

Correspondence Factor Analysis (CFA) of Multivariate Fluid Geochemistry Data from Indian Hot Springs

Amitabha Roy

Ex-Senior Director, Geological Survey of India

Abstract

Correspondence analysis was conducted on a two-way contingency table with 62 rows and 9 columns (Primary dataset: A) representing fluid geochemical data from Indian hot springs. The analysis used the joint probability distribution of two random variables: individual observational samples (rows) and fluid geochemical variables (columns). This data set is really the total of two subsets: (B) Peninsula (25 rows and 9 columns) and (C) Extra-Peninsula (37 rows and 9 columns). The study generated factor loadings for individual samples and variables, which were shown as points on two-dimensional coordinate (factorial) axes with the same origin known as biplots, in order to find discrete geochemical domains defined by natural sample and variable groupings or clusters. The simultaneous changes in trace elements in these three data sets with sample locations appear to reflect broader trends in geothermal evolution in the region.

Keywords: Correspondence factor analysis, Geochemistry, Biplots. Peninsula and Extra-Peninsula, Hot springs of India

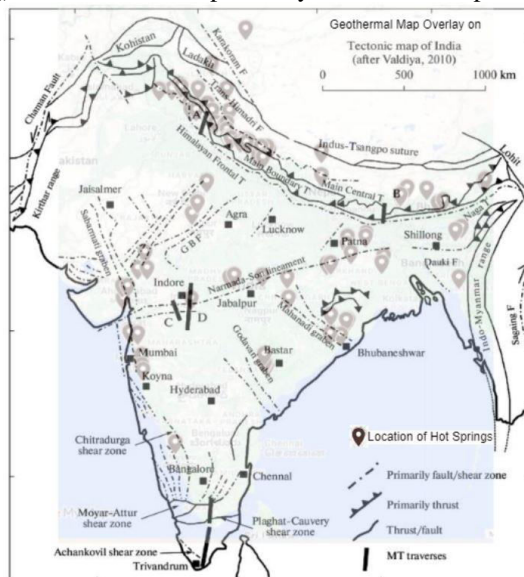
Date of Submission: 04-02-2024

Date of acceptance: 19-02-2024

I. INTRODUCTION

India has about 340 hot springs spread across the peninsular and extra-peninsular regions. The government of India constituted a 'Hot Spring Committee' in 1968 to examine the possibility of developing geothermal plants for power generation. The Central Electricity Authority (CEA) has associated itself with the UNDP geothermal project in India and the Puga and Parvati projects for the utilization of available geothermal resources for power generation (Jonathan Craig, 2013). The Geological Survey of India (GSI) has published a special publication titled "Geothermal Atlas of India" based on data compiled from all sources of information (Ravi Shankar et al. 1991). However, the lack of uniformity in data acquisition practices and manual handling of large amounts of data has made data storage, search, retrieval, and analysis laborious and cumbersome (A.Roy, 1994).

Fig. 1. Geothermal Map overlay on tectonic map of India



Study area and its geologic-tectonic settings

A dataset of 62 samples of multivariate geothermal data spread across two spatially distinct regions of diverse geologic-tectonic settings, one from a 2400 km-long arcuate belt of the tectonically active Extra-Peninsular Himalayan region and the other from Late-Precambrian or Proterozoic mobile belts in the Central Highland in an otherwise stable landmass or shield of Peninsular India, were subjected to robust statistical techniques such as exploratory factor analysis followed by multiple regression analysis to determine the origin of geothermal hot springs (Amitabha Roy, 2023). The model studies distinguish two statistically significant suites of fluid geochemistry: 1. the overall salt assemblage and concentration of Cl-HCO₃-SO₄-Na-F or chloride rich deep seated acidic waters suggestive of the existence of a hydrothermal magmatic system operating in the geotherms of Extra-Peninsular India; and 2. Peninsular springs of K-Na-HCO₃ bicarbonate rich alkaline waters with low SO₄-content and relatively higher contents of HCO₃ compared to other anions SO₄, Cl, and F suggestive of a non-magmatic origin.

Both exploratory factor and multiple regression studies contribute to understanding the origins of these two fluid geochemistry suites. The model studies distinguish two statistically significant suites of fluid geochemistry: 1. the overall salt assemblage and concentration of Cl-HCO₃-SO₄-Na-F or chloride rich deep seated acidic waters suggestive of the existence of a hydrothermal magmatic system operating in the geotherms of Extra-Peninsular India; and 2. Peninsular springs of K-Na-HCO₃ bicarbonate rich alkaline waters with low SO₄-content and relatively higher contents of HCO₃ compared to other anions SO₄, Cl, and F suggestive of a non-magmatic origin.

