

## Combining Precipitation Probabilities

LAWRENCE A. HUGHES AND WAYNE E. SANGSTER

*NOAA National Weather Service, Kansas City, MO 64106*

(Manuscript received 16 June 1978, in final form 16 January 1979)

### ABSTRACT

Two methods are discussed for combining the routine forecasts of the 12 h probability of precipitation made by the National Weather Service, for use when longer period probabilities are desired but cannot be created independently. Both apply a year's forecasts from 28 forecast offices to basic equations of probability to adjust for the obvious dependence of the precipitation events among the forecast periods. Both methods suggest that warm season precipitation events are more independent than cold season ones, as would be expected. One method gave unrealistic results for probability combinations outside the range of those actually used. The other method applied realistic constraints to eliminate this undesirable feature. The largest deviations from probabilities for independent events occurred when combining probabilities of 60%, but the deviations were only about 5% in the warm season and 10% in the cold season. Tables and an equation for combining probabilities are given.

### 1. Introduction

The U.S. National Weather Service (NWS) has been issuing probability of precipitation (PoP) forecasts for a number of years. These consist of the PoP for each of three consecutive 12 h periods, or a 6 h period followed by two 12 h periods. However, some users, especially those in agriculture and forestry, are interested in, say, the PoP for the next 24 or 36 h. NWS forecasters could formulate special forecasts for these longer periods, but unless they are sophisticated in the ways of probability, they might well arrive at unreliable probability forecasts. For example, a forecast of 60% today and 40% tonight might be erroneously converted to a forecast as high as 100% for the next 24 h. Forecasts of 40% today, 40% tonight and 40% tomorrow night might also lead incorrectly to a forecast as high as 100% for the next 36 h. Model Output Statistics (MOS) procedures could be used to generate such forecasts, but unless the demand was sufficiently great, the computer and communications time and extra display space required would not be justified. Therefore, it seems that objective guidelines for combining probabilities are needed and would fill a void. It is the purpose of this paper to provide some insight into this problem and guidelines for the forecaster.

### 2. Brute force method

If the occurrence of precipitation for the first 12 h period (event A) and the occurrence of precipitation for the second 12 h period (event B) are assumed independent (which they are not), we may write

the following [for this and other basic relationships, see, e.g., Hald (1952)]:

$$P(A \cup B) = P(A) + P(B) - P(A)P(B). \quad (1)$$

Extending (1) to include a third 12 h period (event C), we have

$$\begin{aligned} P(A \cup B \cup C) = & P(A) + P(B) + P(C) \\ & - P(A)P(B) - P(A)P(C) - P(B)P(C) \\ & + P(A)P(B)P(C). \quad (2) \end{aligned}$$

Here  $P(A \cup B)$  and  $P(A \cup B \cup C)$  are the PoP's for the 24 and 36 h periods, and  $P(A)$ ,  $P(B)$  and  $P(C)$  are the 12 h PoP's. Our first attempt at combining probabilities was just to assume that  $P(A \cup B \cup C)$  was a linear function of the three 12 h PoP's and all products thereof as in the equation for independent events (2), but with coefficients other than plus and minus one. The coefficients can be determined by usual multiple-regression analysis methods. A constant term was also included because it was computed by the regression program used, and it allows for unreliable low probabilities.

The actual locally made forecasts for 28 stations distributed rather uniformly across the north-central United States for one warm season (April–September) and one cold season (October–March) were treated. Each of the seven terms in (2) was determined from the actual forecasts and then fed into the regression program as predictors, with the observed precipitation as the predictand. The regression estimate of the 36 h probability for the warm season is

$$\hat{P}(A \cup B \cup C) = 0.03 + 1.04P(A) + 0.81P(B) + 0.97P(C) - 0.86P(A)P(B) - 1.36P(A)P(C) - 0.76P(B)P(C) + 1.19P(A)P(B)P(C). \quad (3)$$

Note that the coefficients are all of the same sign as in the equation for independent events (2) and are fairly close to 1 in magnitude. With this equation one can calculate the PoP in a 36 h period for any warm season set of three 12 h PoP's. This was done for all possible combinations of the PoP's in the set used by the NWS. The 36 h PoP from (3) was then compared with the 36 h PoP determined from (2). It was found that in the warm season the PoP's from the two equations were nearly the same (the vast majority of values were within 0.05 of each other over the range of PoP combinations actually used).

