

New operational measure to assess extreme events using site-specific climatology

Michael Sharpe, Clare Bysouth, Philip Gill. Operational Verification Systems and Products,
Met Office UK

IVMW-O, 12th November 2020

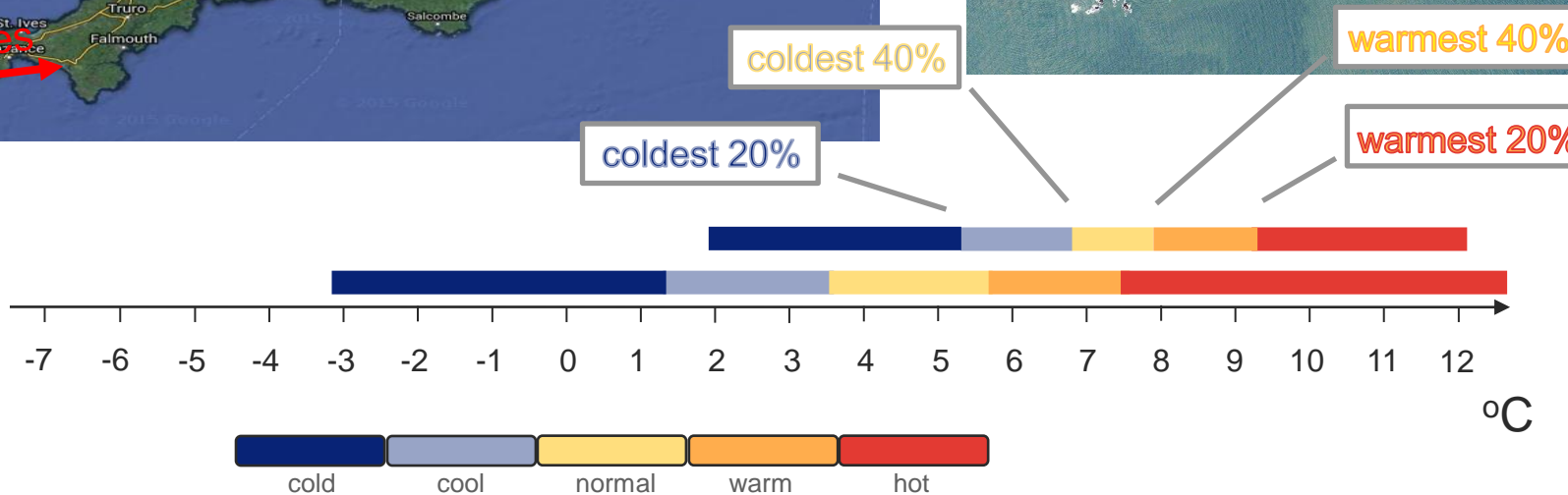
© 2015 Google



© Google



Scilly
Culdrose



Relative-extremes:

define extreme events relative to site-specific climatology

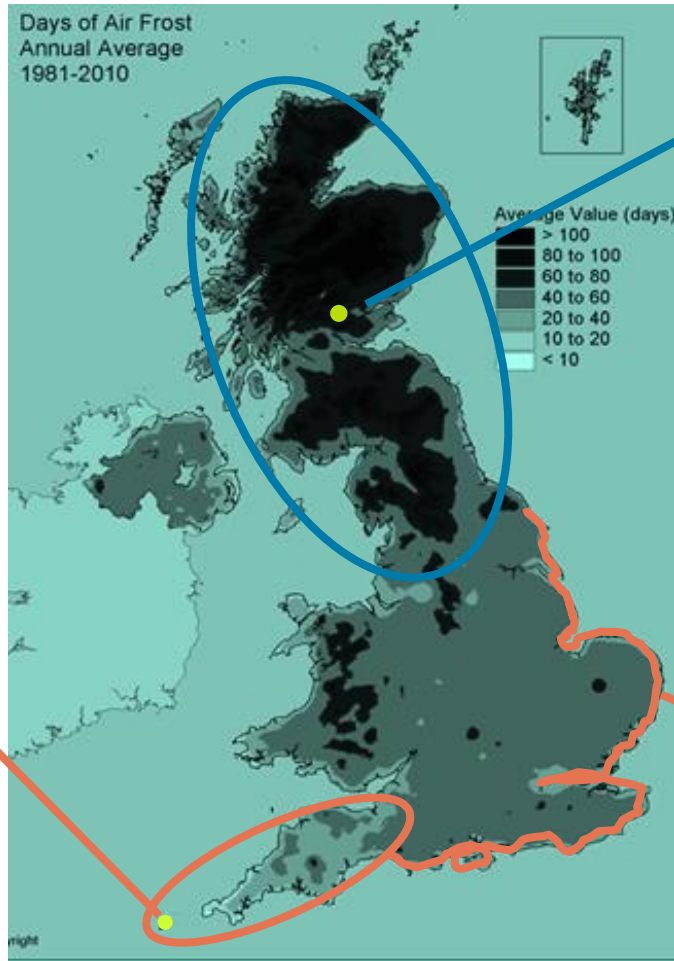
Absolute-extremes:

occur more frequently in some places, and less frequently in others...

...whereas frosts on Scilly are big news...

- burst pipes
- dead plants
- ill pets

...are all likely



Frosts in Gleneagles are common...
...so they have little impact...

...instead it tends to be dominated by the skill in cold areas.

... performance is virtually unaffected by the skill in warmer areas...

Why am I interested in all this...

Public Weather Service requested a

'Key Performance Indicator'

for how well we forecast extreme events.

They want this measure to:

- be easy to explain
- be simple to communicate
- be robust (i.e. no sudden changes)
- directly compare deterministic & probabilistic forecast performance



Any measure of extremes should use local climatology to define an event

Advantages of using local climatology to define extreme events...