Correspondence Factor Analysis

The current study uses correspondence analysis (CA) to investigate the relationship between two variables (rows and columns of a contingency table) as well as the similarities between their categories in the context of Indian hot springs spread across different geologic-tectonic settings. Correspondence Analysis (CA) is often described as an adaptation of PCA for categorical data. A key difference is that the data analyzed in CA is not a covariance/correlation matrix as in PCA. There exists some confusion about multivariate statistical methodologies such as PCA, FA, and CA, all of which have the ultimate goal of identifying hidden patterns inside data structures. PCA and FA are two methods for data reduction. CA is a technique for representing the rows and columns of a non-negative data given in a two-way contingency table as two-dimensional graphical or biplot points (rows in principal coordinates, columns in standard coordinates). Another difference between PCA and CA is that CA calculates two set of factor scores: one for the rows and one for the columns. Thus, the space between rows and column factor scores is interpretable—unlike PCA, where, for comparing observations to variables only the angles and not the distances are compared. Thirdly, PCA uses normal Euclidean distance (Pythagorean distance to best-fit line) to find Inertia, If the coordinates of P and Q such that, (x₁, 0) and (x₂,0), the distance between PQ is given by $D(p,q) = |x_2 - x_1|$; CA uses weighted Euclidean distance (a variant of Euclidean distance) or chi-square distance to measure the proximity and distance among samples in full CA ordination space. The chi-square formula is: $\chi^2 = \sum(O_i - E_i)^2/E_i$, where O_i = observed value (actual value) and E_i = expected value. Correspondence analysis is a way of transforming a contingency table, which shows the frequencies of two or more categorical variables, into a graphical display, which shows the distances and angles between the categories. The graphical display is called a biplot, and it consists of two sets of points: one for the rows and one for the columns of the table. The closer the points are, the more associated they are. The farther the points are, the more dissimilar they are. The angle between the points indicates the nature of the association: acute angles mean positive association, right angles mean independence, and obtuse angles mean negative association. A combined presentation shows the row and column profiles overlaid. This presentation is highly convenient because both row and column points are evenly distributed. The distance between row points (or column points) approximates the inter-row or inter-column Chi-square distance. Alternatively, a row plot only shows the row categories on the biplot, while a column plot only shows the column categories on the biplot; these are useful for highlighting differences among the respective categories with respect to the other variable.

In the present study, correspondence factor analysis, a multivariate statistical technique of proximity and distance measure from the origin of two-dimensional coordinate (factorial) axes, was performed in Excel using XLSTAT software to identify associations or oppositions between observation samples (rows) and multivariate fluid geochemical data (columns), calculating their contribution to total inertia for each factor. The projection of the rows and columns onto the same set of factorial axes with the same origin enables the development of two-dimensional graphs, which aid in the interpretation of the results.

Correspondence factor analysis, a multivariate statistical technique of proximity and distance measure from the origin of two-dimensional coordinate (factorial) axes, was performed in Excel using XLSTAT software

to identify associations or oppositions between observation samples (rows) and multivariate fluid geochemical data (columns), calculating their contribution to total inertia for each factor. The projection of the rows and columns onto the same set of factorial axes with the same origin enables the development of two-dimensional graphs, which aid in the interpretation of the results. Two-way Contingency Table

The raw data has been turned into a two-way contingency table that displays combinations of two category variables for correspondence analysis. Rows represent data or sample points, such as the geographical location of hot springs (qualitative data), whereas columns contain fluid geochemical parameters. A correspondence map is a popular CA output that uses the proximity and distance of variables to measure association or opposition. This concept holds true when comparing rows to rows or columns to columns (a symmetric biplot from CA) in the same space using main coordinates. The biplot, also known as the multivariate variant of the scatterplot, includes more than two axes that are not necessarily perpendicular, but it allows for the graphical depiction of rows (samples) as points and each column (variable) by an axis from the same origin on the same space or plot. Categorical or qualitative data can be saved and identified by names or labels. Numerical or quantitative data is made up of numbers rather than words or descriptions.

Computational strategy

Correspondence analysis is a statistically based geometric technique that displays the rows and columns of a two-way contingency table as points in a two-dimensional vector space (Benzekri, 1973; David et al., 1977; Davis, 1986; Teil, 1975). In this analysis, the contingency table is looked upon as a joint probability distribution of two random variables, namely, individual observations or samples ($i = 1, 2, 3, \dots, N$) and variables ($j = 1, 2, 3, \dots, M$). The raw data matrix (X_{ij}) is converted into a matrix of joint probability (P_{ij}) of occurrence by dividing each cell entry by the sum of the data values in all rows and columns, i.e.

