

Revision
lecture...

University of Glasgow

STATS5099: Data Mining

Revision lecture

Xiaochen Yang
xiaochen.yang@glasgow.ac.uk

University of Glasgow

Data mining methods

- Dimension reduction
 - Principal component analysis (weeks 1-2)
 - Multidimensional scaling (week 2)
- Classification
 - K-nearest neighbours (week 3)
 - Linear discriminant analysis (week 3)
 - Classification trees, bagging, random forest (week 4)
 - Support vector machines (week 5)
 - Neural networks (week 6)
- Clustering
 - Agglomerative hierarchical clustering (week 8)
 - K-means clustering (week 9)
 - K-medoids clustering (week 9)
- Recommendation systems
 - Content-based filtering (week 10)
 - Collaborative filtering (week 10)

200

University of Glasgow

Revision suggestions

- For each method,
 - describe the method
 - explain the procedure of the method, e.g. by using a simple example
 - understand how to apply it in practice, e.g. any parameters
 - understand its applicability, strength and/or limitation
 - implement in R and interpret the results
- Across variants of one method / methods in the same category,
 - understand when to use one or the other

300

University of Glasgow

Principal component analysis

- describe the method
 - PCA is a linear dimensionality reduction method.
 - It finds a small number of uncorrelated linear combinations of the original variables which explain most of the variation in the data.

400

University of Glasgow

Principal component analysis

- explain the procedure of the method
 - The first principal component is the linear combination of original variables with the maximum variance, subject to the normalising constraint:
$$\text{first PC} = \text{linear combination } Y_1 \text{ that maximises } \text{Var}(Y_1) \text{ subject to } \mathbf{a}_1^T \mathbf{a}_1 = 1$$
 - The j^{th} PC is calculated in the same way (i.e. with the objective of maximising variance), with the condition that it is uncorrelated with previous PCs
$$j^{\text{th}} \text{ PC} = \text{linear combination } Y_j \text{ that maximises } \text{Var}(Y_j) \text{ subject to } \mathbf{a}_j^T \mathbf{a}_i = 1 \text{ and } \text{Cov}(Y_j, Y_i) = 0 \text{ for } i < j$$

500

University of Glasgow

Principal component analysis

- explain the procedure of the method
 - The first PC is the eigenvector associated with the largest eigenvalue of the covariance matrix Σ .
 - The j^{th} PC is the eigenvector associated with the j^{th} largest eigenvalue.

600

University of Glasgow

Principal component analysis

- understand how to apply it in practice
 - sample covariance matrix or sample correlation matrix?
 - how many PCs?
 - proportion of Variation
 - Cattell's method
 - Kaiser's method

800

University of Glasgow

Principal component analysis

- understand its applicability, strength and limitation
 - linear dimension reduction for continuous data
 - particularly useful when the original variables are highly correlated
 - no need to make assumptions on model or data distribution
 - has an analytical solution
 - can be severely distorted by outliers (first PC determined by outliers)

700

University of Glasgow

Principal component analysis

- implement in R and interpret the results
 - `princomp(data, cor=T)`
 - `summary(pca.result)`

```
Importance of components:  Comp.1  Comp.2  Comp.3  Comp.4  Comp.5
Standard deviation  2.1861060 1.5894522 1.1601299 0.9617072 0.9282642
Proportion of Variance 0.3676196 0.1946551 0.1056326 0.0714476 0.0668162
Cumulative Proportion 0.3676196 0.5622747 0.6680056 0.7387538 0.8033220

Comp.6  Comp.7  Comp.8  Comp.9  Comp.10
Standard deviation  0.8005460 0.7438816 0.5794282 0.4397897 0.5024590
Proportion of Variance 0.0492982 0.0450719 0.0282554 0.0224132 0.0294184
Cumulative Proportion 0.8523164 0.8948840 0.9207193 0.9431326 0.9625510

Comp.11  Comp.12  Comp.13
Standard deviation  0.4813977 0.4073767 0.2986301
Proportion of Variance 0.0178216 0.0127673 0.0086048
Cumulative Proportion 0.9603736 0.9931366 1.0000000
```

800

University of Glasgow

Principal component analysis

- implement in R and interpret the results
 - `princomp(data, cor=T)`
 - `summary(pca.result)`
 - `pca.result$loadings`
 - `biplot(pca.result)`

Score of each obs.

