

## NOTES

## Retrieval of Water Vapor Profiles via Principal Components: Options and Their Implications

ALAN E. LIPTON AND THOMAS H. VONDER HAAR

*Department of Atmospheric Science, Colorado State University, Fort Collins, CO 80523*

18 August 1985 and 17 December 1986

## ABSTRACT

Principal components have been widely used in regression retrieval of atmospheric parameters, but when applied to water vapor concentrations their use entails special problems. We discuss two of these problems, and present results of retrieval experiments designed to alleviate them. The experiments employed High-resolution Infrared Radiation Sounder satellite data in conjunction with radiosonde observations. We found that mixing ratio is a less appropriate parameter for principal component-based retrieval than is a mean-saturation adjusted mixing ratio. Also, retrieval accuracy was enhanced by identifying the optimum numbers of eigenvectors to use when transforming the water vapor profiles and the satellite brightness temperatures, respectively, into their principal components. In our studies three eigenvectors were optimal for representation of water vapor, implying that HIRS-2 data are capable of retrieving at least third-order vertical resolution in water vapor profiles. In addition, we compared principal component-based retrieval with standard multiple regression and found that a hybrid of the two methods gave the greatest retrieval accuracy.

### 1. Introduction

Regression via principal components has been popularly used for retrieving data on atmospheric properties. One method of this type, often called "eigenvector" retrieval, was applied by Smith and Woolf (1976) to retrieval of temperature and water vapor profiles from satellites, and was designed to minimize interference by random errors.

There are problems particularly relevant to water vapor profile retrieval that were not addressed by Smith and Woolf. We examined two of them:

1) The mixing ratio, as a descriptor of water vapor concentrations, is not well suited to retrieval based on principal components.

2) It is not obvious what is the optimal number of modes (or eigenvectors) to use in approximating the satellite brightness temperatures and in approximating the water vapor profiles. Should both approximations involve the same number of modes?

We experimented with retrieval of water vapor profiles from High-resolution Infrared Radiation Sounder (HIRS-2) data in order to quantify the importance of these problems, and as a by-product, we were able to infer the vertical resolution in water vapor profiles that may be retrieved from HIRS-2 brightness temperatures. Also, principal component-based schemes were compared with a simpler type of multiple regression and also with a hybrid method.

The experimental data consisted of radiosonde observations (RAOB) paired with nearly coincident

brightness temperature measurements from HIRS-2 (on board the NOAA-series satellites). The data included a set of 316 midlatitude (30°–60°N) soundings and a similar set of 191 tropical (30°N–30°S) soundings. All soundings were taken in December 1982 at sites where the surface pressure was at least 100 kPa, to ensure complete statistics at all of the considered levels. The radiosonde data were interpolated to 13 fixed levels from 40 to 100 kPa. The HIRS-2 data had been processed according to the operational procedure at the National Environmental Satellite Data and Information Service (NESDIS), including calibration, scan angle compensation, and determination of the clear column radiance in partly cloudy conditions (Werbowetzki, 1981). No completely cloudy soundings were considered. Eighteen of the twenty HIRS-2 channels were used for retrieval (see Table 1). The ozone channel and the 0.7- $\mu\text{m}$  visible channel were excluded.

The rest of this Introduction employs the notation of Smith and Woolf (1976), with some modifications. A water vapor mixing ratio profile, represented by a vector  $\mathbf{q}$ , is related to its principal component vector  $\mathbf{a}$  by

$$\mathbf{q} = \mathbf{Q}^* \mathbf{a}, \quad (1)$$

where  $\mathbf{Q}^*$  is a matrix whose columns are the eigenvectors ordered such that the one corresponding to the largest eigenvalue is first, and so on to the smallest (Kendall, 1975, chapter 2). Smith and Woolf referred to the principal components as "expansion coefficients".

TABLE 1. HIRS-2 channel characteristics (After Smith et al., 1979).

