

# Application of factor analysis for quantification of climate-forming processes in the eastern part of the Baltic Sea region

Arunas Bukantis\*

Department of Hydrology and Climatology, Vilnius University, M.K. Ciurlionio St. 21, Vilnius 2009, Lithuania

**ABSTRACT:** Factor analysis (FA) is applied to identify climate-forming factors and quantitatively evaluate their importance in the eastern part of the Baltic Sea region. Monthly data from 7 meteorological stations are used: average maximum and average minimum of air temperature ( $T_1$  and  $T_2$ ), atmospheric pressure ( $P$ ), wind speed at a height of 10 m ( $v$ ), monthly precipitation ( $Q$ ), duration of sunshine ( $S$ ) and partial pressure of water vapour ( $e$ ). FA reveals that the variables form 2 main groups (meteorological complexes): a hygrothermal complex— $T_1$ ,  $T_2$  and  $e$ —and a baric-radiational complex— $P$ ,  $Q$  and  $S$ . Both meteorological complexes exist almost independently of each other (in cold seasons in particular), i.e. it is possible to distinguish the 2 most important types of climate-forming processes in the eastern part of the Baltic Sea region. One of them is advection related to the input of various air masses whose features are best reflected by air temperature and humidity. The other process mostly takes place within 1 pressure system (cyclone or anticyclone): air mass transformation, vertical mixing, formation of clouds and related precipitation and input of solar radiation.

**KEY WORDS:** Climate-forming processes · Factor analysis · Baltic Sea region

*Resale or republication not permitted without written consent of the publisher*

## 1. INTRODUCTION

In studying the processes that form climate in any particular region, one usually faces the problem of quantification of the role of separate climatic factors. Today the meteorological observation system provides various types of information about the physical state of weather at a certain moment. This information, i.e. all meteorological parameters, must be systematised and classified to gain information about the climate-forming processes that are themselves not measurable. In other words, the initial meteorological information must be compressed and unified for the convenience of its further scientific interpretation (Razuvaev 1984, Bukantis 1993).

This problem can be efficiently approached by using factor analysis (FA), which enables one to describe the climate-forming factors as general ones whose number

is by far smaller than the number of initial indices. This smaller number of variables can be used to find meaningful structure in the observed variables. This structure will aid in the interpretation and explanation of the process that generated the observations (Überla 1977, Afifi & Clark 1996). The term factor is here applied to a hypothetical latent process related in a certain way to the measurable meteorological elements.

The second reason to carry out an FA is for data reduction. Since the observed variables are represented in terms of a smaller number of unobserved or latent variables, the number of variables in the analysis is reduced and so are the storage requirements. By having a smaller number of factors (vectors of smaller dimension) to work with that capture the essence of the observed variables, only this smaller number of factors needs to be stored. This smaller number of factors can also be used for further analysis, thus reducing computational requirements.

\*E-mail: arunas.bukantis@gf.vu.lt

The main task of FA is not only to check the working hypothesis about the interdependence of the initial variables but also to formulate it (Comrey & Lee 1992). FA aids in the investigation of interdependence of variables without the deductive postulate that these variables fully describe the studied field. Rather, the main task is to determine the quantity and character of linearly independent variables (factors) which would sufficiently precisely express the interdependence of the initial variables. This is the main difference between FA and other statistical methods, such as discriminant analysis and cluster analysis (Überla 1977, Gorsuch 1983, Morrison 1990). The named features of the FA method were also one of the reasons to choose it for the present study. The hypothesis about the structure of climate-forming factors in the eastern part of the Baltic Sea region is formulated on the basis of many initial meteorological data. While searching for hypothetical factors on which meteorological elements depend, we assumed that they could be interpreted and would have equivalents in reality.

