# PCA Revealed

## Part 5: Geometric Approach

**G**aston **S**anchez

August 2014

# Readme

**License:**

Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License
http://creativecommons.org/licenses/by-nc-sa/4.0/

**You are free to:**

**Share** — copy and redistribute the material

**Adapt** — rebuild and transform the material

**Under the following conditions:**

**Attribution** — You must give appropriate credit, provide a link to the license, and indicate if changes were made.

**NonCommercial** — You may not use this work for commercial purposes.

**Share Alike** — If you remix, transform, or build upon this work, you must distribute your contributions under the same license to this one.

## PCA

**Principal Components Analysis** (PCA) allows us to study and explore a set of quantitative variables measured on a set of objects.

## Core Idea

With PCA we seek to reduce the dimensionality (reduce the number of variables) of a data set while retaining as much as possible of the variation present in the data.

# Presentation

## About
In these slides we cover PCA from a geometric perspective.

## Working Principle
The underlying notion of this approach is that of Projected Inertia, and the intensive use of geometric principles.

## Visually Intended
Visualization plays a notable role in the Geometric Approach of PCA. One of the main reasons to reduce dimensions is to obtain graphical representations of data.

# Data considerations

# Data Structure

## Data

The analyzed data takes the form of a table (i.e. `matrix`) $\mathbf{X}$:

$$\mathbf{X}_{n,p} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

- $n$ objects in the rows
- $p$ quantitative variables in the columns

# Data Considerations

### Variables

The $p$ variables in $\mathbf{X}$ are denoted by $X_1, X_2, \ldots, X_p$

### Mean centered

For convenience, we will assume that the data is centered, i.e. Variables with mean $= 0$:

$$\bar{X}_j = \sum_{i=1}^{n} x_{ij} = 0$$
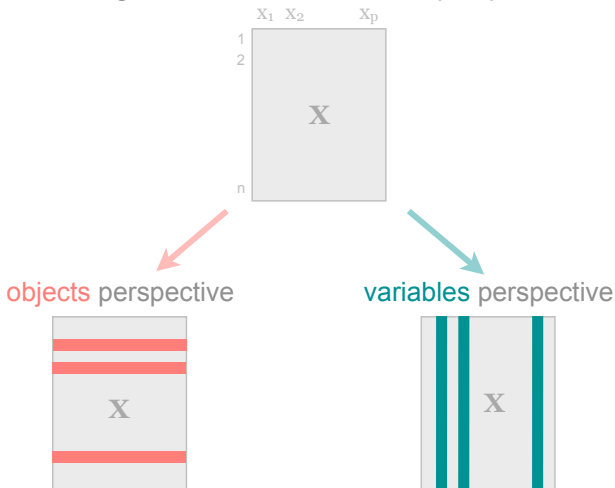
# Data from a geometric perspective

Geometric Frame of Mind

Looking at PCA from a geometric standpoint requires you to think about data in terms of points living in a multidimensional space —both objects and variables—

# Data Perspectives

looking at a data matrix from two perspectives



objects perspective

variables perspective

# Objects and Variables Perspectives

## Data Perspectives

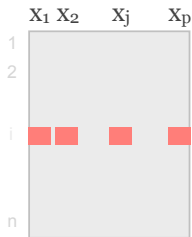We are interested in analyzing a data set from both perspectives: objects and variables

## Main Interests

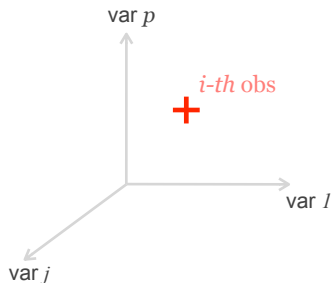At its simplest we are interested in 2 fundamental purposes:

► Study (dis)similarities among objects

► Study relationships among variables

# Objects in Multidimensional Space



each object described
by *p* variables
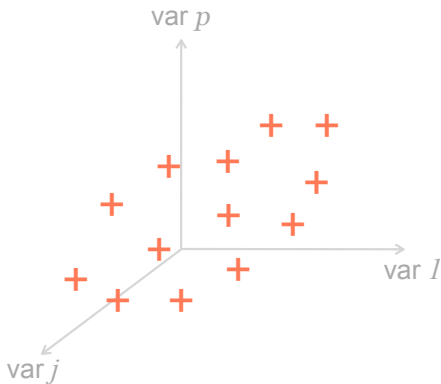
Associated
*p*-dimensional space

