

09/05/15 MLA
EXAM PAPER 2014 Q1

DAVID WEBBRECHT

I A. Usefulness of PCA.

- Dimension reduction - included dimensions are orthogonal and ordered to account for as much variation as possible.
- Visualisation
- Lower Dimensional Summary.
- Helps identify relationships in data
- Potentially assists other clustering or classification algorithms (performance)

B. Lagrange Multiplier

$$P = \mathbf{a}^T \mathbf{z} - \lambda (\mathbf{a}^T \mathbf{a} - 1)$$

Solve $dP/d\mathbf{a}_k = 0$ for $k=1 \dots m$

$$dP/d\mathbf{a} = 0 \Rightarrow \mathbf{a}^T \mathbf{a} = 1 \Rightarrow \text{constraint satisfied}$$

$$\Rightarrow \sum_{i=1}^m \mathbf{a}_i \mathbf{b}_i^T - \lambda \mathbf{a} \mathbf{a}^T = 0$$

$$\Rightarrow \mathbf{a}^T \mathbf{z} = -\lambda \mathbf{a}^T \Rightarrow \mathbf{a} \text{ is an eigenvector of } \mathbf{z}$$

C. Standardisation of data

- If not standardised, one component will account for most variation with figure
- more prominently in the solution, swamping out all other values
- PCA will focus on that variable

ii. Number of PCs

- Could select 2 or 3 3 - account for more variation (96.9%) but 2 dimension is easier to graph/visualise and interpret.

iii. PC1: high score = high birth rate, high death rate, high ID and small LEM and LEF.
Low score = opposite
Interpreted as health care - independent of GNP

PC2: High score = High loading on GNP \Rightarrow wealth
High score for poor countries, low score for rich countries

C. iv. $(-0.3, -1.1, -0.5, 0.8, 0.9, -0.6)$

$$0.43(-0.3) + 0.37(-1.1) + 0.46(-0.5) - 0.47(0.8) - 0.48(0.9) + 0.08(-0.6) = -1.622$$

$$-0.03(-0.3) + 0.14(-1.1) + 0.06(-0.5) - 0.03(0.8) - 0.04(0.9) - 0.99(-0.6) = 0.282$$

$$(-1.622, 0.282)$$

D PCA and MDS

- Same aim - Dimension reduction
- Create configuration of points (map) in lower dimension so that dissimilarities match original dissimilarities as much as possible
- Choice of loss function (sum of squares or stress)
- Choice metric or non-metric scaling
- In classical setting - eigen decomposition of a function of a dissimilarity matrix
- PCA - eigen decomposition of covariance matrix
- Classical setting gives same map as PCA does

PCA and FA

- Same aim - dimensionality reduction
- PCA looks for linear combinations of the data matrix X that are uncorrelated and of high variance, while FA seeks unobserved linear combinations of the variables representing underlying fundamental quantities.
- PCA makes no assumption about the form of the cov matrix, FA assumes that data come from a well-defined model in which specific assumptions hold
eg. $E[FF'] = I$
- PCA: data \Rightarrow PC's FA: Factors \Rightarrow data
- When specific variances are large they are absorbed into the PC's whereas FA model a special provision for them. When the specific variances are small, PCA and FA give similar results.
- Choices:
 - orthogonal / oblique methods
 - Choice of rotation and factor loadings
 - Scale invariant

MLA

2014 Q3 EXAM PAPER

3A. Methods for using Logistic Regression

- Parametric classification technique
- Explanatory model (continuous) without discretisation
- Binary variable
- Interaction
- Variable can be categorical or continuous \Rightarrow Categorical outcome
- Predictor variable in LR does not need to be linearly related, normally distributed or have equal variance within each group

B. Explain Output

- Interpret: expected ~~value~~ mean value of y (leaving car) when all $x=0$
if everything is 0, person will not take leaving cert according to model nor a.

Each unit increase in DVRT implies Δ by 0.06

Sex - male likely to do LC if girl.

