# Demonstrating the Mechanics of Principal Component Analysis via Spreadsheets

## Abstract

Principal component analysis (PCA) is a popular multivariate statistical method that is used for dimensionality reduction. When teaching PCA in a marketing research or business analytics course, the mechanics of the analysis are often not communicated to the students. Students observe computer output that contains information pertaining to eigenvalues, component loadings, and rotated loadings, yet an understanding of how these numbers were obtained is lacking. This paper presents an Excel workbook that demonstrates the mechanics of PCA, which include (1) the construction of the correlation matrix from the raw data, (2) the extraction of eigenvalues and eigenvectors from the correlation matrix and the computation of the component loadings and component scores, and (3) the rotation of the component loadings to improve interpretability.

## Keywords

Spreadsheets, principal component analysis, eigenvalues and eigenvectors, rotation

## 1. Introduction

Spreadsheets have enormous practical value for communicating statistical methods and concepts to students in the fields of business and economics. For example, Barr and Scott [2, 3] described a variety of spreadsheet-based simulation tools and recently reported a specific application in the context of portfolio construction [4]. Kwan [11, 12] has demonstrated approaches for the shrinkage of covariance or correlation matrices to a prespecified target matrix, which also has important applications to portfolio theory. There are many other types of multivariate statistical procedures that center on the analysis of covariance or correlation matrices, such as factor analysis, canonical correlation, and structural equation modeling. Spreadsheet modeling can be particularly useful for communicating a basic understanding of the mechanics of such procedures to business students.

Principal component analysis (PCA) is a well-known multivariate statistical technique that is used for data reduction and an excellent treatment of the topic is provided by Jolliffe [9]. Given a data matrix, $\mathbf{X} = [x_{ij}]$, consisting of measurements for $n$ observations on $p$ variables, the goal of PCA is to select a few linear combinations (i.e., components) of the variables that explain most of the variation in the full data matrix. PCA is closely related to several other multivariate techniques, such a singular-value-decomposition [5, 8], exploratory factor analysis [16], multidimensional scaling [15], correspondence analysis [7], and biplots [6].

Although widely used in many scientific disciplines, the particular focus herein is on marketing applications. PCA has a rich history in the field of marketing research and its importance as a data reduction tool also extends to business analytics. When taught in a marketing research or business analytics course, PCA is generally implemented using a statistical software package such as SPSS, SAS, or R. Unfortunately, although a business student might leave the course with a reasonable understanding of the output produced

by a statistical software package, they seldom have a sufficient grasp of how the results reported in that output were obtained. The goal of this paper is to rectify this problem using a spreadsheet demonstration that highlights the mechanics of PCA, thus providing an original, flexible and educational thrust that the major software packages do not afford.

The Excel workbook for PCA consists of three worksheets. The first worksheet, which is described in Section 2, is used to obtain the correlation matrix from the raw data. The second worksheet, explained in Section 3, is used to extract the eigenvectors from the correlation matrix. This is accomplished using the power method for finding the dominant eigenvector and corresponding eigenvalue of a matrix. Approaches for choosing the number of eigenvectors to retain are also described in this section. The selected eigenvectors and their corresponding eigenvalues are used to compute the principal component loadings (correlations between each variable and each component) and principal component scores for the respondents in the sample are also computed. The third worksheet, described in Section 4, is provided for the rotation of the component loadings to improve their interpretability. Practical experience with the Excel workbook and suggestions for adaptations and extensions are provided in Section 5.


## 2. Obtaining the Correlation Matrix from the Raw Data

The first worksheet 'HSM_Data' of the Excel workbook PCA is shown in Figure 1. The worksheet contains the raw data in the form of an $n \times p$ data matrix, $\mathbf{X} = [x_{ij}]$. In the example, the data are 7-point Likert-scale measurements for $n = 30$ respondents for each of $p = 6$ variables (or scale item statements) pertaining to hedonic shopping motivations. More specifically, the measurements correspond to each respondent's level of agreement with each of the six statements in Figure 1, where the measurements range

from 1= strongly disagree to 7 = strongly agree. The raw data occupy cells B2:C31 and are shaded in yellow in Figure 1. Although the data are synthetic, they are based on actual constructs and questionnaire items associated with an extensive study of hedonic shopping motivations conducted by Arnold and Reynolds [1]. The means ($\bar{x}_j$) and standard deviations ($s_j$) of each of the variables ($1 \leq j \leq p$) are computed (they are displayed in cells B33:G34 and shaded with a tan background in the worksheet). The means and standard deviations are used to transform the raw data to an $n \times p$ matrix of z-scores, $\mathbf{Z} = [z_{ij}]$, which is accomplished as follows:

$$z_{ij} = (x_{ij} - \bar{x}_j)/s_j, \; \forall \; 1 \leq i \leq n \text{ and } 1 \leq j \leq p \tag{1}$$

The z-scores, which occupy cells I2:N31 and are shaded in green in Figure 1, are then used to compute the correlation matrix, $\mathbf{R}$, as follows:

$$\mathbf{R} = (1/(n-1))\mathbf{Z}^T\mathbf{Z}. \tag{2}$$

The correlation matrix (contained in cells I36:N41 and shaded in blue in Figure 1) is computed with the aid of the MMULT function based on the z-score columns. The correlation matrix could have been obtained directly using the Data Analysis Toolpack capabilities of Excel. However, the use of the z-scores and MMULT function to compute $\mathbf{R}$ is concordant with the goals of providing a thorough and detailed presentation of all of the computation aspects of PCA. The correlation matrix $\mathbf{R}$ is copied (cell values only, not the formulas) to the top left corner of the second worksheet 'extract'.
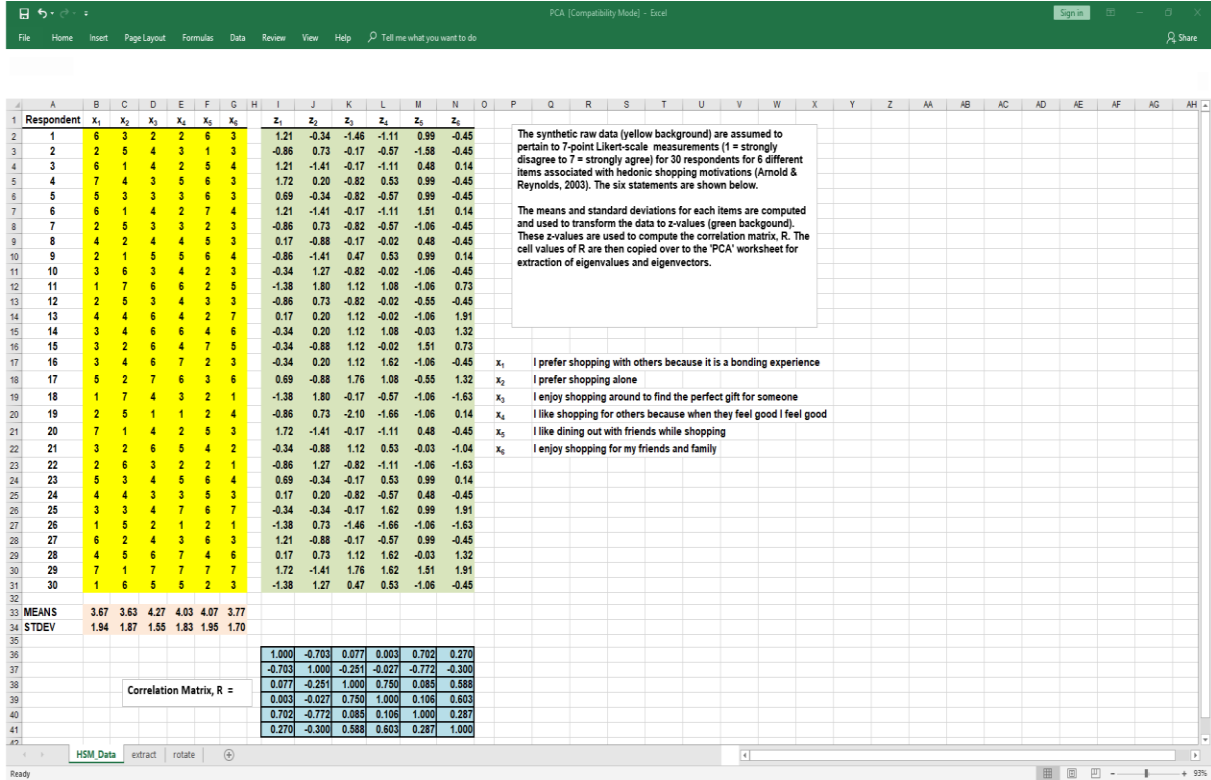
Figure 1: The hedonic shopping motivation data and corresponding correlation matrix.