# *Cloud* of objects
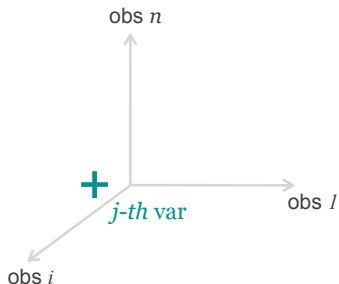
Objects as points in a *p*-dimensional space

# Variables in Multidimensional Space

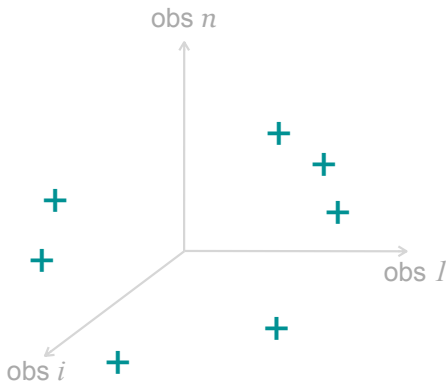each variable described
by *n* observations

Associated
*n*-dimensional space



$X_1$ $X_2$ $X_j$ $X_p$

1
2
i
n

obs *n*

*j-th* var

obs *l*

obs *i*

# *Cloud* of variables

Variables as points in a *n*-dimensional space
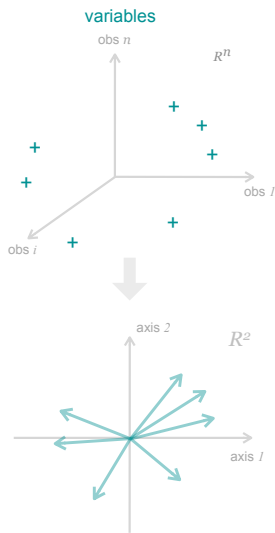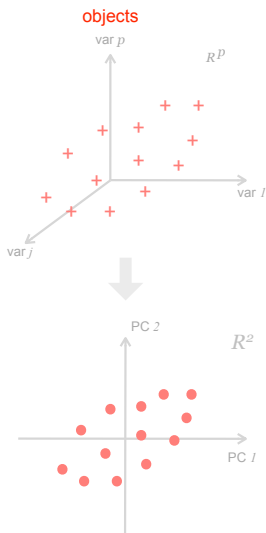
### PCA Visualization

We look for the "best" graphical representation that allows us to visualize the data in a low dimensional space (usually 2-dimensions).

# Best representation in low dimensional space



objects

var $p$

$R^p$

var $l$

var $j$

PC $2$

$R^2$

PC $1$

variables

obs $n$

$R^n$

obs $l$

obs $i$

axis $2$

$R^2$

axis $1$

# Low Dimensional Projections

### Geometric mindset

To help you understand the main idea of PCA from a geometric standpoint, I'd like to begin showing you my *mug-data* toy example.