- a) *Base rate is the same everywhere*
- b) *Simple aggregation gives all-site performance*
- c) *Similar impact everywhere*
- d) *Event defined in terms of a return period*



Disadvantages...

- a) *Can't verify where climatology is unavailable*
- b) *Assimilation affects model climatology (use observed)*
- c) *multiple thresholds → technically more challenging*
- d) *Different to the norm!*





Max-Planck-Institut für Bildungsforschung
Max Planck Institute for Human Development

Extreme Verification ... relatively speaking

Michael Sharpe, Clare Bysouth, Becky Stretton

Summary

- Interest in weather impacts has inevitably led to a desire to examine the ability of weather prediction models to forecast extreme events (Vignati et al., 2015). However, this is particularly difficult because extreme events are rare and consequently many verification metrics are unreliable for their analysis.
- Impact is often related to the chance of occurrence at a particular location therefore it is more appropriate to use relative extreme event thresholds derived from observed site-specific climatology.
- Choosing the 98.5th percentile from each climatological distribution at each site ensures the base rate is the same everywhere so performance is not dominated by locations where the event is more common.
- A 21-month period of Met Office site-specific probabilistic forecasts were analysed and compared with the climate at each site as the reference forecast.
- Individual site scores at a forecast range of 24h show some degree of correlation.
- Two of the three metrics used provide evidence of forecast skill at all forecast ranges.
- Extremes were over-forecast and only predicted with low probability at medium range.

Extreme-event thresholds

The map shows the 98.5th percentile of 24-hour rainfall accumulation at all sites in the UK where 30 consecutive years of climatology are available between 1983 and 2012; this percentile corresponds to the threshold which is exceeded an average of 4 times per year.

Verification methodologies

The Met Office post-processed probabilistic site-specific forecast is a blend of various ensemble models, the output from the model is expressed using 15 quantiles.

- The Symmetric Extremal Dependency Index (SEDI) (Farrs and Stephenson, 2011) is a diagnostic performance measure that is valid when the frequency bias (FB) = 1.
 - SEDI can be used to verify a probabilistic forecast if its quantiles are viewed as a means of calibration; thereby only verifying the quantile for which (FB-1) is increased and FB=1.
- The Continuous Ranked Probability Skill Score (CRPSS) (Murphy, 1971) is a probabilistic performance measure which considers every quantile of the forecast.
 - CRPSS has been evaluated whenever the event was either forecast or observed; the deterministic parallel is the Equitable Threat Score (the proportion correct, adjusted for climatology, given that an event was either observed or forecast).
- The threshold weighted CRPSS (wCRPSS) (Goring and Rayner, 2011) only evaluates the quantiles of the forecast that exceed the event threshold.
 - To ensure wCRPSS is always evaluated the 5th and 100th percentiles of every forecast are set to their site-specific climatological minimum and maximum value.

Results

Figure 1 shows that on day 1 many different probability forecasts appear more reliable. On day 4 events are over-forecast and only low probabilities (from clear bias) are forecast.

Figure 4 shows decreasing skill with increasing forecast range and the bias of the bootstrapped 95% confidence intervals noticeably increases for SEDI but (less noticeably) decreases for wCRPSS.

- SEDI and CRPSS⁹⁵ only forecast skill to T+120
- wCRPSS provides little evidence of forecast skill (relative to site-specific climatology) after T+24

Figure 1. 24h post-processed probabilistic site-specific distribution of 24-hour rainfall accumulation at sites in the UK.

Figure 2. SEDI reveals more discriminatory skill on day 1.

Figure 3. Figure 3 shows that on day 1 many different probability forecasts appear more reliable. On day 4 events are over-forecast and only low probabilities (from clear bias) are forecast.

Figure 4. Figure 4 shows decreasing skill with increasing forecast range and the bias of the bootstrapped 95% confidence intervals noticeably increases for SEDI but (less noticeably) decreases for wCRPSS.

- SEDI and CRPSS⁹⁵ only forecast skill to T+120
- wCRPSS provides little evidence of forecast skill (relative to site-specific climatology) after T+24

Figure 5. SEDI may only be evaluated at sites where a forecast quantile gives FB = 1.

Figure 6. CRPSS⁹⁵ and wCRPSS can only be evaluated when an event is observed or forecast so it is more suited to slightly less extreme events.

Conclusions

- SEDI may only be evaluated at sites where a forecast quantile gives FB = 1.
- CRPSS⁹⁵ can only be evaluated when an event is observed or forecast so it is more suited to slightly less extreme events.
- wCRPSS is dependent on an arbitrary threshold weighting function and assumed 5th and 100th forecast percentiles to ensure it can always be evaluated.

References

Farrs CAT and Stephenson DR, 2011. The skill of probabilistic weather forecasts: A review. *International Journal of Numerical and Analytical Weather Modelling*, 2, 1-10.

Stephenson DR, 2012. *Weather Forecasting: An Introduction to the Theory and Practice of Forecasting*, 2nd edn. Wiley, 304 pp.

Stephenson DR, 2013. *Weather Forecasting: An Introduction to the Theory and Practice of Forecasting*, 2nd edn. Wiley, 304 pp.

Stephenson DR, 2014. *Weather Forecasting: An Introduction to the Theory and Practice of Forecasting*, 2nd edn. Wiley, 304 pp.

Stephenson DR, 2015. *Weather Forecasting: An Introduction to the Theory and Practice of Forecasting*, 2nd edn. Wiley, 304 pp.

How well do Met Office post-processed site-specific probabilistic forecasts predict relative-extreme events?