## 3. Extraction of the Principal Components

In the 'extract' worksheet displayed in Figure 2, the correlation matrix $\mathbf{R}$ is in cells A3:F8. The $p \times p$ correlation matrix is a real symmetric positive semidefinite matrix and, therefore has real nonnegative eigenvalues that sum to the trace of $\mathbf{R}$. Because the main diagonal elements of the correlation matrix are all one, the trace of $\mathbf{R}$ is equal to the number of variables, $p$. Moreover, based on the principles of eigen-decomposition, also known as spectral decomposition, the correlation matrix can be written as a function of its eigenvalues ($\lambda_1,\ldots,\lambda_p$) and eigenvectors ($\mathbf{u}_1,\ldots,\mathbf{u}_p$) as follows:

$$\mathbf{R} = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \lambda_2 \mathbf{u}_2 \mathbf{u}_2^T + \cdots + \lambda_p \mathbf{u}_p \mathbf{u}_p^T = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T \tag{3}$$

where $\boldsymbol{\Lambda}$ is a $p \times p$ diagonal matrix of containing the eigenvalues of $\mathbf{R}$ and $\mathbf{U} = [\mathbf{u}_1,\ldots,\mathbf{u}_p]$ is the $p \times p$ matrix of corresponding eigenvectors. Without loss of generality, we will assume that the eigenvalues are sequenced in nonincreasing order of magnitude (i.e., $\lambda_1$

$\geq \lambda_2,\dots, \geq \lambda_{p-1} \geq \lambda_p$). The eigenvalues are measures of the explained variation in the correlation matrix and, accordingly, the eigenvalue-eigenvector pairs associated with the largest eigenvalues are those that make the greatest contribution to the decomposition of $\mathbf{R}$.

Rather than a complete decomposition of $\mathbf{R}$ as in Equation (3), in PCA we would like a low-rank ($q \ll p$) approximation that explains the greatest amount of variation in $\mathbf{R}$. The underlying optimization problem associated with the first principal component, $\mathbf{u}$, is as follows:

$$\text{Maximize: } \mathbf{u}^T\mathbf{R}\mathbf{u} \tag{4}$$

Subject to: $\qquad\qquad\qquad$ Subject to: $\mathbf{u}^T\mathbf{u} = 1$ $\qquad\qquad\qquad\qquad$ (5)

The Lagrangian function associated with the optimization problem is:

$$\text{Maximize: } L = \mathbf{u}^T\mathbf{R}\mathbf{u} - \lambda(\mathbf{u}^T\mathbf{u} - 1), \tag{6}$$

The first order condition is:

$$\partial L/\partial \mathbf{u} = 2\mathbf{R}\mathbf{u} - 2\lambda\mathbf{u} = \mathbf{0} \text{ or } \mathbf{R}\mathbf{u} = \lambda\mathbf{u} \tag{7}$$

which is an eigen-structure where $\mathbf{0}$ is a $p \times 1$ vector of zeros, the Lagrange multiplier $\lambda$ is the eigenvalue, and $\mathbf{u}$ is its corresponding eigenvector.

Although we know that there are typically $p$ eigenvalues that solve the determinant polynomial corresponding to Equation (7), we also know from the eigen-decomposition in Equation (3) that it is the largest eigenvalue and its corresponding eigenvector that will explain the most variation in $\mathbf{R}$. Therefore, we want to find the largest (or *dominant*) eigenvalue/eigenvector pair associated with Equation (7), that is, $\lambda = \lambda_1$ and $\mathbf{u} = \mathbf{u}_1$. Once this pair is identified, the correlation matrix can be deflated by removing the contribution from $\lambda_1$ and $\mathbf{u}_1$ using Equation (3). We use the notation $\mathbf{R}(q)$ to denote the deflated correlation matrix associated with the elimination of variation stemming from

the first $q$ principal components. After the first eigenvalue/eigenvector pair is extracted, the correlation matrix is deflated via:

$$\mathbf{R}(1) = \mathbf{R} - \lambda \mathbf{u}_1 \mathbf{u}_1^T \tag{8}$$

The second eigenvalue/eigenvector pair ($\lambda_2$, $\mathbf{u}_2$) associated with $\mathbf{R}$ is the dominant eigenvalue associated with $\mathbf{R}(1)$ and would then be extracted in similar fashion. Thus, the process of extracting principal components from $\mathbf{R}$ is sequential [13, p. 98]. After the extraction of ($\lambda_2$, $\mathbf{u}_2$), deflation occurs by setting:

$$\mathbf{R}(2) = \mathbf{R}(1) - \lambda \mathbf{u}_2 \mathbf{u}_2^T, \tag{9}$$

and the third eigenvalue/eigenvector pair ($\lambda_3$, $\mathbf{u}_3$) associated with $\mathbf{R}$ would be extracted as the dominant eigenvalue associated with $\mathbf{R}(2)$. We use this process in our Excel spreadsheet to sequentially extract all of the eigenvalue/eigenvector pairs for $\mathbf{R}$. Each stage of the sequential extraction process is accomplished using successive approximation via the *power method* [14].

The power method is one of the most conceptually straightforward approaches for finding the dominant eigenvector and corresponding eigenvalue of a correlation matrix (as well as other symmetric matrices). Using the power method, the eigenvector at iteration $k+1$ ($\mathbf{u}^{k+1}$) is estimated from the multiplication of the eigenvector at iteration $k$ ($\mathbf{u}^k$) by $\mathbf{R}$. That is, $\mathbf{u}^{k+1} = \mathbf{R}\mathbf{u}^k$. To illustrate why this process will converge to the dominant eigenvector, let us assume that the initial estimate ($\mathbf{u}^0$) is expressed as a linear combination (with coefficients $\alpha_j$, for $1 \leq j \leq p$) function of the $p$ eigenvectors of $\mathbf{R}$ as follows: $\mathbf{u}^0 = \alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \ldots + \alpha_p \mathbf{u}_p$. For the first iteration, we would have $\mathbf{u}^1 = \mathbf{R}\mathbf{u}^0 = \alpha_1 \mathbf{R}\mathbf{u}_1 + \alpha_2 \mathbf{R}\mathbf{u}_2 + \ldots + \alpha_p \mathbf{R}\mathbf{u}_p$ and, because of the eigen-structure relationship $\mathbf{R}\mathbf{u} = \lambda \mathbf{u}$, we can re-write this as $\mathbf{u}^1 = \mathbf{R}\mathbf{u}^0 = \alpha_1 \lambda_1 \mathbf{u}_1 + \alpha_2 \lambda_2 \mathbf{u}_2 + \ldots + \alpha_p \lambda_p \mathbf{u}_p$. Recalling that we have defined $\lambda_1$ as the largest eigenvalue, it is helpful to re-write this as: $\mathbf{u}^1 = \mathbf{R}\mathbf{u}^0 = \lambda_1[\alpha_1 \mathbf{u}_1 + \alpha_2(\lambda_2/\lambda_1)\mathbf{u}_2 + \ldots + \alpha_p(\lambda_p/\lambda_1)\mathbf{u}_p]$. After $m+1$ iterations of the power method, we have: $\mathbf{u}^{m+1} = \mathbf{R}^m \mathbf{u}^0 = (\lambda_1)^m[\alpha_1 \mathbf{u}_1 + \alpha_2(\lambda_2/\lambda_1)^m \mathbf{u}_2 + \ldots + \alpha_p(\lambda_p/\lambda_1)^m \mathbf{u}_p.]$. In the limit, as $m \to \infty$, the terms $(\lambda_2/\lambda_1)^m, \ldots, (\lambda_p/\lambda_1)^m$ go to

zero as long as $\lambda_1$ is strictly larger than all of the other eigenvalues. Assuming $\alpha_1 \neq 0$, the result is convergence to an eigenvector that is a multiple of the dominant eigenvector: that is, $\mathbf{u}^{m+1} = \mathbf{R}^m \mathbf{u}^0 = \alpha_1(\lambda_1)^m \mathbf{u}_1$. Faster convergence is generally achievable if the updated eigenvector is normalized to unit length using the following equation:

$$\mathbf{u}^{k+1} = \frac{\mathbf{R}\mathbf{u}^k}{\sqrt{(\mathbf{R}\mathbf{u}^k)^T (\mathbf{R}\mathbf{u}^k)}}. \tag{10}$$

Thus, assuming $\lambda_1$ is strictly greater than all other eigenvalues, the power method estimates the eigenvector via an iterative process of multiplying by $\mathbf{R}$ and normalizing the result. Upon convergence of the estimation process to the eigenvector $\mathbf{u}$, the corresponding eigenvalue for an eigenvector can be obtained via the Rayleigh quotient as follows:

$$\lambda = \frac{\mathbf{u}^T \mathbf{R} \mathbf{u}}{\mathbf{u}^T \mathbf{u}}, \tag{11}$$

The Rayleigh quotient makes use of the facts that $\mathbf{u}$ is an eigenvector of $\mathbf{R}$ and $\mathbf{R}\mathbf{u} = \lambda\mathbf{u}$. Thus, the numerator of Equation (11) could be re-written as $\lambda\mathbf{u}^T\mathbf{u}$, which implies $\lambda = \lambda$. If the eigenvector is unit length, then the denominator is one and the eigenvalue is equal to the numerator, which is the objective function in Equation (4) and indicates that $\lambda = \mathbf{u}^T\mathbf{R}\mathbf{u}$ is the variance extracted.