900

University of Glasgow

PCA vs MDS

	PCA	metric MDS	nonmetric MDS
data type	continuous variables	proximity data (similarity or dissimilarity matrix)	
data structure	data lies in a linear subspace	—	—
objective	maximise variance	preserve pairwise distances	
computation	analytical solution	iterative	iterative (longer time to converge)
solution	unique (up to sign flip)	sensitive to initial configurations	

900

University of Glasgow

Demo exam Q2

(a) To reduce the dimension of this data set, a researcher has applied principal component analysis (PCA) based on the correlation matrix. Suggest why PCA might have been run on the correlation matrix instead of the covariance matrix? [2 MARKS]

1000

University of Glasgow

Demo exam Q2

(b) Partial output from the principal component analysis is given below.

```
> glass.pca <- princomp(Glass,cor=TRUE)
> glass.pca
Standard deviations:
Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8 Comp.9
1.58 1.43 1.19 1.08 0.96 0.73 0.61 0.28 0.04
> glass.pca$loadings
Loadings:
Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8 Comp.9
RI 0.545 0.286 0.147 0.127 0.115 0.752 0.128 0.312
RA -0.258 0.270 -0.385 0.491 -0.154 -0.558 0.149 0.128 -0.312
Mg 0.111 -0.384 0.379 -0.124 0.306 -0.206 -0.699 0.274 -0.192
Al -0.429 0.295 0.329 -0.138 0.504 -0.274 -0.192 0.380 -0.298
Si -0.229 -0.150 -0.459 -0.653 0.216 0.380 -0.192 0.380 -0.298
K -0.219 -0.154 0.463 0.307 -0.244 0.504 0.110 -0.261 0.251
Ca 0.492 0.345 -0.276 0.188 -0.149 0.216 0.380 -0.192 0.380
Ba -0.250 0.485 -0.133 -0.251 0.657 0.352 -0.145 -0.198 0.579
Fe 0.186 0.284 -0.230 -0.873 -0.243 0.352 -0.145 -0.198 0.579
```

Comment on the loadings of the first principal component. Comment on the role of the variable Fe in the principal component analysis. [4 MARKS]

1000

University of Glasgow

Demo exam Q2

(c) The researcher chose to use the Proportion of Variation approach to determine the number of principal components to be retained. Based on the previous R output, decide how many components should be kept in order to explain 85% of the variability of the data set. [4 MARKS]

1000

University of Glasgow

Support vector machines

- describe the method
 - A classification method
 - The method directly tries to fit decision boundaries in such a way as to maximise the margin or separation between the two classes.
 - Three types: hard-margin linear SVM, soft-margin SVM, nonlinear SVM

1100

University of Glasgow

Support vector machines

- explain the procedure of the method (hard-margin SVM)
 - choose the hyperplane parameters to maximise the margin
 - solve the optimisation problem:
$$\min_{\mathbf{w}, b} \|\mathbf{w}\|$$

subject to $g(\mathbf{w}^T \mathbf{x} + b) \geq 1 \quad \forall i = 1, \dots, n$

1200

University of Glasgow

Support vector machines

- understand how to apply it in practice
 - scale the data before applying SVM
- soft-margin SVM
 - cost parameter *increase/decrease, penalty...*
- nonlinear SVM
 - kernel function *penalties in kernel evaluate the effect*

1300

University of Glasgow

Support vector machines

- understand its applicability, strength and/or limitation
 - designed for binary classification (one vs. all, one vs. one)
 - no model or data assumption
 - very flexible and suitable for nonlinear non-separable data
 - fast inference (even quite large)
 - sensitive to hyperparameters (cost, gamma, kernel)
 - no means of assessing uncertainty
 - low interpretability

