Verification of Numerical Weather Prediction (NWP) Forecasts

Roland Stull

Geophysical Disaster Computational Fluid Dynamics Center Earth, Ocean & Atmos. Science Dept. University of British Columbia 2020-2207 Main mall Vancouver, BC V6T 1Z4 Canada Special Thanks to my NWP Team: Henryk Modzelewski, George Hicks II, Thomas Nipen,

rstull @ eos.ubc.ca

Verification Topics



Dominique Bourdin, Atoossa

Bakhshaii, May Wong, Bruce Thomson, Greg West & many former team members.

- Continuous Variables [T, U, V, P, bounded (RH, Precip)...]
- Categorical (Binary, yes/no) Events
- Probabilistic & Ensemble Forecasts
- Terrain Issues

Verification = a measure of Forecast Quality



- Accuracy (a BAD measure)
- Skill (accuracy relative to some reference such as:
 - climatology (averaged over 30 years)
 - persistence (same as previous weather)
 - random (Monte-Carlo, bootstrapping, etc.)

Continuous Variables



2

4

- Define variables
 - A = initial analysis (based on observations)
 - *V* = verifying analysis (based on later obs.)
 - *F* = deterministic forecast
 - *C* = climatological conditions
 - n = number of grid points being averaged

Anomaly = Difference from climatology

- F C = predicted anomaly
- A C = persistence anomaly
- V C = verifying anomaly





Solved Example (§)

Given the following synthetic analysis (A), NWP forecast (F), verification (V), and climate (C) fields of 50-kPa height (km). Each field represents a weather map (North at top, East at right).

Analysis:				
	5.3	5.3	5.3	5.4
	5.4	5.3	5.4	5.5
	5.5	5.4	5.5	5.6
	5.6	5.5	5.6	5.7
	5.7	5.6	5.7	5.7
Forecast:				
	5.5	5.2	5.2	5.3
	5.6	5.4	5.3	5.4
	5.6	5.5	5.4	5.5
	5.7	5.6	5.5	5.6
	5.7	5.7	5.6	5.6
Verificatio	on:			
	5.4	5.3	5.3	5.3
	5.5	5.4	5.3	5.4
	5.5	5.5	5.4	5.5
	5.6	5.6	5.5	5.6
	5.6	5.7	5.6	5.7
Climate:				
	5.4	5.4	5.4	5.4
	5.4	5.4	5.4	5.4
	5.5	5.5	5.5	5.5
	5.6	5.6	5.6	5.6
	5.7	5.7	5.7	5.7

Find the mean error of the forecast and of persistence. Find the forecast MAE and MSE. Find MSEC and MSESS. Find the correlation coefficient between the forecast and verification. Find the RMS errors and the anomaly correlations for the forecast and persistence.

Solution

Use eq. (20.25): $ME_{forecast} = 0.01 \text{ km} = 10 \text{ m}$ Use eq. (20.26): $ME_{persistence} = 15 \text{ m}$ Use eq. (20.27): MAE = 40 mUse eq. (20.28): $MSE_{forecast} = 0.004 \text{ km}^2 = \frac{4000 \text{ m}^2}{4000 \text{ m}^2}$ Use eq. (20.28): $MSEC = 0.0044 \text{ km}^2 = \frac{4500 \text{ m}^2}{4500 \text{ m}^2}$ Use eq. (20.29): MSESS = 1 - (4000/4500) = 0.11 Use eq. (20.30): $RMSE_{forecast} = 63 \text{ m}$ Use eq. (20.30): $RMSE_{persistence} = 87 \text{ m}$ Use eq. (20.31): r = 0.92 (dimensionless) Use eq. (20.33): forecast anomaly correlation = 81.3% Use eq. (20.34): persist. anomaly correlation = 7.7%

Check: Units OK. Physics OK.

Discussion: Analyze (i.e., draw height contour maps for) the analysis, forecast, verification, and climate fields. The analysis shows a Rossby wave with ridge and trough, and the verification shows this wave moving east. The forecast amplifies the wave too much. The climate field just shows the average of higher heights to the south and lower heights to the north, with all transient Rossby waves averaged out.

Verification Topics



• Continuous Variables [T, U, V, P, bounded (RH, Precip)...]

· Categorical (Binary, yes/no) Events

- Probabilistic & Ensemble Forecasts
- Terrain Issues

Contingency Table

14

Categorical (yes/no)

• True binary:

- Snow / no-snow
- Rain / no-rain
- Sun / shade
- Threshold exceedence
 - T < T_{threshold}
 - Precip > Precip.threshold
 - Wind > Wind, threshold

• etc.



Figure 20.23

Contingency table for a binary (Yes/No) situation. "Yes" means the event occurred or was forecast to occur. (a) Meaning of cells. (b) Counts of occurrences, where a + b + c + d = n. (c) Expense matrix, where C = cost for taking protective action (i.e., for mitigating the loss), and L = loss due to an unmitigated event.

