# Geosoft Technical Note:

# Principal Component Analysis (PCA) and Factor Analysis in OASIS montaj™

**Technical Note Version:**

**OASIS montaj v4.3, SP2, (Dec. 8, 99)**

**GEOSOFT**

**KNOWLEDGE FROM DATA™**

**http://www.geosoft.com**

# Principal Component and Factor Analysis

Principal Component Analysis (PCA) and Factor Analysis are two methods that can help reveal simpler patterns within a complex set of variables. In particular, these methods seek to discover if the observed variables can be explained largely or entirely in terms of a much smaller number of variables called *factors*.

In mineral exploration, the most common application of these multivariate analysis methods is to characterize and map interrelationships within high volume surface geochemistry data sets. Data volumes are a growing problem as geochemists seek to extract more information and knowledge from data sets with up to 50 or more variables.

The desire to simplify processing and analysis of large data sets is renewing interest in PCA and Factor Analysis algorithms and presentation. Requests from major exploration groups for these capabilities led to Geosoft developing PCA and Varimax Fator Analysis in its latest v4.3 release.

This short article reviews some of the considerations in applying PCA and Factor Analysis to exploration problems.

# Quick Review of PCA and Factor Analysis

PCA and Factor Analysis are often misunderstood due to their abstract mathematical nature. Concepts such as eigenvalues, basis vectors (eigenvectors), loadings and scores are difficult to map into the framework of real world geological problems.

Without examining the mathematics in detail, the basic starting point is to compute a correlation matrix using all variables and samples in a verified and properly subsetted geochemical database. This matrix is then used to compute a new set of "artificial" variables called eigenvectors, each with its own distinct eigenvalue. The numerical value of the eigenvalue indicates the contribution of the "artificial" variable or *factor* to the total variation in the data set.

Two other entities also play a role in this mathematical puzzle:

- Loadings. Loadings express the influence of each original variable, such as Au for example, within the factor.

- Scores. Scores are numbers that express the influence of an eigenvector on a specific sample. Scores enable spatial mapping of factors on individual samples.

# Comparison of PCA and Factor Analysis

The main commonality between PCA and Factor Analysis is that they both have eigenvectors, eigenvalues, loadings and scores. Some differences are:

- PCA is often used as a simple starting point in multivariate analysis.

- PCA eigenvectors cumulatively account for all the variability in the data set whereas Factor Analysis results include an unresolved component.

- For this reason, Factor Analysis is often considered to be "statistical" in nature rather than purely mathematical as in PCA.

- Factor Analysis results are often transformed through Varimax and other methods to optimize eigenvectors for interpretation.

## Before Starting…

This technical note is intended to provide a general overview of Principal Component and Factor Analysis in **OASIS montaj**. The interpretation of results, are for demonstration purposes only and may not reflect the true geology of this region. This note describes one possible interpretation of the data. Interpretation of results typically depends on requirements of the project and user's specifications.

Before applying these methods, users should have a solid understanding of the problem they are trying to solve. They must have strategies for determining how they will solve the following questions:

1. Has the data been verified adequately to obtain meaningful results?
2. How many different factors are needed to explain the pattern of relationships among variables?
3. What is the geologic significance of the factors?
4. What proportion of the elements are explained by the most important principal components?
5. What is the spacial relationship of factors (i.e. how do they map)?

**Note:** To perform most analytical methods in the Chimera system, your assay channels **must** be classified as **ASSAY** channels. When importing your data using *ChemImport*, the attribute *Class* for assay channels is defaulted as **ASSAY**. To classify assay channels not imported through *ChemImport*, highlight the assay channel header and right click. Select *Attributes* from the popup menu, the *Assay Information* dialog box will be displayed. In the *Class* box, specify **ASSAY**, click **[OK]** to continue.
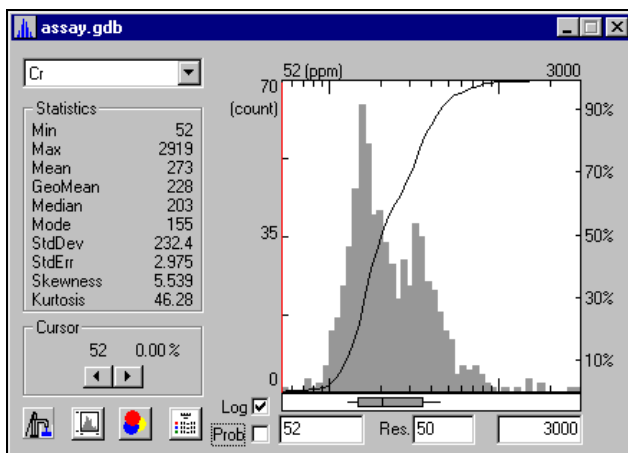
# Step 1: Data Verification

Data verification is absolutely critical in multivariate analysis. Historically, lack of care in preparing data and poor overall results have led to a negative perception of these methods.

