

# Is ETSS really equitable?

L. Kalin

Meteorological and Hydrological Service,  DHMZ  
Croatia

# working title...

- “Defectiveness of equitable skill scores for a multicategorical table”

# Outline

- a brief history of multicategory tables
- some thoughts on ETSS

# ...on multicategory tables...

- verification for standard 2x2 tables goes far into 19th century (“The Finley Affair”: a 1884 paper by J. P. Finley in the *American Meteorological Journal*)
- when it comes to multicategory tables, history of development is not so long

# ...Vernon...

- Vernon, E. M., 1953: A New Concept Of Skill Score For Rating Quantitative Forecasts, *Mon. Wea. Rev.*

## A NEW CONCEPT OF SKILL SCORE FOR RATING QUANTITATIVE FORECASTS

EDWARD M. VERNON

Weather Bureau Forecast Center, San Bruno, Calif.

[Manuscript received November 8, 1951]

### ABSTRACT

Skill scores for rating quantitative forecasts are proposed to take into account the deviations occurring between forecast and observed values. One score, the "deviation" skill score, weights the forecasts linearly according to the deviation; a second score, the "quadratic" skill score weights them according to the square of the deviation. These two scores are compared with the conventional skill score for two sets of forecasts, and for the same forecasts with bias introduced. It is concluded that use of either the deviation or the quadratic skill score is preferable to use of the conventional skill score in rating quantitative forecasts. Examples of the step-by-step computations of the two new scores are given.

### THE DEVIATION SKILL SCORE

The skill score, as first proposed by Heidke [1] and used during recent years for certain forecast verification purposes, may be written

$$S = \frac{R - E}{T - E} \quad (1)$$

where  $S$  is the skill score,  $R$  the number of correct forecasts,  $T$  the total number of forecasts, and  $E$  the number of forecasts expected to be correct on some standard such as chance.

This method of computing a skill score places the same weight on each incorrect forecast regardless of the amount by which the observed condition deviates from the forecast. In other words, a deviation of say 10 class intervals has no more effect on the skill score than one of but 1 class interval. For some purposes it would be advantageous to have the skill score evaluate the actual amount by which forecast and observed conditions differ, i. e., take into account the magnitude of error. To accomplish this end, an analogous equation for skill score may be written

$$S_d = \frac{\sum d_r - \sum d_c}{\sum d_r} \quad (2)$$

where  $S_d$  is the skill score which considers magnitude of deviations, hereafter referred to as the "deviation skill score,"  $\sum d_r$  is the sum of deviations occurring between forecast and observed values, and  $\sum d_c$  is the sum of deviations to be expected on some basis such as chance.

The value of  $\sum d_r$  and  $\sum d_c$  can best be expressed in terms of row, column, and cell totals in the typical contingency table, wherein the frequencies of forecast values are arrayed in columns and of observed values in rows, while a given cell is identified by the row and column to which

it alone is common. When the standard of reference is chance, the summations become

$$\sum d_r = \sum (n_r d_r) \quad (3)$$

$$\sum d_c = \sum \left( \frac{n_r n_c}{T} d_{rc} \right) \quad (4)$$

where  $n_r$  is the number of cases falling in a given row;  $n_c$  is the number of cases falling in a given column;  $n_{rc}$  is the number of cases in the cell at the intersection of row  $r$  and column  $c$ ;  $n_r n_c / T$  is the number which would have fallen by chance in the cell representing the intersection of row  $r$  and column  $c$ ;  $d_{rc}$  is the deviation represented by that cell and is equal to the number of class intervals by which the cell is removed from the perfect hit cell for the same column.

When the standard of reference is climatological expectancy, according to one of the more common definitions of that standard,  $\sum d_r$  remains as expressed in (3) while  $\sum d_c$  becomes

$$\sum d_c = \sum \left( \frac{n_r n_{ec}}{T} d_{rc} \right) \quad (5)$$

where  $n_{ec}$  represents the climatological expectancy for the column, i. e., the number of times which climatological averages would lead one to expect the observed conditions to fall in the particular class interval represented by the column. The other symbols in (5) remain as previously defined in (4) and (1).