(a)	Ē	Ob	servation
		Yes	No
ast	Yes	Hit	False Alarm
Forec	No	Miss	Correct Rejection

b)	l l	Obser	vation	
		Yes	No	_
cast	Yes	а	b	
Fore	No	С	d	

(c)		Observation					
		Yes	No				
cast	Yes	Mitigated Loss (C)	Cost (C)				
Fore	No	Loss (L)	No Cost (0)				

15

Binary Scores The **bias score** *B* indicates over- or under-prediction of the frequency of event occurrence:

$$B = \frac{a+b}{a+c} \tag{20.36}$$

The <u>portion correct PC</u> (also known as portion of forecasts correct PFC) is

$$PC = \frac{a+d}{n} \tag{20.37}$$

But perhaps a portion of PC could have been due to random-chance (dumb but lucky) forecasts. Let E be this "random luck" portion, assuming that you made the same ratio of "YES" to "NO" forecasts:

$$E = \left(\frac{a+b}{n}\right) \cdot \left(\frac{a+c}{n}\right) + \left(\frac{d+b}{n}\right) \cdot \left(\frac{d+c}{n}\right)$$
(20.38)

We can now define the portion of correct forecasts that was actually skillful (i.e., not random chance), which is known as the **Heidke skill score** (*HSS*):

$$HSS = \frac{PC - E}{1 - E} \tag{20.39}$$

Binary Scores

> A <u>true skill score</u> *TSS* (also known as **Peirce's** skill score *PSS*, and as **Hansen and Kuipers'** score) can be defined as

$$TSS = H - F \tag{20.43}$$

which is a measure of how well you can discriminate between an event and a non-event, or a measure of how well you can detect an event.

A <u>critical success index CSI</u> (also known as a threat score TS) is:

$$CSI = \frac{a}{a+b+c}$$
(20.44)

Binary

Scores

The **hit rate** *H* is the portion of actual occurrences (obs. = "YES") that were successfully forecast:

$$H = \frac{a}{a+c} \tag{20.40}$$

It is also known as the **probability of detection** *POD*.

The <u>**false-alarm**</u> rate F is the portion of nonoccurrences (observation = "NO") that were incorrectly forecast:

F

$$=\frac{b}{b+d}$$
 •(20.41)

Don't confuse this with the **false-alarm ratio** *FAR*, which is the portion of "YES" forecasts that were wrong:

$$FAR = \frac{b}{a+b} \tag{20.42}$$

18

Binary Scores Suppose we consider the portion of hits that might have occurred by random chance a_r :

$$a_r = \frac{(a+b)\cdot(a+c)}{n} \tag{20.45}$$

Then we can subtract this from the actual hit count to modify CSS into an **equitable threat score** *ETS*, also known as **Gilbert's skill score** *GSS*:

$$GSS = \frac{a - a_r}{a - a_r + b + c} \tag{20.46}$$

which is also useful for rare events.

For a <u>perfect</u> forecasts (where b = c = 0), the values of these scores are: B = 1, PC = 1, HSS = 1, H = 1, F = 0, FAR = 0, TSS = 1, CSS = 1, GSS = 1. For totally wrong forecasts (where a = d = 0): B = 0 to ∞ , PC = 0, HSS = negative, H = 0, F = 1, FAR = 1, TSS = -1, CSS = 0, GSS = negative.

19

Binary Scores	Solved Example Given the following contingency table, calculate all the binary verification statistics. Observation Yes No Forecast Yes: 90 50 No 75 150	Ve
	Solution: Given: $a = 90$, $b = 50$, $c = 75$, $d = 150$ Find: B , PC , HSS , H , F , FAR , TSS , CSI , GSS First, use eq. (20.35): $n = 90 + 50 + 75 + 150 = 365$ So apparently we have daily observations for a year. Use eq. (20.36): $B = (90 + 50) / (90 + 75) = 0.85$ Use eq. (20.37): $PC = (90 + 150) / 365 = 0.66$ Use eq. (20.38): $E = [(90+50) \cdot (90+75) + (150+50) \cdot (150+75)] / (365^2)$ E = 68100 / 133225 = 0.51 Use eq. (20.39): $HSS = (0.66 - 0.51) / (1 - 0.51) = 0.31$ Use eq. (20.40): $H = 90 / (90 + 75) = 0.25$ Use eq. (20.41): $F = 50 / (50 + 150) = 0.25$ Use eq. (20.42): $FAR = 50 / (90 + 50) = 0.36$ Use eq. (20.43): $TSS = 0.55 - 0.25 = 0.30$ Use eq. (20.44): $CSI = 90 / (90 + 50 + 75) = 0.42$ Use eq. (20.45): $a_r = [(90+50) \cdot (90+75)] / 365 = 63.3$ Use eq. (20.46): GSS = (90-63.3) / (90-63.3+50+75) = 0.18	
	21	