Intercept in \ln by 0.034 by prestige score of person

$$C. \frac{\exp(-8.648 + 0.06(DVRT) + 0.51(\text{Sex}) + 0.034(Pst))}{1 + \exp(-8.648 + 0.06(DVRT) + 0.51(\text{Sex}) + 0.034(Pst))} = P(LC)$$

Female DVRT Sex = 120 Prestige = 32

$$= \frac{\exp(-8.648 + 0.06(120) + 0.51(2) + 0.034(32))}{1 + \exp(-8.648 + 0.06(120) + 0.51(2) + 0.034(32))} = \frac{\exp(0.751)}{1 + \exp(0.751)}$$

$$= 0.67$$

\Rightarrow will do LC.

D. Including interaction adds possible products of the covariates to the model.

- Allows for the effect of one covariate to be altered depending on the value of another covariate

- Interaction would be relevant if or the effect of DVRT had upon LC would differ depending on or not an individual was a girl or boy

For example the ^{DVRT} sex of a person may be considered very important if the person is male but not if boy or female

E Difference between Logistic Regression and LDA

- In LDA the decision boundary between class k and class l given by

$$\log \frac{P(k|x)}{P(l|x)} = \log \frac{\pi_k}{\pi_l} + \log \frac{f(x|k)}{f(x|l)} = 0$$

- In logistic regression model assumption is: $\log \frac{P(k|x)}{P(l|x)} = \beta_0 + \beta^T x$

- Model have same form

- Differ in way the linear coefficients are estimated

Log R - parameters found through MLE - Binary output

LDA - no restriction to which group it belongs

Notes EXAM PAPER 2014 Q2 CLUSTERING

2 A. Hierarchical

- Start with each point as a cluster on its own.
- Clusters are merged by combining cluster 2 at each depth.
- Clusters defined by dissimilarity matrix.
- Continue until a single cluster is formed consisting of all data points.
- Go back and determine when merging should be stopped: $h + 3s_n$ or relatively large jump in next merging distance as seen through dendrogram.

Iterative

- Predefined number of clusters - Specify number of clusters you want.
- Assign data points to a specific cluster.
- Iterate and re-assign membership until convergence.
- Convergence - repeat of previous clustering of data (check that with minimal internal dissimilarity).
- Example includes: K means, partitioning around medoid.

B. Maximum Dissimilarity matrix

- Compare all values in one row for one country against another row from a different country.
- Compute UK(GBR) with USA(GBR) etc for DR, ID and LEM.
- Distance = max of these four distances.

	UK	USA	IRELAND	CHINA
UK	0	3.4	2.4	23.6
USA	3.4	0	1.6	22.9
IRELAND	2.4	1.6	0	24.5
CHINA	23.6	22.9	24.5	0

C. Dendrogram with complete linkage complete linkage $d(A,B) = \max_{x \in A, y \in B} d(x,y)$

Merge USA + Ireland at 1.6
Merge UK and group at 3.4
Merge China and group at 24.5

Choose $h + 3s_n$

We chose 2 group as there is a massive jump in the merging height from two groups to one group.



D Rand Index

Indicates level of agreement between two different clustering methods with possibly different numbers of clusters by considering how frequently the methods agree on when ~~points~~ pairs of data points should be in same cluster.

$$\text{Rand Index} = \frac{\binom{n}{2} + 2 \sum_{i=1}^g \sum_{j=1}^g \binom{n_{ij}}{2}}{\binom{n}{2}} - \left[\sum_{i=1}^g \binom{n_i}{2} + \sum_{j=1}^g \binom{n_j}{2} \right]$$

$$= \frac{\binom{31}{2} + 2 \left[\binom{5}{2} + \binom{6}{2} + \binom{3}{2} + \binom{4}{2} + \binom{8}{2} + \binom{5}{2} \right] - \left[\binom{11}{2} + \binom{7}{2} + \binom{13}{2} + \binom{16}{2} + \binom{12}{2} \right]}{\binom{31}{2}} = 0.49$$

E Dissimilarity from data points Containing binary, numerical and categorical data

- We can work out a dissimilarity for the measured variables (e.g. euclidean).
- Simple matching for binary and proportion of times in agreement in categorical.
- Take a weighted average of these numbers, with weights being proportion of ~~two~~ ~~in agreement~~ variable of that type of data.