### THE QUADRATIC SKILL SCORE

In the foregoing equations all deviations are weighted linearly, a deviation of one class interval scoring as one unit deviation, two class intervals as two unit deviations,

# ...Vernon...

- proposed “deviation skill score” and “quadratic skill score” to take into account the amount of error (“accuracy”)
- such system would discourage forecaster to issue forecasts at extremes
- introduced weights based on chance.

TABLE 3.—Computation of “deviation” skill score,  $S_d$ , for set of forecasts appearing in table 3A below

TABLE 3A.—Array of frequencies of forecast and observed values:  $n_{ij}$ ,  $n_{i.}$ , and  $n_{.j}$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	65	5	2	1	1	74
0.01-0.20.....	9	6	3	3	1	22
0.21-0.50.....	4	5	3	2	1	15
0.51-1.00.....	1	2	3	3	2	10
≥1.01.....	1	1	2	2	2	8
Total.....	80	19	14	11	6	130

TABLE 3B.—Array of deviation represented by each cell:  $d_{ij}$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	0	1	2	3	4	10
0.01-0.20.....	1	0	1	2	3	7
0.21-0.50.....	2	1	0	1	2	6
0.51-1.00.....	3	2	1	0	1	7
≥1.01.....	4	3	2	1	0	10

TABLE 3C.—Array of number of cases expected by chance in each cell:  $\frac{n_{i.}n_{.j}}{T}$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	16.51	10.52	7.87	6.29	4.19	45.38
0.01-0.20.....	13.64	3.22	2.37	1.86	1.02	22.11
0.21-0.50.....	8.80	2.54	1.72	1.35	0.75	15.16
0.51-1.00.....	6.15	1.45	1.06	0.82	0.46	10.04
≥1.01.....	4.02	1.17	0.86	0.68	0.37	7.10

TABLE 3D.—Array of chance frequencies weighted by cell deviation:  $\frac{n_{i.}n_{.j}}{T} d_{ij}$ =(table 3B) (table 3C)

$$\sum \left( \frac{n_{i.}n_{.j}}{T} d_{ij} \right) = \sum d_{ij} = 155.26$$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	0	10.52	15.74	12.75	12.08	51.67
0.01-0.20.....	13.64	0	2.37	3.72	3.06	22.79
0.21-0.50.....	16.70	2.54	0	1.35	1.46	22.05
0.51-1.00.....	18.45	2.92	1.06	0	0.46	22.89
≥1.01.....	16.08	3.51	1.72	0.68	0	22.00

TABLE 3E.—Array of forecast frequencies weighted by cell deviation:  $n_{ij}d_{ij}$ =(table 3A) (table 3B)

$$\sum (n_{ij}d_{ij}) = \sum d_{ij} = 78$$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	0	5	4	3	4	16
0.01-0.20.....	9	0	3	6	3	21
0.21-0.50.....	8	3	0	2	2	15
0.51-1.00.....	3	4	3	0	1	11
≥1.01.....	4	3	4	2	0	13

$$S_d = \frac{\sum d_{ij} - \sum d_{ij}^2}{\sum d_{ij}} = \frac{155.26 - 78}{155.26} = 498$$

TABLE 4.—Computation of “quadratic” skill score,  $S_{d^2}$ , for set of forecasts appearing in table 4A below

TABLE 4A.—Array of frequencies forecast and observed values:  $n_{ij}$ ,  $n_{i.}$ , and  $n_{.j}$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	65	5	2	1	1	74
0.01-0.20.....	9	6	3	3	1	22
0.21-0.50.....	4	5	3	2	1	15
0.51-1.00.....	1	2	3	3	2	10
≥1.01.....	1	1	2	2	2	8
Total.....	80	19	14	11	6	130

TABLE 4B.—Array of squared deviation represented by each cell:  $d_{ij}^2$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	0	1	4	9	16	30
0.01-0.20.....	1	0	1	4	9	15
0.21-0.50.....	4	1	0	1	4	10
0.51-1.00.....	9	4	1	0	1	15
≥1.01.....	16	9	4	1	0	30