Channel	Central wave-length ( $\mu\text{m}$ )	Principal absorbing gas	Channel	Central wave-length ( $\mu\text{m}$ )	Principal absorbing gas
1	15.00	CO <sub>2</sub>	11	7.30	H <sub>2</sub> O
2	14.70	CO <sub>2</sub>	12	6.70	H <sub>2</sub> O
3	14.50	CO <sub>2</sub>	13	4.57	N <sub>2</sub> O
4	14.20	CO <sub>2</sub>	14	4.52	N <sub>2</sub> O
5	14.00	CO <sub>2</sub>	15	4.46	CO <sub>2</sub> /N <sub>2</sub> O
6	13.70	CO <sub>2</sub> /H <sub>2</sub> O	16	4.40	CO <sub>2</sub> /N <sub>2</sub> O
7	13.40	CO <sub>2</sub> /H <sub>2</sub> O	17	4.24	CO <sub>2</sub>
8	11.10	window	18	4.00	window
10	8.30	H <sub>2</sub> O	19	3.70	window

Elimination from (1) of one or more eigenvectors and their corresponding principal components results in the approximation

$$\mathbf{q} \approx \mathbf{Q}_i^* \mathbf{a}_i, \quad (2)$$

where  $i$  is the number of modes retained. The analogous approximation for a corresponding set of brightness temperature observations  $\mathbf{t}_B$  from a satellite is

$$\mathbf{t}_B \approx \mathbf{T}_{Bj}^* \mathbf{b}_j, \quad (3)$$

where  $\mathbf{T}_{Bj}^*$  is the matrix of brightness temperature eigenvectors,  $\mathbf{b}$  is the vector of corresponding principal components, and  $j$  is the number of modes retained. The approximations given by (2) and (3) are useful only if the first eigenvectors correspond to "real" modes

in atmospheric variability and the last eigenvectors correspond primarily to modes of random observational error.

In the Smith and Woolf retrieval method, the vector  $\mathbf{b}_j$  must first be computed from the inverse of (3),  $\mathbf{a}_i$  is determined from  $\mathbf{b}_j$  by regression, and  $\mathbf{a}_i$  is then transformed into  $\mathbf{q}$  by (2). Any modes of  $\mathbf{q}$  or  $\mathbf{t}_B$  that are counterproductive to regression accuracy can be eliminated by this procedure. A commonly used alternative for regression problems (Kendall, 1975, pp. 95-98) is identical to this method, except that only the predictors are transformed. Transformation (2) has also been shown to be useful in nonregression retrieval methods (Weinreb and Crosby, 1977).

## 2. Experiments and results

### a. Adjustment of mixing ratios

Absolute humidity generally decreases rapidly from the earth's surface upward. Accordingly, the variance of mixing ratio values is greatest at the lowest atmospheric levels, and it can be expected that the first eigenvectors in a set will be dominated by these lowest levels. This is illustrated by Fig. 1a, which includes the first four (of thirteen total) eigenvectors of the mid-latitude RAOBs. It is apparent from Fig. 1a that eigenvectors 1 through 4 are dominated by levels below 60 kPa. In fact, the first eigenvector with a contribution greater than 0.3 at the 40 kPa level is number 11 (not shown). Apparently, modes 11 through 13 contain the

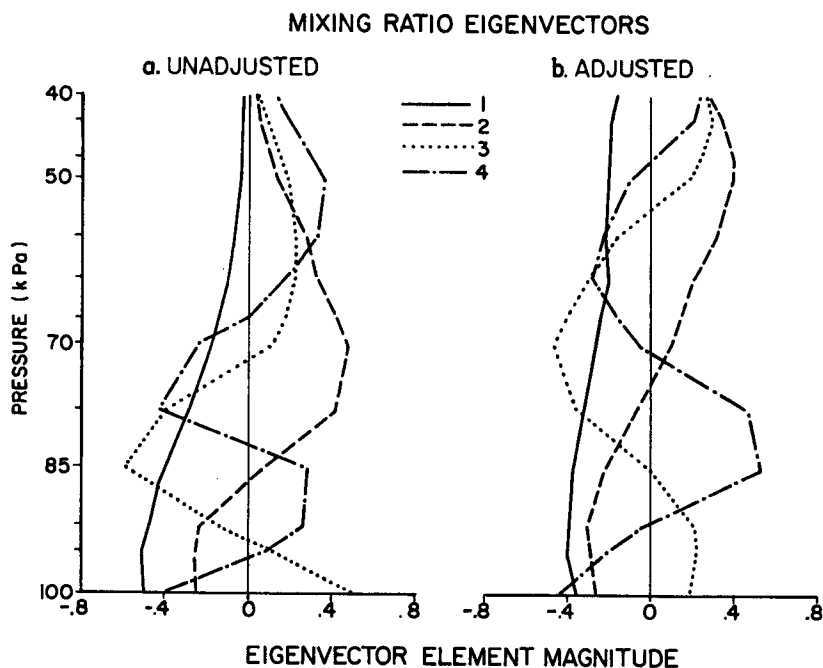


FIG. 1. Eigenvectors of mixing ratio covariance matrices for a) unadjusted and b) adjusted mixing ratios. The first 4 (of 13) eigenvectors are shown for each. Magnitudes are given at each of the 13 levels indicated by tick marks on the vertical axis.

information needed to accurately describe the mixing ratio at 40 kPa.