Brier Skill Score

Probability Fcst Verification

For calibrated probability forecasts, a **Brier skill score** (*BSS*) can be defined relative to climatology as

 $BSS = 1 - \frac{\sum_{k=1}^{N} (p_k - o_k)^2}{\left(\sum_{k=1}^{N} o_k\right) \cdot \left(N - \sum_{k=1}^{N} o_k\right)}$ (20.47)

where p_k is the forecast probability ($0 \le p_k \le 1$) that the threshold will be exceeded (e.g., the probability that the precipitation will exceed a precipitation threshold) for any one forecast k, and N is the number of forecasts. The verifying observation $o_k = 1$ if the observation exceeded the threshold, and is set to zero otherwise.

BSS = 0 for a forecast no better than climatology. BSS = 1 for a perfect deterministic forecast (i.e., the forecast is $p_k = 1$ every time the event happens, and $p_k = 0$ every time it does not). For probabilistic forecasts, $0 \le BSS \le 1$. Larger BSS values are better. erification Topics • Continuous Variables [T, U, V, P, bounded (RH, Precip)...] · Categorical (Binary, yes/no) Events Probabilistic & Ensemble Forecasts Terrain Issues 22 Reliability

Probability Fcst Verification How **reliable** are the probability forecasts? Namely, when we forecast an event with a certain probability, is it observed with the same relative frequency? To determine this, after you make each forecast, sort it into a forecast probability bin (*j*) of probability width Δp , and keep a tally of the number of forecasts (n_i) that fell in this bin, and count how many of the forecasts verified (n_{ojr} for which the corresponding observation satisfied the threshold).

For example, if you use bins of size $\Delta p = 0.1$, then create a table such as:

bin index	bin center	fcst. prob. range	n	noi
j = 0	$p_i = 0$	$0 \leq p_k < 0.05$	no	noo
j = 1	$p_{i} = 0.1$	$0.05 \le p_k < 0.15$	n_1	n_{o1}
<i>j</i> = 2	$p_{j} = 0.2$	$0.15 \leq p_k < 0.25$	n_2	n_{o2}
etc.				
j = 9	$p_{i} = 0.9$	$0.85 \le p_k < 0.95$	n_9	n ₀₉
j = 10 = J	$p_{j} = 1.0$	$0.95 \le p_k \le 1.0$	n ₁₀	<i>n</i> ₀₁₀

A plot of the observed relative frequency (n_{oj}/n_j) on the ordinate vs. the corresponding forecast probability bin center (p_j) on the abscissa is called a **reliability diagram**. For perfect reliability, all the points should be on the 45° diagonal line.



where $BSS_{reliability} = 0$ for a perfect forecast.

Probability Fcst Verification

Reliability Diagram

Ideally, the probabilistic forecast-observation points lie on the diagonal of the reliability diagram, indicating the event is always forecast at the same frequency it is observed. The reliability component of the Brier score in a graphical representation is the weighted, averaged, squared distance between the reliability curve and the 45° diagonal line. If the points lie above (below) the diagonal, it means the event is underforecast (overforecast). Reliability curves with a zigzag shape centered on the diagonal indicate good reliability represented by a small sample size. Poor reliability can be improved substantially by appropriate a posteriori calibration and/or postprocessing of forecasts delivered from an established system, though it is a difficult task to achieve in a real-time operational EPS (Atger 2003).

Solution

Probability Fcst Verification

Solved Example (§)

Given the table below of k = 1 to 31 forecasts of the probability p_k that the temperature will be below threshold 20°C, and the verification $o_k = 1$ if indeed the observed temperature was below the threshold.

(a) Find the Brier skill score. (b) For probability bins of width $\Delta p = 0.2$, plot a reliability diagram, and find the reliability Brier skill score.

k	Pk	ok	BN	j	k	p_k	ok	BN	j
1	0.43	0	0.18	2	16	0.89	1	0.01	4
2	0.98	1	0.00	5	17	0.13	0	0.02	1
3	0.53	1	0.22	3	18	0.92	1	0.01	5
4	0.33	1	0.45	2	19	0.86	1	0.02	4
5	0.50	0	0.25	3	20	0.90	1	0.01	5
6	0.03	0	0.00	0	21	0.83	0	0.69	4
7	0.79	1	0.04	4	22	0.00	0	0.00	0
8	0.23	0	0.05	1	23	1.00	1	0.00	5
9	0.20	1	0.64	1	24	0.69	0	0.48	3
10	0.59	1	0.17	3	25	0.36	0	0.13	2
11	0.26	0	0.07	1	26	0.56	1	0.19	3
12	0.76	1	0.06	4	27	0.46	0	0.21	2
13	0.17	0	0.03	1	28	0.63	0	0.40	3
14	0.30	0	0.09	2	29	0.10	0	0.01	1
15	0.96	1	0.00	5	30	0.40	1	0.36	2
					31	0.73	1	0.07	4

Given: The white portion of the table above. Find: BSS = ?, $BSS_{reliability} = ?$, and plot reliability.