Michael A. Sharpe,* Clare E. Bynott and Rebecca L. Stratten
Weather Science, Met Office, Exeter, UK

ABSTRACT: The Met Office routinely generates post-processed forecasts at sites throughout the United Kingdom; both deterministic and probabilistic products exist and deterministic products are publicly available online. In recent years, provision of weather information has focused upon the impact of events; impact is often related to the frequency of occurrence of an event at a site, which is determined by its climatology. The ability with which a site-specific forecast predicts relative-extremes may be investigated by examining the skill with which these events (defined in terms of a percentile chosen from the climatology at each site) are predicted. The Heaviside, deterministic, website forecast is less likely to forecast extreme events; therefore, the probabilistic forecast product (which does not currently appear on the Met Office website) was evaluated for its ability to predict heavy rainfall (RF_{95}), maximum summer day time temperature (T_{max}), minimum winter night time temperature (T_{min}) and strong winds (WS_{95}) over a 21 month period between December 2013 and August 2015. In this end, four methods of verification are considered: the Symmetric Extremal Dependency Index (SEDI), a threshold-weighted version of the continuous ranked probability skill score (CRPSS) and a conditional version of the CRPSS together with an analysis of the discrimination and reliability. Each method indicates forecast skill, with T_{min} and RF_{95} identified as the most and least skillful respectively and WS_{95} identified as the most reliable. Site-specific values of both versions of the CRPSS appear relatively well correlated and these scores also show correlation with SEDI for WS_{95} .

KEY WORDS extreme; forecasting; verification

Received 14 September 2016; Revised 20 January 2017; Accepted 4 February 2017

1. Introduction

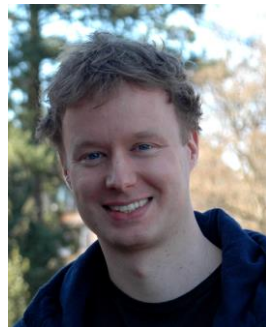
In recent years the general public and meteorological community have become increasingly interested in extreme weather events because of their impact on society and/or infrastructure. This has inevitably led to a desire to examine the ability of weather prediction models to forecast extreme. However, such an examination is particularly difficult because, by their very nature, extreme events are rare.

Before proceeding any further it is necessary to define an extreme event carefully. At most UK locations a rainfall accumulation of 50 mm in 24 h would be considered an extreme event and in some instances this type of 'absolute' definition is appropriate (e.g. road, rail or air transport applications); however, it is often difficult to assess the performance accurately using this type of absolute extreme event threshold because it will be exceeded frequently at some locations but almost never exceeded at others. Consequently, the overall performance can be dominated by locations where the event occurs most frequently and usually has less impact. The impact of an event is often related to its chance of occurrence and, although absolute extreme events are important, when it comes to relating weather events to impacts it is often more appropriate to adopt a relative definition, where event thresholds are derived from the distribution formed by the climate at each location. This may be achieved by choosing the same percentile from the cumulative

distribution function of the climatology (CDF) at each site; consequently, the event base rate is the same at every location, so the overall performance is not dominated by locations where the event is more common. Another advantage is the ability to express events in terms of their expected frequency of occurrence (e.g. the wettest day of the year); such definitions are arguably more meaningful than thresholds such as 50 mm in 24 h. A similar approach was independently taken at the European Centre for Medium-Range Weather Forecasts by Magnusson *et al.* (2014); these authors used the 50th percentile of the observed climatology to define event thresholds which were used to assess these different forecasts using the Symmetric Extremal Dependency Index (SEDI) (Feroz and Stephenson, 2011) and relative economic value (Richardson, 2005). An alternative method of analysis which Magnusson *et al.* (2014) suggest (do not implement) is what they refer to as a 'modified version of the continuous ranked probability score (CRPS), where a function is applied to give more weight to extreme events'. It is likely that this is a reference to the threshold-weighted continuous ranked probability score devised by Gneiting and Ranjan (2011), a methodology which examines the skill of any forecast quantile that predicts an extreme event. One method not considered by the present study is relative economic value because it has been considered in some detail by Magnusson *et al.* (2014).

In the present study, the relative-extreme event definition is adopted to assess the skill with which the probabilistic version of Met Office post-processed site-specific forecasts predict extreme 24 h rainfall accumulation (RF_{95}), maximum summer day time temperature (T_{max}), minimum winter night time temperature (T_{min}) and hourly wind speeds (WS_{95}). Site-specific climatology generation and analysis is discussed in Section 2;

*Correspondence: M. A. Sharpe, Met Office, Fitzroy Road, Exeter, EX1 3PB, UK. E-mail: michael.sharpe@metoffice.gov.uk
This article is published with the permission of the Controller of HMSO and the Queen's Printer for Scotland.



We've tried ...

1. SEDI (deterministic method)

- Consider each (15) forecast quantile separately
- Choose one based on its frequency bias
- Use it as the 'deterministic' forecast

Contrived

2. CRPS (probabilistic method)

- Compare the CDF of the forecast with the CDF of the observation (a Heaviside function)
- Restrict CRPS to only examine extreme events by conditioning or thresholding
- Integrate numerically over forecast percentiles

CRPS^{of} is improper

$$\text{for maximum temperature } O(t) = \begin{cases} 0 & \text{for all } t < ob \\ 1 & \text{for all } t \geq ob \end{cases}$$