For successive extraction of eigenvalues from correlation matrices, the convergence of the power method to the dominant eigenvectors is generally rapid; however, convergence can take more iterations as the ratio of the second largest eigenvalue to the largest eigenvalue approaches one. For example, when computing the dominant eigenvalue for $\mathbf{R}$, the ratio of $\lambda_2/\lambda_1$ may have an effect on the rapidity of convergence. When computing the dominant eigenvalue for $\mathbf{R}(1)$, it is the ratio of $\lambda_3/\lambda_2$ that affects convergence.

**PCA [Compatibility Mode] - Excel**

Numerical Estimation of the First Principal Component via the Power Method for Eigenvalues

R =

| | | | | | | u | Ru | norm(Ru) | |
|---|---|---|---|---|---|---|---|---|---|
| 1.000 | -0.703 | 0.077 | 0.003 | 0.702 | 0.270 | 1.00000 | 1.34779 | 0.28809 | 1.47251 |
| -0.703 | 1.000 | -0.251 | -0.027 | -0.772 | -0.300 | 1.00000 | -1.05279 | -0.22503 | |
| 0.077 | -0.251 | 1.000 | 0.750 | 0.085 | 0.588 | 1.00000 | 2.24882 | 0.48068 | |
| 0.003 | -0.027 | 0.750 | 1.000 | 0.106 | 0.603 | 1.00000 | 2.43578 | 0.52064 | First Eigenvector |
| 0.702 | -0.772 | 0.085 | 0.106 | 1.000 | 0.287 | 1.00000 | 1.40821 | 0.30100 | |
| 0.270 | -0.300 | 0.588 | 0.603 | 0.287 | 1.000 | 1.00000 | 2.44728 | 0.52310 | 34.496 |

Update Eigenvector 1

The First Eigenvalue = 4.67839

% of variance explained by first two components

Numerical Estimation of the Second Principal Component after Deflation of Correlation Matrix

Deflated R =

| | | | | | | u | Ru | norm(Ru) | |
|---|---|---|---|---|---|---|---|---|---|
| 0.878 | -0.608 | -0.127 | -0.218 | 0.574 | 0.048 | 1.00000 | 0.54667 | 0.28809 | 0.59726 |
| -0.608 | 0.925 | -0.092 | 0.146 | -0.672 | -0.127 | 1.00000 | -0.42702 | -0.22503 | |
| -0.127 | -0.092 | 0.660 | 0.382 | -0.128 | 0.218 | 1.00000 | 0.91213 | 0.48068 | |
| -0.218 | 0.146 | 0.382 | 0.601 | -0.125 | 0.202 | 1.00000 | 0.98796 | 0.52064 | Second Eigenvector |
| 0.574 | -0.672 | -0.128 | -0.125 | 0.867 | 0.055 | 1.00000 | 0.57117 | 0.30100 | |
| 0.048 | -0.127 | 0.218 | 0.202 | 0.055 | 0.597 | 1.00000 | 0.99262 | 0.52310 | |

Update Eigenvector 2

The Second Eigenvalue = 1.89757

| | | | | | |
|---|---|---|---|---|---|
| 1.21347 | 0.00000 | 0.42112 | -0.38061 | 0.5110 | -0.2941 |
| 0.00000 | 0.77282 | -0.46564 | 0.33134 | -0.5650 | 0.2561 |
| | | 0.35896 | 0.47791 | 0.4356 | 0.3693 |
| | | 0.31325 | 0.53697 | 0.3801 | 0.4150 |
| | | 0.45010 | -0.35487 | 0.5462 | -0.2743 |
| | | 0.41984 | 0.32044 | 0.5095 | 0.2476 |

This is the "component matrix" and contains elements called "component loadings". These values represent the simple correlations between the variables and the components extracted.
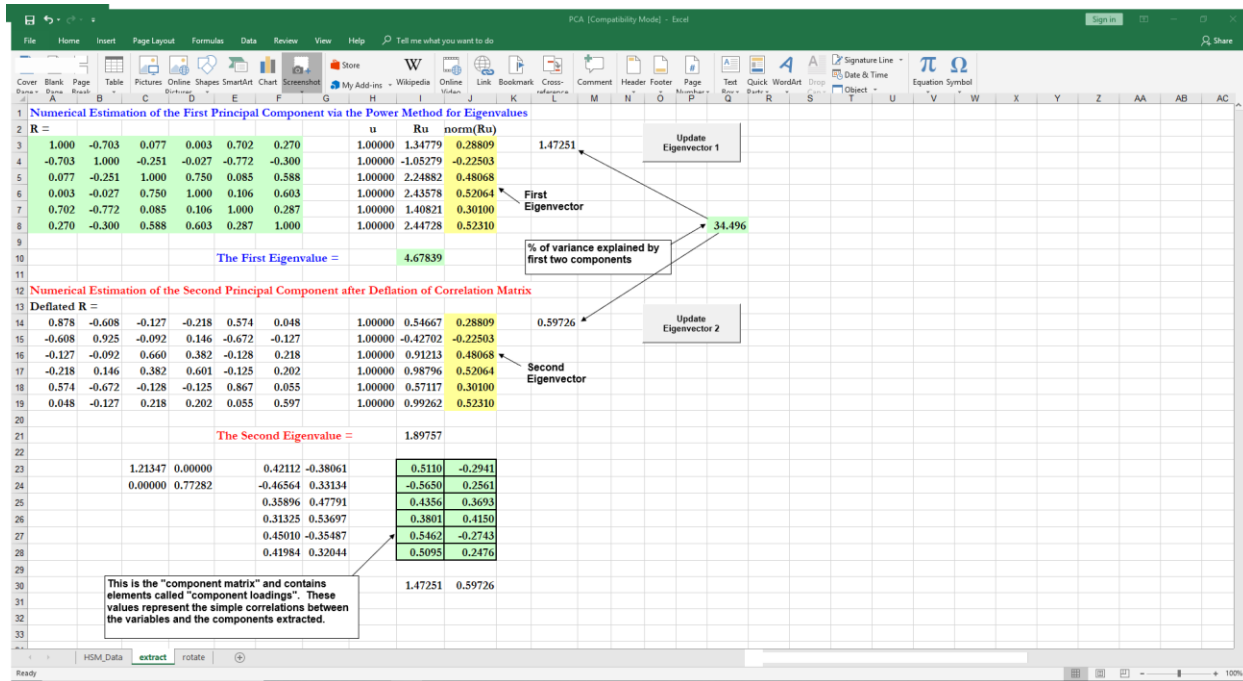
1.47251  0.59726

HSM_Data  extract  rotate

Figure 2: The extraction worksheet with eigenvectors initialized to vectors of ones.

