extended features in the observed and forecast time series allows direct comparison of their timing and magnitude. This is a powerful analytical tool but requires a rigorous a prior definition of an event, which is robust to event-to-event variability, and between observations and forecast. An alternative is to use a scale-selective approach, wherein agreement between forecast and observation is considered at a range of time scales. As the time scales become increasingly coarse, false alarms and missed events increasingly cancel out, reducing the double penalty effect. This allows an assessment of the time scales at which the forecast provides an acceptable level of accuracy.

Part of the job of a metric is to summarize a complex parameter space: different parameter and forecast probability thresholds, different spatial and temporal scales, and different operational sensitivities. The examples shown here consider only the simplest case of solar wind time series. Validation in other domains of the space-weather system also has to deal with intrinsically higher dimensionality. For example, in radiation belt forecasting, in addition to temporal variations, there is a great deal of spatially variability in all three directions (radially from the Earth, and in geomagnetic latitude and magnetic local time), as well as in particle energy space (e.g., Shprits et al., 2015). Often this dimensionality is reduced by averaging over particle drift and bounce motions, but the situation nevertheless remains more complex than a single time series. But the same fundamental issues are still present, just in a more multifarious way.

Finally, it is worth reiterating that these more sophisticated methods of forecast and model validation are intended to complement, not replace, existing metrics. Error functions should undoubtedly continue to be a standard space-weather metric. In additional to continuing the legacy, they are simple to implement and interpret, as well as enabling easy intercomparison of different forecasts and models. But a more diagnostic picture of *why* a forecast is accurate or fails is invaluable too.

**References**

DeForest, C., de Koning, C., & Elliott, H. (2017). 3D polarized imaging of coronal mass ejections: Chirality of a CME. *The Astrophysical Journal*, *850*(2), 130. https://doi.org/10.3847/1538-4357/aa94ca

Dungey, J. W. (1961). Interplanetary magnetic field and the auroral zones. *Physical Review Letters*, *6*(2), 47–48. https://doi.org/10.1103/PhysRevLett.6.47

Ebert, E. E., & Gallus, W. A. Jr. (2009). Toward better understanding of the contiguous rain area (CRA) method for spatial forecast verification. *Weather and Forecasting*, *24*(5), 1401–1415. https://doi.org/10.1175/2009waf2222252.1

Finley, J. P. (1884). Tornado predictions. *American Meteorological Journal*, *1*, 85–88.

Jian, L. K., MacNeice, P. J., Mays, M. L., Taktakishvili, A., Odstrcil, D., Jackson, B., Yu, H.-S., et al. (2016). Validation for global solar wind prediction using Ulysses comparison: Multiple coronal and heliospheric models installed at the Community Coordinated Modeling Center. *Space Weather*, *14*, 592–611. https://doi.org/10.1002/2016SW001435

King, J. H., & Papitashvili, N. E. (2005). Solar wind spatial scales in and comparisons of hourly Wind and ACE plasma and magnetic field data. *Journal of Geophysical Research*, *110*, A02104. https://doi.org/10.1029/2004JA010649

Linker, J., Mikic, Z., Biesecker, D. A., Forsyth, R. J., Gibson, W. E., Lazarus, A. J., et al. (1999). Magnetohydrodynamic modeling of the solar corona during whole Sun month. *Journal of Geophysical Research*, *104*(A5), 9809–9830. https://doi.org/10.1029/1998JA900159

Lockwood, M., Owens, M. J., Barnard, L. A., Bentley, S., Scott, C. J., & Watt, C. E. (2016). On the origins and timescales of geoeffective IMF. *Space Weather*, *14*, 406–432. https://doi.org/10.1002/2016SW001375

MacNeice, P. (2009). Validation of community models: Identifying events in space weather model timelines. *Space Weather*, *7*, S06004. https://doi.org/10.1029/2009sw000463

MacNeice, P., Jian, L., Antiochos, S. K., Arge, C. N., Bussy-Virat, C. D., DeRosa, M. L., et al. (2018). Assessing the quality of models of the ambient solar wind. *Space Weather*. https://doi.org/10.1029/2018SW002040

Mason, I. (1982). A model for assessment of weather forecasts. *Australian Meteorological Magazine*, *30*(4), 291–303.

