

DATA MINING

Dimensionality Reduction

Principal Component Analysis

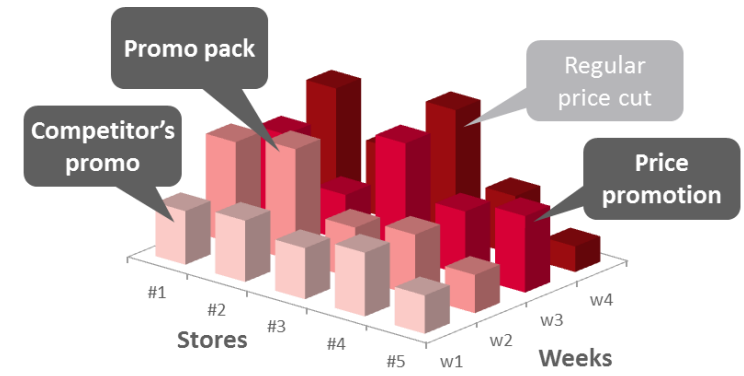
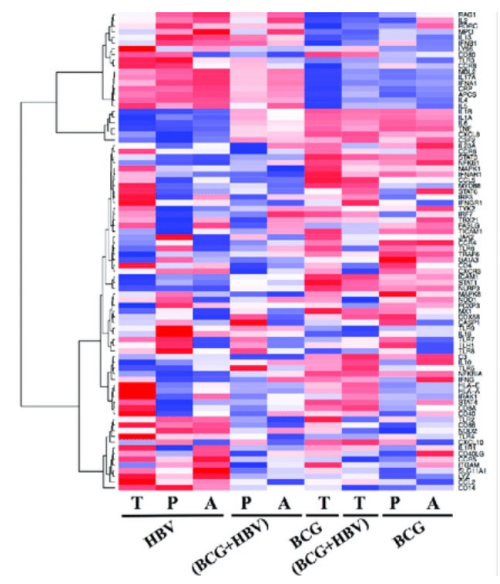
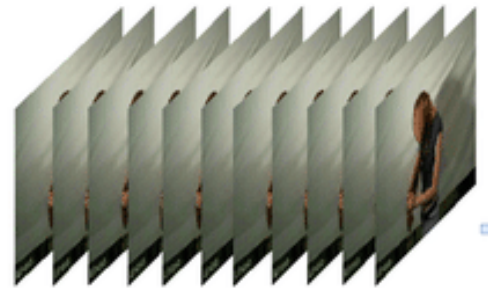
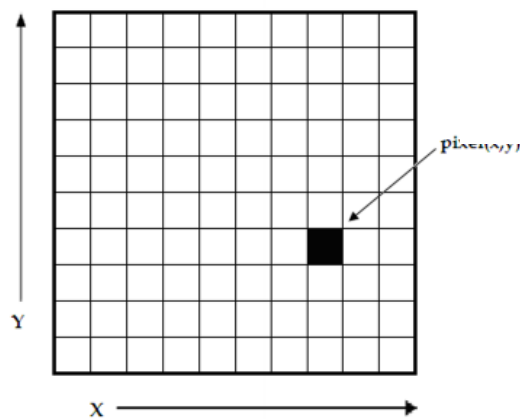
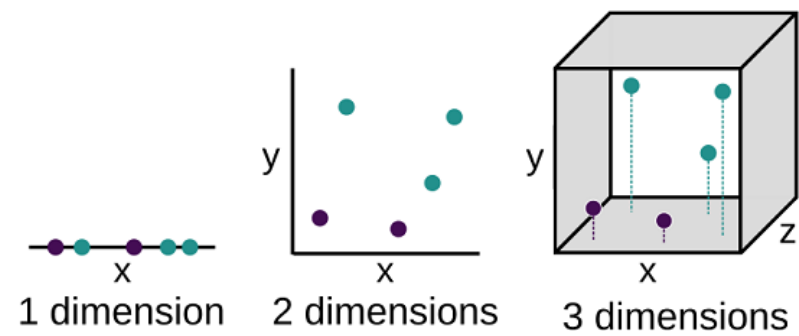
Ashish Pujari

Lecture Outline

- Dimensionality Reduction
- Linear Algebra Review
- Principal Component Analysis
- Linear Discriminant Analysis

DIMENSIONALITY REDUCTION

High Dimensional Data



High Dimensional Data

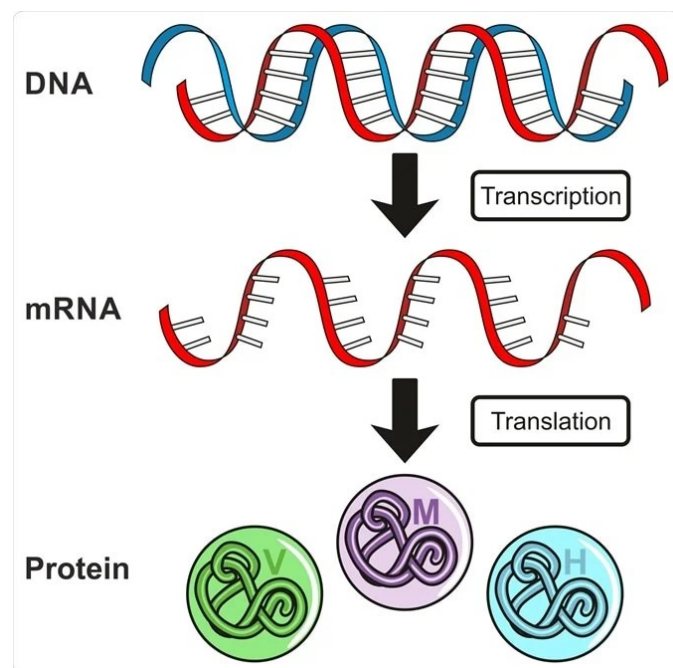
- In complex systems such as neuroscience, photo-science, meteorology and oceanography, etc., the number of variables to measure can be unwieldy and at times deceptive, because the underlying dynamics can often be quite simple.
- Dataset in which the number of features p is larger than the number of observations N , often written as $p \gg N$.

A diagram illustrating a dataset matrix. The matrix is a table with 3 rows and 6 columns. The columns are labeled $p_1, p_2, p_3, p_4, p_5, p_6$ and are grouped under the red label "features" with a red bracket. The rows are labeled n_1, n_2, n_3 and are grouped under the red label "observations" with a red bracket. The matrix is shown with a light gray background and black borders.

