



TutorialSli
de2

STATS5099: Data Mining

PC biplot, regression
and nonlinear dimension reduction

Xiaochen Yang
xiaochen.yang@glasgow.ac.uk

This week's content

- PCA biplots *loadings*
- multidimensional scaling
- principal component regression

2/15

Covariance or correlation-based PCA?

Prefer covariance-based PCA / original variables when:

- variables are comparable, e.g. gene expression from the same platform with similar range and scale
- differences in variation are of interest or meaningful themselves, e.g. rates in USArrests, stock data (high beta stocks will have higher loadings but they probably should)

Prefer correlation-based PCA / standardised variables when:

- difference in the variable variation caused by measurements on different scales, e.g. the trunk diameter in cm, biomass of leaves in kg, the number of branches, overall height in meters
- using correlation will throw out a tremendous amount of information → often need more correlation-PCs than covariance-PCs to retain an equivalent proportion of total variation

The "best" method to use is based on a subjective choice, careful thought and some experience.

<https://stats.stackexchange.com/questions/53/pca-on-correlation-or-covariance>

3/15

Multidimensional scaling (MDS)

- (often) a nonlinear dimension reduction method
- input: dissimilarity matrix $\delta_{ij}, i, j = 1, \dots, n$ *do not need original data*
- output: **configuration points** z_1, \dots, z_n
configuration distances d_{ij} *output configuration*
- objective: preserve pairwise distances δ_{ij} as well as possible

$$\min_{z_1, \dots, z_n} \sum_{i < j} L(\delta_{ij}, d_{ij})$$

The loss function is termed **stress** or **strain**. *disturbance*

4/15

Different types of MDS

- distance scaling vs inner product scaling *configuration distance*
 - distance: given dissimilarities, fit by $d_{ij} = \|z_i - z_j\|$
 - inner product: given similarities, fit by $\langle z_i, z_j \rangle$
- metric scaling vs nonmetric scaling
 - metric: δ_{ij} are interval or ratio, e.g. distance matrix (eurodist), inter-correlations among a set of variables (crimes)
 - nonmetric: δ_{ij} are ordinal; ranks matter more than actual values (wish) *input: ranking*

	distance scaling	inner product scaling
metric scaling	normalised stress Sammon mapping	classical scaling
nonmetric scaling	Stress-1	—

5/15

Implementing MDS and R commands

- compute the dissimilarity matrix
 - convert similarity to dissimilarity: `sim2diss(s, method)`
 - convert data matrix to dissimilarity: `dist`
- apply the appropriate MDS method
 - classical MDS: `cmdscale`
 - metric MDS:
`mds(delta, type=c("ratio", "interval"))`
 - Sammon mapping: `Sammon`
 - metric MDS: `mds(delta, type="ordinal")` *only rank*
- choose the number of dimensions: scree plot (stress vs no. dim) *elbow*
choose between different metric MDS methods: Shepard diagram (fitted distance vs dissimilarity)
- interpretation: configuration points z_i, z_j

6/15

sim2diss

`sim2diss(s, method=c("corr", "z"))`

Method (Argument)	Conversion Formula
Correlation ("corr")	$\delta_{ij} = \sqrt{1 - r_{ij}}$
Reverse ("reverse")	$\delta_{ij} = \min(s_{ij}) + \max(s_{ij}) - s_{ij}$
Reciprocal ("reciprocal")	$\delta_{ij} = 1/s_{ij}$
Membership ("membership")	$\delta_{ij} = 1 - s_{ij}$
Rank orders ("ranks")	$\delta_{ij} = \text{rank}(-s_{ij})$
Exponential ("exp")	$\delta_{ij} = -\log(s_{ij}/\max(s_{ij}))$
Gaussian ("Gaussian")	$\delta_{ij} = \sqrt{-\log(s_{ij}/\max(s_{ij}))}$
Transition frequencies ("transition")	$\delta_{ij} = 1/\sqrt{f_{ij}}$
Co-occurrences ("cooccurrence")	$\delta_{ij} = \left(1 + \frac{f_{ij} \sum_{k \neq i} f_{ik}}{\sum_{k \neq j} f_{kj} \sum_{l \neq j} f_{lj}}\right)^{-1}$
Gravity ("gravity")	$\delta_{ij} = \left(\frac{\sum_{k \neq i} f_{ik} \sum_{l \neq j} f_{lj}}{f_{ij} \sum_{k \neq i} f_{ik} \sum_{l \neq j} f_{lj}}\right)^{-1}$
Confusion proportions ("confusion")	$\delta_{ij} = 1 - p_{ij}$
Probabilities ("probability")	$\delta_{ij} = 1/\sqrt{\arcsin(p_{ij})}$
Integer value z	$\delta_{ij} = z - s_{ij}$