An initial estimate of the first eigenvector (**u**) is placed in cells H3:H8. I usually just start off with all ones in these cells. Cells I3:I8 contain the vector formed by the **Ru** product. The value in cell I10 is the constant $\mathbf{u}^T\mathbf{Ru}$. Again, if **u** is normalized, then this quantity is equal to the eigenvalue estimate. Cells J3:J8 contain **Ru** after it has been normalized to unit length. This is the updated eigenvector. Clicking on the button "Update Eigenvector 1" simply reads the cell values in J3:J8 and re-pastes them in cells H3:H8 to obtain the updated eigenvector for the next iteration. Tapping on this button once (assuming the initial eigenvector of all ones in Figure 2) yields the result in Figure 3. As we continue to tap this button, we can see the eigenvalue get larger. After 15-20 taps, the changes tend to be small. We can continue tapping until the values in H3:H8 and J3:J8 have stabilized (in other words the values in these two cell ranges do not change). The results after convergence are shown in Figure 4. The first eigenvalue is $\lambda = 2.86408$ and its corresponding eigenvector is shown in cells J3:J8.
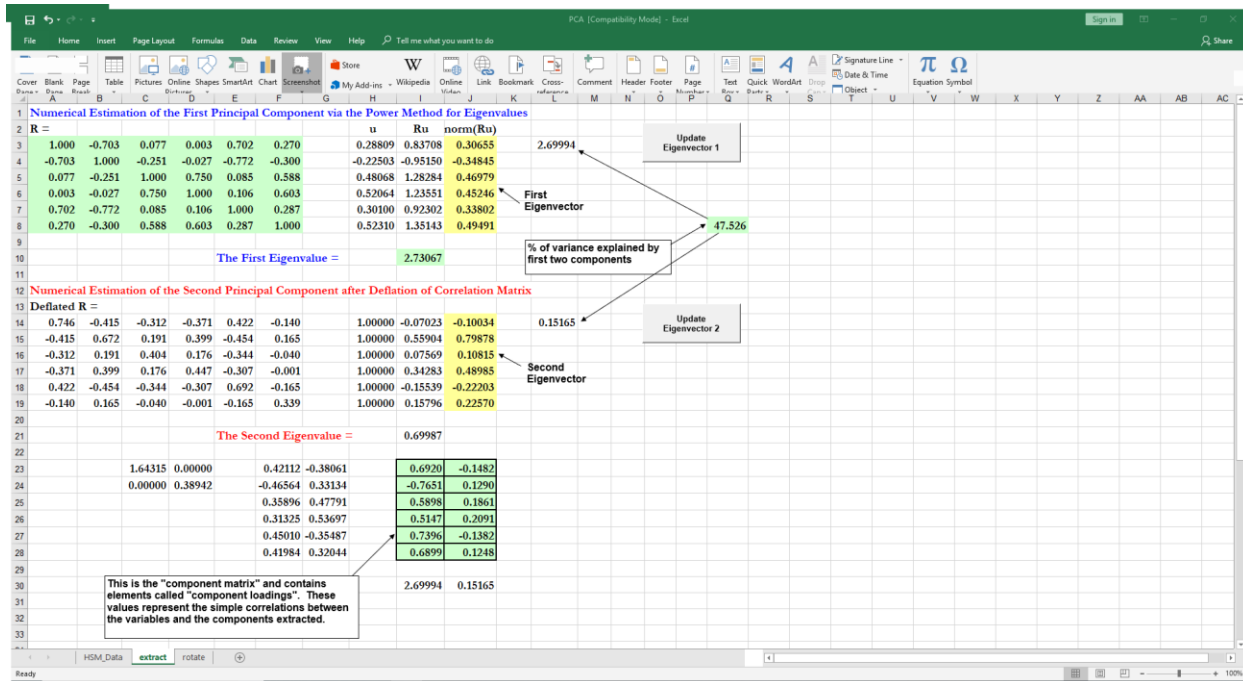
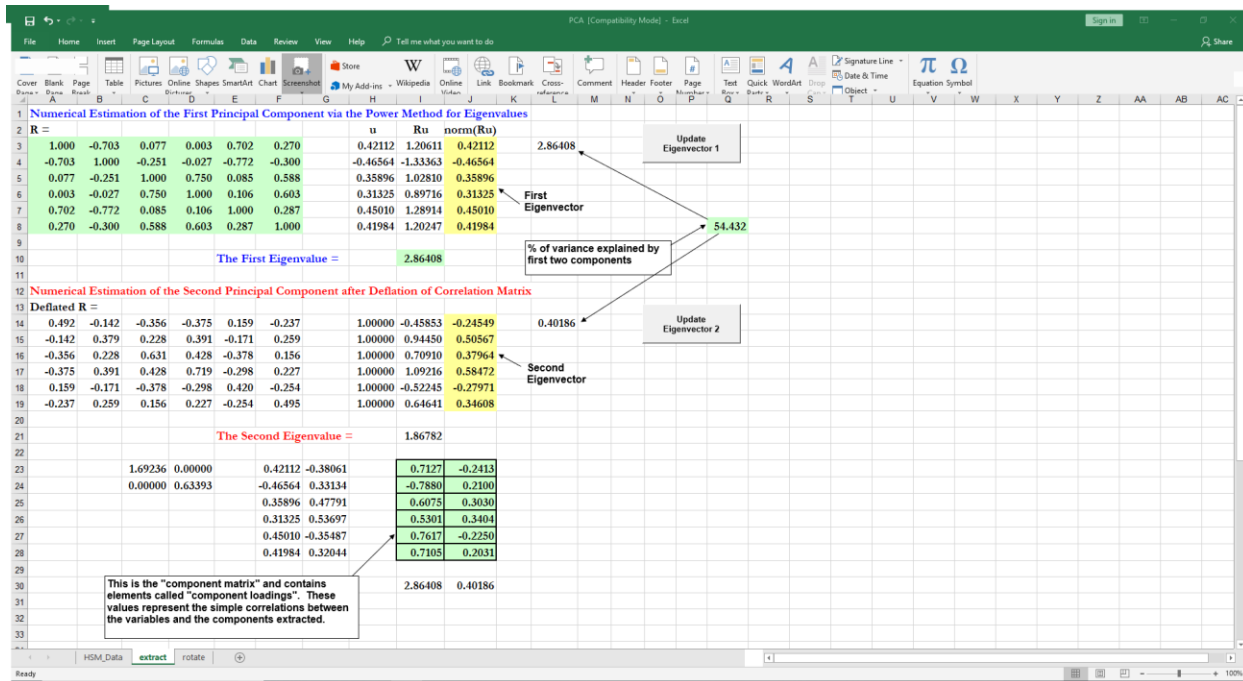Figure 3: The extraction worksheet after tapping the "Update Eigenvector 1" button once.



Figure 4: The extraction worksheet after convergence of the first eigenvector/eigenvalue pair.

Now that the first eigenvalue and eigenvector have been extracted, the correlation matrix is deflated using Equation (8). The deflated correlation matrix is in cells A14:F19. Once again, the initial estimates for the eigenvector are all ones in cells H14:H19.

Tapping the "Update Eigenvector 2" button will lead to convergence of the estimation of the second eigenvalue and eigenvector in a manner similar to the first. The result is shown in Figure 5. The second eigenvalue is $\lambda = 1.93877$ and its corresponding eigenvector is shown in cells J14:J19.



Figure 5: The extraction worksheet after extraction of the first two components.

Continuing in this manner, the remaining four eigenvalue/eigenvector pairs are extracted in rows 34 to 76 of the worksheet. Cells A80:F85 show the fully delated correlation matrix after extraction of all six components. The six eigenvalues are copied into cells U2:U7 so as to facilitate the production of a scree plot, which is shown in Figure 6. The sharp elbow in the plot suggests that two components should be extracted because of the huge drop in the size of the third eigenvalue in comparison to the second. Two components would also be chosen based on the popular default rule (e.g. in SPSS) for correlation matrices of selecting all components with eigenvalues greater than one. The first two components explain just over 80% of the variation in the data set.

Figure 6: A scree plot of the eigenvalues.

