

MLA

CS/IS/IS

EXAM PAPER 202 Q1 CLUSTERING.

- 1 A. We wish to place observations into a group according to their dissimilarities, which has five of the following common properties:
- $d(x,y) \geq 0$ and $d(x,y) = 0$ if $x=y$
 - $d(x,y) = d(y,x)$
 - $d(x,z) \leq d(x,y) + d(y,z)$

Euclidean $\sqrt{\sum_{i=1}^n (x_{ik} - x_{jk})^2}$
 Manhattan (absolute) $\sum_{i=1}^n |x_{ik} - x_{jk}|$
 Maximum $\max_{k=1, \dots, m} |x_{ik} - x_{jk}|$

In HC, aim to find groups of observations st obs within group are similar and different groups are very dissimilar.

Various methods for measuring dissimilarity between 2 groups A and B.

- Single linkage: $d(A,B) = \min_{x \in A, y \in B} d(x,y)$
- Complete linkage: $d(A,B) = \max_{x \in A, y \in B} d(x,y)$
- Average linkage: $d(A,B) = \frac{1}{|A||B|} \sum_{x \in A} \sum_{y \in B} d(x,y)$
- Complete linkage joins first cluster of much larger means of dissimilarity, results in spherical clusters with good internal similarity. Outliers hidden. Invariant under monotonic transform. Small number of large clusters with roughly equal size.
- Single - identifies outliers, invariant under monotonic transform, tends toward chaining effect.
- Average - result in spherical clusters with good internal similarity, invariant under monotonic transform.

B K-Means Clustering

1. Choose number of clusters K and designate cluster centers.
2. Assign each data point to the cluster whose center is closest.
3. For cluster i , calculate its centroid $\bar{x}^i = ((i)_1, (i)_2, \dots, (i)_m)$ where m is dimension (p number of variables) in an observation (these are found by averaging variable s for data points with the cluster).
4. Calculate the sum of squared errors (distance of each object to its cluster centroid):

$$S = \sum_{i=1}^K \sum_{j=1}^n (x_{ij} - (i)_s)^2$$
 Assume N observations. Minimizing S .
5. Re-assign each observation to the cluster whose centroid is closest.
6. Repeat 3-5 until converges.

C Good or finding circular cluster

- k-means is sensitive to starting points, so running the algorithm several times from different starting values help check whether results are robust.

- D. Looking at HC dendrogram we can see that generally see that 2 large groups consisting of several smaller sub-groups.
- As expected the complete linkage joined the 2 groups at a much larger measure of dissimilarity than average linkage
 - Single linkage picked up on another (411) and the chaining effect is particularly noticeable
 - For cluster analysis, a general rule is to plot k vs. S and look for a kink in the curve
 - If no kink, tradeoff between additional complexity by k and better fit by reducing S
 - $k=3$ is where last significant kink appears in curve

E. Standardization

- When constructing a dissimilarity matrix for use in clustering, we need to be aware of how our data is scaled
- Important that different variables are comparably scaled, otherwise variable with greater variance will be more prominent in clustering solution
- These variables are standardized by dividing through by their standard deviation before being used to calculate dissimilarity matrix. Each variable will have variance of 1
- But always outlier may want to give less weight to variable with less info
- PCA can be used to reduce dimensionality. Can see the importance of the components via PCA.
- Ideally take amount of PCs that account for most variance
- Other method may be used - factor analysis and MDS. Cov [4.4]-6

MLA

EXAM PAPER 2012 Q2 FACTOR ANALYSIS

A. Explain Factor Model

- Mathematical approach for attempting to explain the correlation between a large set of variables in terms of a small number of underlying factors
- Main assumption of FA is that it is not possible to observe the underlying factors directly