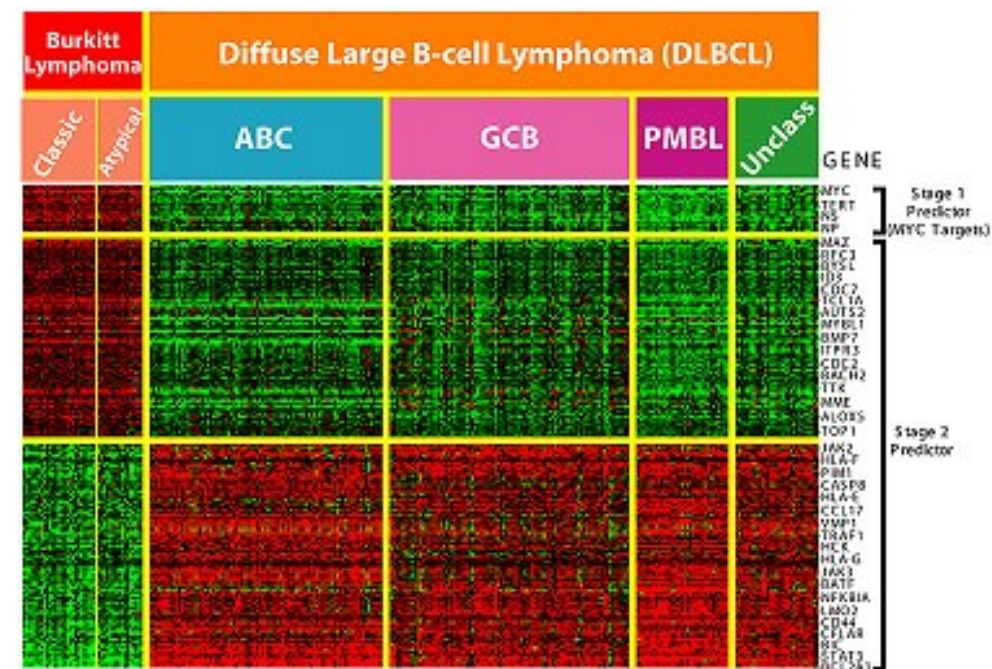
	p_1	p_2	p_3	p_4	p_5	p_6
n_1						
n_2						
n_3						

High Dimensional Data - Genomics

- Gene expression controls the amount and type of proteins that are expressed in a cell at any given point in time. It is essential for understanding protein function, biological pathways, and cellular responses to external and internal stimuli



Source: udaix/Shutterstock.com

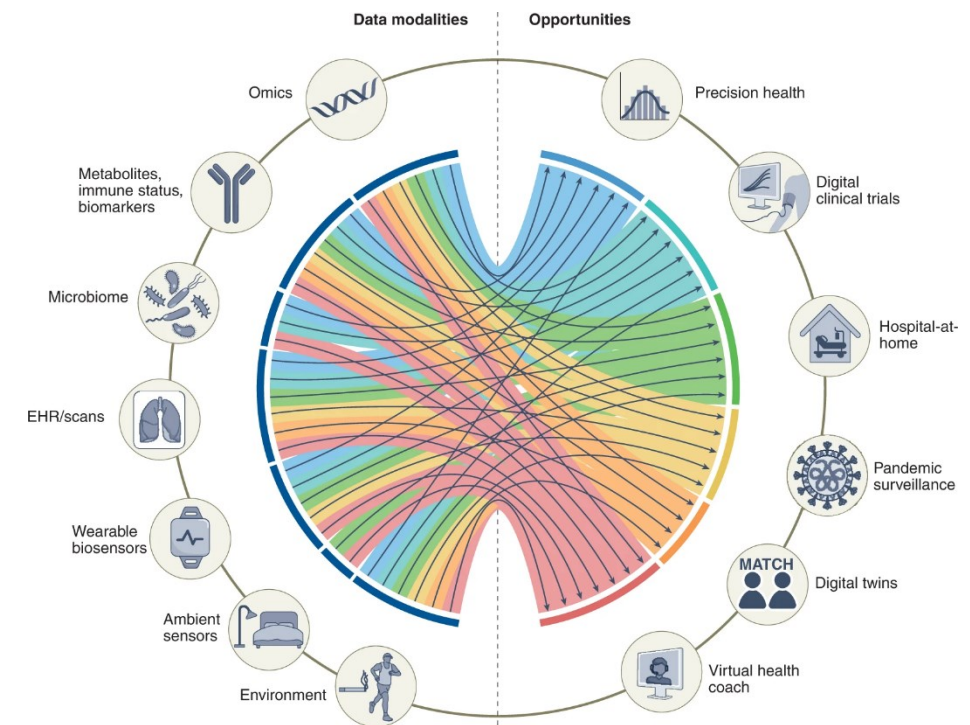
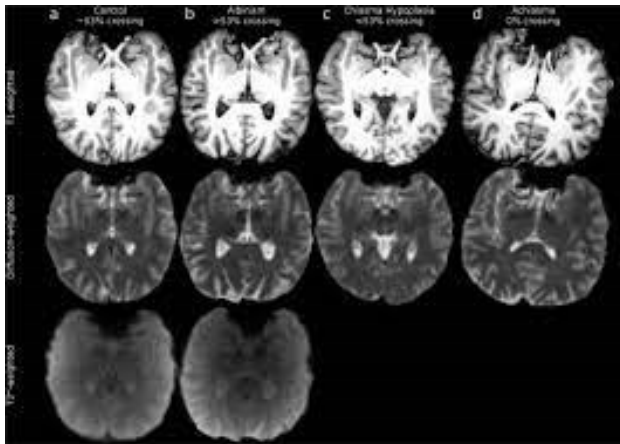


High-dimensional microarray analysis.

<http://archive.ics.uci.edu/ml/datasets/gene+expression+cancer+RNA-Seq>

High Dimensional Data - Healthcare

- High dimensional data is common in healthcare datasets where the number of features for a given individual can be massive
- E.g., MRI, blood pressure, resting heart rate, immune system status, surgery history, height, weight, existing conditions, etc.



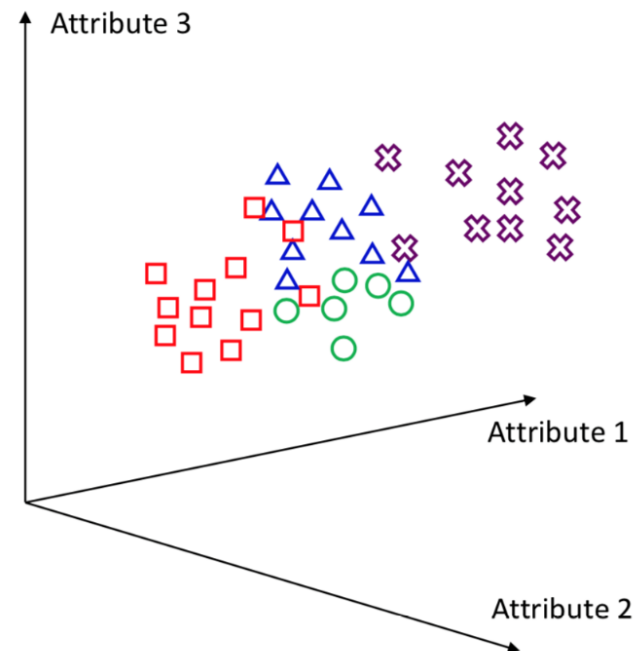
High Dimensional Data - Financial

- High dimensional data is also common in financial datasets where the number of features for a given stock can be quite large (i.e., PE Ratio, Market Cap, Trading Volume, Dividend Rate, etc.)



High Dimensional Data

- It is not easy to figure out the trend in 3D
- Is there a better representation of the data?
- How can we find lower dimensional representation that keeps the most information about the original data?