FA as a generic term includes principal component analysis (PCA), or empirical orthogonal function (EOF) analysis. While the 2 techniques are functionally very similar and are used for the same purpose (data reduction), they are quite different in terms of underlying assumptions. This difference has little to do with the formal definition of methods (Gorsuch 1983, Loehlin 1992). The EOFs are orthogonal spatial patterns that can be thought of as empirically derived basis functions. The low-order EOFs can sometimes be interpreted as natural modes of variation of the observed system. The time coefficients obtained by projecting the observed field onto the EOFs are uncorrelated and represent the variability of the field efficiently (Storch & Zwiers 1999).

FA and PCA (EOF analysis) are similar in the sense that the purpose of both is to reduce the original variables into fewer composite variables, called factors or principal components. However, they are distinct in the sense that the obtained composite variables serve different purposes. In FA, a small number of factors are extracted to account for the intercorrelations among the observed variables—to identify the latent dimensions that explain why the variables are correlated with each other. In PCA, the objective is to account for the maximum portion of the variance present in the original set of variables with a minimum number of composite variables called principal components.

## 2. THE FACTOR ANALYSIS MODEL

FA was performed in a sequence with the following major steps: (1) selecting the variables; (2) computing

the matrix of correlations among the variables; (3) extracting the unrotated factors; (4) rotating the factors orthogonally (Varimax method); and (5) interpreting the rotated factor matrix.

FA customarily uses the matrix of correlations among the selected set of variables to end up with a matrix of factor loadings. This matrix can be interpreted in the orthogonal factor model as a non-zero correlation,  $r$ , between the measured parameters,  $x_j$ , and certain hypothetical constructs, called 'latent factors',  $f_i$ .

There are 2 ways of defining the factor matrix. With communalities in the diagonal cells, the correlation matrix is designated by the symbol  $\mathbf{R}$ . When unities, rather than communalities, are placed in the diagonal cells, the correlation matrix is designated by the symbol  $\mathbf{R}_u$ :

$$\mathbf{R}_u = \mathbf{A}_u \mathbf{A}_u^T \quad (1)$$

where  $\mathbf{R}_u$  is the correlation matrix among data variables with ones in the diagonal cells and  $\mathbf{A}_u$  is the complete matrix of factor loadings, including common, specific, and error factors.  $\mathbf{A}_u$  values also represent correlations between the data variables and the factors. These correlations are called the factor loadings in the orthogonal factor model, which requires all factors to be at right angles to one another, that is, to be uncorrelated.

$\mathbf{R}$ , with communalities in the diagonal cells, may be represented as a product of  $\mathbf{A}$  and its transpose:

$$\mathbf{R} = \mathbf{A} \mathbf{A}^T \quad (2)$$

The transpose of the matrix is obtained by interchanging rows and columns; that is, row 1 of  $\mathbf{A}$  is equivalent to column 1 of  $\mathbf{A}^T$ .  $\mathbf{A}$  consists of only the common factor portion of the factor loadings in  $\mathbf{A}_u$ .

In an FA the communality for a variable has already been defined as the sum of squares of the factor loadings over all the factors. Each observed variable's communality is its estimated squared correlation with its own common portion—that is, the proportion of variance in that variable that is explained by the common factors. Communality values between 1.0 and 0 indicate partial overlapping between the variables and the factors in what they measure.

The matrix of data-variable scores,  $\mathbf{Z}$ , may be obtained by multiplying  $\mathbf{A}_u$  by the matrix of factor scores,  $\mathbf{F}_u$ :

$$\mathbf{Z} = \mathbf{A}_u \mathbf{F}_u \quad (3)$$

As ordinarily applied, FA involves deriving a set of factor loadings from a matrix of correlation coefficients between the data variables. The correlation between a pair of data variables equals the sum of the products of their factor loadings, the  $a$  values from  $\mathbf{A}_u$ , on the common factors.

Eqs. (1) & (2) are 2 forms of a very important theorem in FA which was called by Thurstone (1947) the fundamental equation of FA; they both state that the correlation matrix among the data variables can be decomposed into the product of a factor matrix and its transpose (Comrey & Lee 1992).