$$P_{ij} = X_{ij}/K \text{ where } K = \sum_{i=1}^N \sum_{j=1}^M X_{ij}$$

Sum of the probabilities P_{ij} in all rows and columns is given by

$$\sum_{i=1}^N \sum_{j=1}^M P_{ij} = 1$$

the row-totals of each row

$$P_i = \sum_{j=1}^M P_{ij} = K_i/K$$

and the column-totals of each column

$$P_j = \sum_{i=1}^N P_{ij} = K_j/K$$

give the marginal probabilities of each sample and variable respectively

The joint probability distribution matrix generated from the contingency table (X_{ij}) is then turned into a square, symmetric matrix for the computation of eigenvalues and eigenvectors, from which factor loadings for samples and variables are extracted in the usual way (Davis, 1986). The factor loadings are then utilized to represent samples (I) and variables (J) simultaneously on factorial axes. The correspondence of the i -th sample and the j -th variable on the q -th factorial axis is given by

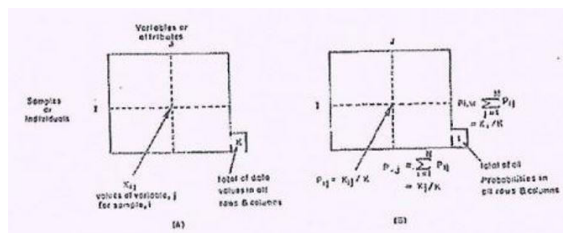


Fig 2. Geometric representation of transformation of (A) two-dimensional contingency table (X_{ij}) having $M(j)$ variables and N observations or samples (i) into (B) joint probability distribution matrix (P_{ij}) by correspondence analysis

$$Sq(i) = \frac{1}{\sqrt{\lambda q}} \sum_{j=1}^M \frac{P_{ij}}{P_i} \quad Vq(j) = \frac{P_{ij}}{P_j}$$

$$Vq(i) = \frac{1}{\sqrt{\lambda q}} \sum_{j=1}^M \frac{P_{ij}}{P_j} \quad Sq(j) = \frac{P_{ij}}{P_j}$$

Where λ_q = eigenvalues or inertia of factorial axis q

$S_{q(i)}$ = abscissae or loadings of i on factorial axis q

$V_{q(j)}$ = abscissae or loadings of j on factorial axis q

$\frac{P_{ij}}{P_i}$

P_i = conditional probability of drawing a variable given that it belongs to sample i

$\frac{P_{ij}}{P_j}$

= conditional probability of drawing a sample given that it belongs to variable of type j

Correspondence analysis dataset (A): a two-way contingency table with 62 rows and 9 columns

The fundamental input data, or in this example, the two-way contingency table, consists of 62 rows or records and 9 columns or categories. Of the 62 rows, 37 represent fluid geochemical data from hot springs in the Extra-Peninsula (Himalayan), whereas 25 reflect data from Indian Peninsula locations.

Distance: Chi-square

Significance level (%): 5

Filter factors Maximum number: 5

Rotation: Varimax (Kaiser normalization) / Based on columns / Number of factors = 2

Test of independence between the rows and the columns:

Chi-square (Observed value)	67098.9188
Chi-square (Critical value)	540.499
DF	488
p-value	<0.0001
alpha	0.05

Test interpretation:

H0: The rows and the columns of the table are independent.

Ha: There is a link between the rows and the columns of the table.

As the computed p-value is lower than the significance level $\alpha=0.05$, one should reject the null hypothesis H0, and accept the alternative hypothesis Ha.

Total inertia: 0.86

Eigenvalues and percentages of inertia:

	F1	F2	F3	F4	F5	F6	F7	F8
Eigenvalue	0.375	0.240	0.160	0.035	0.025	0.013	0.007	0.004
Inertia %	43.624	27.903	18.601	4.061	2.914	1.550	0.837	0.509
Cumulative%	43.624	71.528	90.129	94.190	97.104	98.654	99.491	100.000

Fig. 3. Scree plot showing the percentages of inertia Captured by the new dimensions generated by CA