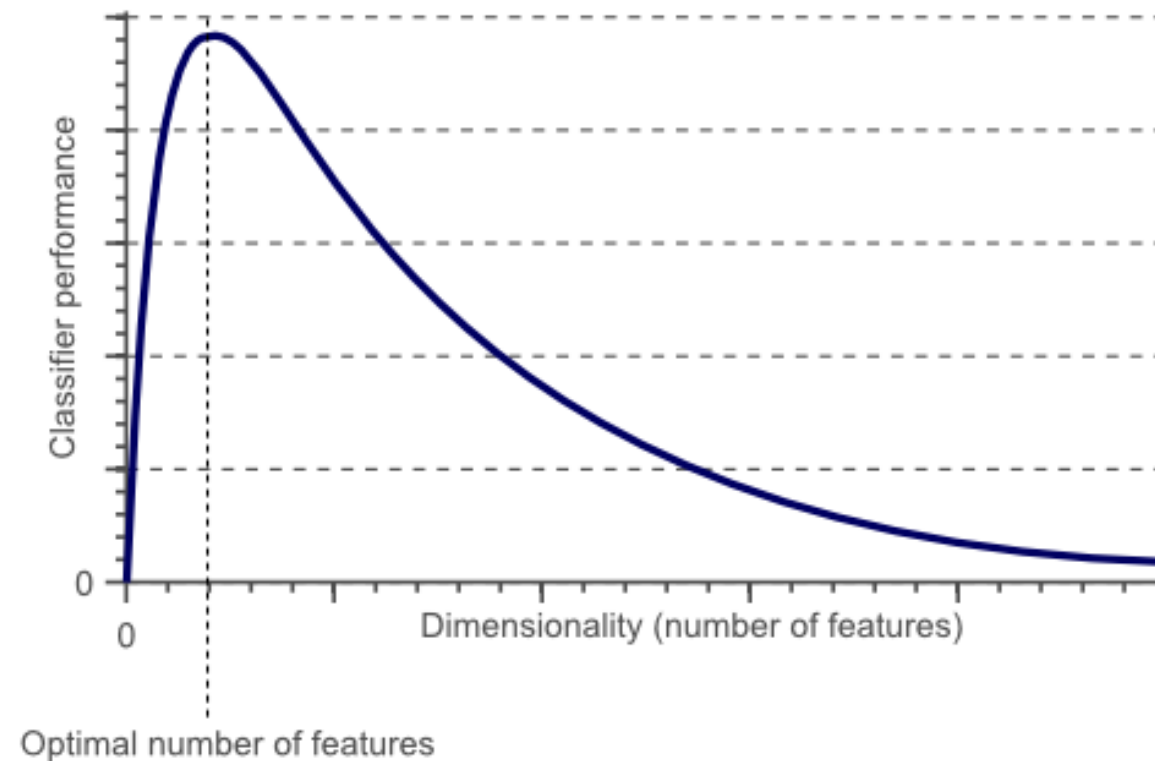
Curse of Dimensionality

As the number of dimensions in the data increases it impacts:

- Sparsity
 - Increases the volume of the feature space, therefore, data becomes more spread out (sparse)
- Sample Size
 - Required sample size n will grow exponentially with data that has d dimensions and might quickly become unmanageable.
- Metrics
 - Distance metrics lose their meaning in higher dimensions
- Performance
 - Machine Learning Model training time increases specially for parametric models.
 - Increase in 'noisy' or irrelevant features could cause degradation in model performance
 - Complex models that are harder to interpret than those with a low number of features.

Curse of Dimensionality

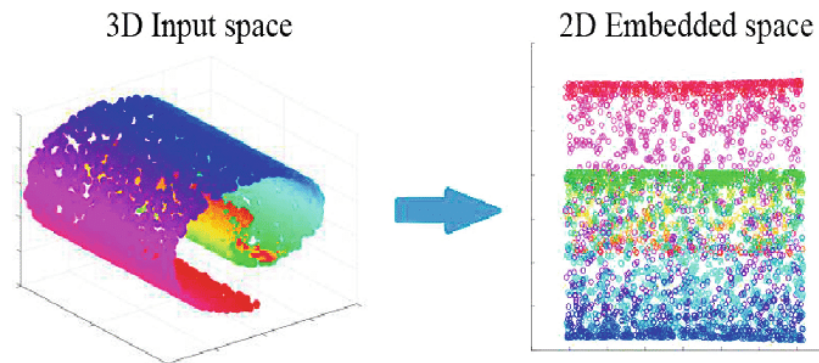
- As the number of features (dimensions) grows the amount of data we need to generalize grows exponentially



Dimensionality Reduction

- Goal
 - Map high dimensional data onto lower-dimensions in a manner that preserves distances / similarities.
 - Simplify understanding of data, either numerically or visually without loss of data integrity.
- Fundamental Idea
 - Exploit redundancy in the data to find a lower dimensional representation.

$$X = \{x_1, x_2, \dots, x_n, \in \mathbb{R}^D\} \rightarrow Y = \{y_1, y_2, \dots, y_n, \in \mathbb{R}^M\}$$



Data Compression

- Consider vectors

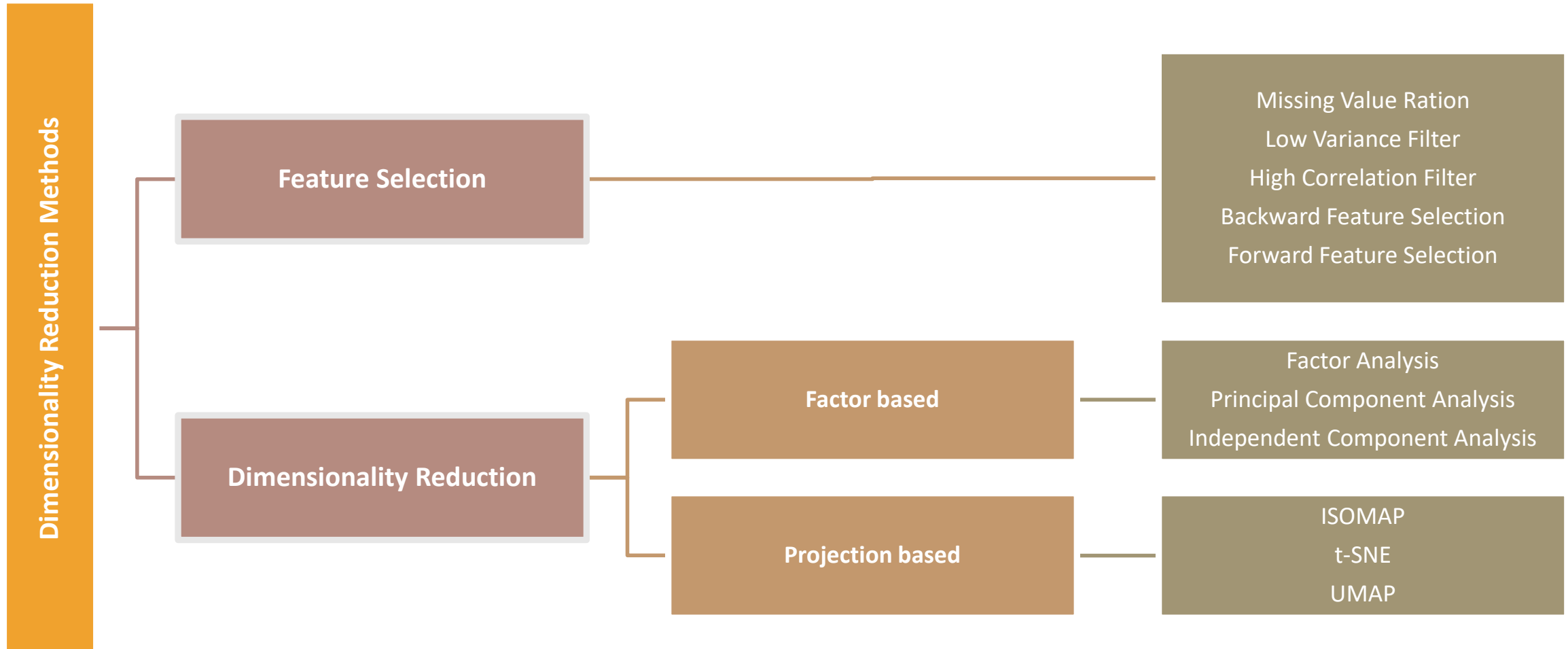
$$v_1 = (1,2,3), v_2 = (2,4,6), v_3 = (3,6,9), v_4 = (4,8,12)$$