In the geometric representation of the factor model, a data variable may be represented as an eigenvector in a space of as many dimensions as there are common factors. In this case, the length of the vector is  $h$ , the square root of  $h^2$ , the communality. As ordinarily carried out, the process of factor extraction starts with a matrix of correlations between data variables with communalities in the diagonals and ends up with a matrix of factor loadings,  $\mathbf{A}$ , such that when multiplied by its transpose,  $\mathbf{A}^T$ , the correlation matrix,  $\mathbf{R}$ , will be reproduced, at least approximately. As the number of monitored meteorological parameters is rather high, it is expedient to extract more than 1 latent factor, i.e. to apply the multi-factor variant of FA. In this work the number of factors was determined by the Bargmann criterion (Überla 1977), suggesting that in the case of 7 data variables each of the factors should have a significant projection to not less than 3 variables. For this reason only 2 factors were extracted.

The initial extracted factor matrix must usually be rotated before the final factor solution is achieved. The orthogonal factors are rotated in the space of input vectors to the position where the groups of variables have maximal loads in concrete factors.

With increasing correlation between  $x_j$  and  $f_i$ , the factor load ( $a_{ji}$ ) of  $x_j$  in  $f_i$  increases. Therefore, having determined which  $x_j$  have the greatest loads in  $f_i$ , we may judge in which variables (meteorological parameters) the manifestation of  $f_i$  is strongest. This provides the possibility of evaluating the latent processes which in FA obtain the shape of the quantitative parameters,  $f_i$ .

Finally, the distribution by factors of the variables' dispersion is determined. The communality of each variable is expanded into parts of dispersion related to discrete factors. These values give an indication of the extent to which the variables overlap with the factor, or more technically, they give the proportion of variance in the variables that can be accounted for by scores in the factors. For a more detailed description of the FA method, see Blahush (1985), Comrey & Lee (1992) and Afifi & Clark (1996).

### 3. THE INPUT METEOROLOGICAL DATA

The results of FA largely depend on the selection of input meteorological parameters. Firstly, they should as fully as possible reflect the object of the study; sec-

ondly, the precision and reliability of measurements should be taken into consideration. In addition, the number of variables (parameters),  $n$ , must be considerably smaller than the number of observations,  $N$ .

On the grounds of the listed criteria, 7 variables,  $x_j$  ( $j = 1, 2, \dots, 7$ ) (average monthly values of meteorological elements), were chosen and the number of observations was  $N = 45$  (1950–1994). The data were taken from the following meteorological stations: St. Petersburg, Kaliningrad (Russia), Riga (Latvia), Birzhai, Vilnius (Lithuania), Minsk and Brest (Belarus). These stations were chosen for homogeneity of data and geographical location (Fig. 1). The first 3 are situated on the Baltic Sea coast; the remaining 4 are situated at a distance of 250 to 400 km from the sea. While choosing the meteorological elements and their time series, the following criteria were taken into consideration: First, the meteorological elements had to be instrumentally measured; they should as fully as possible describe the thermal, humidity and dynamic features of the atmosphere and input of solar radiation. Second, the time series had to be homogeneous and complete, i.e. no missing data.

For a characterisation of the air temperature regime, the average maximal ( $T_1$ , °C) and the average minimal ( $T_2$ , °C) values were used. These 2 parameters provide a more complete view of the meteorological situation than the average temperature. The included indices of atmospheric circulation were atmospheric pressure ( $P$ , hPa) and wind speed at a height of 10 m ( $v$ , m s<sup>-1</sup>). Humidity was described by monthly precipitation ( $Q$ , mm) and the partial pressure of water vapour ( $e$ , hPa), and the input of solar radiation and cloudiness by duration of sunshine ( $S$ , h). We should point out the particular informativeness of  $e$ , which has almost no diurnal variation and is hardly affected by local anthropogenic factors or properties, on a meso- or microscale, of the surface. The variations of water



Fig. 1. Locations of the 7 meteorological stations

vapour pressure include frequent long-term (2 to 2.5 mo) anomalies (Sazonov 1991, Bukantis 1998). This variable, when it is measured at least at 1 point, describes the humidity and genesis of the air mass. Besides, the air humidity is an important factor in atmospheric transparency and in the formation of the radiation balance (Hartmann 1994).