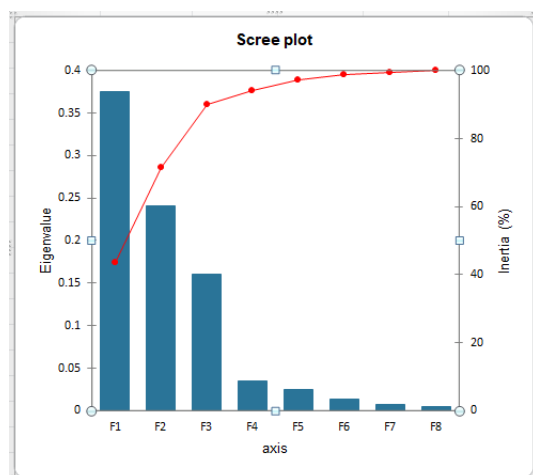


Table 1. A two-way contingency table

	HCO3 mg/L	Cl mg/L	SO4 mg/L	Ca mg/L	Mg mg/L	Na mg/L	K mg/L	F mg/L	B mg/L
1	300.000	163.000	62.000	14.000	5.000	210.000	13.000	12.000	5.000
2	170.000	133.000	36.000	44.000	15.000	88.000	19.000	0.800	33.000
3	490.000	855.000	1244.000	342.000	87.000	600.000	109.000	3.600	138.000
4	210.000	102.000	83.000	30.000	15.000	110.000	19.000	1.200	25.000
5	342.000	232.000	26.000	26.000	1.000	260.000	16.000	10.000	10.000
6	303.000	200.000	340.000	103.000	11.000	260.000	45.000	6.000	13.000
7	173.000	45.000	28.000	13.000	2.000	103.000	5.000	10.000	3.000
8	276.000	170.000	33.000	52.000	12.000	135.000	27.000	3.000	10.000
9	145.000	30.000	55.000	38.000	13.000	30.000	7.000	1.000	3.000
10	15.000	2.000	0.000	3.000	1.000	1.000	0.000	0.200	0.000
11	248.000	72.000	48.000	13.000	2.000	140.000	6.000	5.000	3.000
12	272.000	10.000	14.000	56.000	24.000	8.000	5.000	0.400	0.000
13	445.000	35.000	0.000	50.000	52.000	50.000	10.000	1.200	0.000
14	112.000	1485.000	22.000	70.000	13.000	490.000	37.000	1.600	19.000
15	103.000	8.000	29.000	45.000	44.000	24.000	10.000	0.700	2.000
16	117.000	15.000	30.000	34.000	3.000	30.000	5.000	1.600	0.000
17	861.000	48.000	14.000	14.000	99.000	290.000	43.000	3.000	5.000
18	278.000	12.000	27.000	42.000	26.000	15.000	8.000	0.500	1.000
19	38.000	5.000	0.000	6.000	7.000	2.000	1.000	0.400	0.000
20	953.000	86.000	0.000	0.000	47.000	80.000	83.000	0.000	0.000
21	734.000	12.000	5.000	64.000	10.000	180.000	38.000	2.000	1.000
22	439.000	41.000	21.000	40.000	23.000	163.000	15.000	4.000	2.000
23	254.000	13.000	99.000	13.000	0.000	135.000	6.000	12.500	2.800
24	363.000	17.000	66.000	40.000	8.000	120.000	7.000	10.000	2.000
25	1610.000	85.000	57.000	10.000	2.000	580.000	48.000	10.000	8.000
26	259.000	11.000	1484.000	504.000	22.000	200.000	6.000	2.500	1.000

Correspondence Factor Analysis (CFA) of Multivariate Fluid Geochemistry Data from ..

27	233.000	58.000	383.000	169.000	28.000	10.000	2.000	0.200	0.000
28	32.000	3.000	0.000	9.000	0.000	2.000	0.000	0.400	0.000
29	112.000	30.000	72.000	14.000	1.000	56.000	4.000	6.000	1.000
30	0.000	6.000	12.000	15.000	3.000	9.000	3.000	1.000	0.900
31	0.000	7.000	2.000	27.000	5.000	6.000	2.000	0.200	0.900
32	415.000	596.000	16.000	41.000	21.000	370.000	30.000	3.000	8.000
33	264.000	13.000	10.000	44.000	18.000	19.000	10.000	0.300	0.000
34	49.000	104.000	6.000	7.000	1.000	75.000	3.000	7.000	1.000
35	435.000	10.000	28.000	27.000	11.000	133.000	10.000	2.100	0.000
36	362.000	154.000	370.000	127.000	19.000	150.000	17.000	1.000	0.000
37	353.000	35.000	36.000	54.000	5.000	86.000	9.000	0.000	0.000
38	154.000	1375.000	210.000	204.000	88.000	660.000	18.000	0.700	0.000
39	339.000	165.000	24.000	82.000	16.000	110.000	6.000	0.400	0.900
40	315.000	130.000	33.000	110.000	12.000	70.000	25.000	0.000	0.000
41	390.000	195.000	75.000	65.000	40.000	210.000	5.000	0.000	0.100
42	500.000	140.000	5.000	70.000	40.000	130.000	2.000	1.000	1.200
43	290.000	50.000	5.000	60.000	20.000	30.000	1.000	0.300	1.200
44	190.000	1347.000	5.000	390.000	250.000	6810.000	55.000	0.000	0.000
45	410.000	110.000	25.000	45.000	15.000	95.000	2.000	0.000	0.000
46	150.000	2725.000	10.000	105.000	40.000	1900.000	30.000	0.200	3.000
47	1534.000	2428.000	672.000	9.000	8.000	1167.000	145.000	0.000	0.000
48	195.000	1485.000	0.000	90.000	40.000	875.000	14.000	0.000	0.000
49	183.000	71.000	33.000	40.000	21.000	40.000	2.000	0.000	0.000
50	13.000	4800.000	185.000	186.000	10.000	955.000	13.000	0.000	0.400
51	11.000	850.000	130.000	170.000	0.100	368.000	7.000	2.000	0.400
52	14.000	1210.000	144.000	348.000	0.200	391.000	8.500	7.200	0.000
53	18.000	78.000	242.000	40.000	15.000	155.000	2.000	2.500	0.000
54	71.000	426.000	107.000	32.000	6.000	292.000	4.000	1.500	1.000
55	30.000	375.000	100.000	56.000	1.800	231.000	7.800	4.000	0.400
56	63.000	265.000	108.000	80.000	44.000	148.000	6.000	0.100	0.000
57	177.000	67.000	70.000	3.000	1.000	133.000	0.000	3.000	0.500
58	364.000	30.000	8.000	35.000	3.000	110.000	16.000	0.300	0.000
59	99.000	457.000	128.000	42.000	2.000	360.000	19.000	0.500	0.000
60	366.000	257.000	55.000	96.000	70.000	98.000	15.000	0.200	0.000
61	171.000	50.000	120.600	50.000	7.900	95.000	7.400	4.000	0.000
62	128.600	166.000	182.000	20.000	13.400	208.000	4.000	5.000	0.000