- Need to store 12 integers
- However, they are all related

$$v_1 = 1 \cdot (1,2,3), v_2 = 2 \cdot (1,2,3), v_3 = 3 \cdot (1,2,3), v_4 = 4 \cdot (1,2,3)$$

- We can save $v_0 = (1,2,3)$ and multipliers 1,2,3,4
 - Only 7 integers need to be stored

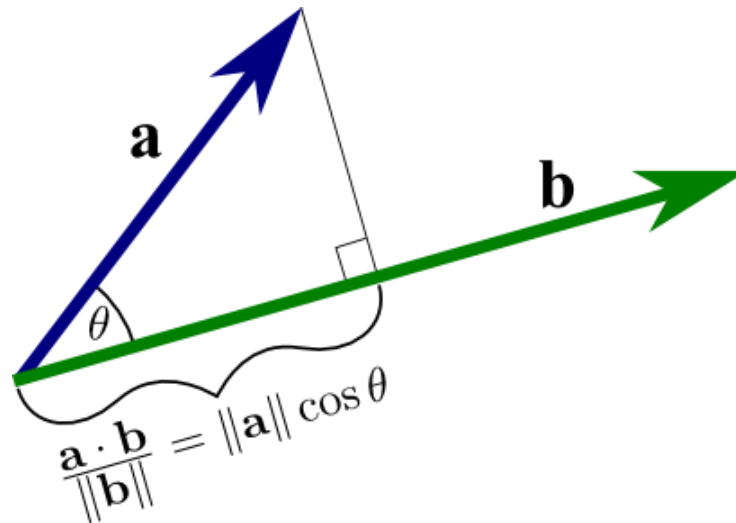
Dimensionality Reduction Methods



LINEAR ALGEBRA REVIEW

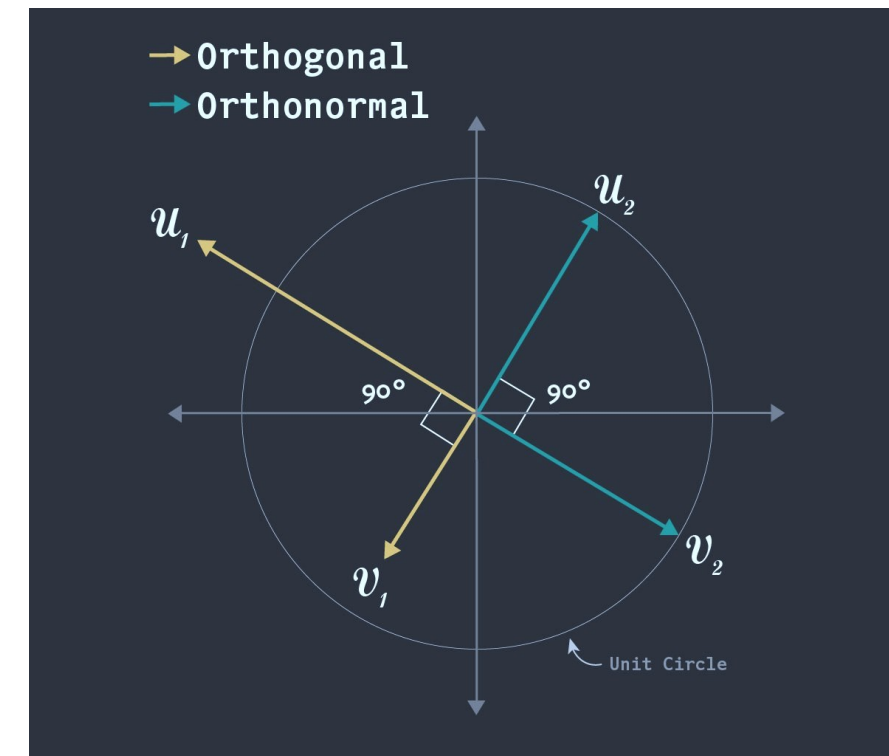
Dot Products

- Dot product between two vectors is based on the projection of one vector onto another
- Output of vector dot products :
 - The angle between the vectors is obtuse if the dot product is < 0
 - The angle between the vectors is acute if the dot product is > 0
 - The vectors are orthogonal (at right angles) if the dot product $= 0$



Orthogonality

- Two vectors u and v are considered to be **orthogonal** $u \perp v$ when the angle between them is 90° i.e. orthogonal vectors are perpendicular to each other.
- The dot product $u \cdot v = 0$
- They are **orthonormal** if they are orthogonal, and each vector has unit length.



Orthogonality (and orthonormality) is necessary to project vectors onto subspaces, find better estimates of nonlinear objects, and measure many properties of vector spaces.

Covariance Matrix

- Variance measures the variation of a single random variable

$$\sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

- Covariance is a measure of how much two random variables vary together

$$\sigma(x, y) = \frac{1}{n-1} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

- Using covariance we can calculate entries of the covariance matrix, which is a square matrix given by $C_{i,j} = \sigma(x_i, x_j)$ where our data set is expressed by the matrix $X \in \mathbb{R}^{n \times d}$

$$C = \frac{1}{n-1} \sum_{i=1}^N (X_i - \bar{X})(X_i - \bar{X})^T$$

Eigenvectors and Eigenvalues

- A vector x of dimension N is an eigenvector of a square $N \times N$ matrix A and λ is the corresponding eigenvalue if

$$Ax = \lambda x$$

Almost all vectors change direction, when they are multiplied by A . Eigenvectors are special unit vectors x that are in the same direction as Ax .