F Real life Clustering Application

- Bank cluster transactions and individual profile - fraud people out similarities
- Detecting terrorist networks - esp. behaviour - some clusters
- Online dating - similar interest / partner email
- Search engine - google suggested paths

3/11/14 DEATH BY MIA Exam Paper 2014

A Usefulness of PCA

- Dimension reduction technique - so included dimensions are orthogonal and ordered to account for as much variation as possible
- Visualisation
- Lower dimensional Summary
- Helps identifying relationships in the data by removing inherent noise
- Potentially assist other clustering or classification algorithms (performance)

$$p = \mathbf{a}^T \mathbf{Z} \mathbf{a} - \lambda (\mathbf{a}^T \mathbf{a} - 1)$$

Solve $\frac{dp}{da_k} = 0$ for $k = 1 \dots m$

$$\frac{dp}{d\lambda} = 0 \Rightarrow \mathbf{a}^T \mathbf{a} - 1 = 0 \Rightarrow \text{constraint satisfied}$$

$$\Rightarrow \sum_{i=1}^m a_i b_{ik} - \lambda a_k = 0$$

$$\Rightarrow \mathbf{a}^T \mathbf{Z} = \lambda \mathbf{a}^T \Rightarrow \mathbf{a} \text{ is an eigen vector of } \mathbf{Z}$$

C.1 One component will account for most of variance of the model, it will dominate the principle component. Will swamp all other values. PCI will focus on that component.

ii. Could be 2 or 3 3 - accounts for more variation.
2 - easier to graph / visualise.

iii. 1 - more people born - more that die, die earlier, life expectancy goes down

PCI: high score - high birth rate, high death rate, high ID, small ZEM and AFF
Low score is the opposite HEALTH CARE
independent of Group

2.

2014 Paper

PC2: High loading on GNP \Rightarrow wealth
 High score for poor country
 Low score for rich country

C10. $(-0.3, -1.1, -0.5, 0.8, 0.9, -0.6)$

$(-1.622, 0.282)$
 coord 1. coord 2

all take by PC1, then all take by PC2
 $= 2 \text{ points}$

$$0.43(-0.3) + (0.37)(-1.1) + \dots + (0.08)(-0.6) = -1.622$$

$$-0.03(-0.3) + (0.14)(-1.1) + \dots + (-0.99)(-0.6) = 0.282$$

D Factor Analysis, multidimensional scaling

MDS

- ~~Representing~~ Create configuration of points (map) in lower dimension so dissimilarities match original dissimilarities as closely as possible
- Choice of loss function (Stress, Sammon Stress)
- whether it is metric or non metric

\hookrightarrow Dissimilarity measure / Standardization

In CLASSICAL SETTING - eigen decomposition of a function of the dissimilarity matrix
 (PCA) \Rightarrow covariance matrix

Classical setting gives same map as PCA does

Factor Analysis

- seek unobserved linear combinations of the variables representing underlying quantities, i.e. the factors
- Tries to model the covariance matrix with a reduced set of parameters under specific assumptions holding true
 factors ind, specific factor $E(\epsilon_i)$ of 0.
- Choices:
 - Orthogonal / oblique method
 - Choice of rotation of factor loadings
 - Scale invariant

1/2/14 MLA 2014 Exam Paper Q2

Q2 A Hierarchical & iterative clustering

Hierarchical

- each point starts at its own cluster
- clusters are merged by combining closest 2 at each depth
- cluster defined by dissimilarity matrix - means for singleton and linkage for distance
- Continue until a single cluster is formed consisting of all data points
- Go back and determine when merging should be stopped e.g. rule average height + 3 s.d of heights or relatively large jump in next merging distance etc. - as seen through a dendrogram