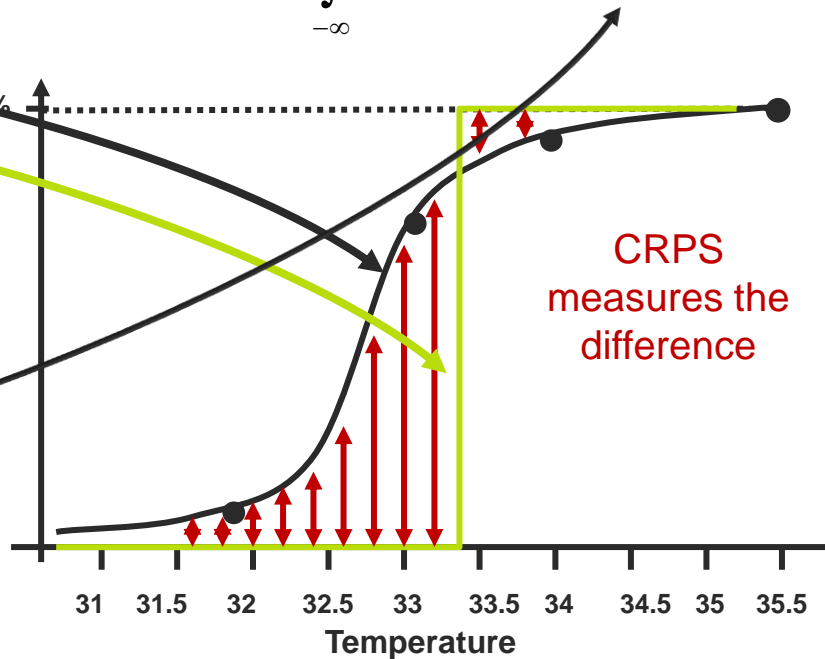
Continuous ranked probability score
 Calculates the difference between
 the forecast CDF
 and the observed CDF

$$tw\ CRPS = \int_{-\infty}^{\infty} (F(t) - O(t))^2 dt$$

We are only interested in extremes...
 ... but the CRPS evaluates everything

Gneiting & Ranjan, 2011:

- Threshold weighted CRPS
- Choose a threshold weighting function that excludes non-extremes



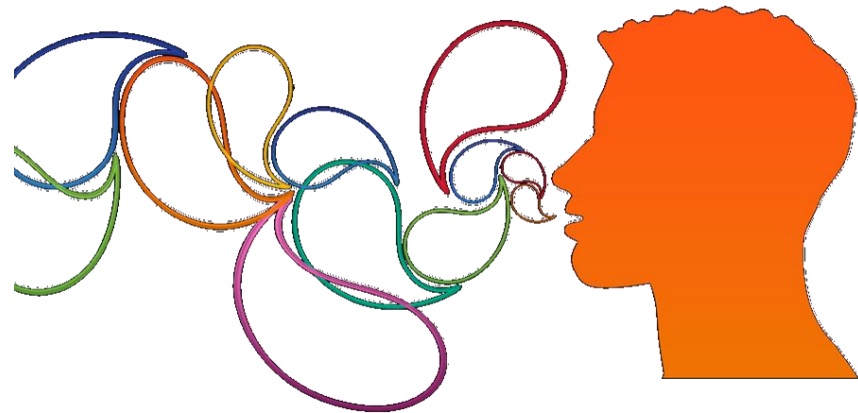
We tried to frame the measure in a way that can be communicated.

So we asked a common question...

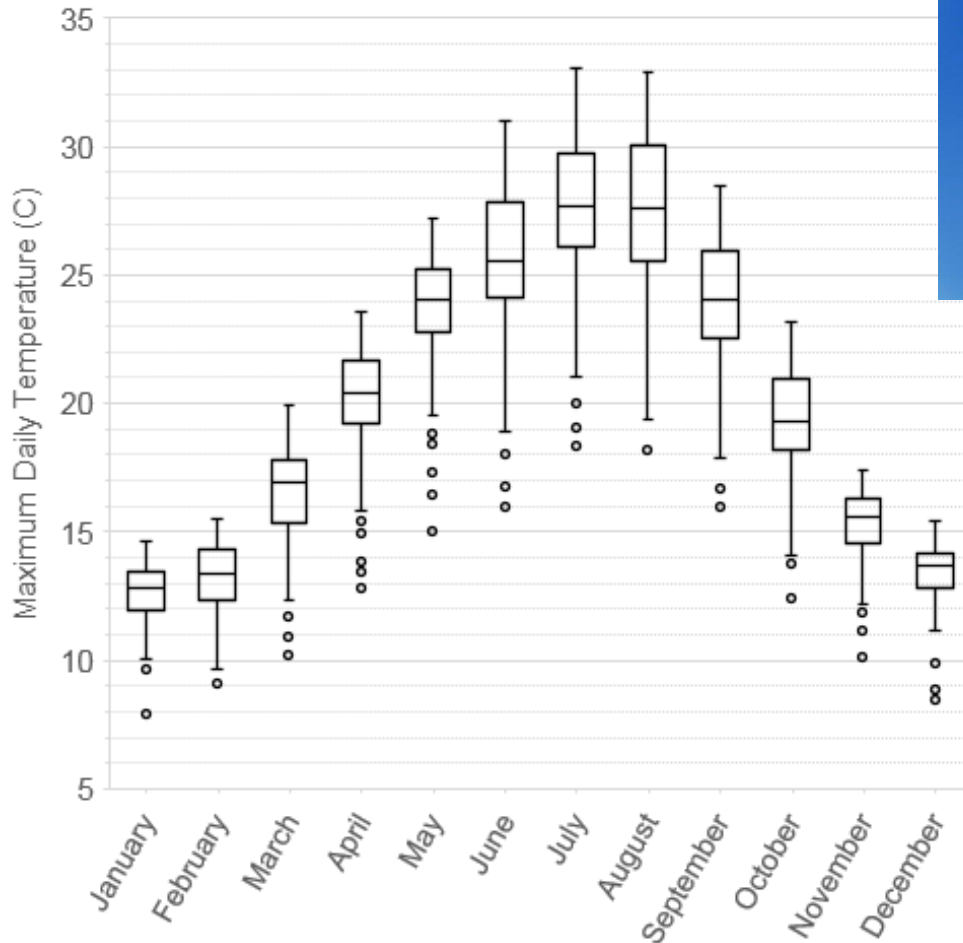
Is today the coldest / warmest / windiest day in 3-years?