#### 4. RESULTS AND DISCUSSION

The FA showed that the 7 variables formed 2 main groups, which will henceforth be referred to as meteorological complexes.

In this work a meteorological complex was taken to be a previously existing one when the factor loads of factor  $f_i$  in 2 to 3 variables exceeded 0.4. This corresponds to 99% statistical significance. For this reason the null hypothesis, that the initial variables are independent of each other and of latent factor, may be rejected. The results obtained by FA revealed that the variation and mutual relations of meteorological parameters in the studied region was influenced throughout the year by 2 main factors,  $f_1$  ( $T_1$ ,  $T_2$ ,  $e$ ) and  $f_2$  ( $P$ ,  $Q$ ,  $S$ ).

This model of 2 factors is characteristic also of other North and East European regions where the climate is humid. In steppes and forest steppes the complex  $T_1$ ,  $T_2$ ,  $P$  occurs in winter and  $T_1$ ,  $Q$  in summer (Sazonov 1991).

Fig. 2 shows the loads ( $a_{ji}$ ) of the 7 meteorological parameters  $x_j$  in factors  $f_1$  and  $f_2$  in Vilnius (at other meteorological stations the distribution pattern of

loads is similar). Positive and negative loads have the same meaning as in correlation—the factor may reduce or increase the numerical value of the variable. The complex  $T_1$ ,  $T_2$ ,  $e$  formed by  $f_1$  may be called hygrothermal, and it is particularly distinct. The factor loads,  $a_{1j}$ , of its variables exceed 0.6 in all seasons, and in the cold season they even reach 0.90 to 0.96. This complex reflects, most likely, the synchronous changes of air humidity and maximal/minimal temperature under varying directions of air advection. In other words,  $f_1$  may be called the factor of advection. The inflow of air from the south (in winter from the west) is usually followed by increases in air temperature and humidity. On the contrary, in winter, when the Baltic region is invaded by dry Arctic or Asian air masses, the  $T_1$  and  $T_2$  fall and  $e$  decreases (Bukantis 1994, Bukantis & Valiushkevichiene 1999). In summer this type of air advection is represented best by minimal temperature (cold nights) and very low air humidity (Bukantis & Valiushkevichiene 2000). Thus even without direct data about air mass trajectories, the probable origin of the air mass and physical processes can be inferred using the hygrothermal complex.

The influence of air mass advection on various meteorological elements has also been investigated by other authors. Keevallik et al. (1999) found that weather in Estonia is strongly related to the 3 general circulation types. Zonal circulation brings to Estonia wet weather that in winter is warmer and in summer cooler than average. Mixed circulation brings cold winters and varying weather in summer. Meridian circulation is responsible for cold winters and hot summers with droughts. Comparison of temperature