In the cold season equation, the coefficients were not as close to 1 as those of (3), and the signs were not always the same as in (2). Using this equation gave differences from (2) significantly larger than those for warm season forecasts.

One problem with (3) is that in order to use tabular data so as to eliminate the need to calculate the probability each time, one would need a table for each possible PoP used in the forecasts—13 in all. A more significant problem with this approach is that some combinations of probabilities were unrealistic. For example, for a given first and second period PoP the 36 h PoP decreased as the third period PoP increased. This effect was observed mainly in the cold season and in PoP combinations that were never used by the forecasters. It was probably caused by extrapolating beyond the range of actual forecast combinations from forecasts that were not perfectly reliable. Nevertheless, the characteristic is undesirable. Another defect of (3) is that it doesn't give zero when all PoP's are zero (the result was 0.03) or 1.0 when one or more probabilities were 100%. Both of these defects are probably the result of slightly unreliable forecasts, but the latter is mainly due to extrapolation from such forecasts to probability combinations not actually used.

### 3. Exponent method

Since the approach used in the previous section has deficiencies, another approach was tried. The addition formula for probabilities *without* the assumption of independence may be written

$$P(A \cup B) = P(A) + P(B) - P(AB). \quad (4)$$

Here  $P(AB)$  is the probability that *both* periods will receive precipitation. At this point we are no better off, since  $P(AB)$  is not known any more than  $P(A \cup B)$ . However, the multiplication formula gives

$$P(AB) = P(A)P(B/A) \quad (5)$$

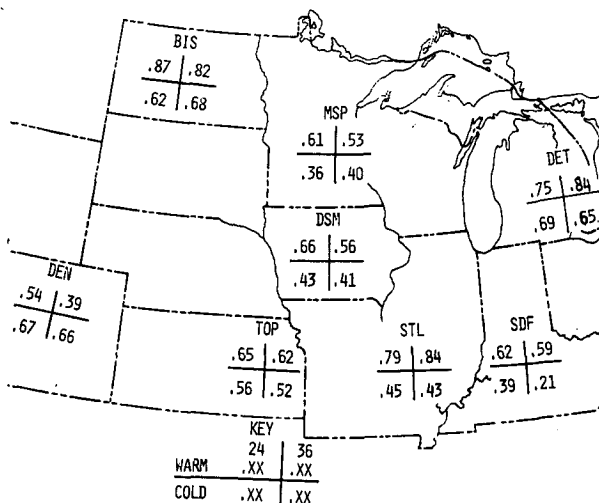


FIG. 1. Values for  $k$  for selected stations.

and

$$P(AB) = P(B)P(A/B), \quad (6)$$

where  $P(B/A)$  and  $P(A/B)$  are conditional probabilities of precipitation in the second period, given that it rains in the first period, and of precipitation in the first period, given that it precipitates in the second period. Although these conditional probabilities are not known, we can set up a statistical model such that estimates of the conditional probabilities are

$$\hat{P}(B/A) = [P(B)]^k \quad (7)$$

or

$$\hat{P}(A/B) = [P(A)]^k. \quad (8)$$

Eqs. (7) and (8) meet the requirements at the extremes (zero and 1) that the conditional probability be equal to the unconditional probability. However, we still must have some rule for determining when to use (7) and when to use (8). We note that if  $P(A) = 1.0$ ,  $P(AB) = P(B)$  and if  $P(B) = 1.0$ ,  $P(AB) = P(A)$ . These conditions are not met if the exponent is attached to the smaller of  $P(A)$  or  $P(B)$ . For example, if  $P(A) = 1.0$  and  $P(B) = 0.5$  the estimate is  $\hat{P}(AB) = 1.0(0.5)^k$ , only when  $k = 1$ , since  $P(AB) = 0.5$ . On the other hand,  $\hat{P}(AB) = 1.0^k(0.5) = 0.5$  for any  $k$ .

Therefore, our final statistical model is

$$\hat{P}(B/A) = [P(B)]^k, \text{ when } P(B) > P(A), \quad (9)$$

$$\hat{P}(A/B) = [P(A)]^k, \text{ when } P(A) > P(B). \quad (10)$$

For  $P(A) = P(B)$  it obviously makes no difference which of (9) or (10) is used. This model gives everything needed to determine  $P(A \cup B)$  from (4), except  $k$ .

