

# Applications of Principle Component Analysis and Linear Discriminant Analysis on Consulting Survey Data

Ruinan(Victor) Zhang

March 24, 2017

## 1 Abstract

This paper address two statistical techniques to analyze surveying data and applications on how to visualize those data for business analytics.

## 2 Motivation

The project is inspired by company, Smari Inc. The challenge was to provide statistical evidence for consulting business decisions by visualizing collected survey data. One specific question is how to evaluate competitors' performance on specific features. For example, cellphone brands may interested to know what do the customers think of them in terms of design, user-friendness, system-smoothness. The goal of this project is to build an open-source online app for consulting companies to upload their dataset and construct visulizations to help them to answer such consulting questions.

## 3 Statistics Concepts Involved

### 3.1 Principal Component Analysis

Principal Component Analysis (PCA for the rest of the document) is a statistical technique to summarize features into algebraic combinations to capture the most correlation between the features and reponses. PCA is usually used to capture the most informations of provided features meanwhile reduce the number of features for modeling complexity.

Suppose I want to do PCA on a set of random vectors  $\vec{X}$ s with size equals to  $p$ . PCA generates a coefficient matrix  $E$  with size typically smaller than the length of  $\vec{X}$  (otherwise, PCA loses its goal of dimentional reduction). Let  $\Sigma$  denotes as variance-covariance matrix of  $\vec{X}$ . Then the eigenvalues of  $\Sigma$ :  $\lambda_1$  through  $\lambda_p$ . These eigenvalues are ordered so that  $\lambda_1$  is the largest eigenvalue and  $\lambda_p$  is the smallest. Then we also calculate the eigenvectors of  $\Sigma$ :  $\vec{e}_1$  through  $\vec{e}_p$ . It turns out the elements for these eigenvectors are the coefficients of principle components. Let the  $i$ th principle component denotes as  $Y_i$ , we can get variance of  $Y_i$  as following:

$$var(Y_i) = var(e_{i1}X_1 + \dots + e_{ip}X_p) = \lambda_i$$

$\lambda_1$  being the biggest eigenvalue indicates that  $Y_i$  captures the most variance. The proportion of variance caputured by  $i$  th principle components can be calculated as  $\lambda_i / \sum \lambda$ . With the coefficients, the  $i$ th principal components can be easily calculated as following:

$$Y_i = e_{i1}X_1 + \dots + e_{ip}X_p$$

In general, to visualize the PCA result on a 2-dimmmension graph, the first two principle components are chosen as the  $x$  and  $y$  axis, and coordinates to represent  $X_i$  can be found by solving the matrix  $Y = \vec{e} \times X$  with  $Y$  already known.

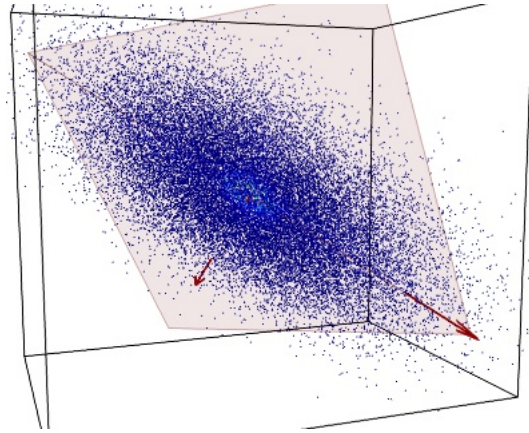


Figure 1: Summarize 3D data into 2 principle components

Figure 1 is a good visualization of summarizing 3-dimmmensional data into 2 principle components:

### 3.2 Linear Discriminant Analysis

Linear discriminant analysis(LDA) is a method used in statistics to find a linear combination of features that characterizes or separates two or more classes of objects or event. The idea of LDA is to find a linear combination of variables to model the difference between the classes of a particular set of variable (usually categorical variables). LDA technique is frequently used as dimensionality reduction before classifications. LDA can be visualized similar to PCA as they both generate set of linear coefficients.

### 3.3 PCA vs. LDA

PCA and LDA are similar in the way they both generate linear combination of variables which best explains data. LDA explicitly attempts to model the difference between the classes of data while PCA models the similarities between

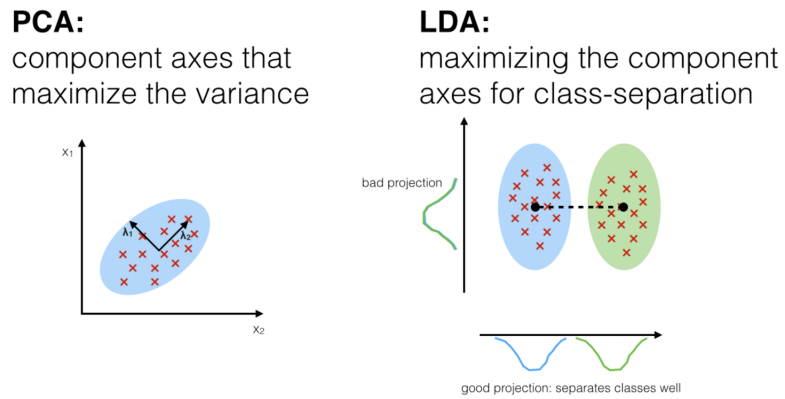


Figure 2: PCA vs. LDA

the classes. Here is a good visulization on difference between PCA and LDA

## 4 Online App on ShinyIO