Eigenvalue λ tells whether the special vector x is stretched or shrunk or reversed or left unchanged when it is multiplied by A

$$Ax = \lambda Ix$$

$$Ax - \lambda Ix = 0$$

$$(A - \lambda I)x = 0$$

$$\det(A - \lambda I) = 0$$

Eigen-decomposition

- Eigen-decomposition is the factorization of a matrix into a canonical form, whereby the matrix is represented in terms of its eigenvalues and eigenvectors
- If a square matrix A is diagonalizable, then there is a matrix P such that

$$A = P D P^{-1}$$

Original Matrix	Eigenvectors Matrix	Eigenvalues Matrix	Inverse of Eigenvectors Matrix
$\begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}$	$= \begin{bmatrix} -1 & -1 \\ 2 & 1 \end{bmatrix}$	$\begin{bmatrix} -2 & 0 \\ 0 & -1 \end{bmatrix}$	$\begin{bmatrix} 1 & 1 \\ -2 & -1 \end{bmatrix}$

- The nondiagonal matrices P and P^{-1} are inverses of each other
- Provides valuable insights into its properties and makes matrix calculations (power, etc.) much easier

Eigen-decomposition

- Eigen-decomposition only exists for square matrices, and even among square matrices sometimes it doesn't exist
- When A is squared, the eigenvectors stay the same; the eigenvalues are squared

Matrix Type	Properties
Square Symmetric	Eigenvalues: always real, non-negative Eigenvectors: always orthogonal
Square Asymmetric	Eigenvalues: can be complex Eigenvectors: don't necessarily exist
Non-square	Eigen decomposition not possible

Singular Value Decomposition (SVD)

- The SVD of a matrix A is a factorization of that matrix into three matrices given by the formula :

$$A = U\Sigma V^T$$

- Vectors in the matrices U and V in the SVD are orthonormal and not necessarily the inverse of one another
- SVD can be used to compute optimal low-rank approximations of arbitrary matrices.
- SVD always exists for any rectangular or square matrix

Singular Value Decomposition (SVD)

The diagram illustrates the Singular Value Decomposition (SVD) of a matrix A . Matrix A is shown as a pink rectangle with dimensions $n \times d$. It is equal to the product of three matrices: U , Σ , and V^T . Matrix U is a pink rectangle with dimensions $n \times n$. Matrix Σ is a blue rectangle with dimensions $n \times d$, containing a pink sub-rectangle $\hat{\Sigma}$ of dimensions $r \times r$. Matrix V^T is a pink rectangle with dimensions $d \times d$, containing a pink sub-rectangle \hat{V}^T of dimensions $r \times d$.

$$\begin{matrix} \boxed{\begin{matrix} A \\ n \times d \end{matrix}} & = & \boxed{\begin{matrix} \hat{U} \\ n \times r \end{matrix}} & \boxed{\begin{matrix} \hat{\Sigma} \\ r \times r \end{matrix}} & \boxed{\begin{matrix} \hat{V}^T \\ r \times d \end{matrix}} \\ & & U & \Sigma & V^T \\ & & n \times n & n \times d & d \times d \end{matrix}$$

U : $n \times n$ matrix of the orthonormal eigenvectors of AA^T .

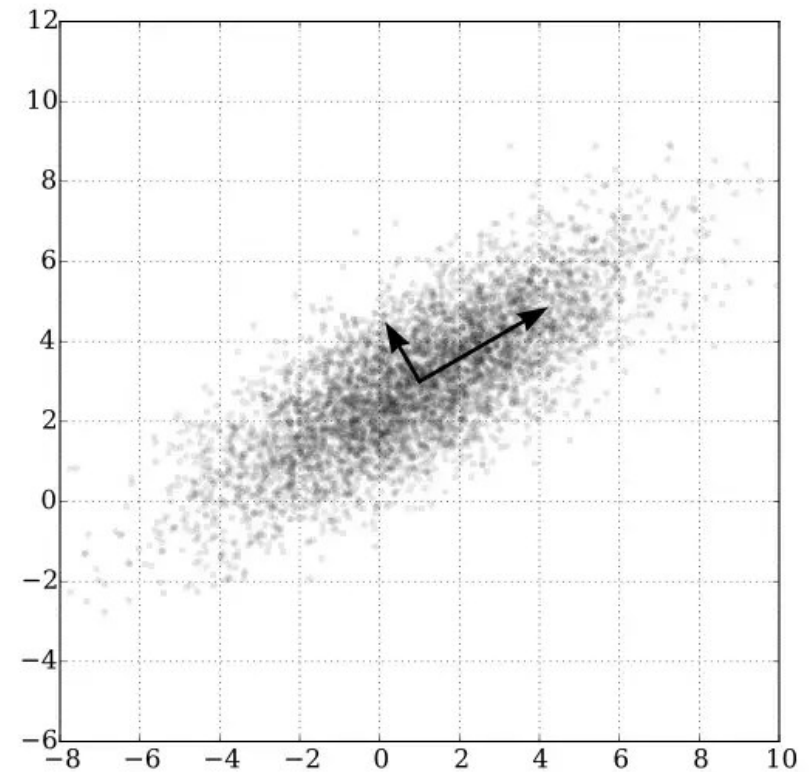
Σ : $n \times d$ diagonal matrix of the singular values of A which are the square roots of the eigenvalues of $A^T A$. The number of non-zero singular values is the rank of A

V^T : transpose of $d \times d$ matrix containing the orthonormal eigenvectors of $A^T A$

PRINCIPAL COMPONENT ANALYSIS (PCA)

Principal Component Analysis (PCA)

- What is PCA
 - PCA finds a lower-dimensional representation of data by constructing new features (Principal Components) which are linear combinations of the original features
- Objectives
 - To reduce the number of variables
 - Examine the relationship between variables
 - Address the problem of multicollinearity
- Assumptions
 - Original variables should be normalized
 - Factors are independent of each other
 - There exist some underlying factors that can describe the original variables



Principal Component Analysis (PCA)

