# Notes on
# Eigenstructures and Factor Analysis

## Eigenstructure

- We know any matrix A can be decomposed (via SVD) as the triple product UDV'.
- When A happens to be square and symmetric (like a correlation matrix or any other cross-products matrix), we will find that U = V, so that A = UDU' or A = VDV'.
- Suppose we compute the cross-products matrix from A. That is, we compute S= A'A. Obviously, we can decompose S into a triple product XGY'. Question is, how does X relate to U, and G to D and Y to V?
- Well, if A = UDV' then A'A = A'UDV' = (UDV')'UDV' = VDU'UDV'. And since U and V are orthogonal (i.e., columns are independent of each other), U'U = I , so VDU'UDV' = $VD^2V'$. So the svd of A'A gets you $VD^2V'$ (and, similarly, the svd of AA' gets you $UD^2U'$)
- We call the svd of a cross-products matrix (such as a correlation matrix) the eigen structure of the matrix. The Us and Vs are called eigenvectors, and the $D^2$s are eigenvalues.

## Eigenvectors

- Since R=A'A = $VD^2V'$, then RV = $D^2V'$. So (simplifying the notation) an eigenvector **v** of a matrix **R** is any vector that satisfies this equation: R**v** = λ**v**. R is a square (normally symmetric) matrix, **v** is the eigenvector, λ is the eigenvalue associated with that eigenvector. The eigenvector is a vector which, if pre-multiplied by a matrix, gets you the vector back again (a property called idempotency).
- Suppose X is a case-by-variable matrix (e.g., the columns of X give responses for each case on a series of attitude questions such as 'Should abortion be legal?' or 'Should citizens be allowed to own guns?') and R is the matrix of correlations among the variables of X. Then the eigenvectors of R (multiplied by their eigenvalues) are known as the factor loadings and are literally the correlations of the each variable in X with an underlying factor or principal component.

## Factor Analysis

- Not a single technique but a family of methods for analyzing a set of observed variables (the data matrix X)
- Two basic branches in family tree: defined factors (aka principal components) and inferred factors (aka common factor analysis or classical factor analysis)
- In principal components, we define new variables (factors), which are linear combinations of our observed variables, that summarize our input data, much like a stock market index summarizes the whole market.
  - the focus is on expressing the new variables (the principal components) as weighted averages of the observed variables
  - the factors (properly called factor scores) have the same order (number of values) as the original variables.
- In common factor analysis, we infer the existence of latent variables that explain the pattern of correlations among our observed variables
  - the focus is on expressing the observed variables as a linear combination of underlying factors.
- In both approaches, the factors are defined as linear combinations of the variables, and the variables are decomposed as linear combinations of the factors. Weird but true.
- Two basic outputs from factor analysis: a set of column (variable) scores called factor loadings (each factor loading has as many values as there are variables in the data matrix), and a set of row (case) scores called factor scores (each factor score has as many values as there are cases in the data matrix).

## Principal Components

- Given as input a rectangular, 2-mode matrix X whose columns are seen as variables, the objective of principal components is to create a new variable (called a factor or principal component) that is a linear combination of the input variables, such that the sum of squared correlations between the principal component (factor) and each of the original variables is maximized.
- Actually, it is to create an ordered set of principal components such that the first principal component explains as much variance in the original variables as possible, and then the second component explains as much of the residual variance not explained by the first as possible, and so on until all variance is accounted for.
    - Since the observed variables may be highly intercorrelated (i.e., share variance), what we are doing is just "collecting together" the shared bits into "components". It doesn't really change anything, it just reallocates things. Think about 10 pieces of luggage of different sizes that are filled with stuff in such a way that each bag is approximately the same weight. Note that if you were to decide to take only 8 bags with you, you would leave behind 20% of your stuff. Now reallocate the contents so that the biggest bags are stuffed to the maximum. This will mean that some of the other bags will be much emptier. Now if you took only 8 bags (the 8 fullest ones) you would leave behind much less than 20% of your stuff.
- We solve this reallocation problem by "factoring" the correlation matrix between the variables. That is, we compute the correlation matrix, and then use SVD to extract the eigenvectors and eigenvalues.
    - the eigenvectors (multiplied by their eigenvalues) are called factor loadings, and these are the correlations of each variable with each factor (principal component)
    - The sum of the squared loadings of each variable with a given factor (the column sum of the squared loadings matrix) will equal the factor's eigenvalue. Hence the eigenvalue summarizes how well the factor correlates with (i.e., summarizes or can stand in for) each of the variables. It is literally the amount of variance accounted for (since correlation squared is variance accounted for).
    - The sum of the squared loadings for each variable across the factors (the row sums of the squared loadings matrix) is defined as the variable's communality. It tells you how much of the variable's variance is captured by the factors. In principal components, the communality of each variable should be 1.0 unless you have chosen to throw away some of the factors (as when we keep only the bigger factors).
    - The loadings can also be seen as a formula for rewriting the variables in terms of the factors: For example if we write $Z1 = b1F1 + b2F2 + b3F3$ ..., i.e., expressing variable Z1 as a linear combination of factors, then the weights b1, b2, b3, etc will turn out to be the factor loadings. There's a different equation for each variable Z.. In matrix form, we can express the collection of equations easily as $X = FB$, where X is the original data (whose columns are the various Z variables), F is the matrix of factor scores (which we haven't discussed yet) and B is the factor loadings matrix (which is just $D^2V$ in the singular value decomposition of the correlation matrix).
    - Since $X=FB=FD^2V$, we can do some simple matrix algebra to obtain a formula for constructing the factor scores. Multiplying both sides by V', we get $XV' = FD^2VV' = FD^2$ and post-multiplying both sides by $D^{-2}$, we get $XV'D^{-2} = F$. This tells us that we can also use the factor loadings to compute the factor scores
- Since a correlation matrix is just a cross-products matrix X'X computed from our original matrix X, where X's columns have been standardized, and since the SVD of X'X gives the same scores as the column scores for the SVD of X (see first section of this handout), it turns that another way to do principal components is to do an SVD of the original matrix X (assuming its columns have been standardized first). If the SVD decomposes X as UDV', then $D^2V$ will be the factor loadings. U will be the factor scores (what we called F above). That is, the columns of U will contain the principal components -- the new variables that summarize X by reallocating the variance so as to load as much as possible on the first few factors. and V will be correspond to factor loadings (actually, factor loadings are the Vs multiplied by the eigenvalues so actually the loadings (called B above) are equal to $D^2V$.

## Common Factor Analysis

- Given as input a rectangular, 2-mode matrix X whose columns are seen as variables, the objective of common factor analysis is to decompose ("factor") the variables in terms of a set of underlying "latent" variables called factors that are inferred from the pattern of correlations among the variables.
- The underlying factors are the "reason for" the observed correlations among the variables. That is, we assume that correlations among the variables are due to the fact that each variable is correlated with the underlying factors. So, if we simplify and assume there is just one underlying factor for a given set of variables, the idea is to see whether we can explain the observed correlation $r(Y,Z)$ between two variables Y and Z as a function of the extent to which each is correlated with an unseen third variable, the factor F. That is, $r(Y,Z) = r(Y,F)*r(Z,F)$. This known as Spearman's fundamental theorem of factor analysis. If there are two underlying factors, then the correlation between two variables is due to their correlations with each of the latent factors, like this: $r(Y,Z) = r(Y,F1)*r(Z,F1) + r(Y,F2)*r(Z,F2)$.