The $2 \times 2$ diagonal matrix (**D**) in cells C23:D24 contains the square roots of the eigenvalues on the main diagonal. The two columns of the $6 \times 2$ matrix (**U₂**) in cells F23:F28 are the two eigenvectors. The matrix product $G = U_2D$ yields the $6 \times 2$ matrix of *component loadings* in cells I23:I28. These (unrotated) loadings, which are interpreted as correlations between the six items and the two components, are visually displayed in the correlation circle plot in Figure 7. Five of the six variables (all but $x_2$) have fairly high positive loadings on component 1. Similarly, four of the six variables (all but $x_1$ and $x_5$) have fairly high positive loadings on component 2. Moreover, the component loadings all fall between 0.4 and 0.8 in absolute value. Ideally, we would like to have a *simple structure* whereby the values in each column are close (in absolute value) to either zero or one. This makes it easy to ascertain which variables correspond most heavily to each component. Therefore, in section 4, we will copy these component loadings into the next worksheet ('rotate') and rotate them to see if we can improve interpretability.
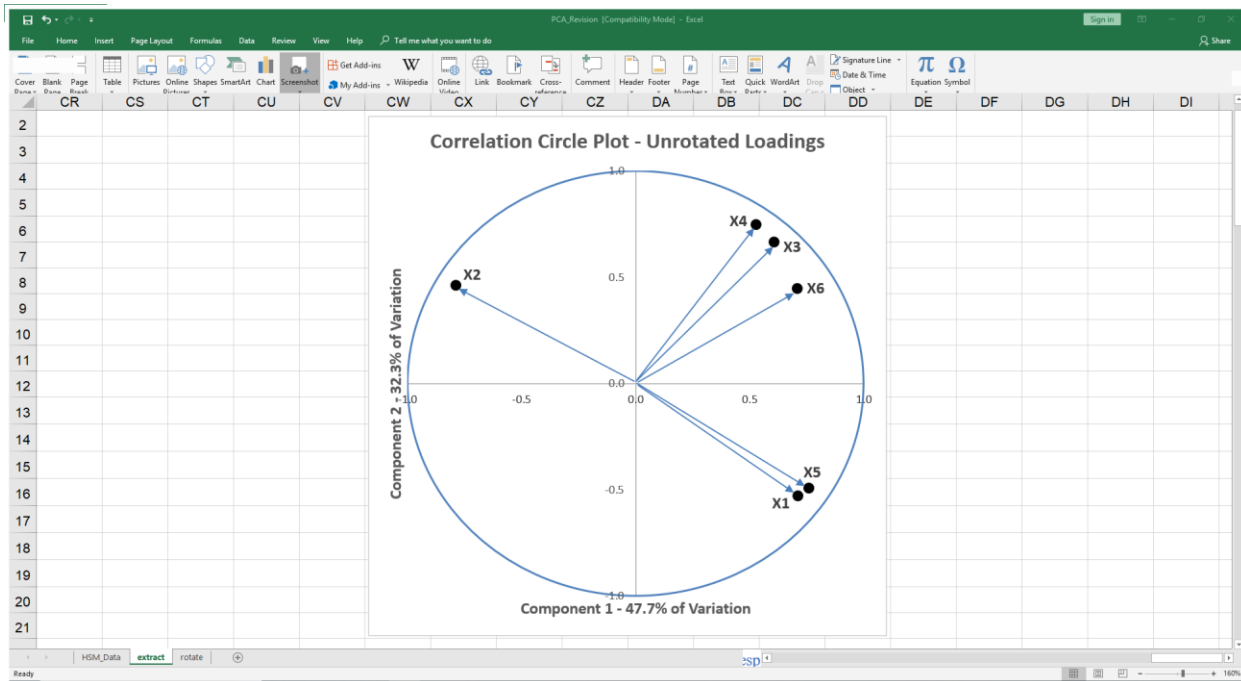
Figure 7: Correlation circle plot for the unrotated component loadings.

Whereas the correlation circle plot in Figure 7 affords a visual display of the relationship between the variables and the two components, Figure 8 is a *component score* plot that provides a visual display of the $n = 30$ respondents in the component space. The component scores are computed by multiplying the $30 \times 6$ matrix of $z$-scores (**Z**) in cells BA4:BF33 by the $6 \times 2$ matrix containing the first two eigenvectors ($\mathbf{U}_2 = [\mathbf{u}_1\ \mathbf{u}_2]$) in cells BH4:BI9. This matrix multiplication is accomplished using the MMULT function and the resulting $30 \times 2$ matrix of raw component scores are contained in cells BK4:BL33. Labels are provided in the plot of these component scores in Figure 8 to identify the respondent associated with each data point. For example, the label 'R29' indicates the position of the component score for respondent #29. This respondent provided extreme answers for all six variables (i.e., a response of 1 for $x_2$ and a response of 7 for the remaining variables and is, therefore, well separated from other respondents in the plot. Ideally, we would like to be able to identify groupings of the respondents (perhaps just by examining the respondents that fall in each quadrant) and characterize

them based on the two dimensions of the plot. Unfortunately, this is difficult because the components are difficult to define based on the loadings.
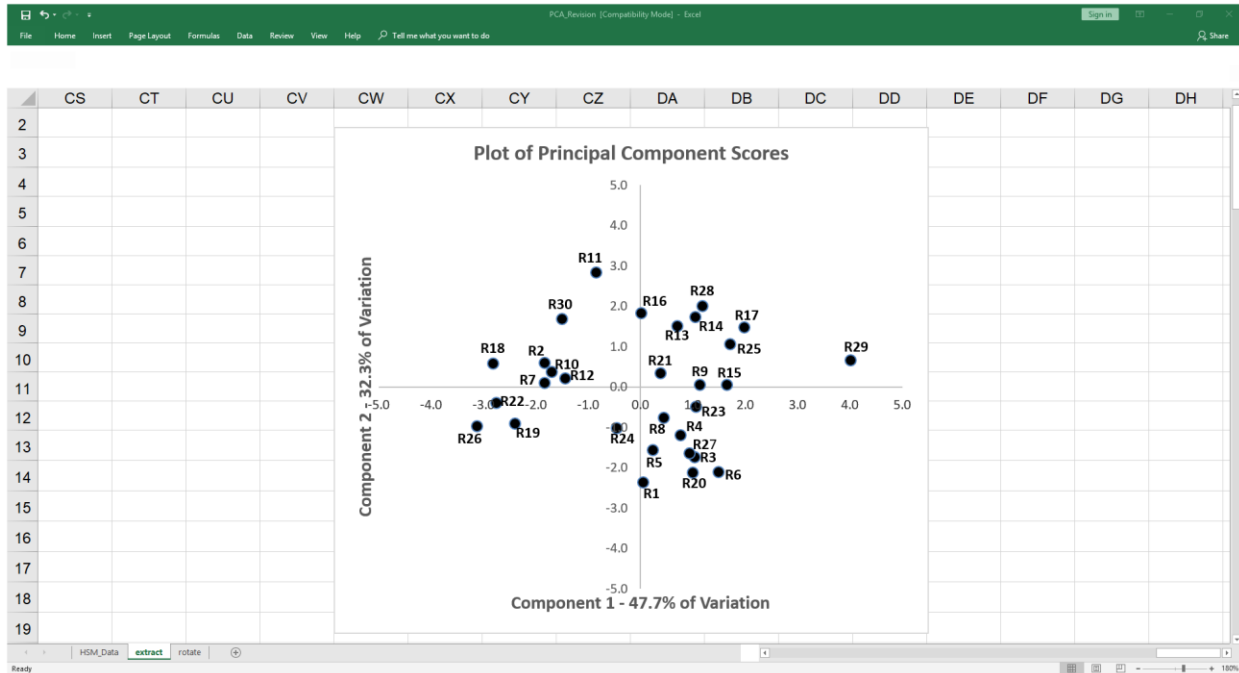


Figure 8. A plot of the raw principal component scores.

We know from the eigenvalues that the first two principal components explain (2.86408 + 1.93877)/6 = 80.048% of the variation in the data. However, we can also see this more concretely from the component scores. To begin the process, we compute the total variation in the data by computing the sum of squares for the $z$-scores in cells BA4:BF33 (note that the mean of the $z$-scores is zero and, therefore, these are squared deviations from the mean). The sum of these squared values, which represents the total variation in the data, is computed to be 174 in cell CD3 using the SUM function. The squared component scores are computed in cells BN4:BO33. The sum of squared values for components 1 and 2 are 83.058 and 56.224, respectively as shown in cells CD6:CD7. The sum for the two components is 139.283 as shown in cell CD9. The percentage of variation in the full data set that is explained by the two components is 139.283/174 = 80.048%. Moreover, the first and second components explain 83.058/174 = 47.7% and 56.224/174 = 32.3% of the variation, respectively.