From the PoP forecasts and precipitation events at eight stations (see Fig. 1),  $k$  was determined for each station for each of the two seasons for the

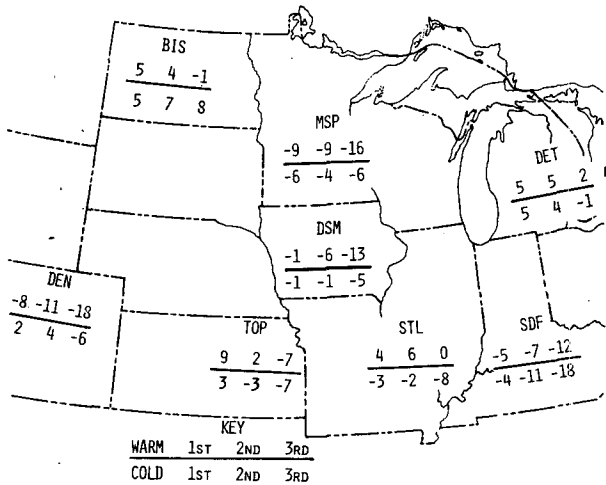


FIG. 2. Bias values for three periods (percent).

period October 1966–September 1973 (seven warm and seven cold seasons). To do this, first the observed relative frequency of precipitation in both periods  $H(AB)$  was determined for each combination of PoP's used by the forecasters at each station. Then, using (9) and (10) to compute  $\hat{P}(AB)$ , a

value of  $k$  was computed by a least-squares method by minimizing the sum of  $[H(AB) - \hat{P}(AB)]^2$  weighted according to the number of forecasts in each combination. This was done simply by choosing a small value of  $k$ , computing the sum of the squares, then increasing  $k$  by 0.01 and repeating the operation. Eventually a minimum value of the sum of the squares was reached, yielding the desired value of  $k$ .

The first pass through the data was used to obtain a 24 h PoP from the first two 12 h PoP's. This 24 h PoP was then rounded to the nearest standard forecast PoP value, as used by the Weather Service. The process was then repeated using this 24 h PoP and the third 12 h period PoP. This procedure resulted in a probability of measurable precipitation sometime in the 36 h period. The values of  $k$  obtained are shown in Fig. 1, with two for the warm season and two for the cold season for each of the eight stations. The values ranged from a low of 0.21 at Louisville for the 36 h cold season to a high of 0.87 at Bismarck for the 24 h warm season.

Since the results obtained are based on the assumption that the forecasts used had no appreciable bias, the forecasts were examined to determine the

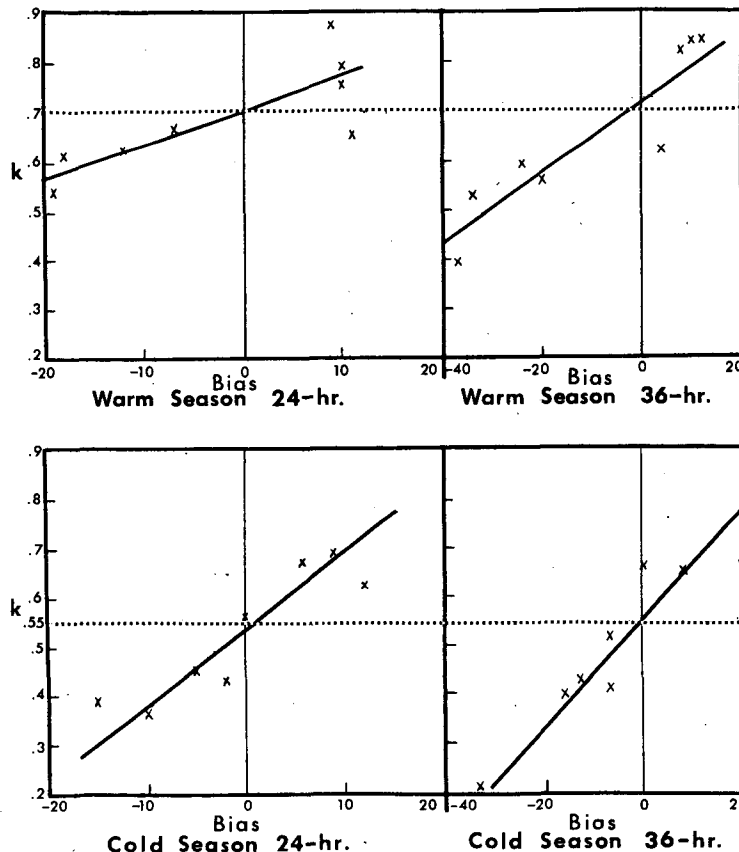


FIG. 3. Graphs of  $k$  versus bias, with first two periods for 24 h and all three periods for 36 h.