Geosoft tools such as Histogram Analysis and the new v4.3 Scatter Analysis tool can help quickly identify data with detection limit problems or outliers. Outliers can be defined as data values outside some statistical range (for example, the 95[th] percentile or the mean $\pm$ 2 standard deviations). After these data or variables are identified, Geosoft's subsetting tools and math expressions can assist in eliminating problematic data or outliers before proceeding with PCA and Factor Analysis.

**Histogram Analysis Tool**

The Histogram Analysis Tool can be used to determine the data distribution and population range. The Histogram tool can also be used to identify outliers by means of the cumulative % option.
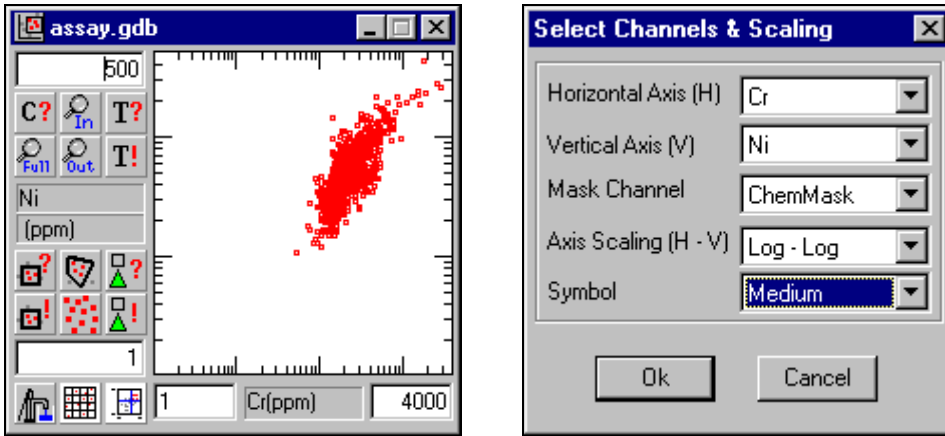


The histogram analysis tool enables you to select the channel you wish to view from within the tool. The default channel is the one selected in your open database. You can select a channel from the drop-down menu in the top left corner of the tool. The original data statistics are displayed for the selected channel.

For more information on the Histogram Analysis Tool press the **[F1]** key.

**Scatter Analysis Tool**

The Scatter Analysis Tool can be used to determine the distribution and population ranges of the data. The Scatter Analysis Tool also enables you to specify a masking channel to record any classification operations performed via other functions in the Scatter Analysis tool. Working with a masking channel is an important part of classification and subsetting since it is through this channel that the system determines which values to select for manipulation and plotting.

The Scatter Analysis tool enables you to plot one channel against another and to interactively interrogate data and information visually (i.e. plotted in an active map on your desktop) and numerically (i.e. contained within the database).

To select channels, click the ( **C?** ) button, the *Select Channels & Scaling* dialog box is displayed. This dialog box enables you to select the channels you wish to plot, the masking channel, the type of scaling (log-log and others) and the symbol size.

For more information on using the Scatter Analysis Tool press the **[F1]** key.

**Math Expressions**

Math expressions enable you to define a range for the data. To apply a math expression highlight an assay channel (click three times on the channel header cell), press ' = ' and then type the function. The following are two math expressions you can use to define a data range:

**WINDOW(exp,Minvalue,Maxvalue)**

This function when applied to a channel, returns a DUMMY value if the specified expression is NOT in the Minimum – Maximum range, otherwise returns the expression value. For example, WINDOW(Au,5,700) will dummy out all the values below (less than) 5 and above (greater than) 700 in the Au channel.