7/15

scree plot

- Look for an elbow in the plot

reduce the feature dimension to 2 *(n-1)*

8/15

Shepard diagram

how well is the conf. points?

metric

fitted

Interval MDS

fitted similar to diss.

non-metric

Ordinal MDS

only rank / step function

*f(sij) disparity
sij input (original)*

along close to the line

9/15

Q&As

MDS for dimension reduction

PCA linear

MDS non-linear

10/15

Tutorial question 1

This confusion matrix shows the frequency with which a stimulus question shown for some milliseconds only is confused with another letter:

Letter	B	C	D	F	G
B	—	—	—	—	—
C	3	—	—	—	—
D	7	5	—	—	—
F	3	5	7	—	—
G	7	12	2	2	—

Task: visualise the letters as points in 1/2/3D and interpret the plot

- Is this a dissimilarity matrix?
- Which MDS method is more appropriate?
- How many dimensions would you keep?
- Do you see any clusters?
- Investigate the sensitivity to different parameters.
- Investigate the sensitivity to initial configurations.

11/15

Principal component regression

Limitations of linear regression

- multicollinearity *(features strongly correlated)*
- high-dimensional problem, i.e. when $p > n$ *features > observations*

→ $X^T X$ is near-singular or singular (recall $\beta_{OLS} = (X^T X)^{-1} X^T y$)

Principal component regression: apply linear regression on PC scores *in higher dim case*

- when keeping all PCs, can get back the original regression coefficients from PC regression coefficients and loadings
- when keeping the first q PCs, reduce variance but increase bias *MSE ↓*

12/15

Bias-variance trade-off

Assume $Y = f(x) + \epsilon$ where $\mathbb{E}(\epsilon) = 0$ and $\text{Var}(\epsilon) = \sigma^2$. The prediction error of a regression fit $\hat{f}(x)$ at an input point $x = x_0$:

predictive error

$$\begin{aligned} \text{PE} &= \mathbb{E}[(Y - \hat{f}(x_0))^2 | x = x_0] \\ &= \sigma^2 + \mathbb{E}[(f(x_0) - \hat{f}(x_0))^2] \\ &= \text{Irreducible error} + \text{MSE} \\ &= \sigma^2 + (f(x_0) - \mathbb{E}[\hat{f}(x_0)])^2 + \mathbb{E}[(\hat{f}(x_0) - \mathbb{E}[\hat{f}(x_0)])^2] \\ &= \text{Irreducible error} + \text{Bias}^2 + \text{Variance} \end{aligned}$$

linear reg → PCA, var ↓ bias ↑

13/15

Bias-variance trade-off

Figure: f (black), \hat{f} : the linear regression line (orange), and two smoothing spline fits (blue and green).

Figure: MSE (red), squared bias (blue), variance (orange), σ^2 (dashed line).

14/15

Bias-variance trade-off

Figure: f (black), \hat{f} : the linear regression line (orange), and two smoothing spline fits (blue and green).

Figure: Training MSE (gray), test MSE (red), and minimum possible test MSE over all methods (dashed line).

14/15

Before leaving the tutorial ...

Check if you feel confident about the learning outcomes:

- Interpreting a biplot.
- Understanding different types of MDS.
- Performing MDS in R and interpreting its output.
- Understanding the bias-variance tradeoff.
- Understanding principal component regression and its advantages (over linear regression).

Quick homework: compare PCA and MDS – what are the differences, pros and cons?

15/15