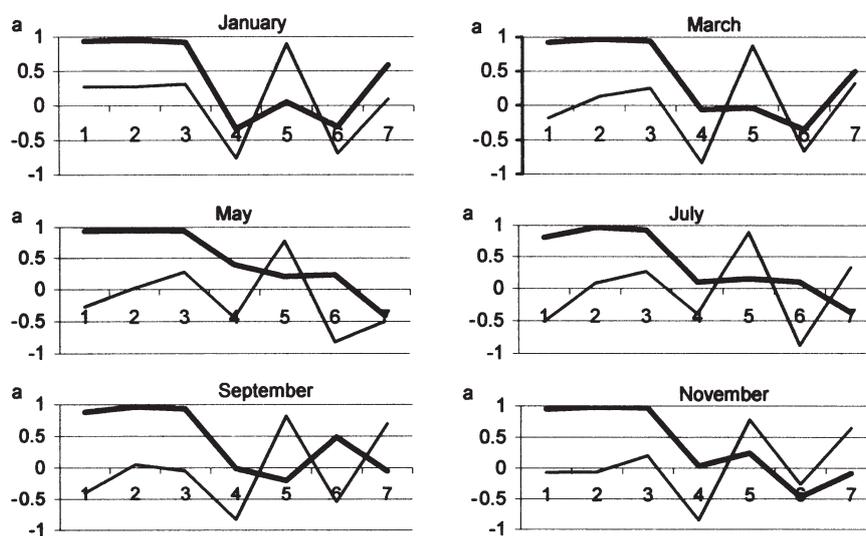


Fig. 2. Loads ( $a$ ) of the 7 meteorological parameters in factors  $f_1$  (thick line) and  $f_2$  (thin line) in Vilnius. 1: average maximal air temperature ( $T_1$ ); 2: average minimal air temperature ( $T_2$ ); 3: partial pressure of water vapour ( $e$ ); 4: atmospheric pressure ( $P$ ); 5: monthly precipitation ( $Q$ ); 6: duration of sunshine ( $S$ ); 7: wind speed at a height of 10 m ( $v$ )

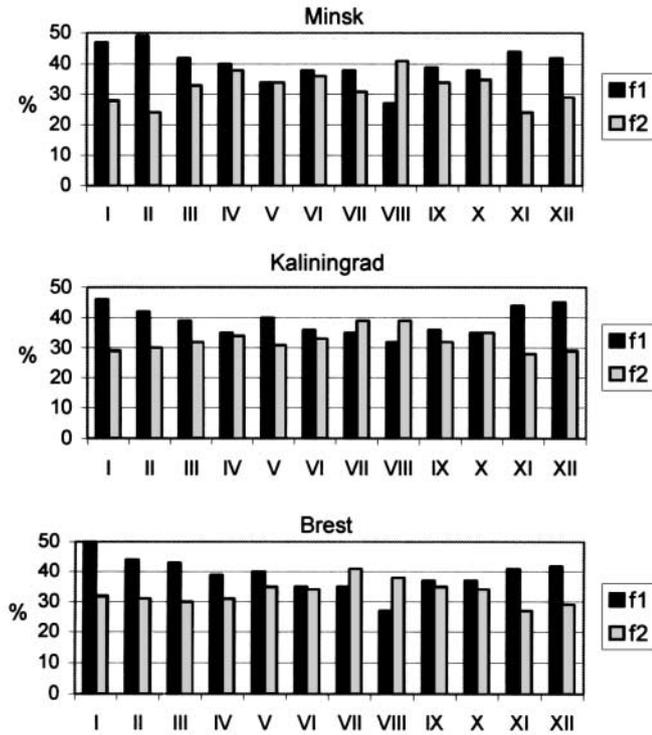


Fig. 3. Proportion of the dispersion of meteorological parameters formed by  $f_1$  and  $f_2$

and precipitation in Estonia and in Central Europe for different patterns shows some principal differences. First of all, the 4 zonal circulation patterns (WA, WZ, WS and WW) that bring similar weather to Central Europe, form 2 different groups for Estonia: WA and WZ bring to Estonia typical maritime weather, while WS and especially WW show continental features. Such differences can also be found amongst the patterns of meridional circulation. These results show that for the eastern part of the Baltic Sea region the grouping of circulation patterns should be somewhat different than for Central Europe.

The second meteorological complex ( $P, Q, S$ ), which may be called the baric-radiation complex, forms as a result of the dominating influence of air pressure on other meteorological parameters. It is also rather distinct throughout the year, but in July–August this complex is joined by  $T_1$  (Fig. 2, factor loads  $a_2$ ): with increasing pressure, the amount of precipitation decreases, sunshine increases, and in summer—due to solar radiation—the maximal day temperatures become higher. In September–December this complex also includes  $v$ , which is related to the intensive cyclonic activity in autumn. When the pressure is low, winds tend to be strong, and cloudiness and precipitation increase. In September–October the region under study has on average 6 to 9  $d\ m^{-1}$  with deep cyclones

(pressure in the centre < 990 hPa), whereas in the remaining months only 1 to 4 such days per month occur (Bukantis 1994). Thus, this complex is formed by the second factor,  $f_2$ , which reflects the cyclonic-anticyclonic activity and vertical instability of the air masses.