**CLIP(exp,Minvalue,Maxvalue)**

This function when applied to a channel, returns the expression value if the expression is in the range Min-Max, otherwise returns the clipped Min or Max value. This functions a little differently than Window (above). Values not in the range are not dummied but held as the min or max value. For example, CLIP(Au,5,700) keeps all the values in the range 5 – 700 as normal values and assign any value not in the range to the nearest min or max value. Simply, any number less than 5 is retained as 5 and any number greater than 700 is kept as 700.

# Step 2: Applying PCA and Factor Analysis in Geosoft

Principal component analysis is a mathematical method designed to reveal the relationships between two or more (often many) variables. Measurements which include many variables are commonly encountered in mineral exploration and geochemistry. PCA and Factor Analysis determines the significance of the correlation of the variables.

The basic problem in PCA and Factor Analysis is determining the number of starting factors. The typical approach is to select an arbitrary, yet standard number, i.e. 5 or 10. In Geosoft, the default is a maximum of 10 components.

Geosoft's implementation of PCA and Factor Analysis is performed via two dialog boxes, *Principal Component Analysis* and *Principal Component Synthesis*.

## Principal Component Analysis

In Geosoft, the *Principal Component Analysis* GX enables you to select the ASSAY channels to include in the analysis, specify the maximum number of components, specify the eigenvalue cutoff limit for Varimax analysis, normalize score values, and save scores as channels in the database.

The following is a brief description of the *Principal Component Analysis* parameters:

**Channels to include**

The channels to include must be members of the **ASSAY** class. You can select either *All ASSAY channels* or *Displayed ASSAY channels*.

**Note:**   To perform the *Principal Component Analysis*, your assay channels **must** be classified as **ASSAY** channels. When importing your data using *ChemImport*, the classification of the channels is accomplished automatically. To classify assay channels not imported through *ChemImport*, highlight the assay channel header and right click. Select *Attributes* from the popup menu, the *Assay Information* dialog box will be displayed. In the *Class* box, specify **ASSAY**, click **[OK]** to continue.

**Maximum # of components for output**

The number specified will determine the number of principal components created. The principal component results will be displayed in the database and in the log file (**princomp.log**).

Prior to the calculation of the principal components, the data is transformed into a condition agreeable to analysis. Depending on whether the assay channel's *Logarithmic Distribution* attribute is set to **Yes** the logarithms of the data are taken. The mean is then removed, and finally the data are normalized through division by the variance (standard deviation).

**Note:**   When importing your data using *ChemImport*, the attribute *Logarithmic Distribution* for assay channels is defaulted as **Yes**. To set the logarithmic distribution for assay channels not imported through *ChemImport*, highlight the assay channel header and right click. Select *Attributes* from the popup menu, the *Assay Information* dialog box will be displayed. In the *Logarithmic Distribution* box, select **Yes**, click **[OK]** to continue.

**Eigenvalue cut-off for Varimax**

A correlation matrix is produced from the transformed data. An eigenvector decomposition is performed to determine the eigenvectors (which are directionally equivalent to the principal components) and eigenvalues. The relative significance of each component is indicated by its eigenvalue. The first principal component will have the largest eigenvalue, and succeeding components will have smaller eigenvalues, as their significance in the data decreases.

The cut-offs determine the number of Varimax factors derived from the Factor Analysis. All components with eigenvalues less than this value are rejected, and the principal component loadings are re-computed using Kaiser's Varimax scheme.

**Save scores as channels?**

*Scores* describe the contribution of each principal component to each data point. A score channel is created for each principal component specified. The score channels are displayed in the database as, SC1, SC2, etc. Score values are commonly presented spatially on a map using colour coded symbols.

**Normalize Scores**

When you select **Yes** to normailize the score values, the values are transformed so that they lie between 0 and 100. If the range of scores values is A to B, then a value X is transformed using the formula:

$$X'=(X-A)*100/(B-A)$$

## Interpreting the PCA and Varimax Log file

When you run the PCA/Varimax dialog, the system generates a **princomp.log** file. The PCA/Varimax log file contains the computed results and is a key interpretation tool.

The following is a partial example of a **princomp.log** file. The complete file can be found in the appendix.