Results after the Varimax rotation (Kaiser normalization):

Rotation matrix:

	D1	D2
D1	-0.931	-0.364
D2	-0.364	0.931

Percentage of variance after Varimax rotation:

	D1	D2	F3	F4	F5
Variability	41.541	29.987	18.601	4.061	2.914
Cumulative	41.541	71.528	90.129	94.190	97.104

Table 2.

Standard coordinates (rows) after varimax rotation			Contribution(rows) after Varimax rotation:		
	D1	D2		D1	D2
1	0.634	-0.214	1	0.004	0.000
2	0.448	0.186	2	0.001	0.000
3	-0.296	1.680	3	0.004	0.140
4	0.601	0.445	4	0.003	0.002
5	0.554	-0.492	5	0.004	0.003
6	0.153	1.178	6	0.000	0.023
7	0.990	-0.198	7	0.005	0.000
8	0.669	-0.162	8	0.004	0.000
9	1.010	0.779	9	0.004	0.003
10	1.820	-0.226	10	0.001	0.000
11	0.950	-0.146	11	0.006	0.000
12	1.989	0.037	12	0.020	0.000
13	1.977	-0.412	13	0.032	0.001
14	-0.978	-0.524	14	0.028	0.008
15	1.152	0.587	15	0.005	0.001
16	1.178	0.521	16	0.004	0.001
17	1.775	-0.631	17	0.056	0.007
18	1.915	0.101	18	0.019	0.000
19	1.851	-0.290	19	0.003	0.000
20	2.197	-0.617	20	0.077	0.006
21	1.969	-0.514	21	0.052	0.004
22	1.541	-0.415	22	0.023	0.002
23	1.096	0.449	23	0.008	0.001
24	1.471	0.103	24	0.018	0.000
25	1.771	-0.616	25	0.097	0.012
26	-0.311	3.596	26	0.003	0.413
27	0.259	2.624	27	0.001	0.078
28	1.829	-0.068	28	0.002	0.000
29	0.671	0.925	29	0.002	0.003
30	-0.371	1.762	30	0.000	0.002
31	-0.260	1.262	31	0.000	0.001
32	0.092	-0.585	32	0.000	0.007
33	1.976	-0.122	33	0.019	0.000
34	-0.183	-0.526	34	0.000	0.001
35	1.791	-0.374	35	0.027	0.001
36	0.349	1.522	36	0.002	0.036
37	1.548	-0.058	37	0.018	0.000

Correspondence Factor Analysis (CFA) of Multivariate Fluid Geochemistry Data from ..

38	-0.781	-0.040	38	0.021	0.000
39	0.879	-0.171	39	0.007	0.000
40	0.953	0.105	40	0.008	0.000
41	0.709	-0.081	41	0.006	0.000
42	1.342	-0.429	42	0.021	0.002
43	1.639	-0.182	43	0.016	0.000
44	-0.482	-1.032	44	0.027	0.123
45	1.356	-0.293	45	0.017	0.001
46	-0.937	-0.771	46	0.056	0.038
47	-0.053	-0.016	47	0.000	0.000
48	-0.785	-0.704	48	0.021	0.017
49	0.992	0.162	49	0.005	0.000
50	-1.316	-0.383	50	0.137	0.012
51	-1.050	0.106	51	0.022	0.000
52	-1.050	0.205	52	0.030	0.001
53	-0.617	2.109	53	0.003	0.032
54	-0.728	0.002	54	0.006	0.000
55	-0.857	0.192	55	0.008	0.000
56	-0.496	0.546	56	0.002	0.003
57	0.632	0.152	57	0.002	0.000
58	1.697	-0.477	58	0.021	0.002
59	-0.624	0.010	59	0.006	0.000
60	0.670	0.009	60	0.006	0.000
61	0.548	1.028	61	0.002	0.007
62	-0.167	0.825	62	0.000	0.006

Table 3.