1400

University of Glasgow

Support vector machines

- implement in R and interpret the results
 - `svm(Y~., data, type="C-classification", cost = C, kernel="radial", ...)`

```
## Self scaling
data(formula = sp ~ ., data = train.data, kernel = "linear",
cost = C.val, type = "C-classification")

Parameters:
SVM-type: C-classification
SVM-kernel: linear, kernel
gamma: 0.5
number of Support Vectors: 55
( 28 27 )
Cost: 0.5
( 28 27 )

Number of Classes: 2
Level: B 0
```

1500

University of Glasgow

Support vector machines

- implement in R and interpret the results
 - `svm(Y~., data, type="C-classification", cost = C, kernel="radial", ...)`
 - `tune.svm` or `tune` *have the pairs in SVM*

```
tune.svm(svm ~ ., variables = data[,c("C", "gamma", "kernel")], kernel="linear",
cost=C.val)
Parameter tuning of 'svm':
- sampling method: 10-fold cross validation
- best parameters:
cost
0.5
- best performance: 0
```

1600

University of Glasgow

Demo exam Q2

See the demo exam.

1600

University of Glasgow

Q&As

- weekly material tasks and tutorial question, which is similar to exam style question
- If we wanted to practice some more exam style questions, would there be another module (past or present) that we can search the past papers portal for, which would provide us with these?
A: multivariate methods (NB: Topics and lecture materials are different)
- an additional drop-in Q&A session at 12pm next Wednesday (28th July)

1700

University of Glasgow

Q&As

- LDA
- biplot
- Tutorial sheet 1 Q2
- Tutorial sheet 1 Q3

1800

University of Glasgow

Q&As: LDA

- Q: The notes say we have to evaluate each group to see which one has the largest LDF, but there is only one LDF for two groups, so what are we comparing the result of the equation with in order to choose which group to assign the data to?
- A: We will assign the object x to the class with the largest linear discriminant function (LDF)
$$LDF_g(x) = x^T \Sigma^{-1} \mu_g - \frac{1}{2} \mu_g^T \Sigma^{-1} \mu_g + \log \pi_g$$

1900

University of Glasgow

Q&As: LDA

Coefficients of linear discriminants:

LD1
balance 2.18221e-03
income 1.056657e-06

create a classification rule, e.g.
LD1 < 0 \Rightarrow group No

calculate the value of LD1 for a new sample x

2000

University of Glasgow

Q&As: LDA

Coefficients of linear discriminants:

LD1
balance 2.18221e-03
income 1.056657e-06

create a classification rule, e.g.
LD1 < 0 \Rightarrow group No

calculate the value of LD1 for a new sample x

2100

University of Glasgow

Q&As: Tutorial 1 Q3

- Show that principal components for standardised variables can be obtained from the eigenvectors of the correlation matrix R of the random vector X .
- Q: The model answer proves that $\text{Cov}(Z) = R$, but I don't understand how that answers the question, which is about the PCs coming from the eigenvalues
- A: PCs for standard variables are eigenvectors of the covariance matrix of the standard variables, which equals to the correlation matrix R of the original variables X

2200

University of Glasgow

Q&As: Tutorial 1 Q3

- Show that principal components for standardised variables can be obtained from the eigenvectors of the correlation matrix R of the random vector X .

Let's standardise original variables $X = (X_1, \dots, X_p)^T$ to $Z = (Z_1, \dots, Z_p)^T$, where

$$Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_i^2}}$$

Using a diagonal matrix $V^{1/2}$ for standard deviations $\sqrt{\sigma_i^2}$, we can write the standardisation in matrix notation

$$Z = (V^{1/2})^{-1}(X - \mu)$$

Clearly, $E(Z) = 0$

$$\text{Cov}(Z) = (V^{1/2})^{-1} \text{Cov}(X) (V^{1/2})^{-1} = (V^{1/2})^{-1} \Sigma (V^{1/2})^{-1} = R$$

2300