Iterative

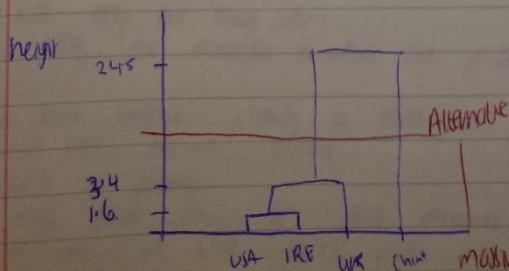
- Pre defined number of clusters
- Assign data points to a specific cluster
- Iterate points between clusters, until convergence
- Convergence - repeat of a previous clustering of the data [choose that with minimum internal dissimilarity]
- Examples include k-means, partitioning around medoids

B Maximum Dissimilarity

	UK	USA	IRELAND	CHINA
UK	0	3.4	2.4	23.6
USA	3.4	0	1.6	22.9
IRELAND	2.4	1.6	0	24.5
CHINA	23.6	22.9	24.5	0

Take max difference between two groups using horizontal values

C.



merge USA and Ireland 1.6
merge UK and China 3.4
merge China and UK 24.5

h + 3 s.d
note 24.5 rule of thumb:
all in one cluster

max jump: 2 groups most appropriate

2

D. Rand Index

Indicates level of agreement between 2 different clustering methods with possibly different numbers of clusters by considering how frequently the methods agree on when pairs of data points should be in the same cluster.

$$= \frac{\binom{31}{2} + 2 \left[\binom{15}{2} + \binom{9}{2} + \binom{7}{2} + \binom{8}{2} + \binom{15}{2} \right] - \left[\binom{11}{2} - \binom{7}{2} + \binom{16}{2} - \binom{15}{2} \right]}{\binom{31}{2}} = 0.49$$

- E. Take each collection of variables of a certain type and calculate dissimilarity for that subset eg euclidean for numeric, simple matching for binary and proportion of time in agreement in 'categorical'
- Take a weighted average of these numbers with weights being proportion of variables of that type of data

F. Steering of students in class

- Bank cluster transactions and individual profile \rightarrow find people doing similarly
- Determining in terrorist networks - same behavior - same cluster
- Online dating - similar interests/interests / partner
- Search engines - google suggested posts/products

1/12/14 PLA Exam Paper

3 (b) Assumptions for LDA

- Know how many groups there are in data
- LDA and QDA assume data within a class are distributed (spread) according to a multivariate normal distribution, which each class having its own mean vector

In LDA a common covariance whilst in QDA, each class can have its own covariance matrix

c Derivation from the notes

$$D \quad \left[\begin{pmatrix} 0.9 \\ 0.6 \end{pmatrix} - \frac{1}{2} \left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right) \right]^T \begin{pmatrix} 0.6 & -0.2 \\ -0.2 & 0.4 \end{pmatrix} \left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right)$$

Group 1 if expression is positive

$$\left[\begin{pmatrix} 0.9 \\ 0.6 \end{pmatrix} - \begin{pmatrix} 1 \\ 0.5 \end{pmatrix} \right]^T \begin{pmatrix} 0.6 & -0.2 \\ -0.2 & 0.4 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$= (-0.1, -0.1) \begin{pmatrix} 0.6 & -0.2 \\ -0.2 & 0.4 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = (-0.29, 0.3) \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0.3$$

Positive \Rightarrow assign to group one

- E Logistic Regression uses a logit link function to link the probability of class assignment to a linear function of data points
- $$\log \frac{P(X=1)}{P(X=0)} = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots$$

Parameters found through MLE

In LDA we have $\log \frac{P(X=1)}{P(X=0)} = \log \left(\frac{P_1}{P_0} \right) + \log \frac{f(x|X=1)}{f(x|X=0)}$

\downarrow ratio of prior probs

Difference is how parameters are obtained

LR only has binary output

LDA, no restriction to which class it belongs to

- F In nearest neighbour - closely to nearest group
- LDA uses maximum distance

LDA is parametric

KNM does not allow for prior distribution of class membership

As number of NN increases, will be changed to class which is most prevalent in data

LDA - separating by hyper plane