

# Predictive Modelling Tutorial 9: PCA

Dr Liew How Hui

Jan 2021

# Tutorial 7: PCA

When variances for features differ a lot,

$$\text{pca\$scale: } \sqrt{\frac{\sum_i (x_{ij} - \bar{x}_{.j})^2}{n-1}}$$

---

Original data:  $X$ ; Data shifted to centre:  $\tilde{X}$

$\text{pca\$center: } \bar{x}_{.j}$

$\text{pca\$sdev: } \sqrt{\lambda_i}$

$\text{pca\$rotation: } [e_1, e_2, \dots]$

$\text{pca\$x: } [\tilde{X}e_1, \tilde{X}e_2, \dots]$

# Tutorial 7

You are given 12 communities that were rated according to transportation and education — the higher the score the better. For example, a better transportation system will higher score. Higher education facilities will score higher as well. The table below shows the score for 12 communities in the two criteria:

# Tutorial 7 (cont)

Obs	Transportation	Education
1	1.1	0.5
2	3.9	3.2
3	1.5	1.6
4	5.6	7.5
5	2.5	1.8
6	7.3	7.2
7	1.4	1.9
8	6.1	7.0
9	1.5	1.6
10	5.1	6.6
11	1.8	2.0
12	7.5	8.0

# Tutorial 7 (cont)

- (a) Plot a scatterplot to visualize your data.
- (b) Generate two principal components for the data.
- (c) Choose one suitable principal component to represent the data.
- (d) Plot your data with the principal component you chose in (c).
- (e) With the eigenvalues computed in (b), calculate the proportion of variance explained by each component and the cumulative proportion.
- (f) With a targeted explained variation of 95%, how many principal components should be considered? State the total variation explained.

# FA May 2020 Q4 (a)

Given the following data with 8 observations in Table 4.1:

Table 4.1: Data with 2 features.

Obs	x	y
A	5.51	5.35
B	20.82	24.03
C	-0.77	-0.57
D	19.30	19.39
E	14.24	12.77
F	9.74	9.68
G	11.59	12.06
H	-6.08	-5.22

Find the first principle component and project the data (5.51, 5.35) to the space span by the first principal component. (4 marks)