## 4. Rotation of the Component Loadings

In the 'rotate' worksheet displayed in Figure 9, the matrix of unrotated component loadings (**G**) has been pasted into cells E6:F11. We seek to find a rotation of this coordinate system to a new set of coordinates, **H** = **GW**, that is more interpretable. The counterclockwise angle of rotation ($\theta$, in degrees) is placed in cell A5. The 2 × 2 rotation matrix, **W**, is in cells A7:B8. To understand the rationale for the rotation matrix, **W**, it is helpful to think of a coordinate pair ($g_1$, $g_2$) for the unrotated loadings in terms of polar coordinates ($d$, $\varphi$), where $d = \sqrt{g_1^2 + g_2^2}$ and $\varphi$ = arccos($d/g_1$) if $g_2 \geq 0$ and $\varphi$ = -arccos($d/g_1$) if $g_2 < 0$. The unrotated component loading coordinates can then be expressed in terms of the polar coordinates as:

$$g_1 = d\cos\varphi \tag{12}$$

$$g_2 = d\sin\varphi \tag{13}$$

Likewise, recognizing that the counterclockwise angle of rotation will be $\theta$, the coordinate pair ($h_1$, $h_2$) for the rotated loadings in terms of polar coordinates will be:

$$h_1 = d\cos(\varphi\text{-}\theta) \tag{14}$$

$$h_2 = d\sin(\varphi\text{-}\theta) \tag{15}$$

Based on the trigonometric rules for compound angles, (14) and (15) can be rewritten as:

$$h_1 = d\cos\varphi\cos\theta + d\sin\varphi\sin\theta \tag{16}$$

$$h_2 = d\sin\varphi\cos\theta - d\cos\varphi\sin\theta \tag{17}$$

Equations (12) and (13) can be used to simplify Equations (16) and (17) as follows:

$$h_1 = g_1\cos\theta + g_2\sin\theta \tag{18}$$

$$h_2 = -g_1\sin\theta + g_2\cos\theta \tag{19}$$

More generally, in matrix notation:

$$[h_1\ h_2] = [g_1\ g_2] \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix} \tag{20}$$

The 2 × 2 matrix on the right side of Equation (20) is the rotation matrix in cells A7:B8. The matrix of rotated component loadings, $\mathbf{H} = \mathbf{GW}$ is computed using this matrix and the values are displayed in cells H6:I11. Because the angle of rotation in cell A5 is set to zero in Figure 9, $\mathbf{H} = \mathbf{G}$. Plots of the unrotated and rotated loadings are displayed at the bottom of Figure 9.

The student can manually change the angle of rotation in cell A5 and the rotated loadings in $\mathbf{H}$ will be updated automatically. Alternatively, an optimal selection of the angle of rotation can be made based on the well-known *varimax* criterion developed by Kaiser [10].

$$V = \frac{1}{p} \sum_{l=1}^{q} \left[ \sum_{j=1}^{p} h_{jl}^4 - \frac{1}{p} \left( \sum_{j=1}^{p} h_{jl}^2 \right)^2 \right], \tag{21}$$

where $q$ is the number of selected components (in this example, $q = 2$). Kaiser [10] developed the *raw varimax criterion* in Equation (21) with the goal of producing an orthogonal rotation of the component axes so as to induce a *simple structure* in the rotated loadings. For this criterion, simplicity is operationalized *as the variance of the squared loadings*. Maximizing this quantity has a propensity to drive the loadings toward zero or one in absolute value, which makes it easier to ascertain the variables that correlate strongly with each component.

Kaiser also proposed a normalized varimax criterion that is particularly relevant in the context of factor analysis. The normalized version adjusts the formula in Equation (21) to account for the *communality* of each variable (i.e., the variance in the variable that is accounted for by the common factor). Because PCA makes no distinction between common and specific factors, communality and the normalized varimax criterion are not relevant in the context of PCA. Although other options for rotation of the loadings are available, varimax is one of the most common and easiest to implement. The values in cells K6:O13 are scratch computations used to facilitate the computation of the value of the varimax criterion, which is in cell J2.
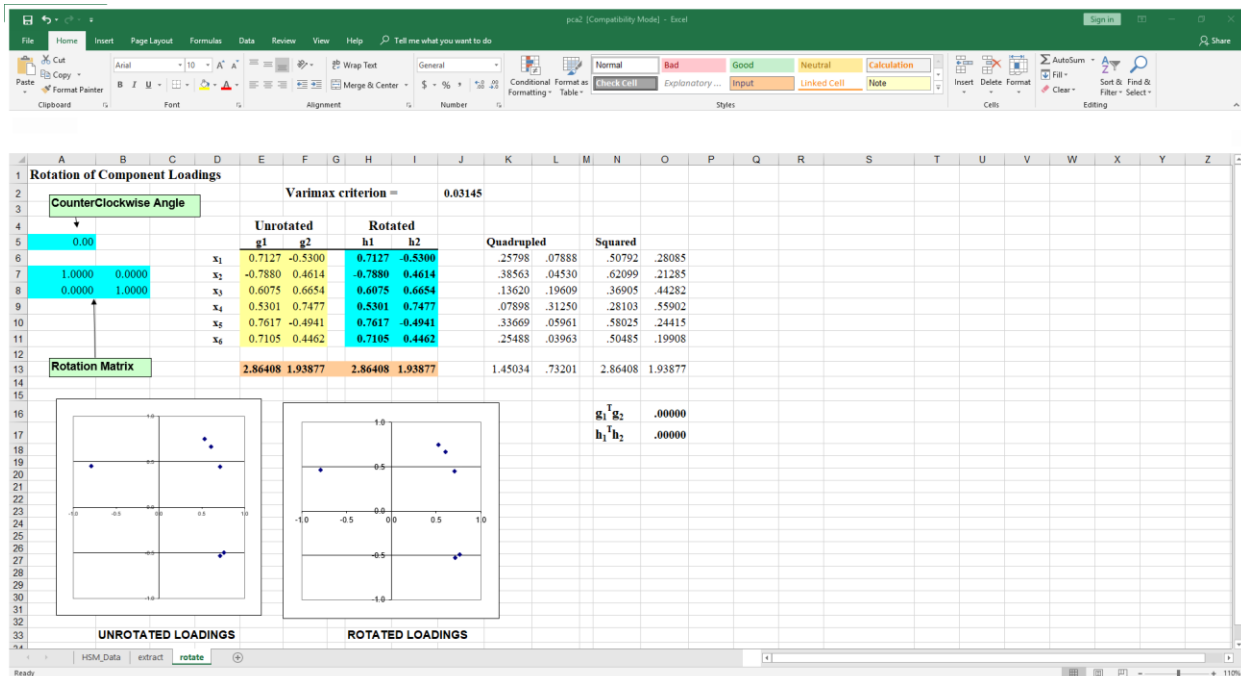
Figure 9: The 'rotate' worksheet prior to varimax rotation (rotation angle $\theta = 0$). The Excel Solver (using $\theta = 0°$ as the initial value in cell A5) was used to find the angle of rotation that maximizes the value of the varimax criterion in cell J2. The results are displayed in Figure 10. The optimal counterclockwise angle of rotation is $\theta = -37.9°$ (i.e., a clockwise rotation of $\theta = 37.9°$). The rotated loadings have changed dramatically and the correlation circle plot of these loadings in Figure 11 now shows that the variables tend to be associated with one of the two components but not both. Cells O16:O17 show that unrotated loading vectors are orthogonal ($\mathbf{g}_1^T\mathbf{g}_2 = 0$) but the rotated loadings are not ($\mathbf{h}_1^T\mathbf{h}_2 \neq 0$). Moreover, it is noteworthy that, although the total variation explained by the two components is unchanged, the relative contribution of the explained variation for the two components is more equally distributed. That is, for the unrotated loadings it is noted that 2.86048 + 1.93877 = 4.80285 and for the rotated loadings 2.51495 + 2.28790 = 4.80285.
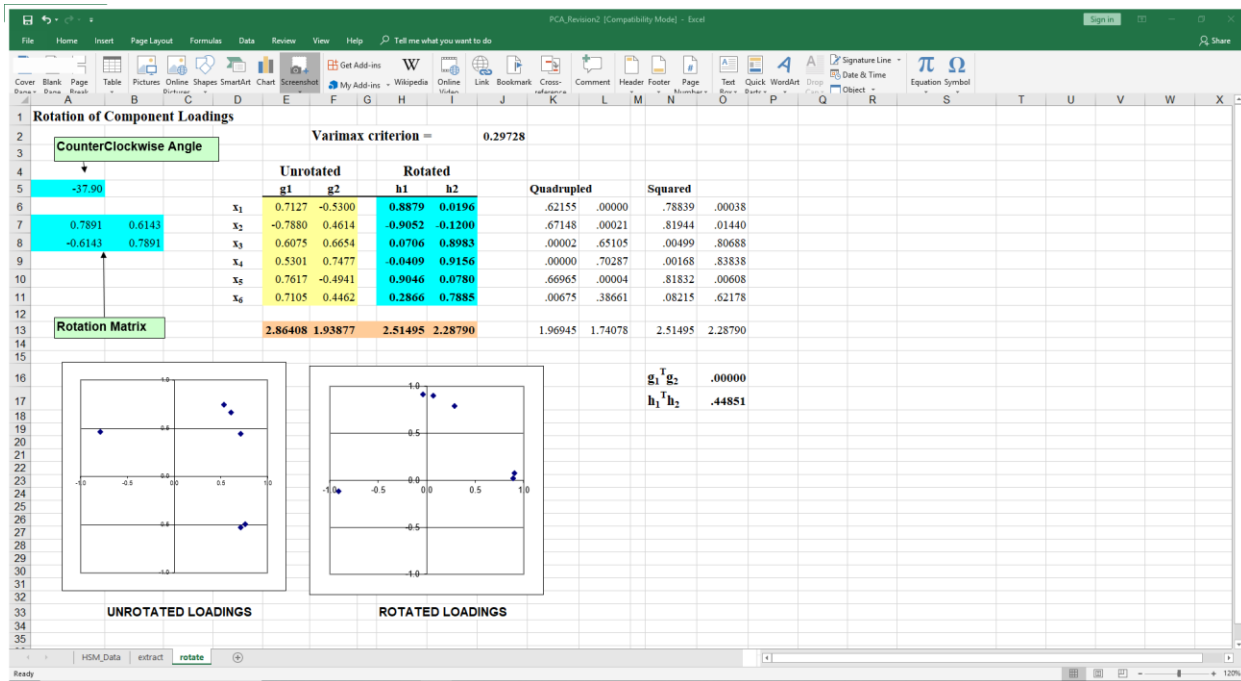
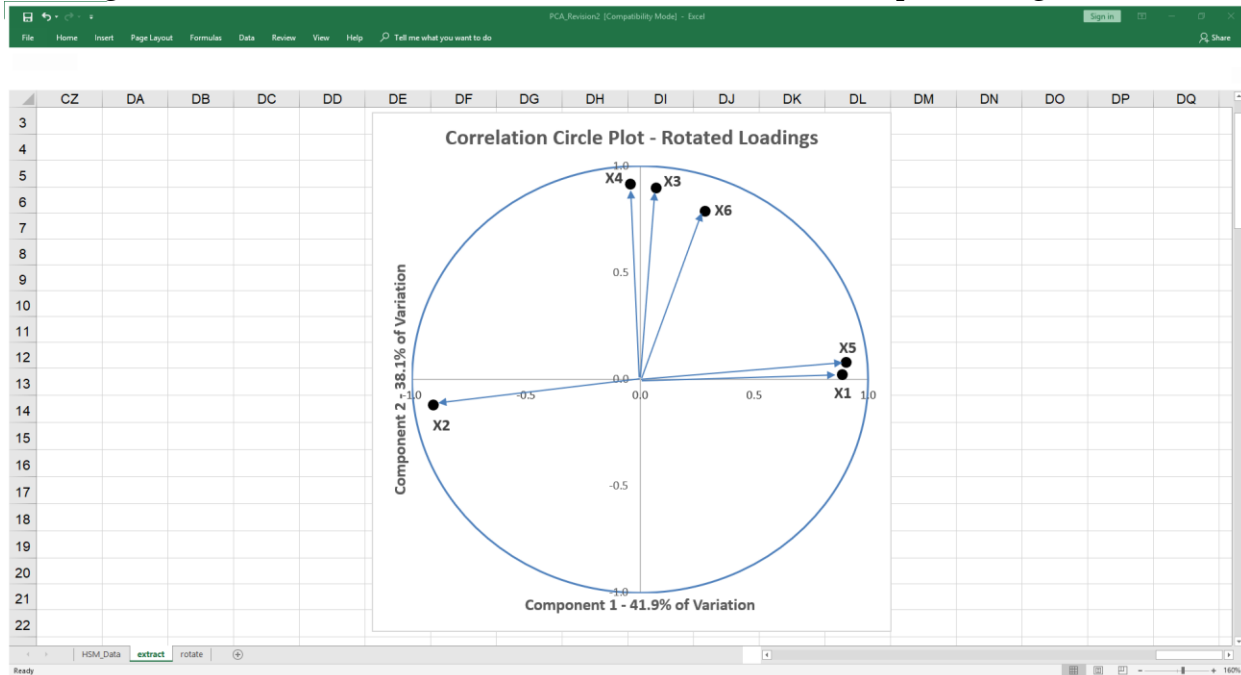Figure 10: The 'rotate' worksheet after varimax rotation (optimal angle -37.9°).