- Dimensionality of matrix can be reduced from m to $m-1$ by expressing the correlated as:

$$X_1 = \lambda_1 F + \epsilon_1$$

$$X_2 = \lambda_2 F + \epsilon_2$$

$$X_m = \lambda_m F + \epsilon_m$$

- The F in these equations is an underlying common factor, the λ_i 's are known (i) factor loadings, whilst the ϵ_i 's are known errors or specific factors
- Common factor can often be given a general interpretation like "general ability"
- The specific factors ϵ_i will have small variance if X_i is closely related to general abilities

- The observable random vector $X^T = (X_1, X_2, \dots, X_m)$ has non-singular covariance matrix Σ
- Factor Model states that X is linearly independent upon a few unobservable random variables F_1, F_2, \dots, F_p called common factors and m additional sources or variation $\epsilon_1, \dots, \epsilon_m$ called specific factors

$$X_i - \mu_i = \lambda_{i1} F_1 + \lambda_{i2} F_2 + \dots + \lambda_{ip} F_p + \epsilon_i$$

$$X_2 - \mu_2 = X_1 - \mu_1$$

$$X_m - \mu_m = \lambda_{m1} F_1 + \lambda_{m2} F_2 + \dots + \lambda_{mp} F_p + \epsilon_m$$

- The λ_{ij} value is called the factor loading of i th variable on j th factor
- Λ is matrix of factor loadings
- Note that i th specific factor ϵ_i is associated only with response X_i
- F_1, F_m and ϵ_1, ϵ_m are all unobservable random variables

62 A Factor Model

- An observable random vector $X = (x_1, \dots, x_m)$ has mean μ and covariance matrix Σ .
- Factor model states that X is linearly dependent upon a few unobservable random variables F_1, \dots, F_p called common factors and an additional source of variation $\epsilon_1, \dots, \epsilon_m$ called specific factors.

$$\text{Here: } X_i - \mu_i = \lambda_{i1}F_1 + \lambda_{i2}F_2 + \dots + \lambda_{ip}F_p + \epsilon_i$$

$$X_m - \mu_m = \lambda_{m1}F_1 + \lambda_{m2}F_2 + \dots + \lambda_{mp}F_p + \epsilon_m$$

- where λ_{ij} is called the factor loading of the i^{th} variable on the j^{th} factor.
- Under this model, it is assumed that:

$$\begin{aligned} E[F] &= 0 & \text{Cov}[F] &= I = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \\ E[\epsilon] &= 0 & \text{Cov}[\epsilon] &= \Psi = \begin{pmatrix} \psi_1 & 0 & \dots & 0 \\ 0 & \psi_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \psi_m \end{pmatrix} \\ \text{Cov}[\epsilon, F] &= 0 \end{aligned}$$

B Variance - Communality and Uniqueness

- The variance of a MV random variable X can be split into 2 parts.
- First portion of the variance for the i^{th} component arises from the m common factors and is referred to as the communality.
- Remainder of the variance for i^{th} component is due to the specific factor, and is referred to as the uniqueness.
- Denote i^{th} communality by h_i^2 :

$$\sigma_i^2 = \lambda_{i1}^2 + \lambda_{i2}^2 + \dots + \lambda_{ip}^2 + \psi_i$$

$$= h_i^2 + \psi_i$$

Sum of Squares of loadings of i^{th} variable on p common factors

$$\text{i.e. Var}[x_i] = \text{communality} + \text{uniqueness}$$

C Factor Rotation and Varimax Rotation

- If initial loadings are subject to an orthogonal transformation (i.e. multiplied by an orthogonal matrix G), the covariance matrix can still be reproduced.
- An orthogonal transformation corresponds to a rigid rotation or reflection of coordinate axes.
- Thus, an orthogonal transformation of the factor loadings (and implied transformation of factors) is called a factor rotation.
- Irrelevant mathematically, but when the objective is statistical interpretation, may be that one rotation is more useful than another.