TABLE 4C.—Array of number of cases expected by chance in each cell:  $\frac{n_{i.}n_{.j}}{T}$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	16.51	10.52	7.87	6.29	4.19	45.38
0.01-0.20.....	13.64	3.22	2.37	1.86	1.02	22.11
0.21-0.50.....	8.80	2.54	1.72	1.35	0.75	15.16
0.51-1.00.....	6.15	1.45	1.06	0.82	0.46	10.04
≥1.01.....	4.02	1.17	0.86	0.68	0.37	7.10

TABLE 4D.—Array of chance frequencies weighted by squared cell deviation:  $\frac{n_{i.}n_{.j}}{T} d_{ij}^2$ =(table 4B) (table 4C)

$$\sum \left( \frac{n_{i.}n_{.j}}{T} d_{ij}^2 \right) = \sum d_{ij}^2 = 387.4$$

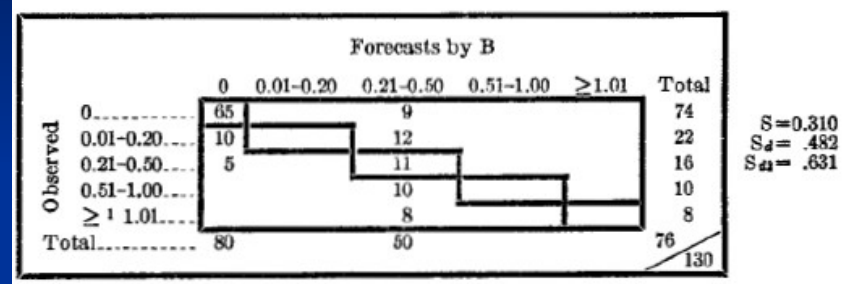
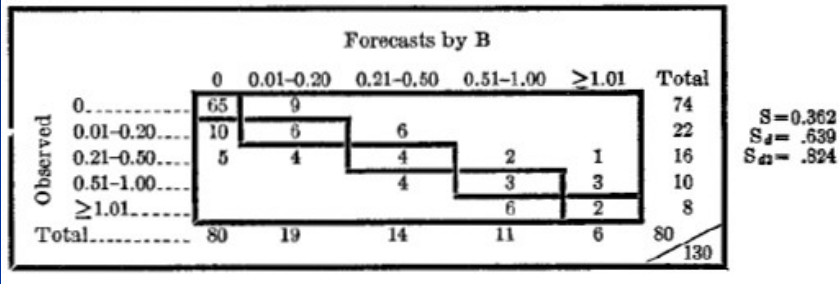
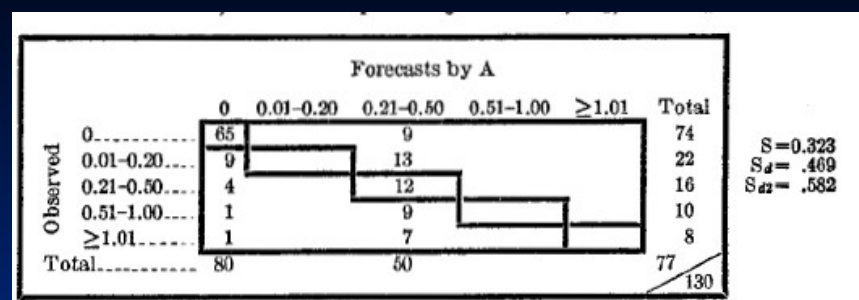
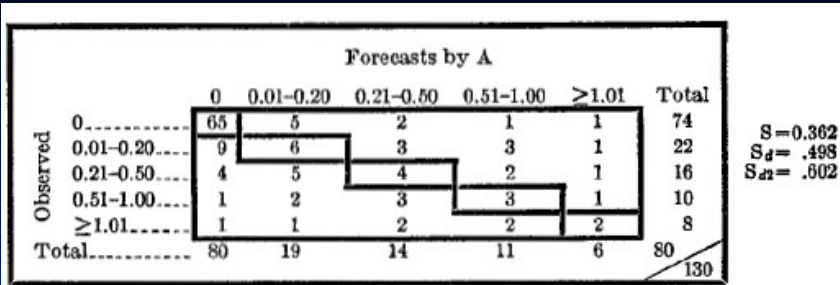
Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	0	10.52	15.74	12.75	12.08	51.67
0.01-0.20.....	13.64	0	2.37	7.44	9.18	32.63
0.21-0.50.....	34.40	2.54	0	1.35	2.92	41.21
0.51-1.00.....	55.35	6.84	1.06	0	0.46	63.71
≥1.01.....	25.72	10.53	3.44	0.68	0	40.37