**Section (1)**

**Section (2)**

**Section (3)**

**Section (4)**

```
Principal Component analysis: .\assay.gdb

Number of channels included: 9
Number of principal components displayed: 9

Data Transformations
--------------------
  Cr  : Logarithmic Normal Distribution
  Co  : Logarithmic Normal Distribution
  Ni  : Logarithmic Normal Distribution
  Cu  : Logarithmic Normal Distribution
  Zn  : Logarithmic Normal Distribution
  As  : Logarithmic Normal Distribution
  Zr  : Logarithmic Normal Distribution
  Mo  : Logarithmic Normal Distribution
  Pb  : Logarithmic Normal Distribution


Correlations of Standardized Data
---------------------------------
          Cr      Co      Ni      Cu      Zn      As      Zr      Mo      Pb

  Cr    1.000   0.339   0.775   0.411   0.273   0.073  -0.082  -0.189   0.078
  Co    0.339   1.000   0.443   0.260  -0.015   0.102  -0.139  -0.122  -0.115
  Ni    0.775   0.443   1.000   0.668   0.316   0.048  -0.426  -0.338  -0.144
  Cu    0.411   0.260   0.668   1.000   0.573  -0.137  -0.339  -0.421   0.016
  Zn    0.273  -0.015   0.316   0.573   1.000  -0.288  -0.007  -0.363   0.494
  As    0.073   0.102   0.048  -0.137  -0.288   1.000  -0.252  -0.113  -0.285
```

The following example describes each section of the **princomp.log** file and gives one possible interpretation of the results.

**Section (1):**

The first comment line in the log file identifies the type of analysis and the name of the current database.

In our example log file, the type of analysis is a **Principal Component Analysis** and the current database is identified as **assay.gdb**.

**Section (2):**

The second comment line in the log file indicates the number of ASSAY channels included in the analysis.

In our example, the number of assay channels included is **9**.

**Section (3):**

The third comment line in the log file displays the number of principal components displayed. Determining how many different factors needed to explain the pattern of relationships among variables is a key task. Strategies include evaluation of the amount of variability reflected in each eigenvector and the acceptable total variance accounted for by a set number of eigenvectors.

In our example, the number of principal components is **9**.

**Section (4):**

This section identifies the transformation performed on each ASSAY to make the data agreeable to analysis.

In our example, each assay channel has been transformed to **Logarithmic Normal Distribution**.

**Section (5):**

This section contains the correlation matrix of the standardized data. To standardize the data, it is first transformed (see above), then the statistical mean value is removed and finally the data is divided by the variance (standard deviation). The maximum positive correlation possible is 1.0, and the maximum negative (or inverse) correlation possible is –1.0.

In our example, the maximum positive correlation is **Cr**, **Ni** with a value of (**0.775**). The maximum inverse correlation is **Cu**, **Mo** with a value of (**-0.440**).

**Section (6):**

This section contains the list of factors, their eigenvalues from the correlation matrix and the cumulative % of each factor. The eigenvalues represent the relative significance of each component.

In our example, the 1$^{st}$ factor has an eigenvalue of **3.388** and represents **37.6%** of the total factors. The first **3** factors have eigenvalues above **1.0** and represent **74.5%** of the total factors.

**Section (7):**

This section contains the eigenvectors derived from the correlation matrix. Eigenvectors are directionally equivalent to the principal components and their eigenvalues.

In our example, for Principal Component 1 (**PC1**) the variable with the greatest eigenvector is **Ni**, followed by **Cu**, **Cr**, **Co** etc.

**Section (8):**

This section contains the Principal Component Loadings. The loadings are the eigenvectors ordered in terms of the size of the eigenvalues and scaled by the square root of the eigenvalues.

In our example, for **PC1**, the largest loading is **Ni**, followed by **Cu**, **Cr**, **Co** etc.

**Section (9):**

This section is displays the proportion of the variables explained by the factors.

In the example log, **Ni** for **Factor 1** has a cumulative component value of (**0.801**), whereas **Zn** only has an cumulative component value of (**0.192**). Note: **Zn** and **As** have no component value in **Factor 9**.

**Section (10):**

This section contains the Varimax Principal Component Loadings (VPCL). The loadings are the eigenvectors ordered in terms of the size of the eigenvalues and scaled by the square root of the eigenvalues.

