# Dimensionality Reduction with Linear Transformations

project update

by
Mingyue Tan
March 17, 2004

---

## Domain and Task
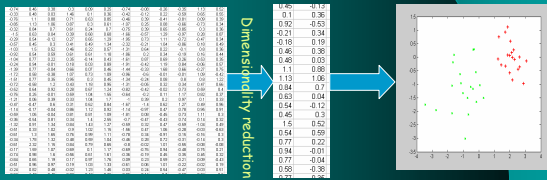
- Questions to answer
  - What's the shape of the clusters?
  - Which clusters are dense/heterogeneous?
  - Which data coordinates account for the decomposition to clusters?
  - Which data points are outliers?
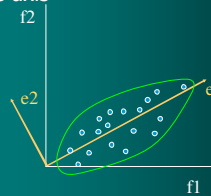
Data are labeled

---

## Solution - Dimension Reduction

1. Project the high-dimensional points in a low dimensional space while preserving the "essence" of the data
   - i.e. distances are preserved as well as possible

2. Solve the problems in low dimensions

Dimensionality reduction

---

## Principal Component Analysis

- Intuition: find the axis that shows the greatest variation, and project all points into this axis

f2

e2        e1

f1

---

## Problem with PCA

- Not robust - sensitive to outliers

- Usually does not show clustering structure

---

## New Approach

- PCA
  - seeks a projection that maximizes the sum

$$\sum_{i \neq j} \left( \mathrm{dist}_{ij}^{p} \right)^{2}$$

- Weighted PCA
  - seeks a projection that maximizes the weighted sum
  - flexibility

$$\sum_{i \neq j} w_{ij} \left( \mathrm{dist}_{ij}^{p} \right)^{2}$$

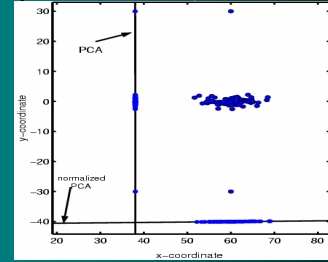Bigger $w_{ij}$ -> More important to put them apart

## Weighted PCA

Varying $w_{ij}$ gives:

- Weights specified by user
- Normalized PCA – robust towards outliers

$$\sum_{i \neq j} w_{ij}\left(\mathrm{dist}_{ij}^{p}\right)^{2} \qquad w_{ij} = \frac{1}{\mathrm{dist}_{ij}}$$
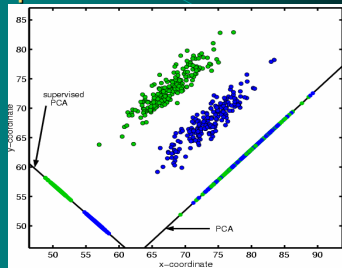
- Supervised PCA – shows cluster structures
  - If $i$ and $j$ belong to the same cluster ➜ set $w_{ij}=0$
  - *Maximize inter-cluster scatter*

---

## Comparison – with outliers



- PCA: Outliers typically govern the projection direction
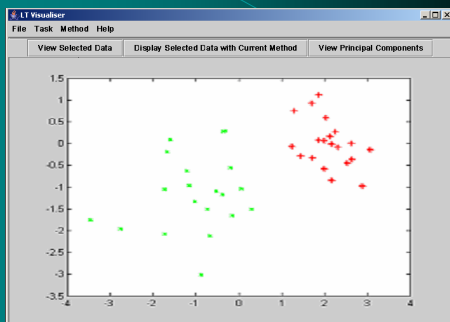
---

## Comparison – cluster structure



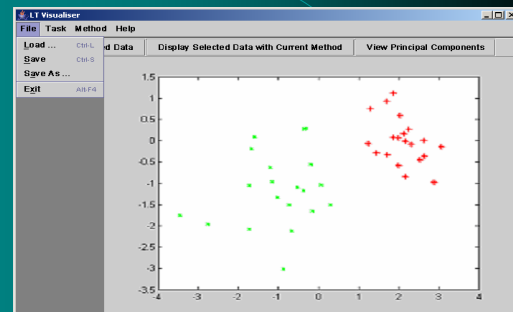- Projections that maximize scatter $\neq$ Projections that separate clusters

---

## Summary

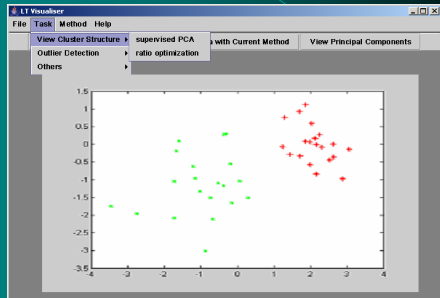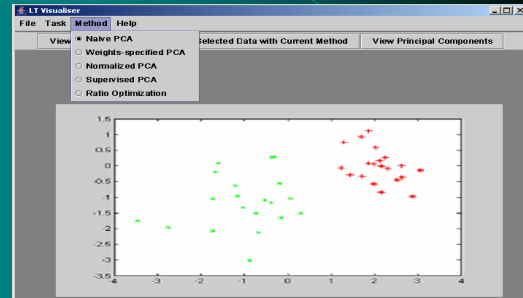| Method | Tasks |
|---|---|
| Naïve PCA | Outlier Detection |
| Weights-specified PCA | General view |
| Normalized PCA | Robustness towards Outliers |
| Supervised PCA | Cluster structure |
| Ratio optimization | Cluster structure (flexibility) |

---

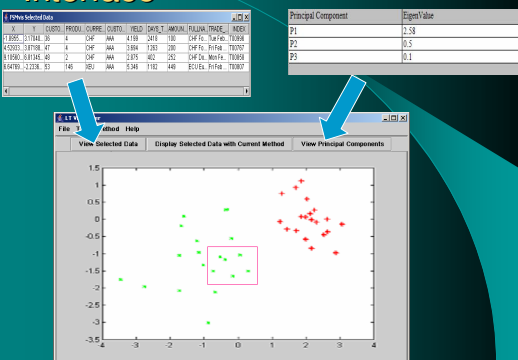## Interface



---

## Interface - File

**Interface - task**

**Interface - method**

**Interface**

## Milestones

- ✓ Dataset Assembled
- - same dataset used in the paper
- ✓ Get familiar with NetBeans
- - implemented preliminary interface (no functionality)
- ● Rewrite PCA in Java (from an existing Matlab implementation) – partially done
- ● Implement four new methods

## Reference

- ● [1] Y. Koren and L. Carmel, "Visualization of Labeled Data Using Linear Transformations", Proc. IEEE Information Visualization (InfoVis?3), IEEE, pp.121-128, 2003.