Figure 11. A correlation circle plot for the rotated loadings.

When examining the rotated loadings, it is evident that the variables $x_1$, $x_2$, and $x_5$ have large absolute values on component 1 and small absolute values on component 2. The reverse is true for the other three variables $x_3$, $x_4$, and $x_6$. Accordingly, the interpretation

of the rotated components is much cleaner. As observed from Figure 1, Likert-scale items $x_1$, $x_2$, and $x_5$ all pertain to shopping *with* other people; whereas items $x_3$, $x_4$, and $x_6$ are related to shopping *for* other people. In the language used by Arnold and Reynolds [1], component 1 would be identified as a *social shopping* construct and component 2 as a *role shopping* construct.

The cell values of the rotation matrix from the 'rotate' worksheet were copied to cells BQ4:BR5 of the 'extract' worksheet to facilitate a rotation of the principal component scores. The rotated component scores in cells BT4:BU33 were obtained using the MMULT function. Cells CD12:CE15 verify that the rotated components still explain 80.048% of the variation in the data, but that the relative contribution of the two components has changed. A plot of the rotated principal scores is displayed in Figure 12. Once again, labels for each respondent are provided in the plot to facilitate interpretation. Respondents with large positive scores on component 1 are those who might generally perceived as social shoppers, whereas those respondents with large negative scores on component 1 have no proclivity for social shopping. Respondents with large positive scores on component 2 are role shoppers whereas those with large negative scores on component 2 are more averse to role shopping.

Respondent 29 (R29 in the plot in Figure 12) has the largest positive score on component 1 and largest positive score on component 2 and, therefore, this person is both a social and role shopper. This customer strongly agreed (Likert scale response of 7) with all of the role shopping statements ($x_3$, $x_4$, $x_6$), as well as the two social shopping statements ($x_1$, $x_5$) that pertained to shopping with others. Concordantly, the respondent strongly disagreed (with a Likert scale response of 1) with the social shopping statement pertaining to shopping alone ($x_2$). At the other corner of the plot, respondent 26 (R26) is someone who is neither a role shopper nor a social shopper. This individual strongly disagreed (Likert scale responses of 1 or 2) with all of the role shopping statements and

the two social shopping statements associated with shopping with others. Respondent 26 generally agreed (Likert scale response of 5) with the statement on shopping alone. Respondent 1 (R1) is a fairly solid social shopper but not a role shopper. This individual strongly agreed (Likert scale response of 6) with the statements pertaining to shopping with others ($x_1$, $x_5$), but generally disagreed (Likert scale response of 3) with the statement pertaining to shopping alone ($x_2$). By contrast, R1 generally disagreed (Likert responses of 2 or 3) with all three of the role shopping statements ($x_3$, $x_4$, $x_6$). Respondent 11 (R11) is a fairly solid role shopper but not a social shopper. This individual agreed (Likert scale responses of 5 or 6) with the statements pertaining to shopping for others ($x_3$, $x_4$, $x_6$). However, R11 disagreed strongly (Likert responses of 1 or 2) with the statements pertaining to shopping with others ($x_1$, $x_5$), but strongly agreed (Likert scale response of 7) with the statement pertaining to shopping alone ($x_2$).