Mason, S. J., & Graham, N. E. (2002). Areas beneath the relative operating characteristics (ROC) and relative operating levels (ROL) curves: Statistical significance and interpretation. *Quarterly Journal of the Royal Meteorological Society*, *128*(584), 2145–2166. https://doi.org/10.1256/003590002320603584

Mays, M. L., Taktakishvili, A., Pulkkinen, A., MacNeice, P. J., Rastätter, L., Odstrcil, D., et al. (2015). Ensemble modeling of CMEs using the WSA–ENLIL+Cone model. *Solar Physics*, *290*(6), 1775–1814. https://doi.org/10.1007/s11207-015-0692-1

McCloskey, A. E., Gallagher, P. T., & Bloomfield, D. S. (2018). Flare forecasting using the evolution of McIntosh sunspot classifications. *Journal of Space Weather and Space Climate*, *8*, A34. https://doi.org/10.1051/swsc/2018022

Murphy, A. H. (1977). The value of climatological, categorical and probabilistic forecasts in the cost-loss ratio situation. *Monthly Weather Review*, *105*(7), 803–816.

Murphy, A. H. (1993). What is a good forecast? An essay on the nature of goodness in weather forecasting. *Weather Forecasting*, *8*(2), 281–293.

Murphy, A. H. (1996). The Finley affair: A signal event in the history of forecast verification. *Weather and Forecasting*, *11*(1), 3–20. https://doi.org/10.1175/1520-0434(1996)011<0003:tfaase>2.0.co;2

Murray, S. A. (2018). The importance of ensemble techniques for operational space weather forecasting. *Space Weather*, *16*, 777–783. https://doi.org/10.1029/2018SW001861

Murray, S. A., Bingham, S., Sharpe, M., & Jackson, D. R. (2017). Flare forecasting at the Met Office Space Weather Operations Centre. *Space Weather*, *15*, 577–588. https://doi.org/10.1002/2016SW001579

Owens, M., & Riley, P. (2017). Probabilistic solar wind forecasting using large ensembles of near-Sun conditions with a simple one-dimensional "upwind" scheme. *Space Weather*, *15*, 1461–1474. https://doi.org/10.1002/2017SW001679

Owens, M. J., Arge, C. N., Spence, H. E., & Pembroke, A. (2005). An event-based approach to validating solar wind speed predictions: High speed enhancements in the Wang-Sheeley-Arge model. *Journal of Geophysical Research*, *110*, A12105. https://doi.org/10.1029/2005JA011343

Owens, M. J., Challen, R., Methven, J., Henley, E., & Jackson, D. R. (2013). A 27 day persistence model of near-Earth solar wind conditions: A long lead-time forecast and a benchmark for dynamical models. *Space Weather*, *11*, 225–236. https://doi.org/10.1002/swe.20040

Owens, M. J., Horbury, T. S., Wicks, R. T., McGregor, S. L., Savani, N. P., & Xiong, M. (2014). Ensemble downscaling in coupled solar wind-magnetosphere modeling for space weather forecasting. *Space Weather*, *12*, 395–405. https://doi.org/10.1002/2014SW001064

Owens, M. J., Riley, P., & Horbury, T. S. (2017). Probabilistic solar wind and geomagnetic forecasting using an analogue ensemble or "similar day" approach. *Solar Physics*, *292*(5), 69. https://doi.org/10.1007/s11207-017-1090-7

Owens, M. J., Spence, H. E., McGregor, S., Hughes, W. J., Quinn, J. M., Arge, C. N., et al. (2008). Metrics for solar wind prediction models: Comparison of empirical, hybrid and physics-based schemes with 8-years of L1 observations. *Space Weather*, *6*, S08001. https://doi.org/10.1029/2007SW000380

Peterson, W., Birdsall, T., & Fox, W. (1954). The theory of signal detectability. *Transactions of the IRE Professional Group on Information Theory*, *4*(4), 171–212. https://doi.org/10.1109/TIT.1954.1057460

Reiss, M. A., Temmer, M., Veronig, A. M., Nikolic, L., Vennerstrom, S., Schöngassner, F., & Hofmeister, S. J. (2016). Verification of high-speed solar wind stream forecasts using operational solar wind models. *Space Weather*, *14*, 495–510. https://doi.org/10.1002/2016SW001390

Richardson, D. S. (2000). Skill and relative economic value of the ECMWF ensemble prediction system. *Quarterly Journal of the Royal Meteorological Society*, *126*(563), 649–667.