To answer that we've looked at:

- all observations at every site from 1987 to 2016
- the most extreme value in 3 randomly chosen years
- **maximum temperature**, minimum temperature & wind speed



The hottest day you should expect every 3 years at each UK site*

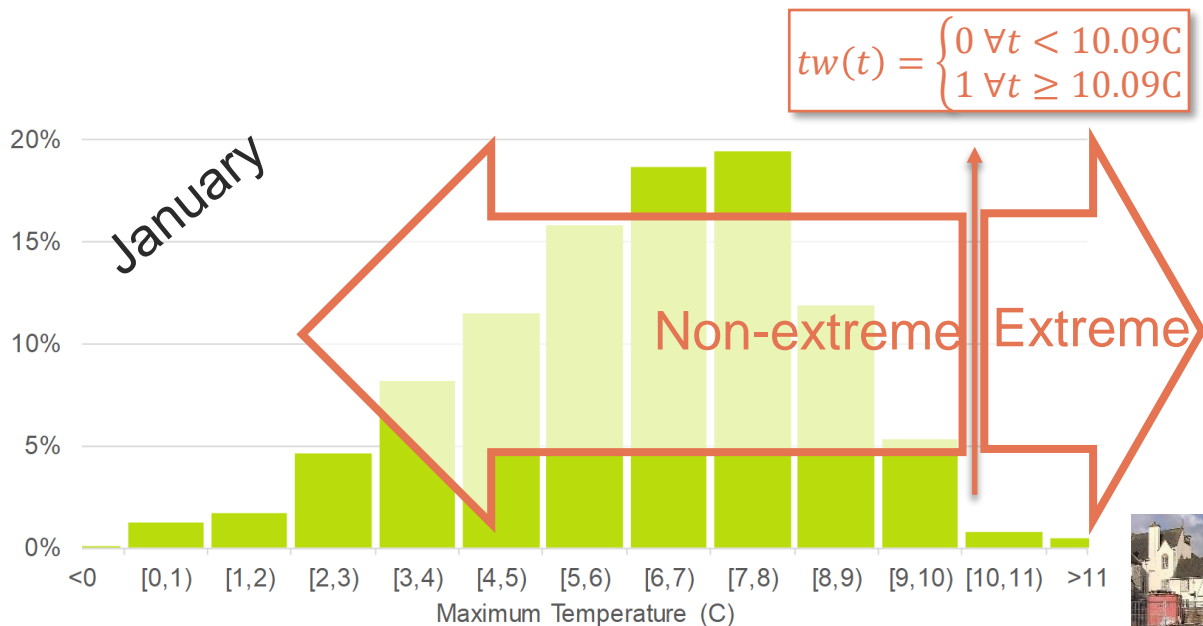


*based on 1987-2016

New operational measure:

- most extreme 3 year event (98.89th percentile)
- 30-year site-specific monthly climatology

$$twCRPS = \int_{10.09}^{\infty} (F(t) - O(t))^2 dt$$



New operational measure:

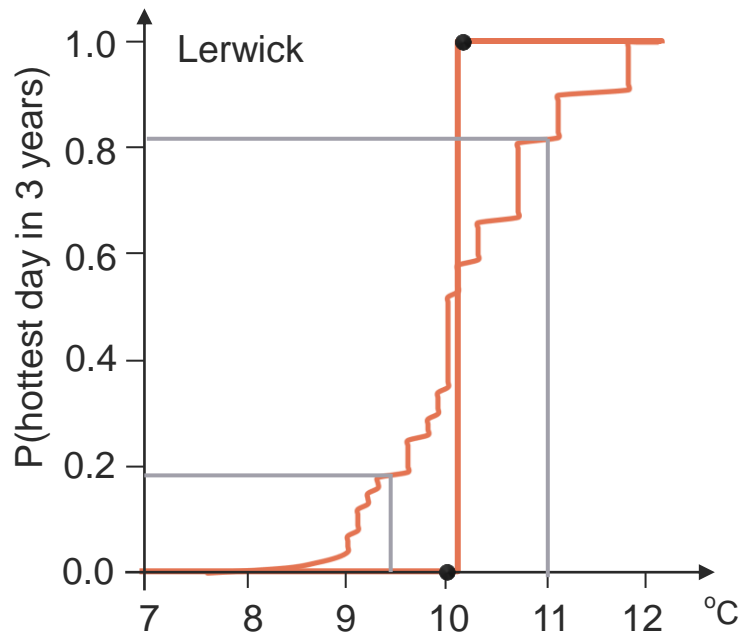
- most extreme 3 year event (98.89th percentile)
- 30-year site-specific monthly climatology
- Monthly sampling gives the probability that t is the hottest day in 3 years

So, for example

- 9.4°C has a 19% chance and...
- 11°C has a 81% chance...

...of being the hottest January day in Lerwick in 3 years

$$twCRPS = \int_{10.09}^{\infty} (F(t) - O(t))^2 dt$$



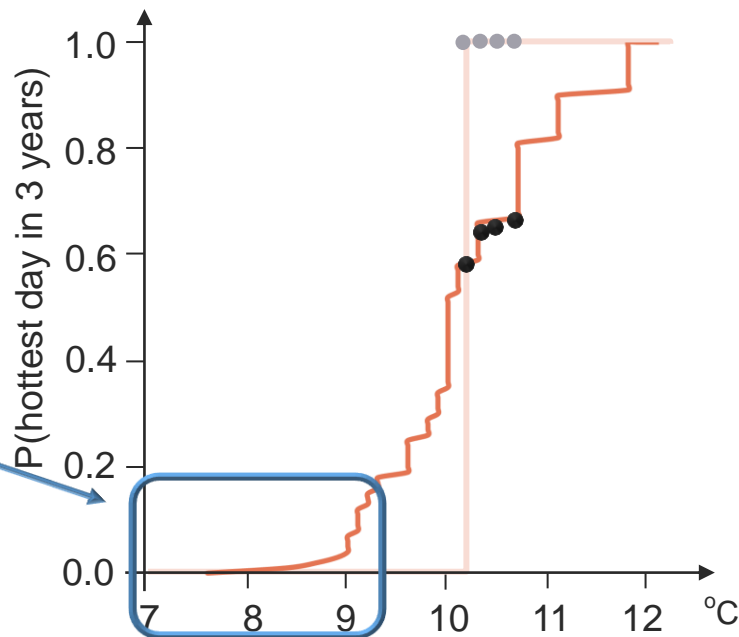
New operational measure:

- most extreme 3 year event (98.89th percentile)
- 30-year site-specific monthly climatology
- Monthly sampling gives the probability that t is the hottest day in 3 years

Advantages:

- Has a real meaning to the public
- Can examine more extreme thresholds due to tail
- Score is less sensitive to small distributional changes
- So annual climate updates unlikely to affect robustness

$$twCRPS = \int_{-\infty}^{\infty} (F(t) - O(t))^2 tw(t) dt$$



Probabilistic measure

$$CRPS = \int_{-\infty}^{\infty} (F(t) - O(t))^2 dt$$

$$twCRPS = \int_{-\infty}^{\infty} (F(t) - O(t))^2 tw(t) dt$$

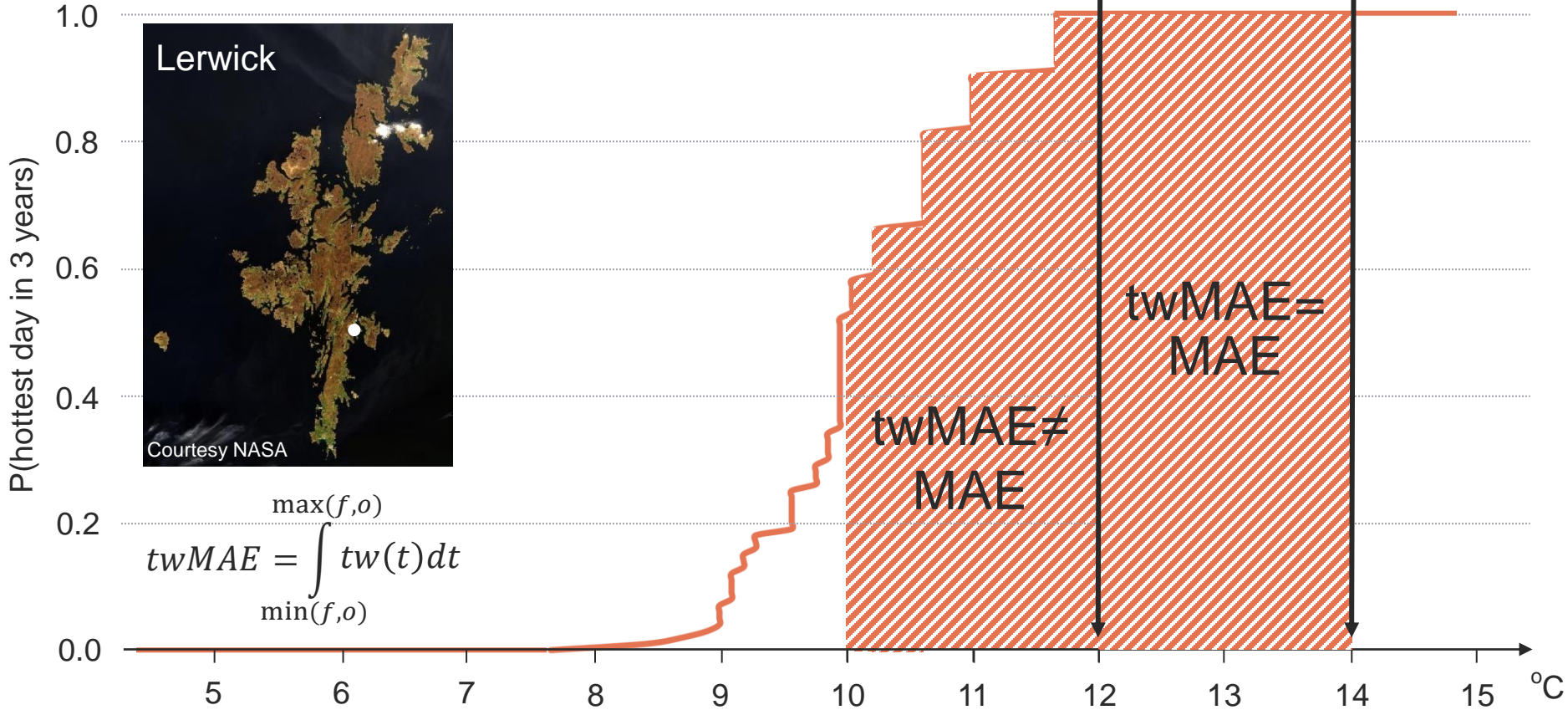
Score range: 0 to ∞

Deterministic equivalent

$$MAE = |f - o| = \int_{\min(f,o)}^{\max(f,o)} 1 dt$$

$$twMAE = \int_{\min(f,o)}^{\max(f,o)} tw(t) dt$$

where 0 is a perfect forecast



But what do twCRPS and twMAE values mean?