Extreme respondents such as R1, R11, R26, and R29 help to characterize the boundaries of the component score plot. It is also interesting to observe those respondents with very similar component scores. For example, {R2, R7, R10, R12} form a tight cluster of similar individuals who are generally not inclined to be either social or role shoppers. For all of these respondents, their Likert scale responses to statements $x_1$, $x_3$, $x_4$, $x_5$, and $x_6$ never exceed four. Moreover, with the exception of statement $x_5$, where R2 responded with a Likert scale measure of 1 and R12 with 3, the difference between the Likert scale responses for all pairs of these four respondents never differed by more than one-unit for any of the other statements.
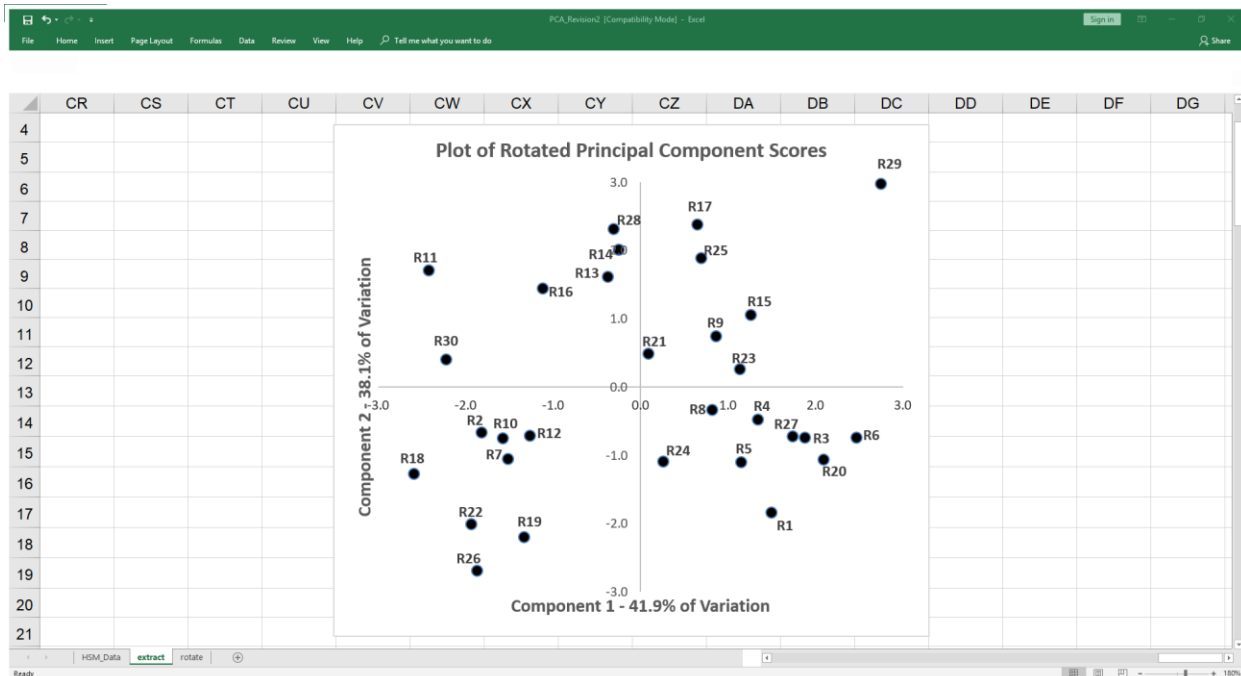
Figure 12: A plot of the rotated principal component scores

## 5. Practical Experience and Extensions

I have used the Excel workbook for PCA for three different courses: (1) an undergraduate marketing research course, (2) a masters level course in marketing analytics, and (3) a Ph.D. seminar in quantitative methods. Naturally, the coverage of the content in the workbook is apt to vary depending on the level of the course. The rigor of the presentation at the masters level typically falls somewhere between the undergraduate and Ph.D. level and is contingent on other factors such as class size. At the undergraduate level, the focus might be restricted to the generation of the correlation matrix, the scree plot of the eigenvalues for choosing the number of components, the correlation circle plots of the loadings, and the plot of the principal component scores. The coverage of the eigen-decomposition would likely be avoided in an undergraduate business course, but not necessarily a course in a more scientific discipline such as engineering. The details of the rotation of the loadings might also be skipped; however, over the years, I have found that many undergraduates are

interested in having at least some rudimentary understanding of how the rotation process works.

At the PhD level, the entire content of the Excel workbook for PCA is presented. The PhD students are expected to have some knowledge of the fundamentals of linear algebra. I have typically required students to do some small matrix analysis examples by hand (i.e., matrix and vector products, inverse, determinants, eigenvalues and eigenvectors) to develop an understanding. The extraction of the eigenvalues and eigenvectors for a six-by-six correlation matrix, however, is virtually impossible by hand. However, the template shows how the power method can do this efficiently and, therefore, the Ph.D. students have some idea as to how eigen-decomposition is done by computer. The feedback that I received has been quite positive, as the students are generally appreciative of gaining a more rigorous understanding of the methodology. Perhaps the most natural extension of the Excel workbook for PCA is to exploratory factor analysis. This would require the incorporation of an approach for computing the communality associated with each variable, which can be accomplished in different ways. Communality may also affect the nature of the varimax rotation. As noted previously, varimax rotation with Kaiser normalization may be preferable for factor analysis and is the default setting in some software packages (e.g., SPSS). Another important extension would be to adapt the workbook to generate biplots. Whereas PCA is based on the eigen-decomposition of the correlation matrix, biplots are grounded in the singular-value-decomposition [5] of the raw $n \times p$ data matrix. Singular-value-decomposition enables the extraction of eigenvector pairs for respondents and variables, thus enabling them to be plotted simultaneously in a two-dimensional space. There have been some important advancements in the design of biplots [6] that should also be taken into consideration when extending the PCA workbook for this interesting topic.

In addition to exploratory factor analysis and biplots, there are several other possible extensions of the Excel workbook for PCA. Many of these possibilities hinge on the fact that the key engine of the workbook is the use of the power method for eigenvalue and eigenvector estimation. Eigen-structure problems also arise in a variety of related multivariate statistical methods, such as multiple discriminant analysis, canonical correlation, and multidimensional scaling. The development of spreadsheet solutions for these and other applications is an interesting avenue for future research.

## References

[1] Arnold, M. J., & Reynolds, K. E. (2003). Hedonic shopping motivations. *Journal of Retailing*, *79*, 77-95.

[2] Barr, G., & Scott, L. (2011). Teaching statistics in a spreadsheet environment using simulation. *eJournal of Spreadsheets in Education*, 4(3), Article 2. Available at: http://epublications.bond.edu.au/ejsie/vol4/iss3/2. [Accessed 17 August 2018].

[3] Barr, G., & Scott, L. (2015). Spreadsheets and simulation for teaching a range of statistical concepts. In: H. MacGillivray, B. Phillips, M. Martin, ed., *Topics from Australian conferences on teaching statistics: OZCOTS 2008-2012*, 1st edition, New York: Springer, pp 99-117.

[4] Barr, G. D., & Scott, L. (2018) An active learning exercise showing some fundamentals of financial portfolio construction, *eJournal of Spreadsheets in Education* Vol. 10(3), Article 5. Available at: https://epublications.bond.edu.au/ejsie/vol10/iss3/5 [Accessed 17 August 2018]

[5] Eckart, C., & Young, G. (1936). The approximation of one matrix by another of lower rank. *Psychometrika*, *1*, 211-218.

[6] Gower, J. C., & Hand, D. J. (1996). *Biplots*. London: Chapman and Hall.

[7] Greenacre, M. J. (1984). *Theory and applications of correspondence analysis*. London: Academic Press.

[8] Johnson, R. A., & Wichern, D. W. (2007). *Applied multivariate statistical analysis* (6th edition). Upper Saddle River, NJ: Pearson Prentice Hall.

[9] Jolliffe, I. T. (2002). *Principal component analysis* (2nd ed.). New York: Springer.

[10] Kaiser, H. F. (1958). The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, *23*, 187-200.

[11] Kwan, C. C. Y. (2011). An introduction to shrinkage estimation of the covariance matrix: A pedagogic illustration. *eJournal of Spreadsheets in Education*, 4(3), Article 6. Available at https://epublications.bond.edu.au/ejsie/vol4/iss3/6/ [Accessed 17 August 2018].

[12] Kwan, C. C. Y. (2017). Shrinkage of the sample correlation matrix of returns towards a constant correlation target: A pedagogic illustration based on Dow Jones stock returns, *eJournal of Spreadsheets in Education*, 10(1), Article 3. Available at http://epublications.bond.edu.au/ejsie/vol10/iss1/3 [Accessed 17 August 2018]

[13] Lattin, J., Carroll, J. D., & Green, P. E. (2003). *Analyzing multivariate data*. Pacific Grove, CA: Thomson.

[14] Lay, D. C., Lay, S. R., & McDonald, J. J. (2015). *Linear algebra and its applications* (5th edition). New York: Pearson.

[15] Okada, A., & Tsurumi, H. (2012). Asymmetric multidimensional scaling of brand switching among margarine brands. *Behaviormetrika*, *30*, 111-126.

[16] Spearman, C. (1904). General intelligence, objectively determined and measured. *American Journal of Psychology*, *15*, 201-292.