Riley, P., Linker, J. A., Lionello, R., & Mikic, Z. (2012). Corotating interaction regions during the recent solar minimum: The power and limitations of global MHD modeling. *Journal of Atmospheric and Solar-Terrestrial Physics*, *83*, 1–10. https://doi.org/10.1016/j.jastp.2011.12.013

Riley, P., Linker, J. A., & Mikić, Z. (2013). On the application of ensemble modeling techniques to improve ambient solar wind models. *Journal of Geophysical Research: Space Physics*, *118*, 600–607. https://doi.org/10.1002/jgra.50156

Riley, P., Linker, J. A., Mikič, Z., Zank, G. P., Borovsky, J., Bruno, R., et al. (2013), Ensemble modeling of the ambient solar wind. Paper presented at AIP conference proceedings. AIP.

Riley, P., Mays, M. L., Andries, J., Amerstorfer, T., Biesecker, D., Delouille, V., et al. (2018). Forecasting the arrival time of coronal mass ejections: Analysis of the CCMC CME scoreboard. *Space Weather*, *16*, 1245–1260. https://doi.org/10.1029/2018SW001962

Roberts, N. (2008). Assessing the spatial and temporal variation in the skill of precipitation forecasts from an NWP model. *Meteorological Applications*, *15*(1), 163–169. https://doi.org/10.1002/met.57

Roberts, N. M., & Lean, H. W. (2008). Scale-selective verification of rainfall accumulations from high-resolution forecasts of convective events. *Monthly Weather Review*, *136*(1), 78–97. https://doi.org/10.1175/2007mwr2123.1

Savani, N. P., Vourlidas, A., Szabo, A., Mays, M. L., Richardson, I. G., Thompson, B. J., et al. (2015). Predicting the magnetic vectors within coronal mass ejections arriving at Earth: 1. Initial architecture. *Space Weather*, *13*, 374–385. https://doi.org/10.1002/2015SW001171

Shprits, Y. Y., Kellerman, A. C., Drozdov, A. Y., Spence, H. E., Reeves, G. D., & Baker, D. N. (2015). Combined convective and diffusive simulations: VERB-4D comparison with 17 March 2013 Van Allen Probes observations. *Geophysical Research Letters*, *42*, 9600–9608. https://doi.org/10.1002/2015GL065230

Siscoe, G. (2007). Space weather forecasting historically viewed through the lens of meteorology. In V. Bothmer, & I. A. Daglis (Eds.), *Space weather—Physics and effects*, (chap. 2, pp. 5–30). Berlin: Springer. https://doi.org/10.1007/978-3-540-34578-7-2

Slingo, J., & Palmer, T. (2011). Uncertainty in weather and climate prediction. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, *369*(1956), 4751–4767. https://doi.org/10.1098/rsta.2011.0161

Spence, H., Baker, D., Burns, A., Guild, T., Huang, C.-L., Siscoe, G., & Weigel, R. (2004). Center for integrated space weather modeling metrics plan and initial model validation results. *Journal of Atmospheric and Solar - Terrestrial Physics*, *66*, 1491–1498.

Taylor, K. E. (2001). Summarizing multiple aspects of model performance in a single diagram. *Journal of Geophysical Research*, *106*(D7), 7183–7192. https://doi.org/10.1029/2000JD900719

Theis, S. E., Hense, A., & Damrath, U. (2005). Probabilistic precipitation forecasts from a deterministic model: A pragmatic approach. *Meteorological Applications*, *12*(3), 257–268. https://doi.org/10.1017/S1350482705001763

Thornes, J. E., & Stephenson, D. B. (2001). How to judge the quality and value of weather forecast products. *Meteorological Applications*, *8*(3), 307–314. https://doi.org/10.1017/S1350482701003061

Tucker-Hood, K., Scott, C., Owens, M., Jackson, D., Barnard, L., Davies, J. A., et al. (2015). Validation of a priori CME arrival predictions made using real-time heliospheric imager observations. *Space Weather*, *13*, 35–48. https://doi.org/10.1002/2014SW001106