It is problematic to attempt to eliminate random error in measurements of  $q$  by transforming to  $a$ , and dropping the last modes. The problem stems from the fact that both the variance and the measurement error in mixing ratios are not distributed uniformly among atmospheric levels. Radiosondes actually measure relative humidity (RH), not mixing ratio, so a 5% RH error translates into a much bigger mixing ratio error at 100 kPa than it does at 40 kPa. Ideally, both the variance and error of the satellite-retrieved water vapor parameter should be uniformly distributed with height. One approach toward that goal would be to retrieve RH, but deriving mixing ratios from RH would introduce errors due to uncertainties in the estimation of the corresponding temperature profile. A compromise was attempted in order to create a parameter independent of the local temperature, yet distributed similarly to RH. The temperature profiles from each of the RAOB sets (midlatitude and tropical) were separately averaged, and the mean profile was used to compute a mean saturation mixing ratio profile  $\bar{q}_s$  by use of the Clausius-Clapeyron equation. All mixing ratios in each set were then adjusted to the mean saturation profile by dividing the elements of  $q$  by their corresponding elements of  $\bar{q}_s$ . The resulting vectors  $q'$  are hereafter referred to as adjusted mixing ratios, and the prime symbol denotes quantities which result from this adjustment.

The  $q'$  vectors of the midlatitude set were subjected to statistical analysis, resulting in the eigenvectors shown in Fig. 1b. The adjusted eigenvectors are much more uniformly distributed with respect to pressure than the unadjusted ones, so  $q'$  should be better suited to principal component-based retrieval than  $q$ .

To determine whether or not retrieval accuracy benefits from adjustment of mixing ratios, retrievals were performed with both adjusted and unadjusted data. We used an averaged explained variance as the measure of success, where the average was computed over the 13 levels with weighting appropriate to the pressure difference between levels. The midlatitude and tropical sets were each randomly divided so that approximately two-thirds were used to compute regression coefficients (dependent data), and the other third was reserved as independent data for testing retrieval skill. (However, these data were not completely independent, since they had been used to compute means and eigenvectors.) The experiments were conducted for a fixed number of mixing ratio modes ( $i = 3$ ), and for a range of brightness temperature modes ( $1 \leq j \leq 18$ ), retaining the modes with the greatest eigenvalues. Figure 2 shows results for principal component-based retrieval for both adjusted and unadjusted data. The choice of limiting retrievals to  $i = 3$  is justified in subsection 2b. Comparison of the explained variance maxima shows 1.7 and 1.9% improvements resulting from adjustment for

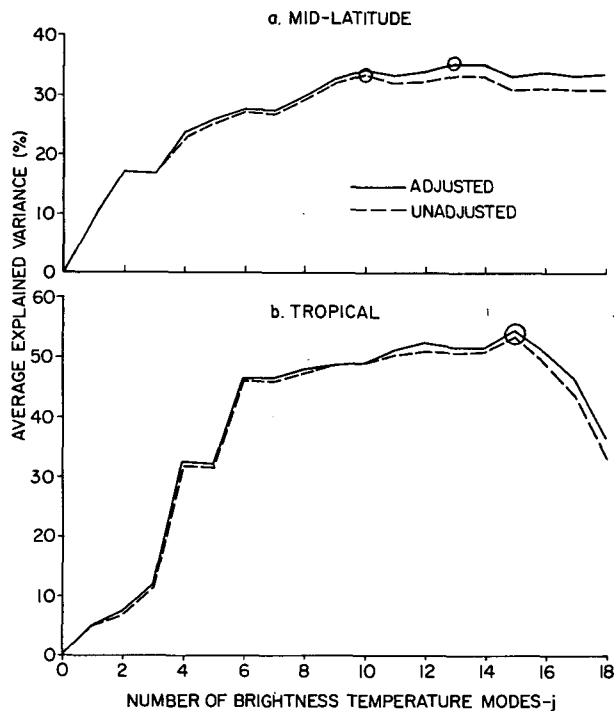


FIG. 2. Comparison of mixing ratio retrieval results for mean-saturation adjusted and unadjusted mixing ratios for the a) midlatitude, and b) tropical data sets. Explained variances are from independent data with  $i = 3$ . The "adjusted" curves are equivalent to cross sections through Fig. 3. The circles indicate where the curves reach their maxima.

the midlatitude and tropical sets, respectively. These improvements are consistent in supporting the hypothesis that adjustment is beneficial to retrieval, although the benefits are minor.

