Principal Components Analysis

Sargur Srihari University at Buffalo

Topics

- Projection Pursuit Methods
- Principal Components
- Examples of using PCA
- Graphical use of PCA
- Multidimensional Scaling

Motivation

- Scatterplots
 - Good for two variables at a time
 - Disadvantage
 - may miss complicated relationships
- PCA is a method to transform into new variables
- Projections along different directions to detect relationships

- Say along direction defined by $2x_1+3x_2+x_3=0$

Projection pursuit methods

- Allow searching for "interesting" directions
- Interesting means maximum variability
- Data in 2-d space projected to 1-d:



Principal Components

- Find linear combinations that maximize variance subject to being uncorrelated with those already selected
- Hopefully there are few such linear combinations-- known as principal components
- Task is to find a *k*-dimensional projection where $0 \le k \le d-1$

Data Matrix Definition

X = n x d data matrix of *n* cases



Assume ${\bf X}$ is mean-centered, so that the value of each variable is subtracted for that variable

Projection Definition

Let **a** be a $p \ge 1$ column vector of projection weights that result in the largest variance when the data X are projected along **a**

Projection of a data vector $\mathbf{x} = (x_1, ..., x_p)^t$ onto $\mathbf{a} = (a_1, ..., a_p)^t$ is the linear combination $\mathbf{a}^t \mathbf{x} = \sum_{j=1}^p a_j x_j$

Projected values of all data vectors in X onto a is Xa -- an *n* X 1 column vector-- a set of scalar values corresponding to *n* projected points Since X is *n* X *p* and a is *p* X 1 Therefore Xa is *n* X 1

Variance along Projection

Variance along a is

$$\sigma_a^{2} = (Xa)^T (Xa)$$
$$= a^T X^t Xa$$
$$= a^T Va$$

where $V = X^{t}X$ is the $p \times p$ covariance matrix of the data since X has zero mean

Thus variance is a function of both the projection line a and the covariance matrix V 8

Maximization of Variance

Maximizing variance along a is not well-defined since we can increase it without limit by increasing the size of the components of a.

Impose a normalization constraint on the a vectors such that $a^{T}a = 1$

Optimization problem is to maximize $u = a^{t}Va - \lambda(a^{t}a - 1)$

Where λ is a Lagrange multiplier. Differentiating wrt a yields

$$\frac{\partial u}{\partial a} = 2Va - 2\lambda a = 0$$

which reduces to
 $(V - \lambda I)a = 0$ Characteristic Equation!

What is the Characteristic Equation?

Given a *d* x *d* matrix V a very important class of linear Equations is of the form $Vx = \lambda x$ dx d dx 1 dx 1which can be rewritten as $(V - \lambda I)x = 0$

If *V* is real and symmetric there are *d* possible solution vectors, called Eigen Vectors, e_1 , e_d and associated Eigen values $\lambda_1, ..., \lambda_d$ Principal Component is obtained from the Covariance Matrix

If the matrix *V* is the Covariance matrix Then its Characteristic Equation is

$$(V - \lambda I)\mathbf{a} = 0$$

Roots are Eigen Values Corresponding Eigen Vectors are principal components

First principal component is the Eigen Vector associated with the largest Eigen value of *V*.

Srihari

Other Principal Components

- Second Principal component is in direction orthogonal to first
- Has second largest Eigen value, etc



Projection into k Eigen Vectors

- Variance of data projected into first k Eigen vectors $e_1, ..., e_k$ is $\sum_{j=1}^{k} \lambda_j$
- Squared error in approximating true data matrix X using only first k Eigen vectors is
 - $\frac{\sum_{j=k+1}^{d} \lambda_{j}}{\sum_{l=1}^{d} \lambda_{l}}$

How to choose k ?

- increase k until squared error is less than a threshold

Usually 5-10 principal components capture 90% variance in data



PCA using correlation matrix and covariance matrix



Graphical Use of PCA

Projection onto first two principal components of six dimensional data 17 pills (data points) Six values are times at which specified proportion of pill has dissolved: 10%, 30%, 50%, 70%, 75%, 90%

N **Principal Component** 17 q 0 10 7 6 Ņ -2 0 **Principal Component 1**

15

Pill 3 is very different

Computational Issue: Scaling with Dimensionality



Can be applied to large numbers of records *n* But does not scale well with dimensionality *d*

Also, appropriate *Scalings of variables* have to be done¹⁷

Multidimensional Scaling

- Using PCA to project on a plane is effective only if data lie on 2-d subspace
- Intrinsic Dimensionality
 - Data may lie on string or surface in *d*-space
 - E.g., when a digit image is translated and rotated
 - Then images in pixel space lie on a 3-dimensional manifold (defined by location and orientation)

Goal of Multidimensional Scaling

- Detecting underlying structure
- Represent data in lower dimensional space so that distances are preserved
 - Distances between data points are mapped to a reduced space
- Typically displayed on a 2-d plot
- Begin with distances and then compute the plot
 - E.g., psychometrics and market research where similarities between objects are given by subjects