(a) Use eq. (20.47). The grey-shaded column labeled BN shows each contribution to the numerator $(p_k - o_k)^2$ in that eq. The sum of BN = 4.86. The sum of $o_k = 16$. Thus, the eq is: $BSS = 1 - [4.86 / \{16 \cdot (31-16)\}] = 0.98$

(b) There are I = 6 bins, with bin centers at $p_i = 0, 0.2$, 0.4, 0.6, 0.8, and 1.0. (Note, the first and last bins are one-sided, half-width relative to the nominal "center" value.) I sorted the forecasts into bins using j =round($p_k/\Delta p$, 0), giving the grey *j* columns above.

For each *j* bin, I counted the number of forecasts n_i falling in that bin, and I counted the portion of those forecasts that verified n_{ai} . See table below:

			-	1					
j	p_i	nj	noj	n_{oj}/n_j	num	-	I		_
0	0	2	0	0	0				
1	0.2	6	1	0.17	0.04	0.0	3		
2	0.4	6	2	0.33	0.16	S.	+	-	-
3	0.6	6	3	0.50	0.36	° c			1
4	0.8	6	5	0.83	0.04	0.4	4	11	1
5	1.0	5	5	1	0	0.3	2	1	
The	obs	erv	ed r	elative	fre-		/		
que	ency	noil	n; p	lotted	agains	t	0 0	2 0	4
Dij	is the	rel	iabi	lity di	iagram	1:	0 0	0.	-

Use eq. (20.48). The contribution to the numerator from each bin is in the num column above, which sums to 0.6. Thus: $BSS_{reliability} = 0.6/ \{16 \cdot (31-16)\} = 0.0025$

b. Continuous ranked probability score

The continuous ranked probability score (CRPS; Hersbach 2000) is a complete generalization of the Brier score. The Brier score (Brier 1950; Atger 2003; McCollor The CRPS can be interpreted as an integral over all

possible Brier scores. A major advantage of the Brier score is that it can be decomposed into a reliability component, a resolution component, and an uncertainty component (Murphy 1973; Toth et al. 2003). In a similar fashion, Hersbach (2000) showed how the CRPS can be decomposed into the same three components:

CRPS = Reli - Resol + Unc.





a reliability C

Probability Fcst Verification

Continuous Ranked Probability Skill Score

The CRPS components can be converted to a positively oriented skill score (<u>CRPSS</u>) in the same manner that the Brier score components are converted to a skill score (BSS; McCollor and Stull 2008b):

 $CRPSS = \frac{Resol}{Unc} - \frac{Reli}{Unc}$ = relative resolution - relative reliability = CRPS_{RelResol} - CRPS_{RelReli}.





Linear Error in Probability Space (LEPS)

Linear error in probability space (LEPS) is defined as the mean absolute difference between the cumulative frequency of the forecasts and the cumulative frequency of the observations (Déqué 2003). LEPS ensures that error in the center of the distribution is treated with more importance than error found in the extreme tail of the distribution. A LEPS skill score can be defined with the climatological median as a reference (Wilks 1995).

ROC diagram

ROC Diagram

A <u>Relative Operating Characteristic</u> (ROC) diagram shows how well a probabilistic forecast can **discriminate** between an event and a non-event. For example, an event could be heavy rain that causes flooding, or cold temperatures that cause crops to freeze. The probabilistic forecast could come from an ensemble forecast, as illustrated next.

33

Probability Fcst Verification

ROC diagram

On Day 2, three of the 10 models forecast 10 mm or more of precipitation, hence the forecast probability is $p_2 = 3/10 = 30\%$. On this day precipitation did NOT exceed 10 mm, so the observation flag is set to zero: $o_2 = 0$. Similarly, for Day 3 suppose the forecast probability is $p_3 = 10\%$, but heavy rain was observed, so $o_3 = 1$. After making ensemble forecasts every day for a month, suppose the results are as listed in the left three columns of Table 20-4.

Probability Fcst Verification

ROC diagram

34

Suppose that the individual NWP models of an N = 10 member ensemble made the following forecasts of 24-h accumulated rainfall *R* for Day 1:

NWP model	<u>R (mm)</u>	NWP model	R (mm)
Model 1	8	Model 6	4
Model 2	10	Model 7	20
Model 3	6	Model 8	9
Model 4	12	Model 9	5
Model 5	11	Model 10	7

Consider a precipitation threshold of 10 mm. The ensemble above has 4 models that forecast 10 mm or more, hence the forecast probability is $p_1 = 4/N = 4/10 = 40\%$. Supposed that 10 mm or more of precipitation was indeed observed, so the observation flag is set to one: $o_1 = 1$.