Principal coordinates (columns) after varimax rotation			Contributions (columns) after Varimax rotation:			Squared cosines (rows) after Varimax rotation:		
	D1	D2		D1	D2		D1	D2
HCO3 mg/	1.000	-0.029	HCO3 mg/	0.643	0.001	HCO3 mg/	0.960	0.001
Cl mg/L	-0.585	-0.206	Cl mg/L	0.296	0.051	Cl mg/L	0.596	0.074
SO4 mg/L	-0.065	1.388	SO4 mg/L	0.001	0.712	SO4 mg/L	0.002	0.959
Ca mg/L	0.016	0.590	Ca mg/L	0.000	0.080	Ca mg/L	0.000	0.460
Mg mg/L	0.438	-0.023	Mg mg/L	0.010	0.000	Mg mg/L	0.154	0.000
Na mg/L	-0.225	-0.376	Na mg/L	0.038	0.146	Na mg/L	0.097	0.271
K mg/L	0.484	0.045	K mg/L	0.009	0.000	K mg/L	0.301	0.003
F mg/L	0.612	0.175	F mg/L	0.002	0.000	F mg/L	0.086	0.007
B mg/L	0.115	0.781	B mg/L	0.000	0.009	B mg/L	0.002	0.097

Fig. 4. Primary Dataset (A): Biplot showing the row, column and supplementary variables in two-dimensional space: Combined Peninsula and Extra-Peninsula/62 rows and 9 columns

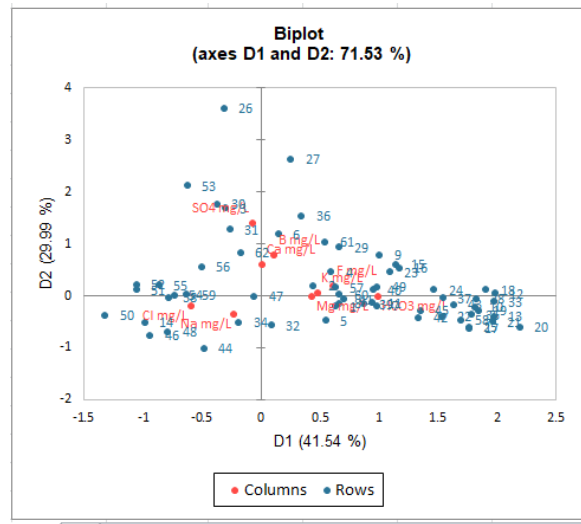


Fig. 5. Biplots: Primary Dataset (A) – 62 rows (Left) and 9 columns (Right)

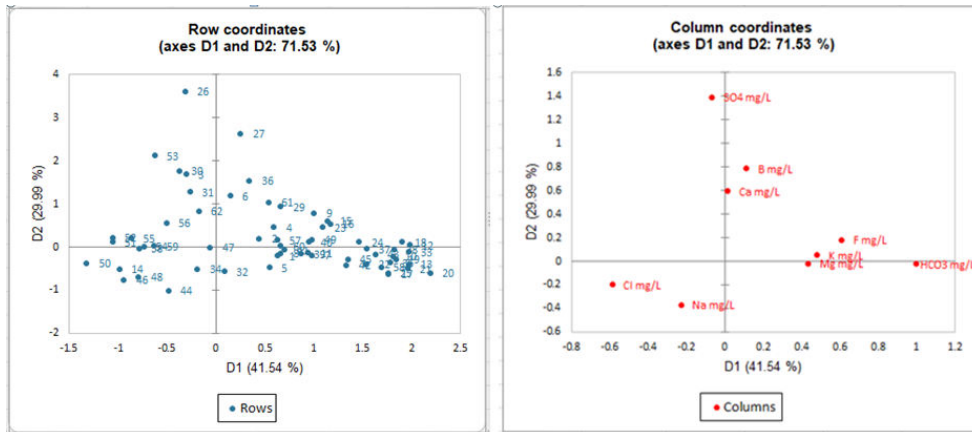


Fig. 6. Biplot showing the row, column and supplementary variables in two-dimensional space: Left: Peninsula/ 25 rows and 9 columns (B) and Right: Extra-Peninsula/ 37 rows and 9 columns (C)