### Key Message

The "name of the game" is **projection**: PCA involves projecting the data onto a low-dimensional space that best captures the original dispersion in the data.
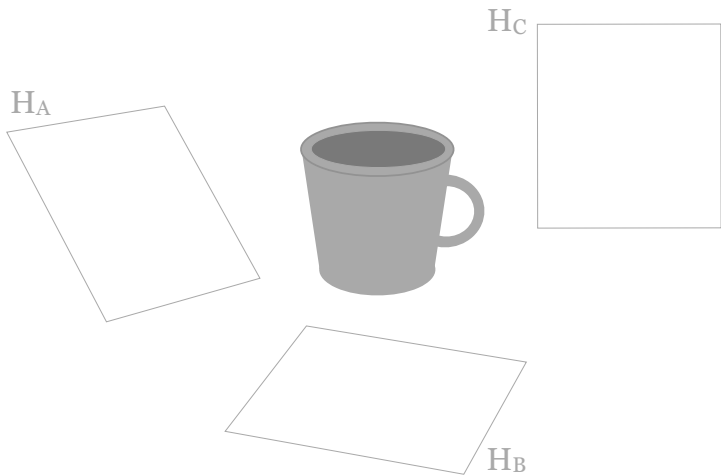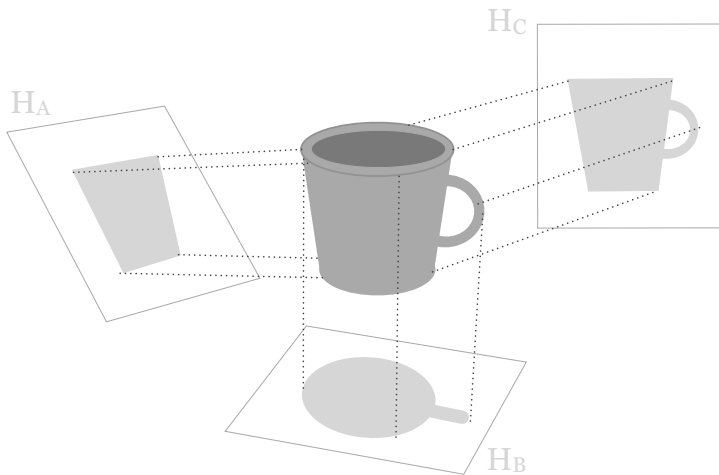
## Example

Imagine we have some data in a "high-dimensional space"

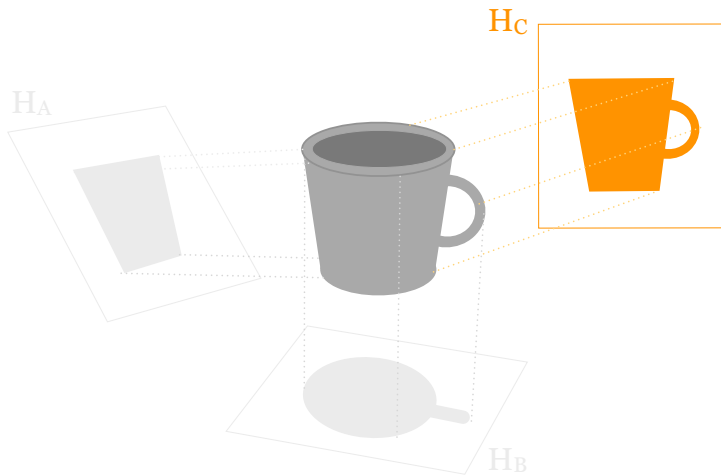# We are looking for Candidate Subspaces



H$_A$

H$_B$

H$_C$

# with the best low-dimensional representation



$H_A$

$H_B$

$H_C$

# Best low-dimensional projection



$H_C$

$H_A$

$H_B$

# Projections!!!

## Projection

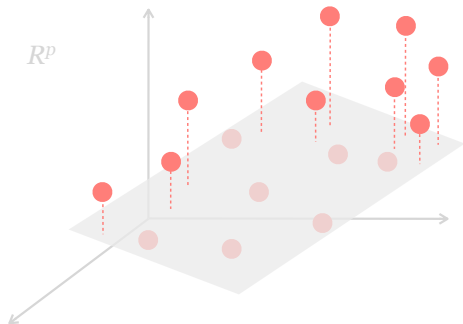We want to find a subspace that provides us the best projection of the data

