

Principal Components Analysis

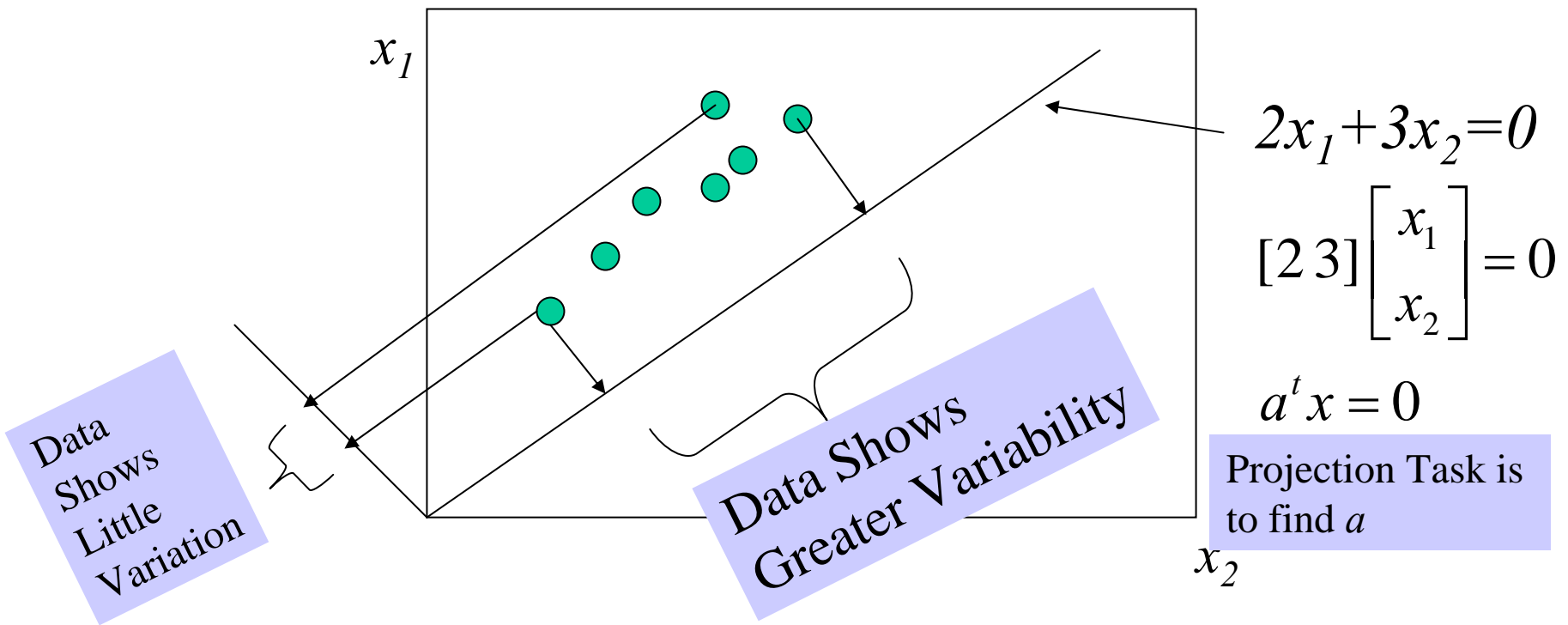
Transformation of data into new variables

Motivation

- Scatterplots
 - Good for two variables at a time
 - Disdvantage: may miss complicated relationships
- Projections along different directions may detect them
 - 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 $0 \leq k \leq d-1$

Principal Components Analysis

$\mathbf{X} = n \times p$ data matrix of n cases

p variables

$\mathbf{x}(1)$

$\mathbf{x}(i)$

$\mathbf{x}(n)$

$\mathbf{x}(i)$ is a $p \times 1$ column vector

Each row of matrix is of the form $\mathbf{x}(i)^T$

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

Projection

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

Projection of a data vector $\mathbf{x} = (x_1, \dots, x_p)^t$ onto $\mathbf{a} = (a_1, \dots, a_p)^t$ is the linear combination

$$\mathbf{a}^t \mathbf{x} = \sum_{j=1}^p a_j x_j$$

Projected values of all data vectors in \mathbf{X} onto \mathbf{a} is $\mathbf{X}\mathbf{a}$ -- an $n \times 1$ column vector-- a set of scalar values corresponding to n projected points