In the example log, there are only **3 VPCL** displayed. The calculation of the **VPC** considers the *Eigenvalue cutoff for Varimax* specified in the *Principal Component Analysis* dialog box. In this example the eigenvalue cutoff was specified as **1.0**, therefore components with eigenvalues less than **1.0** have not been considered in this calculation. For the VPCL, the greatest loading is **Ni**, followed by **Cr**, **Co**, and **Cu** etc.
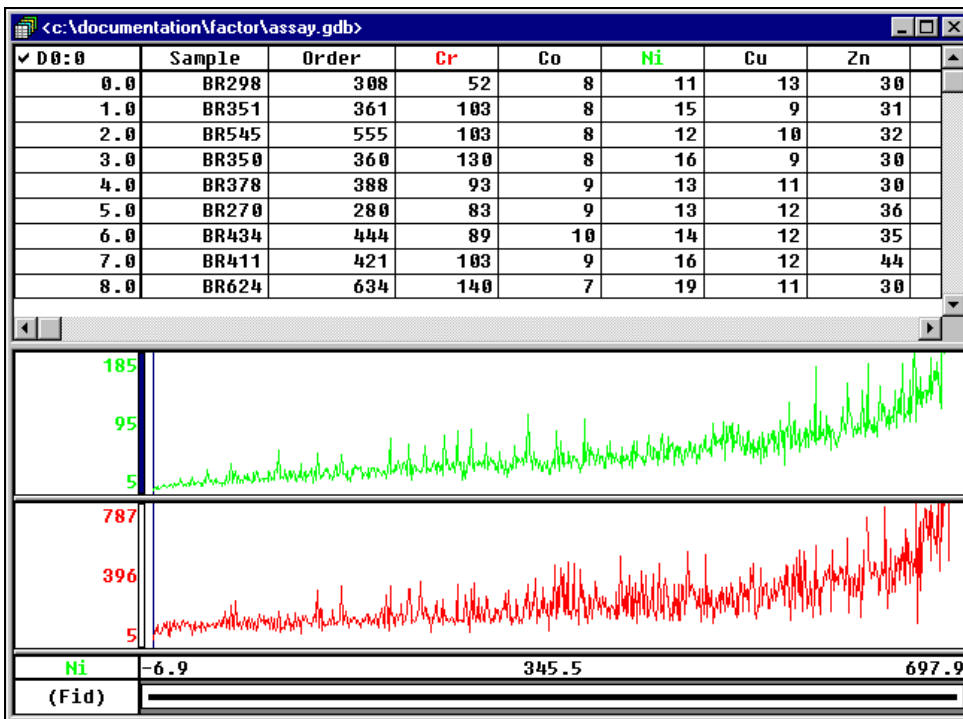
**Section (11):**

This section contains the proportion of variables explained by Varimax factors. This list displays the cumulative factors of each component for each element.

In the example, **Ni** for **Factor 1** has a cumulative component value of (**0.792**), whereas **Zn** only has an cumulative component value of (**0.062**). Note: as only **3** Varimax Principal Components are displayed the cumulative components do not total **1.0**.

# Step 3: Determining Geologic Significance

Determining the geologic significance of the factors requires an understanding of the target, lithogeochemistry and alteration patterns in the area under investigation. Typically, senior geochemical expertise and the knowledge of the exploration project are pooled to interpret results.

An absolutely critical step is to compare the hypothetical factors with the observed results. One simple method to do this is to plot the profiles of the variables that weight most strongly (positively or negatively) in eigenvectors and to sort them using the Varimax (VSC) fields that has been added to the original geochemical database. This type of presentation can also help determine how much purely random or unique variance each observed variable includes.



The database above was sorted by one channel, using the VCS1 as the reference channel. Then the variables that show the strongest positive correlation Cr and Ni are displayed in profile format. This visualisation technique clearly shows the relationship between the variables.

# Step4: How Variables Relate to Principal Components

The ability to synthesise the Principal Components back to the original data units may help in verifying your results. By synthesising the data using only the most important components you can determine what proportion of the variables (elements) are explained by the most important principal components. One way to evaluate the contribution of the principal components is to remove the influence of the less important components in order to bring out more clearly the more important components. Geosoft has implemented a synthesis option that enables you to synthesis a specified number of principal components.