$R^p$

# We look for a subspace such that



$R^p$

# the projection of points on it



$R^p$
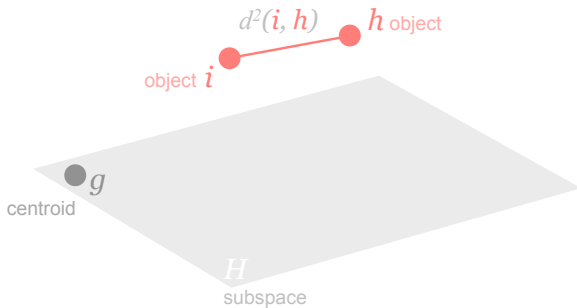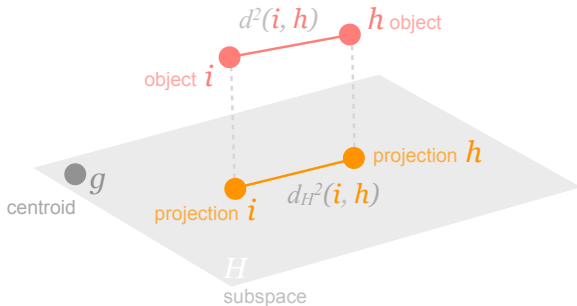
## Best Projection?

Looking for the best low-dimensional projection means that we want to find a subspace in which the projected distances among points are as much similar as possible to the original distances.

# Focus on distances between objects



$d^2(i, h)$

$h$ object

object $i$

$g$

centroid

$H$

subspace

# We want projected dists to preserve original dists

# Distances and Dispersion

## Dispersion of Data

Focusing on distances among all pairs of objects implicitly entails taking into account the **dispersion** (i.e. variation) of the data.

## Data Configuration

The reason to pay attention to distances and dispersion is to summarize in a quantitative way the original configuration of the data points.

## Pair-wise Square distances

One way to consider the dispersion of data (in a mathematical form) is by adding the square distances among all pairs of points.
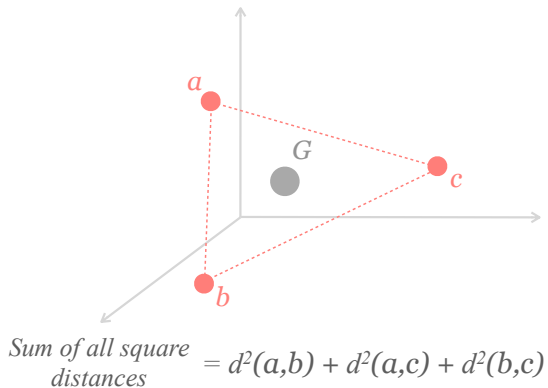
## Square distances from centroid

Another way to measure the dispersion of data is by considering the square distances of all points around the center of gravity (i.e. centroid)

# Imagine 3 points and its centroid

Sum of all square
distances $= d^2(a,b) + d^2(a,c) + d^2(b,c)$

*Sum of all square distances* $= 2d^2(a,G) + 2d^2(b,G) + 2d^2(c,G)$

# Inertia

## Inertia

To better take into account the dispersion of the data we must use the concept of **Inertia**.

## Idea

Simply put, we use the term Inertia to convey the idea of dispersion or *information* (variation) contained in the data.

## Moment of Inertia

Inertia is a term borrowed from the *moment of inertia* in mechanics.

## Inertia in Multivariate Methods

In multivariate methods, the term **Inertia generalizes the notion of variance**. Think of Inertia as a "multidimensional variance"

# Formula of Total Inertia

### Formula

The Total Inertia, $I$, is a weighted sum of square distances among all pairs of objects:

$$I = \frac{1}{2n^2} \sum_{i=1}^{n} \sum_{h=1}^{n} d^2(i, h)$$
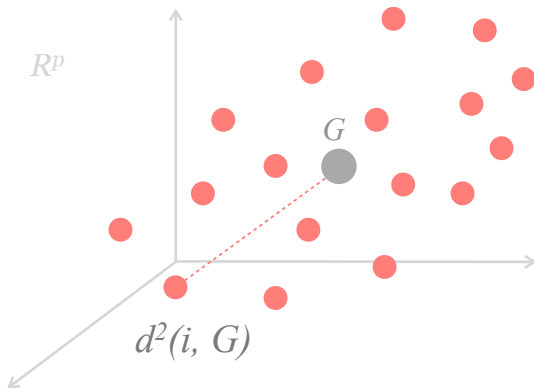
### Equivalent Formula