Since \mathbf{X} is $n \times p$ and

\mathbf{a} is $p \times 1$

Therefore $\mathbf{X}\mathbf{a}$ is $n \times 1$

Variance along Projection

Variance along \mathbf{a} is

$$\sigma_{\mathbf{a}}^2 = (\mathbf{X}\mathbf{a})^T (\mathbf{X}\mathbf{a})$$

$$= \mathbf{a}^T \mathbf{X}^t \mathbf{X} \mathbf{a}$$

$$= \mathbf{a}^T \mathbf{V} \mathbf{a}$$

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

Thus variance is a function of both the projection line \mathbf{a} and the covariance matrix \mathbf{V}

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 $u = a^T V a - \lambda (a^T a - 1)$

Where lambda is a Lagrange multiplier.

Differentiating wrt a yields

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

which reduces to

$$(V - \lambda I) a = 0$$

← Characteristic Equation!

What is the Characteristic Equation?

Given a $d \times d$ matrix M , a very important class of linear Equations is of the form $Mx = \lambda x$

$$\begin{matrix} dx & dx \\ dx & dx \end{matrix}$$

which can be rewritten as $(M - \lambda I)x = 0$

If M is real and symmetric there are d possible solution Vectors, called Eigen Vectors, e_1, e_d and associated Eigen values

$$\lambda_1, \dots, \lambda_d$$

Principal Component is obtained from the Covariance Matrix

If the matrix V is the Covariance matrix

Then its Characteristic Equation is

$$(V - \lambda I)a = 0$$

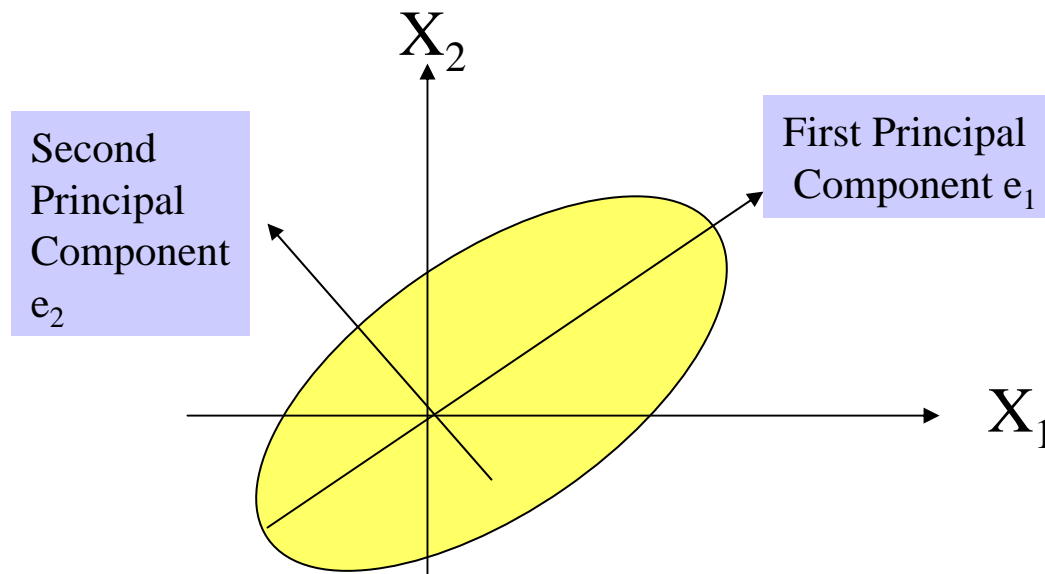
Roots are Eigen Values

Corresponding Eigen Vectors are principal components

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

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, \dots, e_k is $\sum_{j=1}^k \lambda_j$
- Squared error in approximating true data matrix X using only first k Eigen vectors is $\frac{\sum_{j=k+1}^p \lambda_j}{\sum_{l=1}^p \lambda_l}$
- How to choose k ?
 - increase k until squared error is less than a threshold