X1	X2	X3



PC1	PC2	PC3

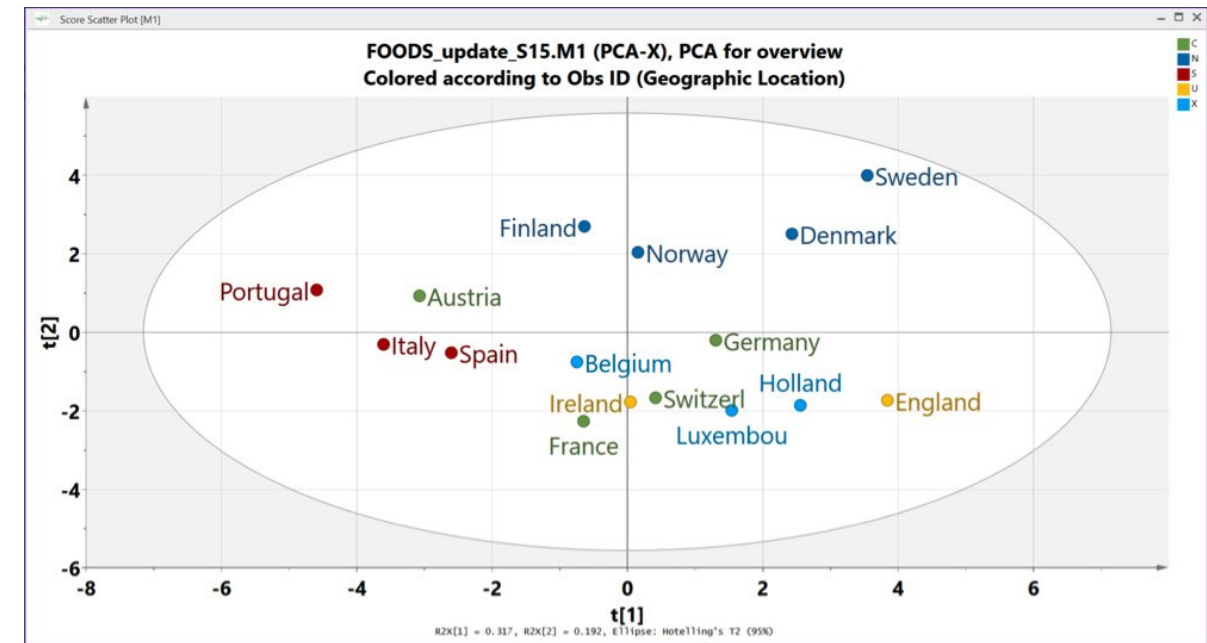
$$PC1 = a_1x_1 + a_2x_2 + a_3x_3$$

$$PC2 = b_1x_1 + b_2x_2 + b_3x_3$$

$$PC3 = c_1x_1 + c_2x_2 + c_3x_3$$

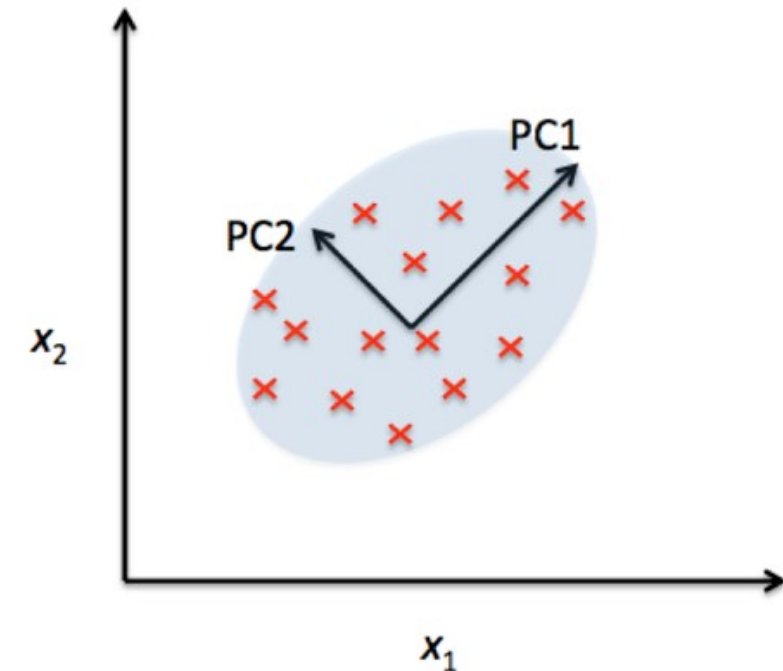
PCA Applications

- Data visualization
- Data compression (Lossy)
- Noise reduction
- Factor analysis
- Feature extraction
 - High dimensionality of the input features
 - Applied to data having multi-collinearity between the features/variables



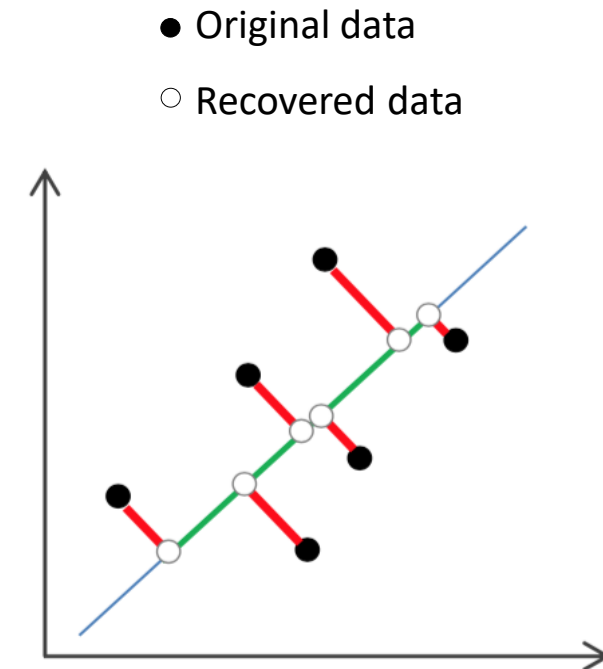
PCA Solution Approach

- PCA seeks a space of lower dimensionality known as principal subspace, such that the orthogonal projection of data points into this subspace maximizes the variance of the projected points
- This is done by projecting (dot product) the original data into the reduced PCA space using the eigenvectors of the covariance matrix (i.e Principal Components)
- The resulting projected data are linear combinations of the original data capturing most of the variance in the data
- The first component explains the most variance in the data with each subsequent component explaining less



PCA: Solution Approach

- Let's say we have an i.i.d dataset $X = \{x_1, x_2, \dots, x_n, \in \mathbb{R}^D\}$ with D dimensions and a mean value of 0
- Goal is to find projections that are as similar to the original data as possible but have lower dimensionality ($M < D$).
- There are two approaches:
 1. Maximum variance
 2. Minimum error



PCA: Solution Approach

$$\underbrace{\text{Variance of data}}_{\text{fixed}} = \underbrace{\text{captured variance}}_{\text{maximize}} + \underbrace{\text{reconstruction error}}_{\text{minimize}}$$

- Maximum variance formulation

- Find a low-dimensional representation which maximizes the variance of the projected data.

$$\max_{X \in \mathbb{R}^{m \times p}} \|AX\|^2 \text{ subject to } X^T X = I$$

- Minimum error formulation

- Find a low-dimensional representation which minimizes the average reconstruction error between the original data and the reconstructed data.

$$\min_{X \in \mathbb{R}^{m \times p}} \|A - AXX^T\|^2 \text{ subject to } X^T X = I$$

PCA: Maximum Variance

- Consider the simplest case $M = 1$. We define a vector $\mathbf{w}_1 \in \mathbf{R}^D$ as the direction of the lower dimensional space.
- Since we are only interested in the direction, we set \mathbf{w}_1 to be of unit length. i.e.

$$\mathbf{w}_1^T \mathbf{w}_1 = 1$$

- Then the data observations \mathbf{x}_n can be projected onto this new space as

$$\hat{\mathbf{x}}_n = \mathbf{w}_1^T \mathbf{x}_n$$

- If $\bar{\mathbf{x}}$ is the mean of the data observations in the original space, then the mean of the samples in the projected space is given by

$$\hat{\bar{\mathbf{x}}} = \mathbf{w}_1^T \bar{\mathbf{x}}$$

PCA: Maximum Variance

- Variance of the projected data can be derived as shown on the right:
- Where \mathbf{S} is the covariance matrix of the observed data in the original high dimensional space.