In the final stage of FA the distribution by factors of the dispersion of the meteorological elements is determined. The portion of dispersion formed by both factors in the variation of meteorological elements makes up on average 70 to 85% of their total dispersion (Fig. 3). In cold seasons (November to February)  $f_1$  is dominant, forming 42 to 54% of the total dispersion. The value of  $f_2$  in these months is 1.5 to 2 times lower. The influence of  $f_2$  increases in March–April, and in July–October it is almost equal with the influence of  $f_1$  (forms 30 to 40% of dispersion in the variables). In July–August  $f_2$  becomes dominant in the southern part of the studied region. It was observed that the influence of  $f_1$  is for most of the year (except for February–April) stronger in the northern than in the southern part of the investigation area, accounting for 2 to 3% more of the dispersion in the variables in the north than in the south. In contrast,  $f_2$  accounts for 2 to 9% more of the dispersion in the southern than in the northern part of the region (Fig. 4). Similar patterns of the influence of  $f_1$  and  $f_2$  have been observed in other regions of Eastern and Central Europe (Sazonov 1991).

There is a question as to why the influence of  $f_2$  decreases in the beginning of the warm season (May–June) when due to intensive solar radiation the influence of this factor should increase. Studies of

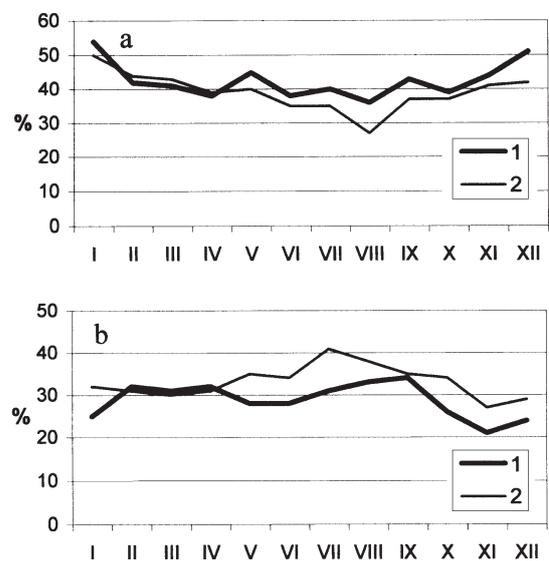


Fig. 4. Proportion of the dispersion of meteorological parameters by (a)  $f_1$  and (b)  $f_2$  in the north (line 1, St. Petersburg) and in the south of the region (line 2, Brest)

atmospheric circulation have revealed a lower recurrence of strong anticyclones over the Baltic region in May–June and an increased recurrence of other anticyclonic formations (weak anticyclones and ridges) generated in the North Atlantic, the Norwegian Sea and the Barents Sea (Bukantis 1994). For this reason the link between the duration of sunshine, cloudiness and precipitation (and also the input of solar radiation) and atmospheric pressure weakens. It remains strong only in the strong anticyclones. Sunshine duration is closely related to the origin of the air mass (Russak 1990, Weber 1990). The most favourable insolation conditions during the whole year are accompanied by north-easterly anticyclonic circulation. Less favourable conditions for direct radiation are provided by westerly anticyclonic circulation. A comparably low sunshine duration is characteristic for westerly and north-westerly cyclonic circulation (Chodakova 1980, Dubicka & Karal 1998). Besides, the temperature contrast between the sea and land surfaces remains large in the beginning of summer, which makes the meteorological conditions in the region sensitive to the direction of air advection, i.e. to  $f_1$ .

The remaining 15 to 25 % of dispersion in the meteorological elements is due, presumably, to innumerable, highly complicated and often interrelated local and macro factors.

## 5. CONCLUSIONS

FA is an effective mathematical method. It is helpful in finding out the existing natural meteorological complexes which under the effect of a certain factor join a few meteorological elements. At the same time, FA considerably reduces the contribution of casual fluctuations and errors.