#### b. Number of modes used

Principal component-based retrievals were done to evaluate how many modes of  $q$  and  $t_B$  should be used for optimal results. For these tests only adjusted data were used, and both  $i$  and  $j$  were varied throughout their respective ranges. Here the included modes were the ones with the greatest eigenvalues, although it is not assured that the *first*  $j$  modes of  $t_B$  are the *best*  $j$  modes for determining  $q$ . Note that when  $i = 13$  and  $j = 18$ , this method is identical to standard linear regression.

Results for the independent data are shown in Fig. 3. For the midlatitude set, the explained variance values were highest for  $i = 4$  and  $j = 13$ , while for the tropical set the values peaked at  $i = 3$  and  $j = 15$ . However, the midlatitude value for  $i = 4$ ,  $j = 13$  was not significantly greater than the values for  $i = 3$ ,  $j = 13$  and  $i = 3$ ,  $j = 14$ . One implication is that three modes of  $q'$  are enough to represent all of the retrievable information in water vapor profiles, based on regression on HIRS-2 brightness temperatures. The retrieved product

### a MIDLATITUDE

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
13	9.8	17.1	16.7	23.7	26.1	28.1	27.9	30.5	33.2	34.6	33.5	34.2	35.2	34.9	33.1	33.5	33.2	33.0
10	9.8	17.1	16.7	23.7	26.1	28.0	27.9	30.4	33.2	34.6	33.5	34.2	35.2	34.9	33.1	33.5	33.2	33.0
5	9.7	17.1	16.8	23.8	26.2	28.1	28.0	30.5	33.3	34.7	33.7	34.4	35.3	35.1	33.3	33.8	33.4	33.3
4	9.7	17.2	16.8	23.8	26.4	28.3	28.1	30.7	33.4	34.9	33.9	34.5	35.4	35.2	33.4	33.9	33.5	33.3
3	9.6	17.1	16.7	23.8	26.2	27.8	27.5	30.1	32.9	34.4	33.4	34.4	35.4	35.4	33.5	34.0	33.6	33.4
2	8.9	16.6	16.2	23.3	25.2	26.8	26.3	28.8	31.3	32.6	31.5	32.2	33.3	33.1	32.1	32.5	32.3	32.3
1	8.8	16.2	15.8	21.0	22.6	24.3	23.7	26.5	28.8	29.0	28.3	28.5	29.2	29.1	28.4	28.4	28.5	28.4

### NUMBER OF WATER VAPOR MODES - i

### b TROPICAL

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
13	4.6	8.9	15.6	32.3	31.2	45.6	45.7	47.0	48.2	48.3	49.3	50.2	48.8	48.3	51.2	46.6	41.2	31.1
8	4.6	9.0	15.6	32.4	31.3	45.7	45.8	47.2	48.3	48.4	49.5	50.4	49.1	48.5	51.5	46.9	41.7	31.5
5	4.7	9.0	15.7	32.4	31.7	46.3	46.2	47.6	48.7	48.8	50.0	51.0	50.1	49.7	52.8	48.3	43.1	33.0
4	4.7	9.0	15.7	32.3	31.7	46.2	46.0	47.4	48.5	48.7	49.9	51.0	50.3	50.1	53.1	48.6	43.5	33.5
3	4.6	9.0	15.6	32.4	31.9	46.4	46.3	48.0	48.7	49.0	51.2	52.4	51.8	51.5	54.6	50.8	46.2	36.6
2	4.6	8.9	15.5	31.8	31.8	46.1	46.1	47.6	47.8	48.1	49.8	50.5	49.0	47.9	52.1	50.9	46.3	38.4
1	1.5	5.7	11.2	26.5	26.2	39.2	39.4	40.2	40.6	40.9	42.7	43.3	41.9	41.2	45.2	44.2	39.9	31.7

### NUMBER OF BRIGHTNESS TEMPERATURE MODES - j

FIG. 3. Mixing ratio retrieval results for various numbers of (i) mixing ratio and (j) brightness temperature modes. Averaged explained variances (percent) are given for a) midlatitude and b) tropical sounding sets, with independent data, and mean-saturation adjusted mixing ratios. The boxes indicate the maxima, and the contours are at 10% intervals. Note that the results for i = 13 are the same as would be computed if only the brightness temperatures were transformed into principal components (e.g., Kendal, 1975, pp. 95-98).

could thus be described as having third-order vertical resolution (Hillger et al., 1984).