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T$$

$$\begin{aligned} \sigma^2(\hat{\mathbf{x}}) &= \frac{1}{N} \sum_{n=1}^N (\hat{\mathbf{x}}_n - \hat{\bar{\mathbf{x}}})^2 \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}_1^T \mathbf{x}_n - \mathbf{w}_1^T \bar{\mathbf{x}})^2 \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}_1^T \mathbf{x}_n - \mathbf{w}_1^T \bar{\mathbf{x}})(\mathbf{w}_1^T \mathbf{x}_n - \mathbf{w}_1^T \bar{\mathbf{x}})^T \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}_1^T \mathbf{x}_n - \mathbf{w}_1^T \bar{\mathbf{x}})(\mathbf{x}_n^T \mathbf{w}_1 - \bar{\mathbf{x}}^T \mathbf{w}_1) \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{w}_1^T (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n^T - \bar{\mathbf{x}}^T) \mathbf{w}_1 \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{w}_1^T (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{w}_1 \\ &= \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 \end{aligned}$$

PCA: Maximum Variance

- We need to maximize the variance of \hat{x} : $\sigma^2(\hat{x}) = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1$
- To prevent the trivial solution where $\|\mathbf{w}_1\| \rightarrow \infty$, we make use of the unit norm constraint we set earlier $\mathbf{w}_1^T \mathbf{w}_1 = 1$
- We introduce a Lagrange multiplier λ_1 and formulate our optimization objective as :

$$J(\mathbf{w}_1) = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 + \lambda_1 (1 - \mathbf{w}_1^T \mathbf{w}_1)$$

$$\frac{\partial J(\mathbf{w}_1)}{\partial \mathbf{w}_1} = 2\mathbf{S}\mathbf{w}_1 - 2\lambda_1 \mathbf{w}_1$$

- This shows that at the stationary point, \mathbf{w}_1 must be an eigenvector of \mathbf{S} and λ_1 the eigenvalue. corresponding to the eigenvector \mathbf{w}_1

$$\mathbf{S}\mathbf{w}_1 = \lambda_1 \mathbf{w}_1$$

PCA: Maximum Variance

- Left-multiplying with \mathbf{w}_1^T

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \mathbf{w}_1^T \lambda_1 \mathbf{w}_1$$

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \lambda_1 \mathbf{w}_1^T \mathbf{w}_1$$

$$\mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \lambda_1$$

- The maximum variance in the lower dimensional space is equal to the eigenvalue corresponding to eigenvector \mathbf{w}_1

PCA: Maximum Variance

- Identifying additional components
 - We can identify additional principal components by choosing directions that maximize variance while being orthogonal to the existing ones.

$$\mathbf{w}_2 \perp \mathbf{w}_1, \mathbf{w}_2^T \mathbf{w}_2 = 1$$

- General case of a lower dimensional space
 - with M dimensions with $M < D$, the principal components are the eigenvectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m$ corresponding to the M largest eigenvalues: $\lambda_1, \lambda_2, \dots, \lambda_m$

Relation between PCA and SVD

- Applying SVD to the data matrix \mathbf{X}

$$\begin{aligned}\mathbf{X} &= \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \\ \frac{1}{n-1}\mathbf{X}\mathbf{X}^T &= \frac{1}{n-1}(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T \\ &= \frac{1}{n-1}(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)(\mathbf{V}\mathbf{\Sigma}\mathbf{U}^T)\end{aligned}$$

- As \mathbf{V} is an orthogonal matrix $\mathbf{V}^T\mathbf{V} = \mathbf{I}$

$$\frac{1}{n-1}\mathbf{X}\mathbf{X}^T = \frac{1}{n-1}(\mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^T)$$

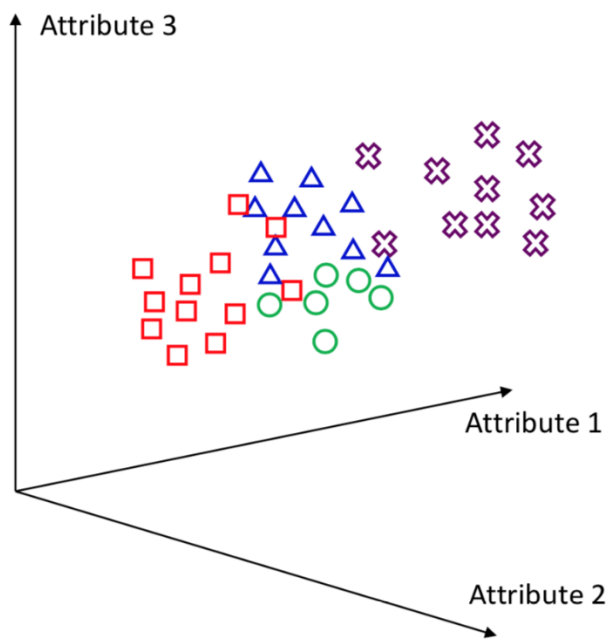
- The square roots of the eigenvalues of $\mathbf{X}\mathbf{X}^T$ are the singular values of \mathbf{X} .
- SVD can therefore be used to perform PCA without forming the covariance matrix as $\mathbf{X}\mathbf{X}^T$ can cause loss of numerical precision.

PCA: Steps

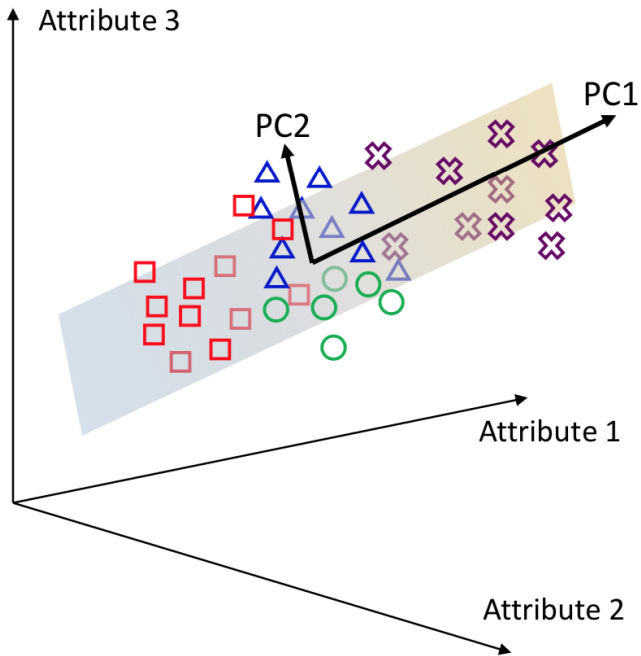
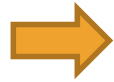
Let X be a matrix containing the original data with shape $[n_samples, n_features]$

1. Input variables of X are z-scored such each original variable has zero mean and unit standard deviation.
2. Construction and eigendecomposition of the covariance matrix. Covariance is equal to the correlation matrix for z-scored data.
3. Eigenvalues are then sorted in a decreasing order representing decreasing variance in the data (the eigenvalues are equal to the variance).
4. Finally, the projection of the original normalized data onto the reduced PCA space is obtained by multiplying (dot product) the originally normalized data by the leading eigenvectors of the covariance matrix i.e. the PCs.
5. To visualize the projected data as well as the contribution of the original variables, in a joint plot, we can use the biplot.