Usually 5 to 10 principal components capture 90% of variance in the data

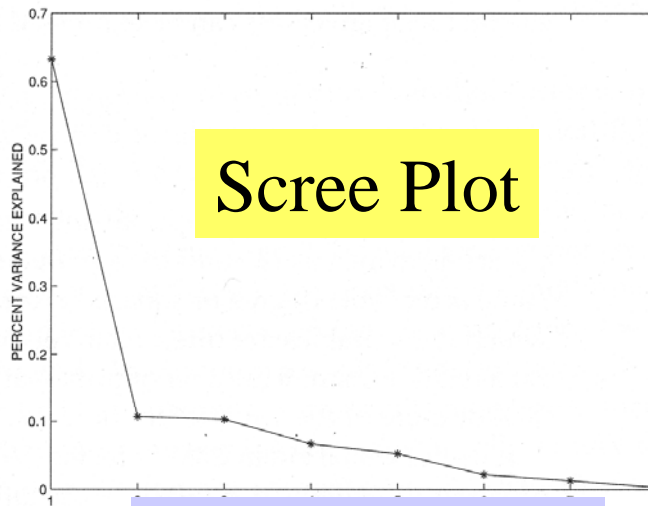
Example of PCA

CPU data

Eigen values of
Correlation
Matrix

CPU data
8 Eigen values:
63.26
10.70
10.30
6.68
5.23
2.18
1.31
0.34

Percent
Variance Explained



Amount of variance
explained by each
consecutive
Eigen value

Eigen Value number



Weights put by
first component e_1

on eight variables are:

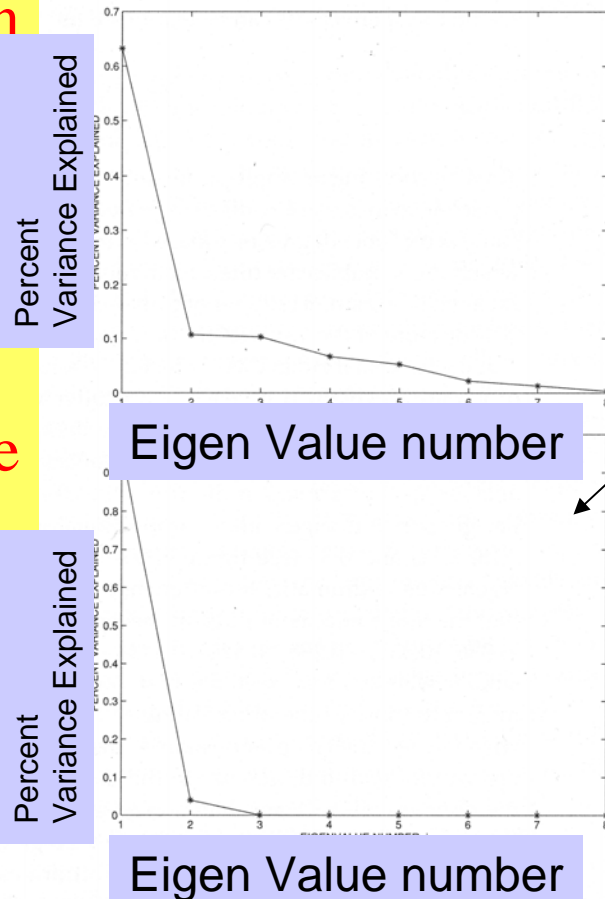
0.199
-0.365
-0.399
-0.336
-0.331
-0.298
-0.421
-0.423

An example
Eigen Vector

PCA using correlation matrix and covariance matrix

Scree Plot
Correlation
Matrix

Scree Plot
Covariance
Matrix



Proportions of variation attributable to different components:

96.02

3.93

0.04

0.01

0

0

0

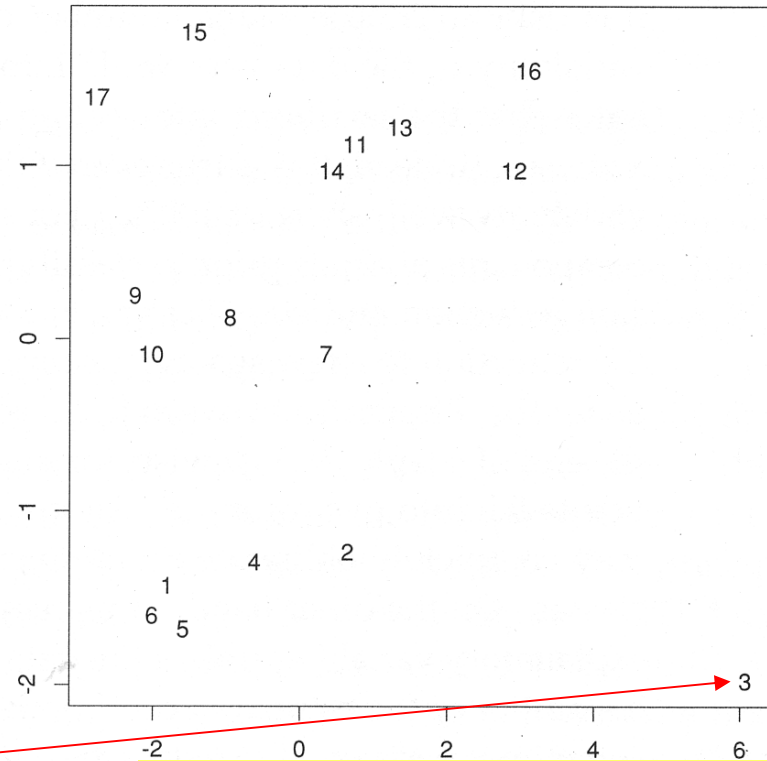
0

Figure 3.17 Scree plots for the computer CPU data set. The upper plot displays the eigenvalues from the correlation matrix, and the lower plot is for the covariance matrix.

Graphical Use of PCA

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%

Principal Component 2



Pill 3 is very different

Principal Component 1

Figure 3.18 Projection onto the first two principal components.

Computational Issues

- $O(np^2 + p^3)$

To calculate V

Solve Eigen value equations
for the $p \times p$ matrix

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

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

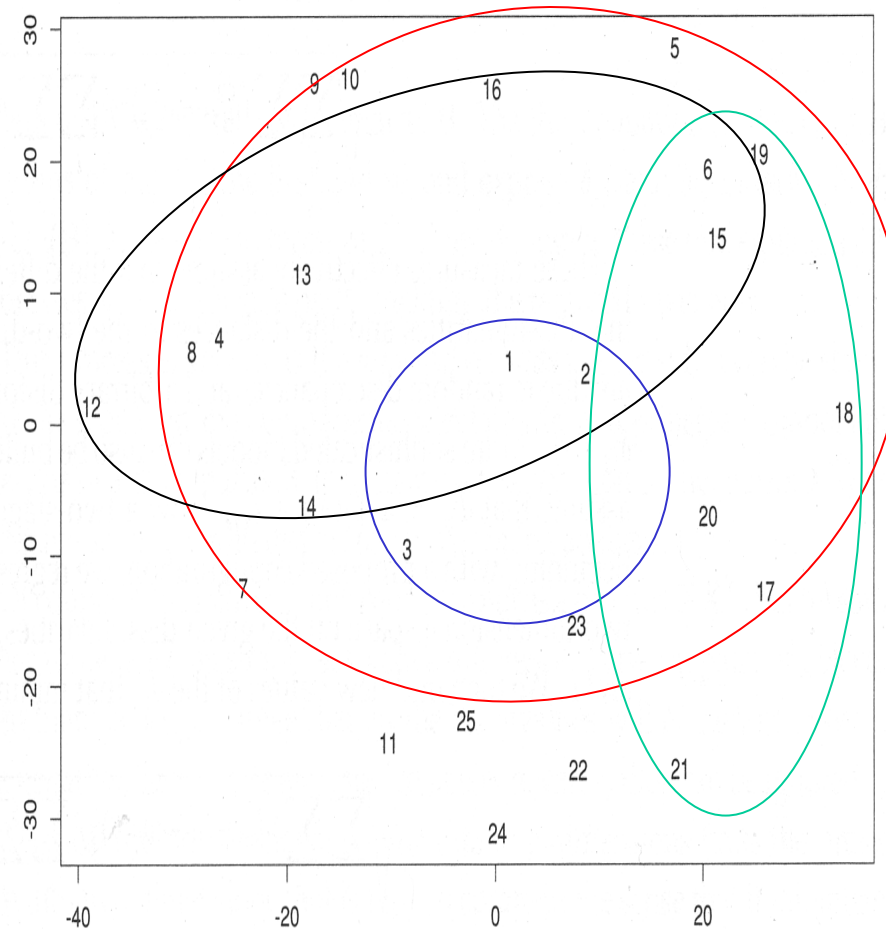
Multidimensional Scaling

- Detecting underlying structure
 - e.g., data lies on a string or surface in multi-dimensional space
- Distances are preserved
 - Distances between data points are mapped to a reduced space
- Typically displayed on a two-dimensional plot