- The varimax procedure selects the orthogonal transformation G maximizing V , the sum of the column variances across all factors.
- Maximizing V corresponds to 'spreading out' the squares of the loadings on each factor as much as possible.
- Hence groups of large and negligible coefficients are found in any column of Λ^* i.e. it aims for a rotation that makes the squared loadings Λ^* either large or small i.e. few medium sized values.

D Output Interpretation

- Communalities = Sum of squared loadings in each row of loading matrix
- Communalities + Uniqueness for each variable = total variance for that variable
- SS-loadings = Sum of squared loadings for that factor
- Proportion variance = $SS/m = 55/6$ in this case
- The two factors account for 59.7% of total variance
- Factor 1: literacy: 36% variance. Correlates to good reading abilities
- Factor 2: Spatial awareness: Focuses on visual spatial
- Reading seems to be better explained than small uniqueness - smaller the uniqueness, larger the variance in the model

E PCA, MDS, and FA

- PCA find linear combination of the variables in the data which capture most of the variation in the original data
- MDS seeks to produce lower dimensional summary of data, such that distances between i and j are in the representation as close to dissimilarities between their points as dis_{ij} & (i, j)
- Take a set of dissimilarities and returns a set of dimensional points s.t. the distance between points is dissimilarities
- Classical MDS performs eigen decomposition of data dist matrix to find low-dim summary such that distances are preserved as closely as possible in a lower space
- FA share same obs as PCA but is much more elegant
- Not affected by rescaling unlike PCA
- PCA, MDS no data assumption FA assumed data come from well defined matrix in which assumption holds i.e. ECPJ=0
- PCA: data \Rightarrow PCs FA: obs \Rightarrow factors
- PCA - when specific variables are large, often absorbed into PCs. Where FA must special provision for them. When small PCA and FA produce similar results

- Assumed that:

$$E[F] = 0$$

$$\text{Cov}[F] = I = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$E[E] = 0$$

$$\text{Cov}[E] = \Phi = \begin{pmatrix} \psi_1 & 0 & 0 \\ 0 & \psi_2 & 0 \\ 0 & 0 & \psi_3 \end{pmatrix}$$

F and E are independent so $\text{Cov}[F, E] = 0$

B Variance, Communality and Uniqueness

- Variance of X can be split into 2 parts

- The first portion of variance for i^{th} component (score) from k in common factors and is referred to as the i^{th} communality

- Remainder of variance for i^{th} component is due to specific factor, referred to as uniqueness

- Denoting the i^{th} communality h_i^2 then:

$$\sigma_i^2 = \lambda_{i1}^2 + \lambda_{i2}^2 + \dots + \lambda_{ip}^2 + \psi_i$$

$$= h_i^2 + \psi_i^2$$

$$\text{Var}[X_i] = \text{communality} + \text{uniqueness}$$

- i^{th} communality is the sum of squares of the loadings of i^{th} variable on p common factors

C Factor Rotation and Varimax Rotation

- If initial loading of subject is an orthogonal transformation (ie multiplied by an orthogonal matrix G), the covariance matrix Σ can still be reproduced

- An orthogonal transformation corresponds to a rigid rotation or reflection of coordinate axes

- Hence, the orthogonal transformation of factor loadings (and the implied transformation of the factors) is called a factor rotation

- VARIMAX - Note that the squared loading λ_{ij}^2 is the proportion of the variance of variable i attributable to common factor j :

$$\text{Var}[X_i] = \lambda_{i1}^2 + \lambda_{i2}^2 + \dots + \lambda_{ip}^2 + \psi_i = h_i^2 + \psi_i$$

- We aim for a rotation which makes the loadings either large or small (ie. few medium sized values)

- Varimax procedure selects the orthogonal transformation G maximizing the sum of column variances and all factors $j=1 \dots p$

MLA

EXAM PAPER Q2 FACTOR ANALYSIS

Maximizing V corresponds to 'spreading out' the squares of the loadings on each factor as much as possible. Hence, both groups of large and negligible coefficients are found in any column of Λ^* .