D		(0/)				Prob	ability T	hreshold	1 p _{threshol}	d (%)			
Day	0	p (%)	0	10	20	30	40	50	60	70	80	90	100
1	1	40	1	1	1	1	1	0	0	0	0	0	0
2	0	30	1	1	1	1	0	0	0	0	0	0	0
3	1	10	1	1	0	0	0	0	0	0	0	0	0
4	1	50	1	1	1	1	1	1	0	0	0	0	0
5	0	60	1	1	1	1	1	1	1	0	0	0	0
6	0	30	1	1	1	1	0	0	0	0	0	0	0
7	0	40	1	1	1	1	1	0	0	0	0	0	0
8	1	80	1	1	1	1	1	1	1	1	1	0	0
9	0	50	1	1	1	1	1	1	0	0	0	0	0
10	1	20	1	1	1	0	0	0	0	0	0	0	0
11	1	90	1	1	1	1	1	1	1	1	1	1	0
12	0	20	1	1	1	0	0	0	0	0	0	0	0
13	0	10	1	1	0	0	0	0	0	0	0	0	0
14	0	10	1	1	0	0	0	0	0	0	0	0	0
15	1	70	1	1	1	1	1	1	1	1	0	0	0
16	0	70	1	1	1	1	1	1	1	1	0	0	0
17	1	60	1	1	1	1	1	1	1	0	0	0	0
18	1	90	1	1	1	1	1	1	1	1	1	1	0
19	1	80	1	1	1	1	1	1	1	1	1	0	0
20	0	80	1	1	1	1	1	1	1	1	1	0	0
21	0	20	1	1	1	0	0	0	0	0	0	0	0
22	0	10	1	1	0	0	0	0	0	0	0	0	0
23	0	0	1	0	0	0	0	0	0	0	0	0	0
24	0	0	1	0	0	0	0	0	0	0	0	0	0
25	1	70	1	1	1	1	1	1	1	1	0	0	0
26	0	10	1	1	0	0	0	0	0	0	0	0	0
27	0	0	1	0	0	0	0	0	0	0	0	0	0
28	1	90	1	1	1	1	1	1	1	1	1	1	0
29	0	20	1	1	1	0	0	0	0	0	0	0	0
30	1	80	1	1	1	1	1	1	1	1	1	0	0
		a =	13	13	12	11	11	10	9	8	6	3	0
Contin	gency	b = 1	17	14	10	7	5	4	3	2	1	0	0
Table V	/alues	C =	0	0	1	2	2	3	4	5	7	10	13
		d =	0	3	7	10	12	13	14	15	16	17	17

ROC diagram

An end user might need to make a decision to take action. Based on various economic or political reasons, the user decides to use a probability threshold of $p_{threshold} = 40\%$; namely, if the ensemble model forecasts a 40% or greater chance of daily rain exceeding 10 mm, then the user will take action. So we can set forecast flag f = 1 for each day that the ensemble predicted 40% or more probability, and f = 0 for the other days. These forecast flags are shown in Table 20-4 under the $p_{threshold} = 40\%$ column.

Probability Fcst
Verification

ROC diagram

37

39

Other users might have other decision thresholds, so we can find the forecast flags for all the other probability thresholds, as given in Table 20-4. For an *N* member ensemble, there are only (100/N) + 1 discrete probabilities that are possible. For our example with *N* = 10 members, we can consider only 11 different probability thresholds: 0% (when no members exceed the rain threshold), 10% (when 1 out of the 10 members exceeds the threshold), 20% (etc.), . . . 90%, 100%.