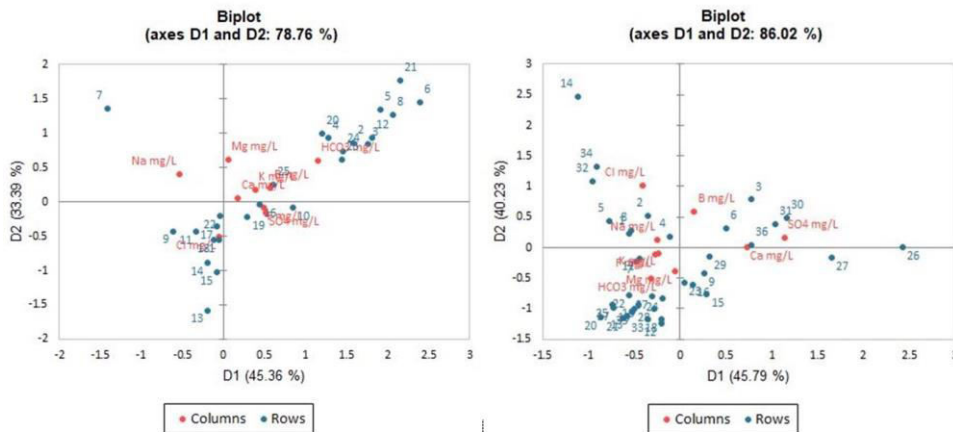


Fig. 7. Biplots: Subset (B) - Peninsula/ 25 rows (Left) and 9 columns (Right)

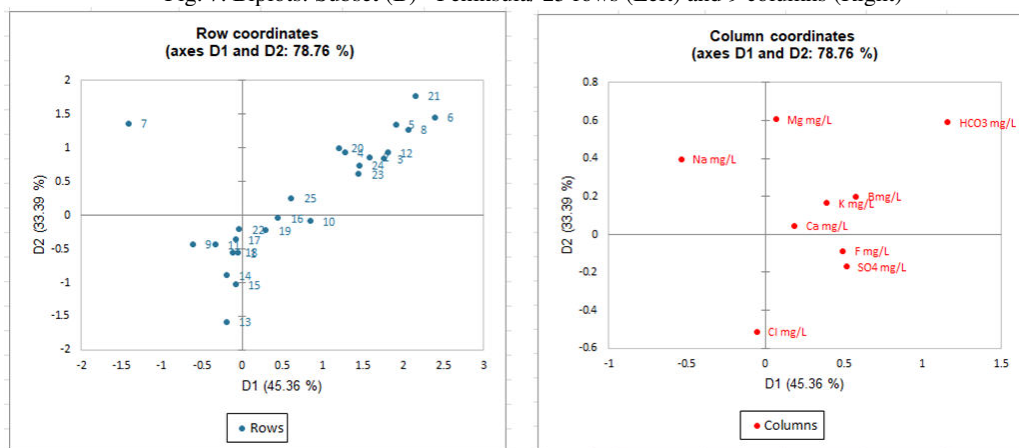
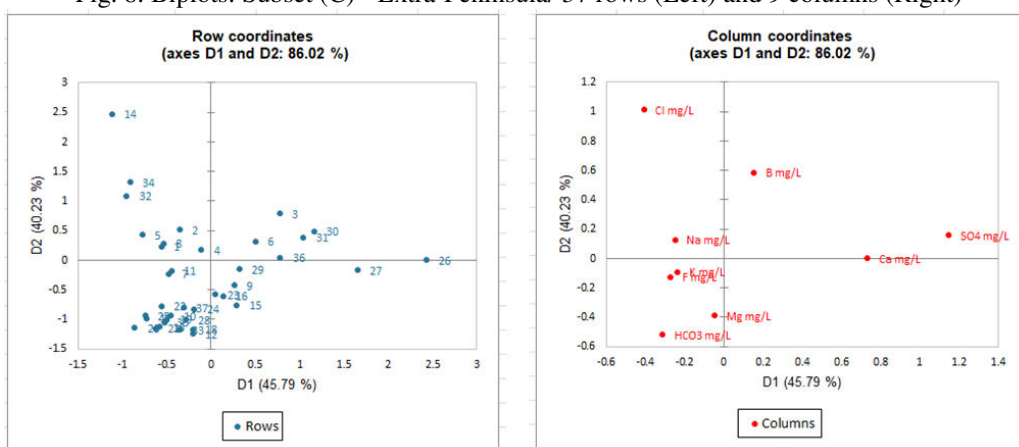


Fig. 8. Biplots: Subset (C) - Extra-Peninsula/ 37 rows (Left) and 9 columns (Right)



Interpreting the results

The findings of correspondence analysis are visually interpreted using scree plots and biplots, as well as statistically using output statistics.