Equivalently, the Total Inertia can be calculated in terms of the center of gravity $G$:

$$I = \frac{1}{n} \sum_{i=1}^{n} d^2(i, G)$$

The Inertia is an average sum of square distances around the centroid $G$

# Data Points with their Centroid



$R^p$

$G$

$d^2(i, G)$

# Inertia around the Centroid



$R^p$

$G$

Inertia = $(1/n) \sum d^2(i, G)$

**Notation**

$x_{i.}$ $i$-th object $(i = 1, \ldots, n)$

$m_i$ mass of $i$-th object (usually $m_i = \frac{1}{n}$)

$G$ center of gravity (if data is mean-centered then $G = 0$)

$d^2(i, G)$ distance between $i$-th object and centroid $G$

# Computing Inertia

## Inertia Formula

$$Inertia = \sum_{i=1}^{n} m_i d^2(i, G) \tag{1}$$

$$= \sum_{i=1}^{n} \frac{1}{n} (x_{i.} - G)'(x_{i.} - G) \tag{2}$$

$$= \frac{1}{n} tr(X'X) \tag{3}$$
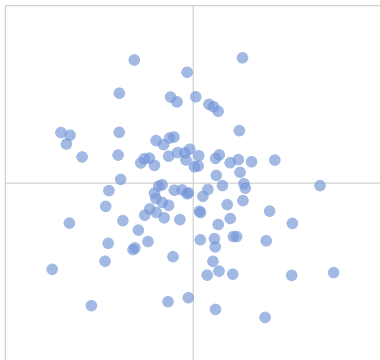
$$= \frac{1}{n} tr(XX') \tag{4}$$

## What's Important?

Two data sets can have the same inertia. The amount of dispersion is important, but it is also important the shape-form of that dispersion.

# Two data sets with similar inertia but different shape



Inertia =  2.02

Inertia =  2

# Projected Inertia and Dimension Reduction

# Inertia Concept

## Inertia and PCA

In PCA we look for a low-dimensional subspace having Projected Inertia as close as possible to the Original Inertia.

## Criterion

The criterion used for dimensionality reduction implies that the inertia of a cloud of points in the optimal subspace is maximum, but that would still be less than that in the true space.

# Subspace Preserving Information

### Global Criterion

We want the subspace that maximices the total dispersion, that is, the total inertia (maximize projected inertia)

### Interpretation

This means that we want to find the subspace that keeps the maximum amount of information of the original configuration

# Criterion

## Maximize Projected Inertia

We want to maximize the Projected Inertia on subspace $H$:

$$max \text{ projected} \sum_i d_H^2(i, G)$$

## Axis of Inertia

To find the subspace $H$ we can look for each of its dimensions (i.e. axes) $\mathbf{u_1}, \mathbf{u_2}, \ldots, \mathbf{u_k}$

# Projection 1

# Projection

# Projection



$$\text{Projection of obj } i \text{ on axis } \mathbf{u_1} = \boldsymbol{x_i}' \mathbf{u_1} = \sum_{j=1}^{p} x_{ij} \, u_{1j}$$

# Projection



axis
$u_1$

$R^p$

$G$

Projection of all points $= \mathbf{X}\, \mathbf{u_1}$

## PCs as Data Projections

It turns out that:

Looking for all the projections $Z_1 = \mathbf{Xu_1}, Z_2 = \mathbf{Xu_2}, \ldots, Z_k = \mathbf{Xu_k}$ implies looking for the Principal Components.

# Finding PCs

### How to find the 1st PC

In order to find the first principal component $Z_1 = \mathbf{X}\mathbf{u_1}$, we need to find $\mathbf{u_1}$ such that

$$\max_{\mathbf{u_1}} \ \mathbf{u_1'}\mathbf{X'}\mathbf{X}\mathbf{u_1}$$

subject to $\mathbf{u_1'}\mathbf{u_1} = 1$

### What to do?

Being a maximization problem, the typical procedure to find the solution is by using the Lagrangian multiplier method.