### c. Comparison with standard regression and hybrid

To further assess the applicability of principal component-based regression to water vapor profile retrieval, we compared results of that method with those of standard regression and a hybrid method. For standard regression we determined  $q'$  directly from  $t_B$ , skipping transformations of the forms (2) and (3). For hybrid regression we retained the transformation of  $q'$  into  $a'_i$  and retrieved  $a'_i$  from  $t_B$ , with  $i = 3$ . For standard and hybrid regression the brightness temperatures comprising  $t_B$  were the available predictors, while the elements of  $b$  were the available predictors in principal component-based retrieval.

For the comparisons, stepwise predictor selection was applied to all three regression methods because it is a simple and commonly practiced way to improve regression performance. However, there is no guarantee that stepwise selection identifies the best possible combination of predictors. The forward stepwise procedure (Kendall, 1975, p. 98) was modified in that the stepping was not halted until all predictors were included, and the explained variance values for the independent data were used to identify which combination of predictors was best.

The maximum explained variance values for independent data are shown in Table 2 for standard, hybrid, and full principal component stepwise regression.<sup>1</sup> The hybrid method produced the highest values for both the midlatitude and tropical data sets. Apparently, transformation of brightness temperatures did not succeed in concentrating the information relevant to water vapor detection into a few modes. The transformation actually seemed to muddle the information in the midlatitude case. The muddling could be due to nonrepresentativeness of the statistics, or due to the possibility that individual brightness temperatures are better than their principal components at determining mixing ratios.

### 3. Conclusion

We studied two features of water vapor profile retrieval by regression via principal components. Exper-

<sup>1</sup> Stepwise selection clearly failed to identify the best combination of predictors for the full principal component method, as is seen by comparing the last column of Table 2 with the boxed numbers in Fig. 3. However, stepwise selection put all three methods on an equal basis, and better combinations of predictors may exist for standard and hybrid regression also.

TABLE 2. Comparison of regression methods.

Sounding set	Explained variance (%)		
	Standard	Hybrid	Principal component
Midlatitude	36.6	36.7	33.7
Tropical	51.3	53.6	53.6

iments on HIRS-2 data indicated that adjustment of mixing ratios by mean saturation mixing ratios slightly enhanced the accuracy of retrievals. This type of regression performed at its best when a different number of modes was used to represent mixing ratios than was used to represent brightness temperatures.

Tests with both midlatitude and tropical data showed that the optimal number of mixing ratio profile modes to use in retrieval is three, which indicates that HIRS-2 data contain information relevant to at least three statistically independent modes of a water vapor profile.

The principal component-based retrieval method was compared with standard regression and a hybrid method. The hybrid retrievals were the most accurate.

*Acknowledgments.* We thank Dr. Larry McMillin of NESDIS for providing the satellite and radiosonde data used in this study, and Dr. Donald Hillger for contributing helpful comments. Mrs. Loretta Wilson typed the manuscript, and Ms. Judy Sorbie drafted the figures. The research reported in this paper was sponsored by NOAA Contract NA84AA-H-00020 to the Cooperative Institute for Research in the Atmosphere (CIRA).

### REFERENCES

- Hillger, D. W., A. E. Lipton and T. H. Vonder Haar, 1984: Vertical resolution in water vapor soundings from satellites: Goals and limitations. *IRS '84: Current Problems in Atmospheric Radiation*, G. Fiocco, Ed., *Proc. of the International Radiation Symposium*, Perugia, A. Deepak., 438 pp.
- Kendall, M., 1975: *Multivariate Analysis*. Hafner Press, 210 pp.
- Smith, W. L., and H. M. Woolf, 1976: The use of eigenvectors of statistical covariance matrices for interpreting satellite sounder radiometer observations. *J. Atmos. Sci.*, **33**, 1127-1140.
- , C. M. Hayden, D. Q. Wark and L. M. McMillin, 1979: The TIROS-N operational vertical sounder. *Bull. Amer. Meteor. Soc.*, **60**, 1177-1187.
- Weinreb, M. P., and D. S. Crosby, 1977: A technique for estimating atmospheric moisture profiles from satellite measurements. *J. Appl. Meteor.*, **16**, 1214-1218.
- Werbowetzi, A., Ed., 1981: *Atmospheric Sounding User's Guide*. NOAA Tech. Rep. NNESS83, 82 pp. [NTIS PB81 230476].