

Methods for Dimensionality Reduction

Ramón Roales-Welsch

11 Januar 2018

Contents

Introduction	2
Overview	2
Motivation, Goal and Selection Criteria	2
Multidimensional Scaling (MDS)	3
Classical MDS	3
Metric MDS	3
Non-metric MDS	3
Principal Components Analysis (PCA) - Slides 81-88	4
Projection to m -Dimensions	4
Maximizing Variance	4
Eliminate Correlation between Features	4
Feature Selection	5
Overview (Slide 89/90)	5
Outlier and Duplicates	5
Appendix: Technical Implementation	6
Data Sample	6
R	6
Python (sklearn - package)	6
Appendix: Code examples	7
Example: Table	7
Example: Plot	7
Appendix: Literature Sources	8
Appendix: Covariation Matrix Calculation (More formulas)	8

Introduction

Overview

1. Multidimensional Scaling (MDS)
 - Dimension reduction by using feature-distances between observations
2. Principal Components Analysis (PCA)
 - Reduce dimensions while maximizing variance of features
 - Slide 88: PCA is not applicable in Big Data (*FastMap* method is proposed.)
3. Feature Selection
 - Drop weak features and select features with strong explanatory power

Motivation, Goal and Selection Criteria

Wikipedia names three main motivations of dimensionality reduction: Reduce costs (time and storage), reduce multicollinearity and provide visualisation opportunities. In the following the goals and selection criteria for mentioned methods are discussed:

1. MDS
 - Goal: Visualize your dataset in a 2-dimensional space and to show the similarities between objects(observations).
 - Alternative to Factor Analysis: Focus on dissimilarities (distances) between objects rather than similarities between features (via correlation matrices).
2. PCA
 - Goal: Dimension reduction and eliminaitaion of correlation between features. (*Bishop 2006*: Dimensionality reduction, lossy data compression, feature extraction and data visualitation.)
 - Projects d features on m features with $m < d$.
3. Feature Selection
 - Goal: Reduce overfitting, improve generalization and speed up computation.
 - Either drop weak features or select strong features.

Multidimensional Scaling (MDS)

Classical MDS

To build up the 2-dimensional matrix $D = [d_{i,j}]$, you need to compute the single components of the features x and y . Therefore, you use the subsequent formula according to Wikipedia:

$$d_{i,j} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

Lecture Slide 76 gives the formula for $d - dimensions$:

$$d(x_i, x_j)^2 = \sum_{k=1}^d (x_{i,k} - x_{j,k})^2$$

Not found on slides In order to evaluate, whether the MDS provides a useful dimensionality reduction, one need to examine the Stress/Strain (referring to the information loss) of the reduction.

Metric MDS

$$Stress_D(x_1, x_2, \dots, x_N) = \left(\frac{\sum_{i,j} (d_{i,j} - \|x_i - x_j\|)^2}{\sum_{i,j} d_{i,j}^2} \right)^{1/2}$$

Non-metric MDS

$$Stress = \sqrt{\frac{\sum (f(x) - d)^2}{\sum d^2}}$$

Principal Components Analysis (PCA) - Slides 81-88

Projection to m -Dimensions

PCA projects $x \in \mathbb{R}^d$ to $z \in \mathbb{R}^m$. When $\|w\| = 1$ holds, the following transformation can be applied:

$$z = w^T x$$

Maximizing Variance

When we maximize the variance of our projection z and therefore minimize the information loss of the transformation, $\|w_1\| = w_1^T * w_1 = 1$ must hold. Otherwise $\|w_1\| \rightarrow \infty$ applies and the variables get inflated to ∞ . When we project x_i with $z_i = w_1^T x_i$, the mean $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$ of the projection then is $\bar{z} = w_1^T \bar{x}$. From this follows the variance definition, which we want to maximize (Slide 83):

$$\frac{1}{N} \sum_{i=1}^N (w_1^T x_i - w_1^T \bar{x})^2 = \frac{1}{N} \sum_{i=1}^N (w_1^T (x_i - \bar{x}))^2 = w_1^T \left[\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x}) * (x_i - \bar{x})^T \right] w_1 = w_1^T \Sigma w_1$$

Σ is the $m \times m$ covariance matrix and further referred to as $\Sigma^{(m)}$ to indicate the m -dimensions. Using the langrange multiplier $\lambda_1 \in \mathbb{R}_{\leq 0}$, we have the following optimization problem (Slide 84):

$$\underset{w_1}{\text{maximize}} \quad f(w_1) = w_1^T \Sigma w_1 + \lambda_1 (1 - w_1^T w_1)$$

$$f'(w_1) \rightarrow \Sigma w_1 - \lambda_1 w_1 \stackrel{!}{=} 0$$

Following Slide 85, two observations can be made:

1. With $\Sigma w_1 = \lambda_1 w_1$ or $w_1^T \Sigma w_1 = \lambda_1$, w_1 have to be an ‘eigenvector’ of Σ , which is scaling the projection Σw_1 .
2. The eigenvector with the highest eigenvalue of w_1 maximizes the variance of the projection.