Specify the number of components to re-synthesis as less than your original starting components. The system will calculate the data and display the results in the database. The resultant channels will contain the same name as the original Assay channels with "_#" on the end (the # will reflect the number of components used to re-synthesis the data).\

The following is a brief description of the *Principal Component Synthesis* parameters:

**Channels to include**

The channels to include must be members of the **ASSAY** class. You can select either *All ASSAY channels* or *Displayed ASSAY channels*.

**Note:**   To perform the Principal Component Analysis, your assay channels **must** be classified as **ASSAY** channels. When importing your data using *ChemImport*, the classification of the channels is accomplished automatically. To classify assay channels not imported through *ChemImport*, highlight the assay channel header and right click. Select *Attributes* from the popup menu, the *Assay Information* dialog box will be displayed. In the *Class* box, specify **ASSAY**, click **[OK]** to continue.

**# of components in synthesis**

The number of principal components to use to re-synthesis the data. This number must be less than the number of original components.

| Element | Original Data | 7 Components in Re-Synthesis | 4 Components in Re-Synthesis |
|---------|---------------|------------------------------|------------------------------|
| **Cr** | 52 | 48 | 55 |
| **Ni** | 11 | 13 | 13 |
| **Co** | 8 | 8 | 6 |

The table above compares the values of the variables Cr, Ni, Co from the original data with the results derived from re-synthesising the data using 7 components and 4 components.

# Step 5: Mapping Factor Analysis Results

After results are verified, a common method of presenting the data spatially is to display the results in a map using proportional, colour-coded symbols. Because the VSC results give the most clearly differentiated geologically interpretable results, mapping these values and integrating the results with a geology grid of the area may assist in showing the trends in your data.



In this example, Factor 1 (Cr, Ni, Co) mapped the Kv – Bv boundary as identified on the geological map. Other factor scores could also be mapped to confirm their relationship to the geology.

For more information on Factor Analysis
http://www.yorku.ca/dept/psych/lab/psy6140/fa/factorbi.htm

# Appendix: (princcomp.log)

**Section (1)** ——— Principal Component analysis: .\assay.gdb

**Section (2)** ——— Number of channels included: 9

**Section (3)** ——— Number of principal components displayed: 9

**Section (4)**

Data Transformations
--------------------
```
 Cr  : Logarithmic Normal Distribution
 Co  : Logarithmic Normal Distribution
 Ni  : Logarithmic Normal Distribution
 Cu  : Logarithmic Normal Distribution
 Zn  : Logarithmic Normal Distribution
 As  : Logarithmic Normal Distribution
 Zr  : Logarithmic Normal Distribution
 Mo  : Logarithmic Normal Distribution
 Pb  : Logarithmic Normal Distribution
```

**Section (5)**

Correlations of Standardized Data
---------------------------------
```
         Cr      Co      Ni      Cu      Zn      As      Zr      Mo      Pb

 Cr    1.000   0.445   0.775   0.412   0.273   0.087  -0.082  -0.191   0.011
 Co    0.445   1.000   0.580   0.378   0.088   0.104  -0.205  -0.255  -0.137
 Ni    0.775   0.580   1.000   0.668   0.316   0.083  -0.426  -0.403  -0.173
 Cu    0.412   0.378   0.668   1.000   0.572  -0.078  -0.338  -0.440   0.031
 Zn    0.273   0.088   0.316   0.572   1.000  -0.243  -0.007  -0.285   0.539
 As    0.087   0.104   0.083  -0.078  -0.243   1.000  -0.307  -0.275  -0.331
 Zr   -0.082  -0.205  -0.426  -0.338  -0.007  -0.307   1.000   0.771   0.405
 Mo   -0.191  -0.255  -0.403  -0.440  -0.285  -0.275   0.771   1.000   0.088
 Pb    0.011  -0.137  -0.173   0.031   0.539  -0.331   0.405   0.088   1.000
```

**Section (6)**

Eigenvalues of correlation matrix
---------------------------------
```
 Factor  Eigenvalue  cum. %
 --------------------------
    1       3.388     37.6
    2       2.054     60.5
    3       1.267     74.5
    4       0.774     83.1
    5       0.577     89.6
    6       0.456     94.6
    7       0.221     97.1
    8       0.178     99.0
    9       0.086    100.0
```