Defining the *B* Matrix

- For an *n* x *d* data matrix X we could compute *n* x *n* matrix $B = XX^t$
- We will see (next slide) that the Euclidean distance between the *ith* and *jth* objects is given by

$$d_{ij}^{2} = b_{ii} + b_{jj} - 2b_{ij}$$

• Matrices *XX*^{*t*} and *X*^{*t*}*X* are both meaningful

X^tX versus XX^t

- If X is $n \times d$ d=4 $\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix}$ Covariance Matrix • X¹X is $d \times d$ $d \times n$ $n \times d$ $d \times d \swarrow$ $\begin{bmatrix} x_{11} & x_{21} & \vdots & x_{n1} \\ x_{12} & x_{22} & \vdots & x_{n2} \\ x_{13} & x_{23} & \vdots & x_{n3} \\ x_{14} & x_{24} & \vdots & x_{n4} \end{bmatrix} \times \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{n} x_{i1}^{2} & \sum_{i=1}^{n} x_{i1}x_{i2} \\ \sum_{i=1}^{n} x_{i1}x_{i2} \\ \sum_{i=1}^{n} x_{i1}x_{i2} \end{bmatrix}$
- $B = XX^t$ is $n \times n$

$$n \times d \qquad d \times n$$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix} \times \begin{bmatrix} x_{11} & x_{21} & \vdots & x_{n1} \\ x_{12} & x_{22} & \vdots & x_{n2} \\ x_{13} & x_{23} & \vdots & x_{n3} \\ x_{14} & x_{24} & \vdots & x_{n4} \end{bmatrix} = \begin{bmatrix} b_{11} = \sum_{j=1}^{4} x_{1j}^{2} & b_{12} \sum_{i=1}^{4} x_{1j} x_{2j} & \vdots & b_{1n} \\ b_{21} & b_{22} & \vdots & b_{2n} \\ b_{31} & b_{32} & \vdots & b_{3n} \\ b_{41} & b_{42} & \vdots & b_{4n} \end{bmatrix}$$
B Matrix contains distance information

Factorizing the *B* matrix

- Given a matrix of distances D
 - Derived from original data by computing n(n-1)/2 distances
 - Compute elements of B by inverting

$$d_{ij}^{2} = b_{ii} + b_{jj} - 2b_{ij}$$

- Factorize B
 - in terms of *eigen vectors* to yield coordinates of points
 - Two largest eigen values would give 2-d representation



Thus expressing b_{ij} as a function of d_{ij}^2 Method is known as *Principal Coordinates Method*

Criterion for Multidimensional Scaling

Find projection into two dimensions to minimize

Srihari

$$\sum_{i}\sum_{j}(\delta_{ij}-d_{ij})^{2}$$

Observed distance between points *i* and *j* in *d*-space

Distance between the points in two-dimensional space

Criterion is invariant wrt rotations and translations. However it is not invariant to scaling Better criterion is $\sum_{i} \sum_{j} (\delta_{ij} - d_{ij})^2$ or $\overline{\sum d_{ij}^{2}}$

$$\sqrt{\frac{\sum_{i} \sum_{j} (\delta_{ij} - d_{ij})^{2}}{\sum_{i,j} d_{ij}^{2}}} \leftarrow \begin{array}{c} \text{Called} \\ \text{stress} \\ 24 \end{array}$$

Algorithm for Multidimensional Scaling

- Two stage procedure
- Assume that $d_{ij} = a + b \delta_{ij} + e_{ij}$

Original dissimilarities

- Regression in 2-D on given dissimilarities yielding estimates for *a* and *b*
- Find new values of d_{ij} that minimize the stress
- Repeat until convergence

Multidimensional Scaling Plot: Dialect Similarities

Numerical codes of villages and their counties

1	1	North Wheatley	Nottinghamshire
	2	South Clifton	Nottinghamshire
	3	Oxton	Nottinghamshire
	4	Eastoft	Lincolnshire
	5	Keelby	Lincolnshire
	6	Wiloughton	Lincolnshire
	7	Wragby	Lincolnshire
	8	Old Bolingbroke	Lincolnshire
	9	Fulbeck	Lincolnshire
	10	Sutterton	Lincolnshire
	11	Swinstead	Lincolnshire
	12	Crowland	Lincolnshire
	13	Harby	Leicestershire
	14	Packington	Leicestershire
	15	Goadby	Leicestershire
	16	Ullesthorpe	Leicestershire
	17	Empingham	Rutland
	18	Warmington	Northamptonshire
	19	Little Harrowden	Northamptonshire
	20	Kislingbury	Northamptonshire
	21	Sulgrave	Northamptonshire
	22	Warboys	Huntingdonshire
	23	Little Downham	Cambridgeshire
	24	Tingewick	Buckinghamshire
	25	Turvey	Bedfordshire



Each Pair of villages rated by percentage of 60 items for which villagers used different words 26

We are able to visualize 625 distances intuitively

Variations of Multidimensional Scaling

- Above methods are called metric methods
- Sometimes precise similarities may not be known– only rank orderings
- Also may not be able to assume a particular form of relationship between d_{ij} and δ_{ij}
 - Requires a two-stage approach
 - Replace simple linear regression with monotonic regression

Multidimensional Scaling: Disadvantages

- When there are too many data points structure becomes obscured
- Highly sophisticated transformations of the data (compared to scatter lots and PCA)
 - Possibility of introducing artifacts
 - Dissimilarities can be more accurately determined when they are similar than when they are very dissimilar
 - Horseshoe effect when objects manufactured in a short time span differ greatly from objects separated by greater time gap
- Biplots show both data points and variables