Multidimensional Scaling Plot of Dialect Similarities

Numerical codes of villages and their counties

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

We are able to visualize 625 distances intuitively

Mathematical Formulation of Multidimensional Scaling

- For an $n \times p$ data matrix X we could compute $n \times n$ matrix $B = XX^t$ (*see next slide*)
 - Not practical since n is large
- The distance between the i^{th} and j^{th} objects is given by

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

X^tX versus XX^t

- If X is $n \times p$ $p=4$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix}$$

Covariance Matrix
- X^tX is $p \times p$ $p \times n$ $n \times p$ $p \times p$

↙

$$\begin{bmatrix} x_{11} & x_{21} & \cdot & x_{n1} \\ x_{12} & x_{22} & \cdot & x_{n2} \\ x_{13} & x_{23} & \cdot & x_{n3} \\ x_{14} & x_{24} & \cdot & x_{n4} \end{bmatrix} \times \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n x_{i1}^2 & \sum_{i=1}^n x_{i1}x_{i2} \\ \cdot & \cdot \\ \cdot & \cdot \\ \sum_{i=1}^n x_{in}^2 & \sum_{i=1}^n x_{in}x_{i1} \end{bmatrix}$$

- XX^t is $n \times n$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ \cdot & \cdot & \cdot & \cdot \\ x_{n1} & x_{n2} & x_{n3} & x_{n4} \end{bmatrix} \times \begin{bmatrix} x_{11} & x_{21} & \cdot & x_{n1} \\ x_{12} & x_{22} & \cdot & x_{n2} \\ x_{13} & x_{23} & \cdot & x_{n3} \\ x_{14} & x_{24} & \cdot & x_{n4} \end{bmatrix} = \begin{bmatrix} b_{11} = \sum_{j=1}^4 x_{1j}^2 & b_{12} \sum_{i=1}^4 x_{1j}x_{2j} & \cdot & b_{1n} \\ b_{21} & b_{22} & \cdot & b_{2n} \\ \cdot & \cdot & \cdot & \cdot \\ b_{41} & b_{42} & \cdot & b_{4n} \end{bmatrix}$$

B Matrix contains distance information

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

Multidimensional Scaling Procedure

- Given a matrix of distances D
 - (Derived from the original data by computing $n(n-1)/2$ distances)
 - we can compute the elements of B by inverting
$$d_{ij}^2 = b_{ii} + b_{jj} - 2b_{ij}$$
- We could factorize B
 - in terms of *eigen vectors* to yield coordinates of points
 - Two largest eigen values would give 2-D representation

Inverting distances to get B

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

- Summing over i

$$\sum_i d_{ij}^2 = tr(B) + nb_{jj}$$

Can obtain b_{jj}

- Summing over j

$$\sum_j d_{ij}^2 = tr(B) + nb_{ii}$$

Can obtain b_{ii}

- Summing over i and j

$$\sum_{i,j} d_{ij}^2 = 2ntr(B)$$

Can obtain $tr(B)$

Thus expressing b_{ij} as a function of d_{ij}^2

Method is known as *Principal Coordinates Method*

Classical Multidimensional Scaling into Two Dimensions

- Minimize

$$\sum_i \sum_j (\delta_{ij} - d_{ij})^2$$

Observed distance between points i and j in p -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


$$\frac{\sum_i \sum_j (\delta_{ij} - d_{ij})^2}{\sum_{i,j} d_{ij}^2}$$

or

$$\sqrt{\frac{\sum_i \sum_j (\delta_{ij} - d_{ij})^2}{\sum_{i,j} d_{ij}^2}}$$

Called stress

Algorithm for Multidimensional Scaling

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

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

Disadvantages of Multidimensional Scaling

- When there are too many data points structure becomes obscured
- Highly sophisticated transformations of the data (compared to scatter plots 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