TABLE 1. Probabilities (percent) for  $k$  of 0.70 (warm season).

PROB1	PROB2												
	0	2	5	10	20	30	40	50	60	70	80	90	100
0	0	2	5	10	20	30	40	50	60	70	80	90	100
2	2	4	7	12	21	31	41	51	61	70	80	90	100
5	5	7	9	14	23	33	42	52	62	71	81	90	100
10	10	12	14	18	27	36	45	54	63	72	81	91	100
20	20	21	23	27	34	41	49	58	66	74	83	91	100
30	30	31	33	36	41	47	54	62	69	77	84	92	100
40	40	41	42	45	49	54	59	65	72	79	86	93	100
50	50	51	52	54	58	62	65	69	75	81	87	94	100
60	60	61	62	63	66	69	72	75	78	83	89	94	100
70	70	70	71	72	74	77	79	81	83	85	90	95	100
80	80	80	81	83	84	86	87	89	90	92	96	100	100
90	90	90	90	91	91	92	93	94	94	95	96	96	100
100	100	100	100	100	100	100	100	100	100	100	100	100	100

bias. The bias is the average forecast probability minus the observed precipitation frequency. This was expressed as a percentage of the observed frequency, and is shown in Fig. 2 for each season for the first, second and third periods of the forecast. An inspection of the values of the bias and  $k$  shows that the low values of  $k$  for the cold season at Louisville are associated with negative bias (underforecasting), especially in the second and third periods of the forecast, and the high values of  $k$  at Bismarck were associated with positive bias (overforecasting) in the first and second periods.

A plot of the 24 and 36 h  $k$  values versus the bias figures is shown in Fig. 3. It is apparent that there is a strong correlation between bias and  $k$ . While  $k$  undoubtedly varies at least somewhat from one place to another and from one time to another due to factors other than bias, one could, as at least a first approximation, say that if future forecasts are assumed to be unbiased, both the 24 and 36 h  $k$ 's should be about 0.70 in the warm season and 0.55 in the cold season.

Tables of probabilities for  $k$  of 0.70 and 0.55 are shown as Tables 1 and 2. To use the tables, take

TABLE 2. Probabilities (percent) for  $k$  of 0.55 (cold season).

PROB1	PROB2												
	0	2	5	10	20	30	40	50	60	70	80	90	100
0	0	2	5	10	20	30	40	50	60	70	80	90	100
2	2	4	7	11	21	31	41	51	60	70	80	90	100
5	5	7	9	14	23	32	42	52	61	71	81	90	100
10	10	11	14	17	26	35	44	53	62	72	81	91	100
20	20	21	23	26	32	40	48	56	65	74	82	91	100
30	30	31	32	35	40	45	52	60	67	75	83	92	100
40	40	41	42	44	48	52	56	63	70	77	85	92	100
50	50	51	52	53	56	60	63	66	72	79	86	93	100
60	60	60	61	62	65	67	70	72	75	81	87	93	100
70	70	70	71	72	74	75	77	79	81	82	88	94	100
80	80	80	81	81	82	83	85	86	87	88	89	95	100
90	90	90	90	91	91	92	92	93	93	94	95	95	100
100	100	100	100	100	100	100	100	100	100	100	100	100	100

TABLE 3. Probabilities from Table 1 minus those for independent events.

PROB1	PROB2												
	0	2	5	10	20	30	40	50	60	70	80	90	100
0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	-0	-0	-0	-0	-0	-0	-0	-0	-0	-0	-0	0
5	0	-0	-0	-0	-1	-1	-1	-1	-1	-0	-0	-0	0
10	0	-0	-0	-1	-1	-1	-1	-1	-1	-1	-1	-1	0
20	0	-0	-1	-1	-2	-3	-3	-2	-2	-2	-1	-1	0
30	0	-0	-1	-1	-3	-4	-4	-3	-3	-2	-2	-1	0
40	0	-0	-1	-1	-3	-4	-5	-5	-4	-3	-2	-1	0
50	0	-0	-1	-1	-2	-3	-5	-6	-5	-4	-3	-1	0
60	0	-0	-0	-1	-2	-3	-4	-5	-6	-5	-3	-2	0
70	0	-0	-0	-1	-2	-2	-3	-4	-5	-6	-4	-2	0
80	0	-0	-0	-1	-1	-2	-2	-3	-3	-4	-4	-2	0
90	0	-0	-0	-0	-1	-1	-1	-1	-2	-2	-2	-3	0
100	0	0	0	0	0	0	0	0	0	0	0	0	0