# Lagrangian Multiplier

### Finding 1st PC

Using Lagrange multipliers we get:

$$\mathbf{u_1'X'Xu_1} - \lambda(\mathbf{u_1'u_1} - 1)$$

Differentiation with respect to $\mathbf{u_1}$ gives:

$$\mathbf{X'Xu_1} - \lambda_1\mathbf{u_1} = \mathbf{0}$$

Rearranging some terms we get:

$$\mathbf{X'Xu_1} = \lambda_1\mathbf{u_1}$$

What does this mean?

$$\mathbf{X'Xu_1} = \lambda_1 \mathbf{u_1}$$

It means that

- $\lambda_1$ is an eigenvalue of $\mathbf{X'X}$
- and $\mathbf{u_1}$ is the corresponding eigenvector

### How to find the 2nd PC

In order to find the second principal component $Z_2 = \mathbf{X}\mathbf{u_2}$, we need to find $\mathbf{u_2}$ such that

$$\max_{\mathbf{u_2}} \ \mathbf{u_2'}\mathbf{X'}\mathbf{X}\mathbf{u_2}$$

subject to $\quad \|\mathbf{u_2}\| = 1 \quad$ and $\quad Z_1' Z_2 = 0$

# Finding 2nd PC

### Another eigenvalue-eigenvector pair

Applying the Lagrange multipliers, it can be shown that the desired $\mathbf{u_2}$ is such that

$$\mathbf{X'Xu_2} = \lambda_2 \mathbf{u_2}$$

### In other words

- $\lambda_2$ is an eigenvalue of $\mathbf{X'X}$
- and $\mathbf{u_2}$ is the corresponding eigenvector

# Finding all PCs

## Diagonalization

All Principal Components can be found simultaneously by **diagonalizing** $\mathbf{X'X}$

## Eigenvalue Decomposition (EVD)

Diagonalizing a matrix is nothing more than obtaining its eigenvalue decomposition (a.k.a. spectral decomposition)

# Variables Perspectives

## PCA model

A similar Inertia approach can be conceived for the variables. In this case, we look for vectors $\mathbf{v_j}$ such that the projections $\mathbf{X'v_j}$ have maximum inertia.

## It turns out that:

The solution consists in **diagonalizing $\mathbf{XX'}$**

# Algebraic Perspective

# Data Decomposition

## Algebraically

PCA involves a **Singular Value Decomposition** (SVD) of the data matrix $\mathbf{X}$, that is:

$$\mathbf{X} = \mathbf{UDV}'$$

- $\mathbf{U}$ is orthonormal (i.e. $\mathbf{U}'\mathbf{U} = \mathbf{I}$)
- $\mathbf{D}$ is a diagonal matrix
- $\mathbf{V}$ is orthonormal (i.e. $\mathbf{V}'\mathbf{V} = \mathbf{I}$)

# SVD Decomposition

## PCA by SVD

The data matrix $\mathbf{X}$ is decomposed as a product of three simpler matrices $\mathbf{U}$, $\mathbf{D}$ and $\mathbf{V}$

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}'$$

- $\mathbf{U}_{n,k}$ (information about the basic structure)
- $\mathbf{D}_{k,k}$ (information about the scale)
- $\mathbf{V}_{p,k}$ (information about orientation or correlations)

# SVD Approach

## SVD and PCA

The relationship between **SVD** and **PCA** is:

$$\mathbf{X} = \mathbf{UDV}' = \mathbf{ZP}'$$

where:

$\mathbf{Z} = \mathbf{UD}$ (PCs or scores)

$\mathbf{P} = \mathbf{V}$ (Loadings)

Note that:

$\mathbf{W} = \mathbf{V}$ since $\mathbf{XV} = \mathbf{UDV}'\mathbf{V}$

i.e. $\mathbf{Z} = \mathbf{XW} = \mathbf{XV}' = \mathbf{UD}$

# PCA and Data Decomposition

## Computation of all PCs

We can obtain as many PCs as the rank of $\mathbf{X}$ (i.e. $k = rank(\mathbf{X})$)

$$\mathbf{X}_{n,p} = \mathbf{Z}_{n,k}\mathbf{P'}_{k,p}$$

## Keeping just a few PCs

But usually we will only retain just a few PCs (i.e. $m \ll p$)

$$\mathbf{X}_{n,p} \approx \mathbf{Z}_{n,m}\mathbf{P'}_{m,p} + Residual$$

(just a few PCs will *optimally* summarize the main structure of the data)