TABLE 4E.—Array of forecast frequencies weighted by square of cell deviation:  $n_{ij}d_{ij}^2$ =(table 4A) (table 4B)

$$\sum (n_{ij}d_{ij}^2) = \sum d_{ij}^2 = 154$$

Observed	Forecast					Total
	0	0.01-0.20	0.21-0.50	0.51-1.00	>1.00	
0.....	0	5	8	12	9	34
0.01-0.20.....	9	0	3	12	9	33
0.21-0.50.....	16	3	0	2	4	25
0.51-1.00.....	9	8	3	0	1	21
≥1.01.....	16	9	8	2	0	35

$$S_{d^2} = \frac{\sum d_{ij}^2 - \sum d_{ij}^2}{\sum d_{ij}^2} = \frac{387.4 - 154}{387.4} = 602$$



- possibility that rating on the basis of the size of the deviations will lead forecasters to bias their forecasts towards the middle class interval (where is a minimum possible deviation), rather than rather than trying to catch extreme

# ...Bryan...

- Bryan, J. G., and I. Enger, 1967: Use of probability forecasts to maximize various skill scores. *J. Appl. Meteor.*, 6, 762-769

768

JOURNAL OF APPLIED METEOROLOGY

VOLUME 6

by applying Eq. (3.11) to the same probability distribution used in Section 2C.

$v=0.0$			$v=0.2$			$v=0.4$			
1.00	0.00	-1.00	0.80	-0.20	-1.20	0.60	-0.40	-1.40	
-0.46	0.54	-0.46	-0.57	0.43	-0.57	-0.68	0.32	-0.68	
-1.00	0.00	1.00	-1.20	-0.20	0.80	-1.40	-0.40	0.60	
$v=0.6$			$v=0.8$			$v=1.0$			
0.40	-0.60	-1.60	0.20	-0.80	-1.80	0.00	-1.00	-2.00	
-0.78	0.22	-0.78	-0.89	0.11	-0.89	-1.00	0.00	-1.00	
-1.60	-0.60	0.40	-1.80	-0.80	0.20	-2.00	-1.00	0.00	
$v=0.2$			$v=0.5$						
1.60	0.60	-0.40	-1.40	-2.40	1.00	0.00	-1.00	-2.00	-3.00
-0.03	0.97	-0.03	-1.03	-2.03	-0.39	0.61	-0.39	-1.39	-2.39
-1.27	-0.27	0.73	-0.27	-1.27	-1.54	-0.54	0.46	-0.54	-1.54
-2.03	-1.03	-0.03	0.97	-0.03	-2.39	-1.39	-0.39	0.61	-0.39
-2.40	-1.40	-0.40	0.60	1.60	-3.00	-2.00	-1.00	0.00	1.00
$v=0.8$									
	0.40	-0.60	-1.60	-2.60	-3.60				
	-0.76	0.24	-0.76	-1.76	-2.76				
	-1.82	-0.82	0.18	-0.82	-1.82				
	-2.76	-1.76	-0.76	0.24	-0.76				
	-3.60	-2.60	-1.60	-0.60	0.40				



# Gringorten, (1967)

- Gringorten score awards forecasts of extremes (too much?)
- compared to Bryan, bad forecasts are not equivalently punished

TABLE 3. Example of scores for the forecast of mutually exclusive events  $X_0, X_1, X_2, X_3$  whose climatic frequencies are  $cP_0=0.05, cP_1=0.10, cP_2=0.10, cP_3=0.75$ .

