

Homework 3

Sunday, April 30, 2017 8:33 PM

PCA

Logistics

- Submit your solutions on Canvas, this applies to each and all sections:
 - Code in one or more Matlab (.m) files
 - Report and figures as a single MS Word or pdf file
 - Put everything into a single ZIP file
- You must use Octave or Matlab code/instructions to obtain the results
 - Submit any scripts you used
 - You can use the code provided in the PCA in-class exercise
- Besides just finding the results, you must discuss your findings, explain the why's and hows, and give reasonable interpretations.
- Don't forget to include references in your report!
- List any assumptions when needed to disambiguate

Goal

Practice PCA and interpret results. In this homework we will use PCA to investigate potential clustering in the data.

Dataset

The following dataset is also commonly used to understand and explore the potential and application of PCA.



wine_beer_...

This is a csv file, but you can open it like a txt file.

The "Wine Beer" dataset, contains the following information regarding alcohol consumption in various countries and a couple of health-concerning metrics:

	Liquor consumption liter/year	wine consumption liter/year	beer consumption liter/year	life expectancy years	hear disease rate cases/10e5/year
France	2.5	63.5	40.1	78	61.1
Italy	0.9	58	25.1	78	94.1
Switzerland	1.7	46	65	78	106.4
Australia	1.2	15.7	102.1	78	173
Great Britain	1.5	12.2	100	77	199.7
United States	2	8.9	87.8	76	176
Russia	3.8	2.7	17.1	69	373.6
Czech Republic	1	1.7	140	73	283.7
Japan	2.1	1	55	79	34.7

Mexico	0.8	0.2	50.4	73	36.4
--------	-----	-----	------	----	------

There is only one class in this dataset. Each country is an observation, the columns are the features. We will use PCA to see if the data clusters in any way.

Part 1 Initial Statistics

As usual, obtain relevant statistics: means and standard deviations for each of the features.

Obtain covariance matrix

Do a summary of 2D scatter plots like those in Homework 2 Part 0 using the original data points (observations)

Show and discuss results in your report.

Part 2 Preprocessing

- Preprocess your data like in Homework 2 (mean center and scale).
- Discuss whether there are any "useless" or "noisy" features

Part 3 SVD

Perform SVD on the preprocessed data. Do not use the built-in PCA function, use the SVD function

Obtain the Scree plots and discuss the results: identify relevant principal components (dimensions)

Obtain and plot loading vectors, discuss results: what features are positively/negatively correlated and why? Are there any redundant features and why? What features are relevant and which ones seem irrelevant and why?

Based on the relevant principal components (new dimensions) identified with the help of the scree plot, create 2D scatter plots and discuss clusterings and location of the data points in the original and transformed spaces

Rubric

Criteria	Points	Notes
Code	30p	<p>Code should be well commented and structured</p> <p>Must include all the following implementations:</p> <ul style="list-style-type: none"> • Center and scale: 2pt • Svd: 2pt • Scree plots • Loading Vectors • Scatter: before and after transformation <p>NOT: Must NOT use Matlab's pre-built PCA function</p>
Analysis	<p>Outstanding: 70pt</p> <p>Sufficient: 65pt</p> <p>Poor: 40pt</p> <p>Missing: 0pt</p>	<p>Include and discuss statistics and plots.</p> <ul style="list-style-type: none"> • The discussion must include interpretations for each of the following: • Scree Plots: e.g., discussion of how many components are useful • Loading Vectors: should have four plots (one per vector), discussion about relevant /irrelevant features, correlations, etc. • Scatter: 2D plots in summary form, for data before and after svd <p>Discuss your results, putting them in context of the results (plots and statistics)</p> <p>Focus on the WHY not just the What's.</p> <p>Remember, your interpretation of the results is what matters the most!</p>

		Points will be deducted for lousy redaction, spelling/grammar mistakes, typos -- professionalism is expected.
--	--	---