D	225					Prob	ability 1	hreshold	1 pthreshol	Id (%)			
Day	0	p (%)	0	10	20	30	40	50	60	70	80	90	100
1	1	40	1	1	1	1	1	0	0	0	0	0	0
2	0	30	1	1	1	1	0	0	0	0	0	0	0
3	1	10	1	1	0	0	0	0	0	0	0	0	0
4	1	50	1	1	1	1	1	1	0	0	0	0	0
5	0	60	1	1	1	1	1	1	1	0	0	0	0
6	0	30	1	1	1	1	0	0	0	0	0	0	0
7	0	40	1	1	1	1	1	0	0	0	0	0	0
8	1	80	1	1	1	1	1	1	1	1	1	0	0
9	0	50	1	1	1	1	1	1	0	0	0	0	0
10	1	20	1	1	1	0	Ő	0	0	0	0	0	0
11	1	90	1	1	1	1	1	1	1	1	1	1	0
12	Ô	20	1	1	1	Ô	Ô	Ô	0	Ô	Ô	Ô	0
13	Ő	10	1	1	Ô	0	Ő	Ő	0	0	0	0	Ő
14	0	10	1	1	0	0	0	0	0	0	0	0	0
15	1	70	1	1	1	1	1	1	1	1	0	0	0
16	Ô	70	1	1	1	1	1	1	1	1	0	0	0
17	1	60	1	1	1	1	1	1	1	Ô	0	0	0
18	1	90	1	1	1	1	1	1	1	1	1	1	0
19	1	80	1	1	1	1	î	1	1	1	1	Ô	Ő
20	Ô	80	1	1	1	1	1	1	1	1	1	0	Ő
21	Ő	20	1	1	1	Ô	Ô	Ô	0	Ô	Ô	0	0
22	Ő	10	1	1	Ô	0	Ő	0	0	0	0	0	0
23	Ő	0	1	Ô	0	0	Ő	0	0	0	0	Ő	0
24	Ő	0	1	Ő	Ő	0	Ő	Ő	0	Ő	0	Ő	Ő
25	1	70	1	1	1	1	1	1	1	1	0	0	0
26	Ô	10	1	1	0	0	0	0	0	0	0	Ő	0
27	Ő	0	1	Ô	0	0	0	0	0	0	0	0	Ő
28	1	90	1	1	1	1	1	1	1	1	1	1	0
29	Ó	20	1	1	1	Ô	Ô	Ô	Ô	Ô	Ô	Ô	0
30	1	80	1	1	1	1	1	1	1	1	1	0	0
									_				
022 020		<i>a</i> =	13	13	12	11	11	10	9	8	6	3	0
Contin	gency	<i>b</i> =	17	14	10	7	5	4	3	2	1	0	0
Table V	/alues	C =	0	0	1	2	2	3	4	5	7	10	13
		d =	0	3	7	10	12	13	14	15	16	17	17

Dave	0		Probability Threshold p _{threshold} (%)												
Day	0	p (%)	0	10	20	30	40	50	60	70	80	90	100		
1	1	40	1	1	1	1	1	0	0	0	0	0	0		
2	0	30	1	1	1	1	0	0	0	0	0	0	0		
3	1	10	1	1	0	0	0	0	0	0	0	0	0		
4	1	50	1	1	1	1	1	1	0	0	0	0	0		
5	0	60	1	1	1	1	1	1	1	0	0	0	0		
6	0	30	1	1	1	1	0	0	0	0	0	0	0		
7	0	40	1	1	1	1	1	0	0	0	0	0	0		
8	1	80	1	1	1	1	1	1	1	1	1	0	0		
9	0	50	1	1	1	1	1	1	0	0	0	0	0		
10	1	20	1	1	1	0	0	0	0	0	0	0	0		
11	1	90	1	1	1	1	1	1	1	1	1	1	0		
12	0	20	1	1	1	0	0	0	0	0	0	0	0		
13	0	10	1	1	0	0	0	0	0	0	0	0	0		
14	0	10	1	1	0	0	0	0	0	0	0	0	0		
15	1	70	1	1	1	1	1	1	1	1	0	0	0		
16	0	70	1	1	1	1	1	1	1	1	0	0	0		
17	1	60	1	1	1	1	1	1	1	0	0	0	0		
18	1	90	1	1	1	1	1	1	1	1	1	1	0		
19	1	80	1	1	1	1	1	1	1	1	1	0	0		
20	0	80	1	1	1	1	1	1	1	1	1	0	0		
21	0	20	1	1	1	0	0	0	0	0	0	0	0		
22	0	10	1	1	0	0	0	0	0	0	0	0	0		
23	0	0	1	0	0	0	0	0	0	0	0	0	0		
24	0	0	1	0	0	0	0	0	0	0	0	0	0		
25	1	70	1	1	1	1	1	1	1	1	0	0	0		
26	0	10	1	1	0	0	0	0	0	0	0	0	0		
27	0	0	1	0	0	0	0	0	0	0	0	0	0		
28	1	90	1	1	1	1	1	1	1	1	1	1	0		
29	0	20	1	1	1	0	0	0	0	0	0	0	0		
30	1	80	1	1	1	1	1	1	1	1	1	0	0		
			6							-					
200 - 10 M		<i>a</i> =	13	13	12	11	11	10	9	8	6	3	0		
Contin	gency	b =	17	14	10	7	5	4	3	2	1	0	0		
Table V	Values	c =	0	0	1	2	2	3	4	5	7	10	13		
		d =	0	3	7	10	12	13	14	15	16	17	17		

ROC diagram

For each probability threshold, create a 2x2 contingency table with the elements *a*, *b*, *c*, and *d* as shown in Fig. 20.23b. For example, for any pair of observation and forecast flags (o_j, f_j) for Day *j*, use *a* = count of days with hits $(o_j, f_j) = (1, 1)$. *b* = count of days with false alarms $(o_j, f_j) = (0, 1)$. *c* = count of days with misses $(o_j, f_j) = (1, 0)$. *d* = count of days: correct rejection $(o_j, f_j) = (0, 0)$. For our illustrative case, these contingency-table elements are shown near the bottom of Table 20-4.