Note that the 2 meteorological complexes mentioned are formed (as revealed by FA) by 2 orthogonal factors, that is, they are uncorrelated with each other and exist almost independently of each other (particularly in cold seasons). Thus, it is possible to distinguish 2 main types of climate-forming processes in the eastern part of the Baltic Sea region. One of them is represented by the advection process (factor  $f_1$ ) related with the inflow of different air masses whose properties are best represented by air temperature and humidity. The other process (factor  $f_2$ ) takes place mostly within 1 pressure system (cyclone or anticyclone): air mass transformation, vertical mixing, formation of clouds and related precipitation, and input of solar radiation. Thus FA does not confine itself only to the systematisation of input data. It is helpful in objectivization of climate-forming peculiarities and, in this particular case, in confirmation of the theory of climate-forming pro-

cesses in the eastern part of the Baltic Sea region. It also makes it possible to quantify the importance of climate-forming processes.

## LITERATURE CITED

- Afifi A, Clark VA (1996) Computer-aided multivariate analysis, 3rd edn. Chapman and Hall, London
- Blahush P (1985) Faktorova analiza a jeji zobecneni. SNTL, Prague
- Bukantis A (1993) Application of factor analysis for quantification of climate factors. *Geography. Sci Ann (Vilnius)* 29: 44–46
- Bukantis A (1994) Climate of Lithuania. Vilnius University Press, Vilnius
- Bukantis A (1998) Extreme cold and warm winters in the Baltic Sea area. In: Lemmela R, Heleniusla N (eds) Proceedings of the Second International Conference on Climate and Water, Espoo, Finland, 17–20 August 1998. Helsinki University of Technology, p 468–478
- Bukantis A, Valiushkevichiene L (1999) Air temperature anomalies in Lithuania: (I) cold seasons. *Geographical Yearbook XXXII*, Institute of Geography and Lithuanian Geographical Society, Vilnius, p 57–64
- Bukantis A, Valiushkevichiene L (2000) Air temperature anomalies in Lithuania: (II) warm seasons. *Geographical Yearbook XXXIII*, Institute of Geography and Lithuanian Geographical Society, Vilnius, p 45–55
- Chodakova WP (1980) About sunshine duration in the territory of Western Europe. *Ann Glavnoi Geofizicheskoi Observatorii (Leningrad)* 441:44–49 (in Russian)
- Comrey AL, Lee HB (1992) A first course in factor analysis. Academic Press, San Diego, and University of California
- Dubicka M, Karal J (1998) Sunshine duration on Szerenica Mt. and its relation to the atmospheric circulation. *Acta Univ Wratislaviensis No. 1590*, Ser C, T1:35–43
- Gorsuch RL (1983) Factor analysis. Erlbaum Associates, Hillsdale, NJ
- Hartmann DL (1994) Global physical climatology. In: Dmowska R, Holton JR (eds) International geophysics series, Vol 56. Academic Press, New York
- Keevallik S, Post P, Tuulik J (1999) European circulation patterns and meteorological situation in Estonia. *Theor Appl Climatol* 63(1/2):117–127
- Loehlin JC (1992) Latent variable models. Erlbaum Associates, Hillsdale, NJ
- Morrison DF (1990) Multivariate statistical methods. McGraw-Hill, New York
- Razuvaev VN (1984) Statistical methods of climatological information analysis. Factor analysis. *Gidrometeoizdat, Obninsk* (in Russian)
- Russak V (1990) Trends of solar radiation, cloudiness and atmospheric transparency during recent decades in Estonia. *Tellus* 42B(2):206–210
- Sazonov BI (1991) Heavy winters and droughts. *Gidrometeoizdat, St Petersburg* (in Russian)
- Storch H, Zwiers FW (1999) Statistical analysis in climate research. Cambridge University Press, Cambridge
- Thurstone LL (1947) Multiple factor analysis. University of Chicago Press, Chicago
- Überla K (1977) Faktorenanalyse. Springer-Verlag, Berlin
- Weber GR (1990) Spatial and temporal variation of sunshine in the Federal Republic of Germany. *Theor Appl Climatol* 41(1/2):1–9