D Interpretation of output.

- Communalities are the sum of squared loadings in each row of loading matrix
- The communality for each variable added to the uniqueness for each variable is the total variance for that variable. (equal to 1 if data is standardized)
- The SS loadings are the sum of squared loadings for each factor (column).
- Cumulative = 59.7% of variance of data accounted for by 2 factors

Factor 1: Highest loading on reading, vocab and general, this could be interpreted as English Skill/Knowledge factor.

Factor 2: High loading on blocks, picture and general and more, could be interpreted as visualisation in the mind factor or problem solving factor.

E. Other Dimension Reduction techniques

PCA v FA

- PCA looks for linear combinations of the data matrix X that are uncorrelated and of high variance, whilst FA seeks uncorrelated linear combinations of the variables representing underlying

fundamental quantities

- PCA makes no assumption about form of covariance matrix, whilst FA assumes data comes from a well defined model in which specific assumptions hold e.g. $\mathbb{E}[E_i] = 0$, $\mathbb{E}[E_i^2] = 1$

- PCA: data \Rightarrow PCs. FA: factors \Rightarrow data

- When specific variables or large bias are absorbed into the PCs whereas FA model special provision for them. When specific variables are small, PCA and FA give similar results

- Two analysis often performed together. Example \rightarrow can conduct a PCA to determine number of factors to extract in FA

FA v MDS

- FA requires that the underlying data are distributed as MVN and that the relationships are linear. MDS implies no such restriction.
- As long as rank ordering of distances (or similarities) in matrix 1) meaningful, MDS can be used.

- In terms of robustness differences, FA tends to extract more factors (dimensions) than MDS, as a result MDS often yields more readily interpretable solutions.

- MDS can be applied to any kind of distance or similarity, while FA requires us to first compute a correlation matrix.

- MDS can be based on subjects' direct assessment of similarity between stimuli, while FA requires subjects to rate those stimuli on some list of attributes.

MLA

EXAM PAPER 2012 Q1

1 A Dissimilarity Measure

Euclidean $\sqrt{\sum_{k=1}^m (x_{ik} - x_{jk})^2}$

Manhattan Distance (Manhattan): $\sum_{k=1}^m |x_{ik} - x_{jk}|$

Maximum: $\max_{k \in \{1, 2, \dots, m\}} |x_{ik} - x_{jk}|$

LINKAGE

Single: $d(A, B) = \min_{x \in A, y \in B} d(x, y)$ result in spherical cluster, good internal symmetry, display outlier, invariant under monotonic transform

Complete: $d(A, B) = \max_{x \in A, y \in B} d(x, y)$ Join cluster at each layer instead of dist result in spherical cluster, good internal symmetry, how outlier transform, smaller number of large cluster

Average Linkage: $d(A, B) = \frac{1}{|A||B|} \sum_{x \in A} \sum_{y \in B} d(x, y)$ variant under monotonic transform

B K-Means Clustering Algorithm

1. Choose number of clusters k and designate cluster centers.
2. Assign each data point to the cluster whose center is closest.
3. For cluster i , calculate (i) centroid $c_i = (c_{i1}, c_{i2}, \dots, c_{im})$ where m denotes the number of variables in an observation.
4. Calculate the sum of squared distances of each object to its cluster centroid:
$$SS = \sum_{i=1}^n \sum_{j=1}^m (x_{ij} - c_{ij})^2$$

Assume total of n observations. Want to minimize SS

5. Re-assign each observation to the cluster whose centroid is closest.
6. Repeat (3)-(5) until convergence.

C. Different Start points may give different results
Want to find convergence

K-means clustering works for circular shapes

D 2-groups from some plot and graph

E Standardization prior variable with larger variance appear prominently in the clustering result.
- Great extra weight to small variables etc

Dimension reduction will make a easier solution more interpretable, will remove non-significant variables from analysis