Day	2.2	p (%)	Probability Threshold <i>p</i> _{threshold} (%)												
	0		0	10	20	30	40	50	60		70	80	90	100	
1	1	40	1	1	1	1	1	0	0	(b)	Ť	0	bservation		
2	0	30	1	1	1	1	0	0	0			Yes		No	
3	1	10	1	1	0	0	0	0	0						
4	1	50	1	1	1	1	1	1	0	ast	Yes	a		b	
5	0	60	1	1	1	1	1	1	1	Forec	1000			24	
6	0	30	1	1	1	1	0	0	0		No	C	d		
7	0	40	1	1	1	1	1	0	0	-	-		-		
8	1	80	1	1	1	1	1	1	1	-	1	1	0	0	
9	0	50	1	1	1	1	1	1	0		0	0	0	0	
10	1	20	1	1	1	0	0	0	0		0	0	0	0	
11	1	90	1	1	1	1	1	1	1		1	1	1	0	
12	0	20	1	1	1	0	0	0	0		0	0	0	0	
13	0	10	1	1	0	0	0	0	0		0	0	0	0	
14	0	10	1	1	0	0	0	0	0		0	0	0	0	
15	1	70	1	1	1	1	1	1	1		1	0	0	0	
16	0	70	1	1	1	1	1	1	1		1	0	0	C	
17	1	60	1	1	1	1	1	1	1		0	0	0	0	
18	1	90	1	1	1	1	1	1	1		1	1	1	0	
19	1	80	1	1	1	1	1	1	1		1	1	0	0	
20	0	80	1	1	1	1	1	1	1		1	1	0	C	
21	0	20	1	1	1	0	0	0	0		0	0	0	0	
22	0	10	1	1	0	0	0	0	0		0	0	0	(
23	0	0	1	0	0	0	0	0	0		0	0	0	0	
24	0	0	1	0	0	0	0	0	0		0	0	0	0	
25	1	70	1	1	1	1	1	1	1		1	0	0	C	
26	0	10	1	1	0	0	0	0	0		0	0	0	0	
27	0	0	1	0	0	0	0	0	0		0	0	0	C	
28	1	90	1	1	1	1	1	1	1		1	1	1	0	
29	0	20	1	1	1	0	0	0	0		0	0	0	0	
30	1	80	1	1	1	1	1	1	1		1	1	0	0	
		<i>a</i> =	13	13	12	11	11	10	9		8	6	3	0	
Contin	gency	b = 1	17	14	10	7	5	4	3		2	1	0	0	
Table V	alues	c =	0	0	1	2	2	3	4		5	7	10	13	
		d =	0	3	7	10	12	13	14		15	16	17	17	

Probability Fcst Verification

ROC diagram

Next, for each probability threshold, calculate the hit rate H = a/(a+c) and false alarm rate F = b/(b+d), as defined earlier in this chapter. These are shown in the last two rows of Table 20-4 for our example. When each (*F*, *H*) pair is plotted as a point on a graph, the result is called a **ROC diagram** (Fig. 20.24).