Skill Scores are more intuitive...



$$\text{Skill Score} = \frac{100 \times (\text{score}_{\text{climate}} - \text{score}_{\text{forecast}})}{\text{score}_{\text{forecast}}}$$

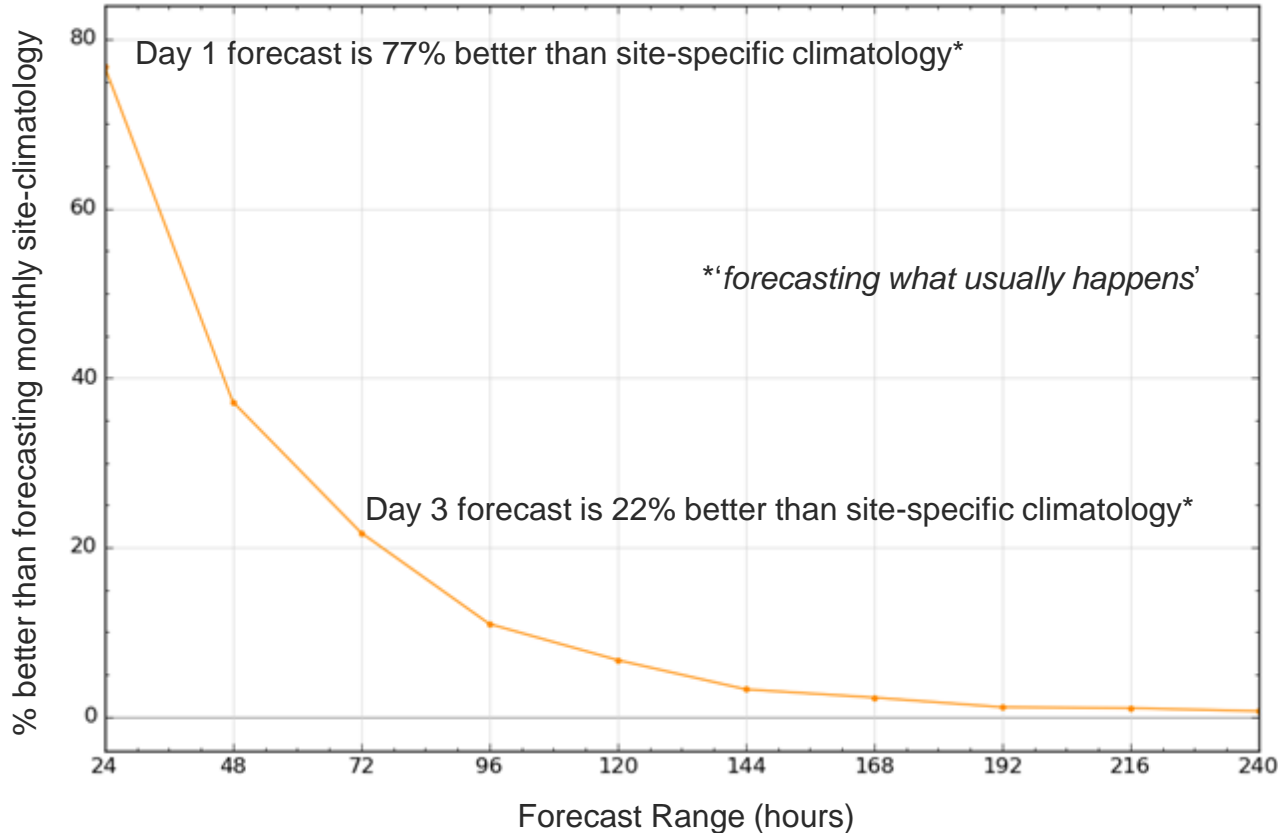
Reference:

- Site-specific climatology

Range: $-\infty$ to 1

Transform to % better than reference

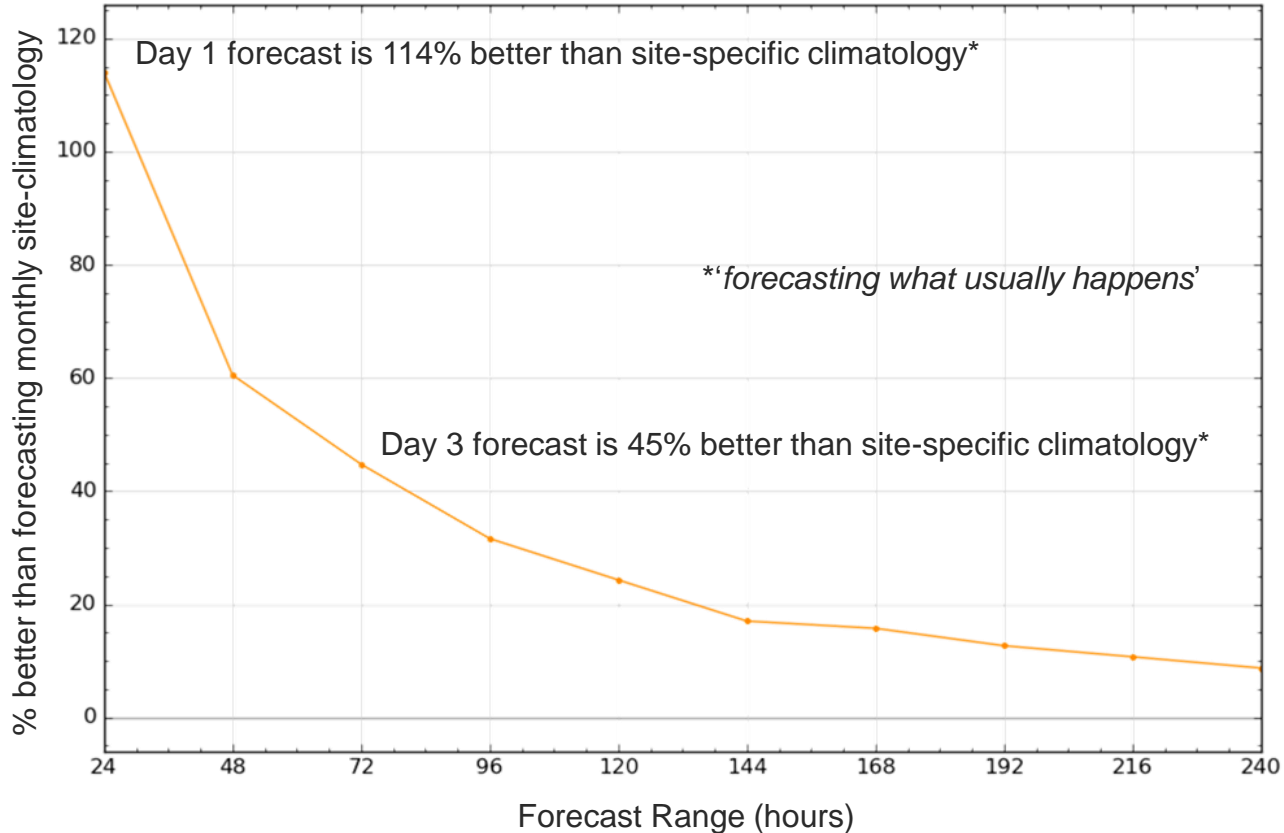
Maximum daily temperature



Forecast: deterministic
Reference: site climatology
Extreme event: hottest day expected every 3-years*?
Period: Jan-17 to Dec-19

*98.89th percentile

Maximum daily temperature



Forecast: Probabilistic

Reference: site climatology

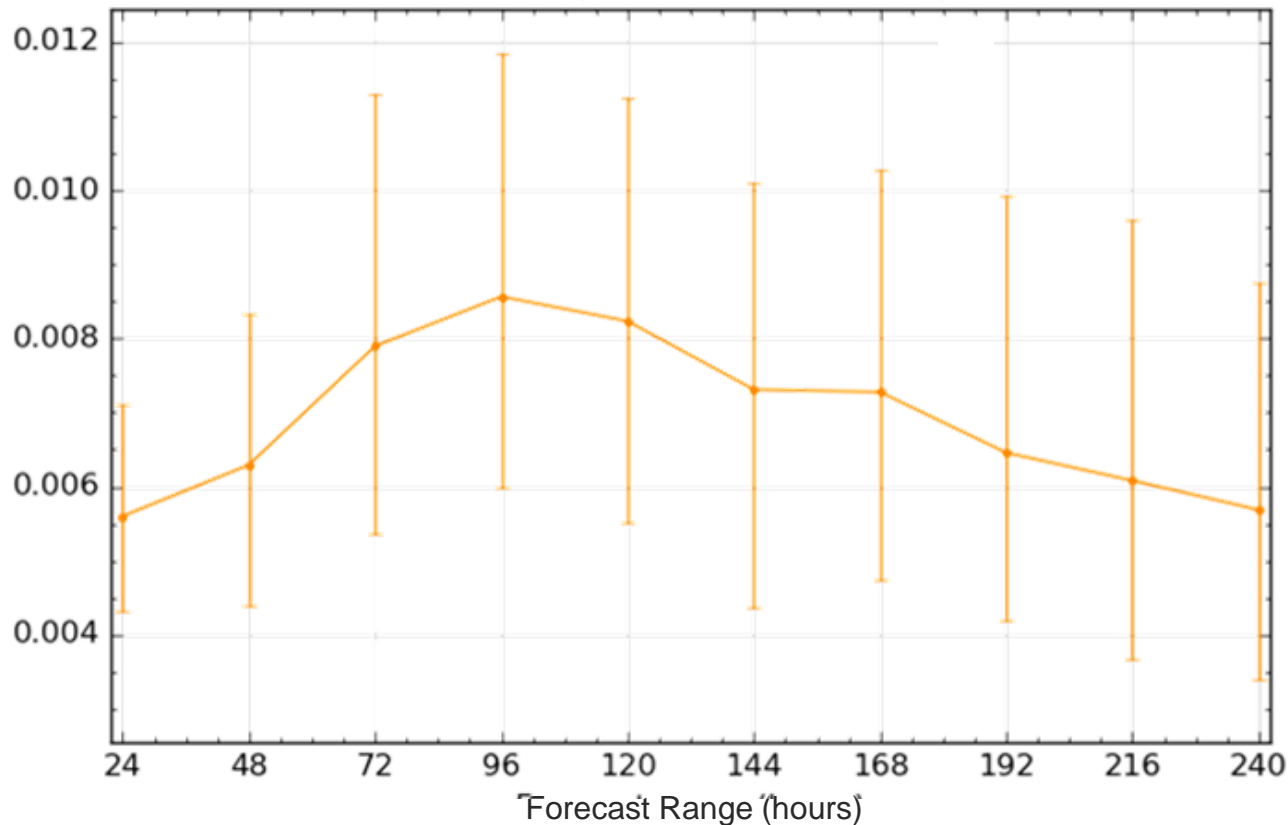
Extreme event: hottest day expected every 3-years*?

Period: Jan-17 to Dec-19

*98.89th percentile

Maximum daily temperature

Difference (twMAE- twCRPS) 99% Confidence Interval



Strong evidence that probabilistic forecast better than deterministic forecast at predicting the hottest day expected every 3-years*?

*98.89th percentile

Requirements for Public Weather Service

'Key Performance Indicator'

to measure how well we forecast extreme events:

- be easy to explain
forecast error for the most extreme value that's likely to occur every 3-years
- be simple to communicate
how much better is the forecast at predicting extremes than forecasting what usually happens
- be robust (i.e. no sudden changes)
by calculating an all-site rolling 3-year performance using a rolling 30-year climatological CDF
- directly compare deterministic & probabilistic forecast performance
yes, because $twMAE = twCRPS$



Thank you

Michael Sharpe, Clare Bysouth, Philip Gill. Operational Verification Systems and Products,

Met Office UK

IVMW-O, 12th November 2020