II. Discussion

A few crucial factors to remember while analyzing correspondence analysis include: 1) Use raw data to verify findings, 2) The further things are from their origin, the more discriminating they are, 3) The closer anything is to its origin, the less distinct it is, 4) The more variance explained, the less insights will be lost, 5) Proximity between row labels suggests similarity, 6) Proximity between column labels shows similarity; and 7) If a tiny angle connects a row and column label to the origin, they are most likely related,

8) If a row and column label form a 90-degree angle with the origin, they are most likely unrelated, 9) If a row and column label are on opposing sides of the origin, they are most likely negatively connected, and 10) The further a point from the origin, the greater its positive or negative relationship.

The primary dataset (A), the two-way contingency table (Table-1), contains 62 rows and 9 columns. It is divided into two subsets: 25 rows and 9 columns for Peninsular India (B), and 37 rows and 9 columns for Extra-Peninsular India (C). While the findings of the correspondence analysis of the main dataset has been replicated in their entirety, only the results of biplots for the two subsets have been provided with the goal of matching the overall results.

A) Primary dataset: contingency table with 62 rows and 9 columns: Looking at the biplots, three clusters are discernible: a) predominantly alkaline HCO₃-Mg-K-F further right from the origin in close proximity and aligning with the factorial D1 axis; b) acidic to neutral group SO₄-B-Ca further north of the origin aligning with the vertical D2 factorial axis and forming a 90-degree angle in respect of group (a) with the origin, thus uncorrelated; and c) acidic to neutral group Cl-Na making a tiny angle that connects a row and column label to the the left of the origin close to the factorial D1 axis.

B) Subset Peninsula with 25 rows and 9 columns: Biplots show two different opposing groups north and south of the origin: a) the generally alkaline group HCO₃-Mg-Ca-Na-K-B, which is negatively related with b) the mostly acidic group Cl-SO₄-F, south of the origin.

C) Subset Extra-Peninsula with 37 rows and 9 columns: Biplots show two opposing groups north and south of the origin: a) the typically acidic Cl-SO₄-B-Na north of the origin, which is negatively connected to b) the largely alkaline HCO₃-Mg-Ca-K-F south of the origin. Biplots of subsets (B) and (C) swap places across the origin of two-dimensional coordinate (factorial) axes.

Conclusion

To interpret a correspondence analysis, key output includes principal components, inertia, proportion of inertia, quality, mass, and several graphs (Tables 1-3; Figs. 1-9)

1. The first component (D1) accounts for 41.541% of the inertia and the second component (D2) accounts for 29.987% of the inertia. Together, these 2 components account for 71.528% of the total inertia (Cumulative = 0.71.528). Therefore, specifying 2 components for the analysis may be sufficient.
2. In the squared cosines (rows) after varimax rotation table, the highest quality values occur for HCO₃ (0.960) and SO₄ (0.959). Therefore, these two rows are best represented by the two components.
3. The proximity of the alkaline group (HCO₃, F, K, and Mg) in feature space shows positive associations, as does the tighter angle between the weak acidic and neutral group (Cl and Na) with regard to the origin. They have a lower correlation.
4. The row plot shows the row principal coordinates. Component D1, which best explains Cl and HCO₃, shows these two fields farthest from the origin, but with opposite signs. Component D1 contrasts the Cl and Na with HCO₃, F, K, Mg. Component D2 contrasts SO₄, B and Ca with Cl and Na.
5. Alkaline group HCO₃ and acidic group SO₄ representing the principal components D1 and D2 respectively form a 90-degree angle with the origin, they are most unrelated.
6. The quality values determined by the proportion of the row inertia or column inertia for the rows and/or columns can help to interpret the components. Quality is always a number between 0 and 1. Larger quality values indicate that the row or column is well represented by the components. Lower values indicate poorer representation.

References

- [1]. Amitabha Roy, 2023. Geostatistics As Applied To The Fluid Geochemistry Of Indian Hot Springs. Jour. Applied Geology and Geophysics, V. 11, Issue 4, Ser. II, pp 1-37
- [2]. A.Roy, 1994. Gthermis – An Information Management And Analysis System For Geothermal Data of India, A Field Season Report (1993-94).
- [3]. Benzekri, J.P, 1973. L'Analyss des donnaes, Tome II, Paris, pp. 619
- [4]. David, M. et al., 1977. Correspondence analysis. Quart. Colorado Sch Mines, V.72, pp. 11-57
- [5]. Davis, J.C., 1986. Statistical data analysis in geology. John Wiley, N.Y., pp. 646
- [6]. Jonathan Craig, 2013. Hot Springs And The Geothermal Energy Potential Of Jammu & Kashmir State, N.W. Himalaya, India.
- [7]. Kaiser, H.F. (1958). The varimax criterion for analytic rotation in factor analysis. Psychometrika, 23, pp.187-200.
- [8]. Ravi Shankar et al. Geothermal Atlas of India. GSI Spec Publ, (1991)
- [9]. Teil, H, 1975. Correspondence factor analysis:an outlining of method. Mathematical Geology, 7, pp. 3-12