Day	о	(1)	Probability Threshold pthreshold (%)											
		p (%)	0	10	20	30	40	50	60	70	80	90	10	
1	1	40	1	1	1	1	1	0	0	0	0	0	0	
2	0	30	1	1	1	1	0	0	0	0	0	0	0	
3	1	10	1	1	0	0	0	0	0	0	0	0	0	
4	1	50	1	1	1	1	1	1	0	0	0	0	0	
5	0	60	1	1	1	1	1	1	1	0	0	0	0	
6	0	30	1	1	1	1	0	0	0	0	0	0	0	
7	0	40	1	1	1	1	1	0	0	0	0	0	0	
8	1	80	1	1	1	1	1	1	1	1	1	0	0	
9	0	50	1	1	1	1	1	1	0	0	0	0	0	
10	1	20	1	1	1	0	0	0	0	0	0	0	0	
11	1	90	1	1	1	1	1	1	1	1	1	1	0	
12	0	20	1	1	1	0	0	0	0	0	0	0	0	
13	0	10	1	1	0	0	0	0	0	0	0	0	0	
14	0	10	1	1	0	0	0	0	0	0	0	0	0	
15	1	70	1	1	1	1	1	1	1	1	0	0	0	
16	0	70	1	1	1	1	1	1	1	1	0	0	0	
17	1	60	1	1	1	1	1	1	1	0	0	0	0	
18	1	90	1	1	1	1	1	1	1	1	1	1	0	
19	1	80	1	1	1	1	1	1	1	1	1	0	0	
20	0	80	1	1	1	1	1	1	1	1	1	0	0	
21	0	20	1	1	1	0	0	0	0	0	0	0	0	
22	0	10	1	1	0	0	0	0	0	0	0	0	0	
23	0	0	1	0	0	0	0	0	0	0	0	0	0	
24	0	0	1	0	0	0	0	0	0	0	0	0	0	
25	1	70	1	1	1	1	1	1	1	1	0	0	0	
26	0	10	1	1	0	0	0	0	0	0	0	0	0	
27	0	0	1	0	0	0	0	0	0	0	0	0	0	
28	1	90	1	1	1	1	1	1	1	1	1	1	0	
29	0	20	1	1	1	0	0	0	0	0	0	0	0	
30	1	80	1	1	1	1	1	1	1	1	1	0	0	
		<i>a</i> =	13	13	12	11	11	10	9	8	6	3	0	
Contin	Contingency Table Values		17	14	10	7	5	4	3	2	1	0	0	
Table \			0	0	1	2	2	3	4	5	7	10	13	
		<i>d</i> =	0	3	7	10	12	13	14	15	16	17	17	
Hit Rate: $H = a/(a+c) =$		1.00	1.00	0.92	0.85	0.85	0.77	0.69	0.62	0.46	0.23	0.0		
lse Alarm Rate: $F=b/(b+d)$			1.00	0.82	0.59	0.41	0.29	0.26	0.18	012	0.06	0.00	0.0	

44





Sharpness

Associated with reliability is sharpness, which characterizes the relative frequency of occurrence of the forecast probabilities. Sharpness is often depicted in a histogram indicating the relative occurrence of each forecast probability category. If forecast probabilities are frequently near 0 or near 1, then the forecasts are sharp, indicating the forecasts deviate significantly from the climatological mean, a positive attribute of an ensemble forecast system. Sharpness measures the variability of the forecasts alone, without regard to their corresponding observations; hence, it is not a verification measure in itself. In a perfectly reliable forecast system, sharpness is identical to resolution.

Ensemble Fcst Verification

Prob. Integral Transform (PIT) Let $f_t(x) =$ **forecast** prob. density for the value ft x of any variable (e.g., Temperature) forecast prob. distr. Let $F_t(x)$ denote the forecasted cumulative distribution function (CDF) given by $F_t(x) = \int_{-\infty}^{\infty} f_t(s) \, ds.$ x (°C) (1) F_t^1 cum, prob. distr. In addition, let x, denote the observed state of X at time t. Let p_t denote the CDF value corresponding to the of fcst. observed state: pt (2) $p_t = F_t(x_t).$ \cap x (°C) Often, p_t is called the probability integral transform observed value (PIT value) corresponding to the observation. value of x Good fcst if $p_t = p$, where p = observed cum. freq at x_t .

Ensemble Fcst Verification

Resolution

j. Resolution

A useful probabilistic forecast system must be able to a priori differentiate future weather outcomes, so that differing forecasts are, in fact, associated with distinct verifying observations. This is the most important attribute of a forecast system (Toth et al. 2003) and is called *resolution*.

Resolution cannot be improved through simple adjustment of probability values or statistical postprocessing. Resolution can be gained only by improving the actual forecast model engine that produces the forecasts.

Prob. Fcst Verif.

Ignorance Score (IGN)

$$\operatorname{IGN}(f) = \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} -\log_2[f_t(x_t)].$$
(21)

IGN rewards forecasts that place high confidence in the value where the observation falls. Low ignorance scores are desired.

where $f_t(x_t) =$ **forecast** prob. density for the **observed** value x_t of any variable (e.g., Temperature)

55

53



More so for coarser grids.

Thus, the elevation of a verifying obs is often different from the elevation in the NWP model.

Thus, the model is **not representative** of the observation.

Verification Topics



• Continuous Variables [T, U, V, P, bounded (RH, Precip)...]

- Categorical (Binary, yes/no) Events
- Probabilistic & Ensemble Forecasts
- Terrain Issues

Canadian Terrain Elevation



• West-East terrain cross section through Whistler (50.12°N)

59



How Fine is Fine Enough? Many valleys are narrower than 1 km



An Obvious Trick: Use finer horizontal grid size

Example: If grid size is $\Delta x = 7$ km -

Then the modeled terrain is closer to the actual terrain. **Good.** And the modeled slopes become steeper (closer to real). **Difficult**.



+ West-East terrain cross section through Whistler (50.12°N), where 0.1° lon \approx 7 km.

66

Summary



- Continuous Variables [T, U, V, P, bounded (RH, Precip)...]
- Categorical (Binary, yes/no) Events
- Probabilistic & Ensemble Forecasts
- Terrain Issues

Roland Stull, UBC