**Section (7)**

Eigenvectors of correlation matrix
----------------------------------
```
         PC1     PC2     PC3     PC4     PC5     PC6     PC7     PC8     PC9

 Cr     0.367   0.148   0.450   0.272   0.369   0.410  -0.049  -0.259  -0.440
 Co     0.341  -0.022   0.395  -0.056  -0.829  -0.081  -0.124   0.015  -0.119
 Ni     0.486   0.045   0.266  -0.071   0.217   0.099   0.088   0.398   0.680
 Cu     0.424   0.214  -0.109  -0.234   0.141  -0.596   0.499  -0.206  -0.202
 Zn     0.238   0.506  -0.312   0.151   0.063  -0.249  -0.690   0.149  -0.047
 As     0.105  -0.445  -0.005   0.791   0.029  -0.388   0.050   0.105   0.017
 Zr    -0.339   0.369   0.401   0.194  -0.013  -0.257  -0.039  -0.537   0.441
 Mo    -0.379   0.167   0.502  -0.085   0.166  -0.312   0.023   0.593  -0.303
 Pb    -0.080   0.558  -0.220   0.412  -0.277   0.290   0.495   0.243  -0.009
```

**Section (8)**

```
Principal component loadings
----------------------------
          PC1      PC2      PC3      PC4      PC5      PC6      PC7      PC8      PC9

  Cr     0.676    0.213    0.507    0.239    0.280    0.277   -0.023   -0.109   -0.129
  Co     0.628   -0.032    0.445   -0.049   -0.630   -0.055   -0.058    0.006   -0.035
  Ni     0.895    0.064    0.300   -0.063    0.165    0.067    0.041    0.168    0.199
  Cu     0.780    0.307   -0.123   -0.206    0.107   -0.402    0.235   -0.087   -0.059
  Zn     0.438    0.726   -0.352    0.133    0.048   -0.168   -0.324    0.063   -0.014
  As     0.194   -0.638   -0.006    0.695    0.022   -0.262    0.024    0.044    0.005
  Zr    -0.624    0.529    0.451    0.171   -0.010   -0.174   -0.018   -0.226    0.129
  Mo    -0.698    0.240    0.565   -0.075    0.126   -0.211    0.011    0.250   -0.089
  Pb    -0.147    0.799   -0.247    0.362   -0.210    0.195    0.233    0.102   -0.003
```

**Section (9)**

```
Proportion of variables explained by factors
--------------------------------------------
            No. of factors
  Variable    1        2        3        4        5        6        7        8        9

  Cr        0.457    0.502    0.759    0.816    0.895    0.971    0.972    0.983    1.000
  Co        0.394    0.395    0.593    0.596    0.992    0.995    0.999    0.999    1.000
  Ni        0.801    0.805    0.895    0.899    0.926    0.930    0.932    0.960    1.000
  Cu        0.609    0.703    0.718    0.760    0.772    0.934    0.989    0.997    1.000
  Zn        0.192    0.719    0.842    0.860    0.862    0.891    0.996    1.000    1.000
  As        0.038    0.445    0.445    0.928    0.929    0.997    0.998    1.000    1.000
  Zr        0.389    0.669    0.872    0.901    0.902    0.932    0.932    0.983    1.000
  Mo        0.487    0.544    0.864    0.869    0.885    0.929    0.930    0.992    1.000
  Pb        0.022    0.660    0.722    0.853    0.897    0.935    0.990    1.000    1.000
```

**Section (10)**

```
Varimax Principal component loadings
------------------------------------
          PC1      PC2      PC3

  Cr     0.868    0.071    0.011
  Co     0.753   -0.139   -0.082
  Ni     0.890    0.021   -0.320
  Cu     0.586    0.388   -0.474
  Zn     0.249    0.834   -0.292
  As     0.043   -0.579   -0.328
  Zr    -0.128    0.288    0.879
  Mo    -0.164   -0.027    0.914
  Pb    -0.133    0.822    0.170
```

**Section (11)**

```
Proportion of variables explained by varimax factors
-----------------------------------------------------
            No. of factors
  Variable    1        2        3

  Cr        0.754    0.759    0.759
  Co        0.567    0.586    0.593
  Ni        0.792    0.792    0.895
  Cu        0.343    0.493    0.718
  Zn        0.062    0.757    0.842
  As        0.002    0.337    0.445
  Zr        0.016    0.099    0.872
  Mo        0.027    0.028    0.864
  Pb        0.018    0.693    0.722
```