Forecast event	Bryan Observed event				Gringorten Observed event			
	$X_0$	$X_1$	$X_2$	$X_3$	$X_0$	$X_1$	$X_2$	$X_3$
$X_0$	95	-5	-5	-5	20	0	0	0
$X_1$	-10	90	-10	-10	0	10	0	0
$X_2$	-10	-10	90	-10	0	0	10	0
$X_3$	-75	-75	-75	25	0	0	0	1.33

# LEPS

- Linear Error in Probability Space
- originally introduced by Ward and Folland (1991)

$$\text{LEPS} = (1/n) \sum |Pv -$$

$Pf |$

- corresponds to MAE transformed into probability space
- Revised, normalised LEPS (Potts et al., 1996)

# LEPS...

- encourages forecasting of extremes of the climatological distribution (not as much as Gringorten)

TABLE 2. LEPS  $S''$  scores for quintiles.

Forecast	Observation				
	Q1	Q2	Q3	Q4	Q5
Q1	1.28	0.52	-0.20	-0.68	-0.92
Q2	0.52	0.56	0.04	-0.44	-0.68
Q3	-0.20	0.04	0.32	0.04	-0.20
Q4	-0.68	-0.44	0.04	0.56	0.52
Q5	-0.92	-0.68	-0.20	0.52	1.28

# ETSS

- proposed by Gandin and Murphy (1992)
- “equitability”
  - equitable score takes zero value (“no skill”) for random forecasts and for unvarying forecasts of a constant category
  - takes value 1 for perfect forecast

# ETSS

- Gerrity (1992) showed that, for a multcategory table  $N \times N$ , ETSS can be calculated as mean of  $N-1$  values of Pierce Skill Score of collapsed  $2 \times 2$  tables
- since PSS has some deficiencies, ETSS will probably inherit them?

# ETSS

- ECMWF 24-hour precipitation forecast (Green Book, 2007, Hungary)
- 4 categories
  - C0 less than 0.1 mm
  - C1 0.1 - 2 mm
  - C2 2 - 10 mm
  - C3 bigger than 10 mm

# ETSS

F\O	C0	C1	C2	C3
C0	5812	179	46	5
C1	3292	1165	383	56
C2	608	735	971	279
C3	15	45	247	326



- ETSS = .561

# ETSS - modified

F\O	C0	C1	C2	C3
C0	5812	179	46	5
C1	3292	1165	383	56
C2	0	0	0	0
C3	623	780	1218	605



- ETSS = .643



# ...scoring matrices...

- Based on marginal probabilities one can calculate the scoring matrix

$$P_0=0.687$$

$$P_1=0.150$$

$$P_2=0.116$$

$$P_3=0.047$$

F\O	C0	C1	C2	C3
C0	.233	-.251	-.650	-1
C1	-.252	.812	.414	.064
C2	-.650	.414	2.455	2.105
C3	-1	.064	2.105	9.194

# ...scoring matrices...

- Based on marginal probabilities one can calculate the scoring matrix

$$P_0=0.687$$

$$P_1=0.150$$

$$P_2=0.116$$

$$P_3=0.047$$

F\O	C0	C1	C2	C3
C0	.233	-.251	-.650	-1
C1	-.252	.812	.414	.064
C2	-.650	.414	2.455	2.105
C3	-1	.064	2.105	9.194



# ETSS

- dependent on the number of categories
- if a  $N \times N$  contingency table is reduced to  
 $(N-1) \times (N-1) \dots$   
...and finally to  $2 \times 2$   
ETSS is constantly decreasing...

# ETSS

- favorizing rare events...
- dependent on bias, which may lead to “hedging”
- overforecasting of extreme (and rare) categories leads to increase of the score compared to the original forecast
- equitable
- proper?
- suitable for “hedging”

# ...on hedging...