Eliminate Correlation between Features

We seek an orthogonal relation between projected feature z_1 and z_j , which means the correlation between the features is equal to zero: $cor(z_1, z_2) = 0$. This lead to the condition $w_2^T w_1 = 0$ and the following optimization problem:

$$\underset{w_2}{\text{maximize}} \quad f(w_2) = w_1^T \Sigma w_1 + \lambda_2 (1 - w_1^T w_1) + \beta (0 - w_2^T w_1)$$

The first-order derivation is set to zero: $f'(w_2) \rightarrow 2\Sigma w_2 - 2\lambda_2 w_2 - \beta w_1 \stackrel{!}{=} 0$ When we multiply this term from the left with w_1 we obtain the following expression:

$$2w_1^T \Sigma w_2 - 2\lambda_2 w_1^T w_2 - \beta w_1^T w_1 \stackrel{!}{=} 0$$

As $w_1^T w_1 = 1$, then $\beta = 0$. We conditioned $w_1^T w_2 = 0$ and observe $w_1^T \Sigma w_2 = w_2^T \Sigma w_1 = w_1^T \lambda_1 w_2 = 0$.

Finally, we have a result, which is analogous to w_1 :

$$\Sigma w_2 = \lambda_2 w_2$$

Feature Selection

Overview (Slide 89/90)

1. Evaluate the ‘quality’ of a feature.
2. Find best subset of features using the quality criterium.
 - highest quality of features
 - Lowest multicollinearity in subset

Use greedy algorithm with either bottom-up or top-down approach:

- bottom-up: Start with no feature and add new feature with every iteration. Evaluation via an error function and repeated until the required dimension is achieved.
- top-down: Start with all feature and remove the least explanatory with each step under consideration of an error function and until the required dimension is achieved.

Outlier and Duplicates

Outlier (Slide 101/102)

- A boxplot can reveal outlier.
- Also mean and standard deviation (sd) can be used. Usually outlier can be identified by $mean \pm sd$.
- Use cluster, every observation, which does not belong to a cluster, is a outlier.
- Parametrized processes, for instance mixture of gaussians.

Duplicates (Slide 103)

Use distances to find duplicates or observations which are very similar (very small distance). Reduce the number of observations and only keep one of them.

Appendix: Technical Implementation

Data Sample

```
# Code for Data Sample
```

Mirror the code here. It is executed in a separate document.

R

Python (sklearn - package)

Eigenvalue and Eigenvector (Numpy)

You can calculate 'eigenvalue' and 'eigenvector' with the function *numpy.linalg.eig()*.

Useful function/package

```
sklearn.manifold.MDS(n_components=2, metric=True, n_init=4, max_iter=300, verbose=0, eps=0.001,  
n_jobs=1, random_state=None, dissimilarity='euclidean')
```

<http://scikit-learn.org/stable/modules/generated/sklearn.manifold.MDS.html>

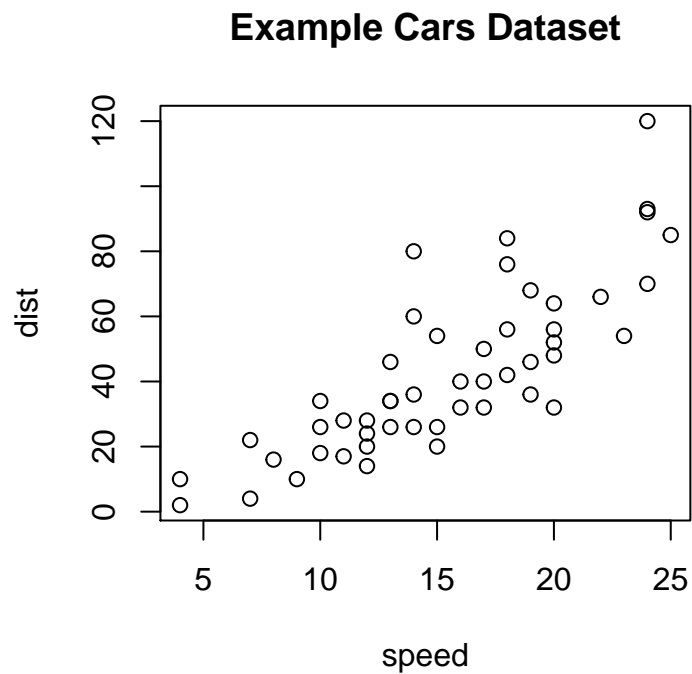
Appendix: Code examples

Example: Table

Table 1: Example Summary Cars Dataset

speed	dist
Min. : 4.0	Min. : 2.00
1st Qu.:12.0	1st Qu.: 26.00
Median :15.0	Median : 36.00
Mean :15.4	Mean : 42.98
3rd Qu.:19.0	3rd Qu.: 56.00
Max. :25.0	Max. :120.00

Example: Plot



Appendix: Literature Sources

- Lecture Slides from Machine Learning, Dr. Thomas Fober
- Bishop, C. (2007). *Pattern Recognition and Machine Learning* (Information Science and Statistics), 1st edn. 2006. corr. 2nd printing edn. Springer, New York.
- Wikipedia
 - https://en.wikipedia.org/wiki/Multidimensional_scaling
 - https://en.wikipedia.org/wiki/Principal_component_analysis
 - https://en.wikipedia.org/wiki/Feature_selection
- STUFF
 - STUFF

Appendix: Covariation Matrix Calculation (More formulas)

$$\Sigma^{(m)} = E \left[(X^{(m)} - E[X^{(m)}])(X^{(m)} - E[X^{(m)}])^T \right]$$

$$\Sigma^{(m)} = E \left[\left(X^{(m)} - \frac{1}{N} \sum_{i=1}^N X_i^{(m)} \right) \left(X^{(m)} - \frac{1}{N} \sum_{i=1}^N X_i^{(m)} \right)^T \right]$$

$$\Sigma^{(m)} = \sum_{i=1}^N \frac{1}{N} \left[(x_i^{(m)} - \bar{x}^{(m)})(x_i^{(m)} - \bar{x}^{(m)})^T \right]$$