the first 12 h PoP as PROB1 and the second 12 h probability as PROB2 to find the 24 h PoP. For combining more than two PoP's, use the table iteratively. For example, for a 36 h PoP, use the 24 h PoP as PROB1 (interpolating as needed) and the third 12 h PoP as PROB2. Actually, the order of using the probabilities is immaterial.

If one of the periods is a 6 h period, there is slightly more dependence among the periods than if all three contiguous periods were 12 h. Thus the  $k$  value would be slightly lower than the values indicated and the resultant 30 h probability obtained by use of Tables 1 and 2 would be slightly too high, but probably only by 1% or so as a maximum. This would occur when both probabilities were middle values.

Tables 3 and 4 give the values in Tables 1 and 2 minus the values that would be obtained for independent events. Notice that the maximum difference is in the middle probability values (actually centered on 60, 60), with no difference on the borders of the table. The latter result is due to the method of determining  $P(AB)$ , but would be true for any reliable forecasts.

TABLE 4. Probabilities from Table 2 minus those for independent events.

PROB1	PROB2												
	0	2	5	10	20	30	40	50	60	70	80	90	100
0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	-0	-0	-0	-0	-0	-0	-0	-0	-0	-0	-0	0
5	0	-0	-1	-1	-1	-1	-1	-1	-1	-1	-1	-0	0
10	0	-0	-1	-2	-2	-2	-2	-2	-2	-2	-1	-1	0
20	0	-0	-1	-2	-4	-4	-4	-4	-3	-2	-2	-1	0
30	0	-0	-1	-2	-4	-6	-6	-5	-5	-4	-3	-1	0
40	0	-0	-1	-2	-4	-6	-8	-7	-6	-5	-3	-2	0
50	0	-0	-1	-2	-4	-5	-7	-9	-8	-6	-4	-2	0
60	0	-0	-1	-2	-3	-5	-6	-8	-9	-7	-5	-3	0
70	0	-0	-1	-2	-4	-5	-6	-7	-9	-6	-3	0	0
80	0	-0	-0	-1	-2	-3	-3	-4	-5	-6	-7	-3	0
90	0	-0	-0	-0	-1	-1	-2	-2	-3	-3	-3	-4	0
100	0	0	0	0	0	0	0	0	0	0	0	0	0

The difference would get smaller in absolute value for  $k > 0.70$  (less dependence) and larger in absolute value for  $k < 0.55$ . For example, the largest absolute value of the difference for a  $k$  value of 0.40 is 13%. The difference would always be negative or zero for  $k$  values  $< 1.00$ .

A seasonal contrast is shown in the 24 and 36 h probabilities in that the maximum difference in the cold season ( $k = 0.55$ ) reveals more dependence. But note that the differences are always negative, i.e., the 36 h probability for dependent events is less than that for independent events, and are quite small, especially for warm season forecasts even in the mid-probability range. The seasonal contrast is less pronounced in the results obtained by the second method of determining the probability, in that the difference resulting from the second method is a bit smaller in absolute value in the cold season and a bit larger in the warm season than that obtained by the first method.

The method discussed in this section is valid anywhere, but the specific values of  $k$  might change significantly in radically different precipitation climates such as the west coast of the United States.

#### 4. Conclusions

Either method discussed here could be used with probability forecasts from any part of the world to create equations or tables for any climatology. However, the equation and tables given here should be reasonably applicable in the United States east of the Rockies and in similar climates elsewhere. The tables derived from the second method should be more universal because they are reasonable outside the range of the forecast values actually used. The first method gave a larger difference between the two seasons, and had smaller warm season differences from the equation for independent events.

The best way to create 24 and 36 h probabilities would be directly by MOS or by forecasters (after a suitable training and verification period). That time is not near. In the meantime, either of the above methods will prevent major errors by forecasters responding to occasional requests by giving a reasonable estimate to a longer term probability.

#### REFERENCE

- Hald, A., 1952: *Statistical Theory with Engineering Applications*. Wiley 783 pp. (see pp. 12–15).