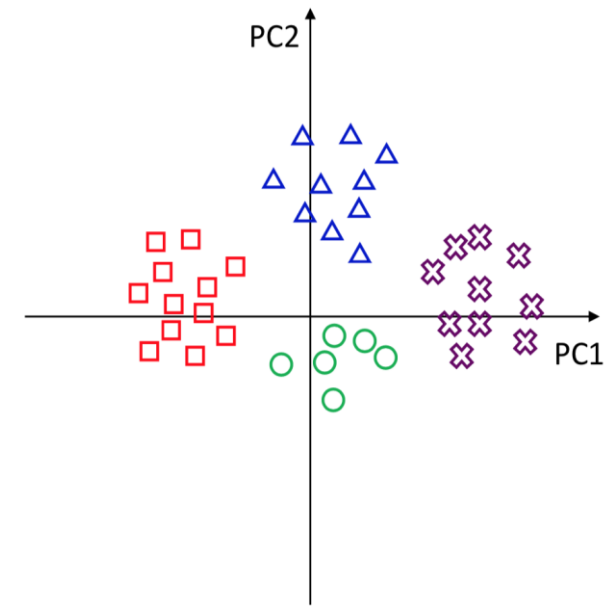
PCA: Data Visualization



1 Data in 3D



2 Principal Components



3 Data Visualization

Number of Principal Components

- For N original dimensions, sample covariance matrix is $N \times N$, and has up to N eigenvectors. So, N PCs.
- We can ignore the components of lesser significance
- Some information is lost, but if the eigenvalues are small, you don't lose much
 - N dimensions in original data
 - Calculate N eigenvectors and eigenvalues
 - Choose only the first D eigenvectors, based on their eigenvalues
 - Final data set has only D dimensions

Example 1

- Covariance matrix

$$\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

- Eigenvalue-eigenvector pairs are

$$\lambda_1 = 5.83, \quad x_1 = (0.383, -0.924, 0)$$

$$\lambda_2 = 2, \quad x_2 = (0, 0, 1)$$

$$\lambda_3 = 0.17, \quad x_3 = (0.924, 0.383, 0)$$

- $\lambda_1 > \lambda_2 > \lambda_3 \Rightarrow$ order of importance: x_1, x_2, x_3
- Dimensionality reduction
 - Pick eigenvectors (PC) with largest p eigenvalues
 - If $p = 1$, pick x_1
 - If $p = 2$, pick x_1 and x_2

Example 1 - How Many PCs?

- Comparison of recovered matrices with $p = 1, 2$ and original matrix

Original data

0.8	4.4	-0.9
5.8	12.4	6.1
-3.2	-14.6	-5.9
-6.2	-15.6	-1.9
-2.2	-8.6	2.1
1.8	8.4	0.1
4.8	8.4	5.1
1.8	13.4	6.1
-3.2	-12.6	-6.9
-4.2	-10.6	-7.9
2.8	19.4	6.1
-0.2	1.4	2.1
1.8	-5.6	-3.9

Recovered data with $p = 2$

1.2	4.3	-0.9
3.4	13.0	6.0
-3.8	-14.4	-5.9
-4.4	-16.1	-1.9
-2.4	-8.6	2.1
2.3	8.2	0.1
2.4	9.1	5.0
3.4	12.9	6.1
-3.3	-12.6	-6.9
-2.8	-11.0	-7.9
5.0	18.8	6.1
0.3	1.2	2.1
-1.2	-4.8	-4.0

Recovered data with $p = 1$

0.9	3.5	1.3
3.6	13.4	5.0
-3.9	-14.6	-5.5
-4.0	-14.9	-5.6
-1.9	-7.0	-2.6
2.0	7.3	2.7
2.5	9.5	3.6
3.6	13.3	5.0
-3.5	-13.3	-5.0
-3.2	-12.1	-4.5
4.9	18.5	6.9
0.5	1.7	0.6
-1.5	-5.5	-2.0

Example 1 - How Many PCs?

- Comparison of recovered matrices with $p = 1, 2$ and original matrix

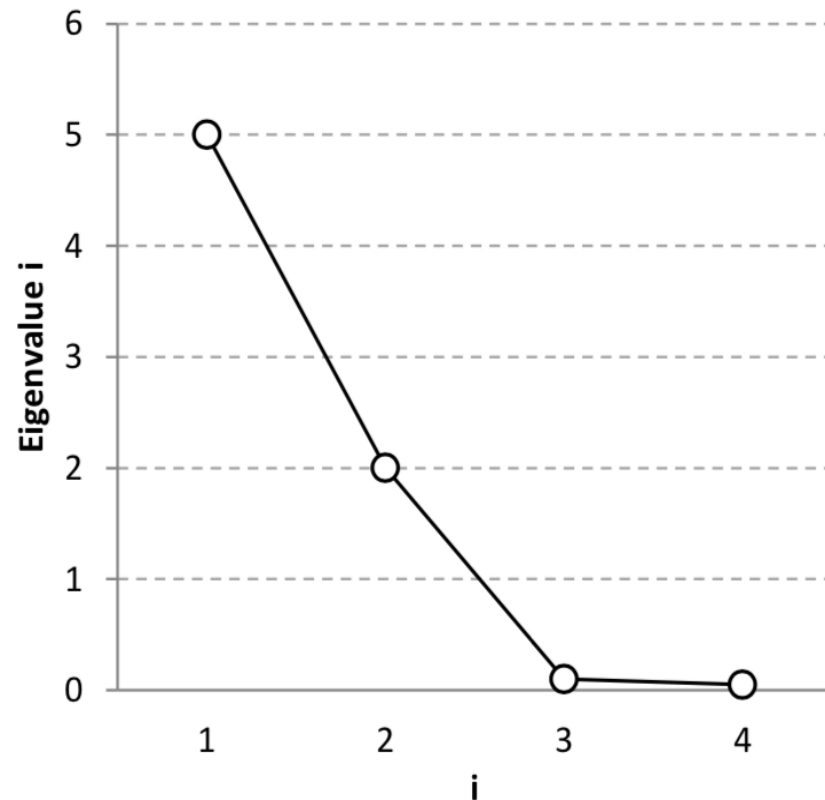
Original data			Error matrix with $p = 2$			Error with $p = 1$		
0.8	4.4	-0.9	-0.4	0.1	0.0	0.2	-0.9	2.2
5.8	12.4	6.1	2.3	-0.6	0.1	-2.2	1.0	-1.1
-3.2	-14.6	-5.9	0.6	-0.2	0.0	-0.7	0.0	0.5
-6.2	-15.6	-1.9	-1.8	0.5	-0.1	2.3	0.8	-3.6
-2.2	-8.6	2.1	0.2	-0.1	0.0	0.4	1.7	-4.7
1.8	8.4	0.1	-0.5	0.1	0.0	0.2	-1.0	2.7
4.8	8.4	5.1	2.4	-0.7	0.1	-2.2	1.2	-1.5
1.8	13.4	6.1	-1.6	0.5	0.0	1.8	-0.1	-1.1
-3.2	-12.6	-6.9	0.1	0.0	0.0	-0.3	-0.6	2.0
-4.2	-10.6	-7.9	-1.4	0.4	0.0	1.0	-1.5	3.4
2.8	19.4	6.1	-2.3	0.6	-0.1	2.2	-0.9	0.8
-0.2	1.4	2.1	-0.5	0.1	0.0	0.7	0.3	-1.4
1.8	-5.6	-3.9	3.0	-0.8	0.1	-3.2	0.2	1.9

- Errors greater than 1 are in red
- Sum of squared errors are 120.2 and 36.9 with $p = 1$ and 2, respectively

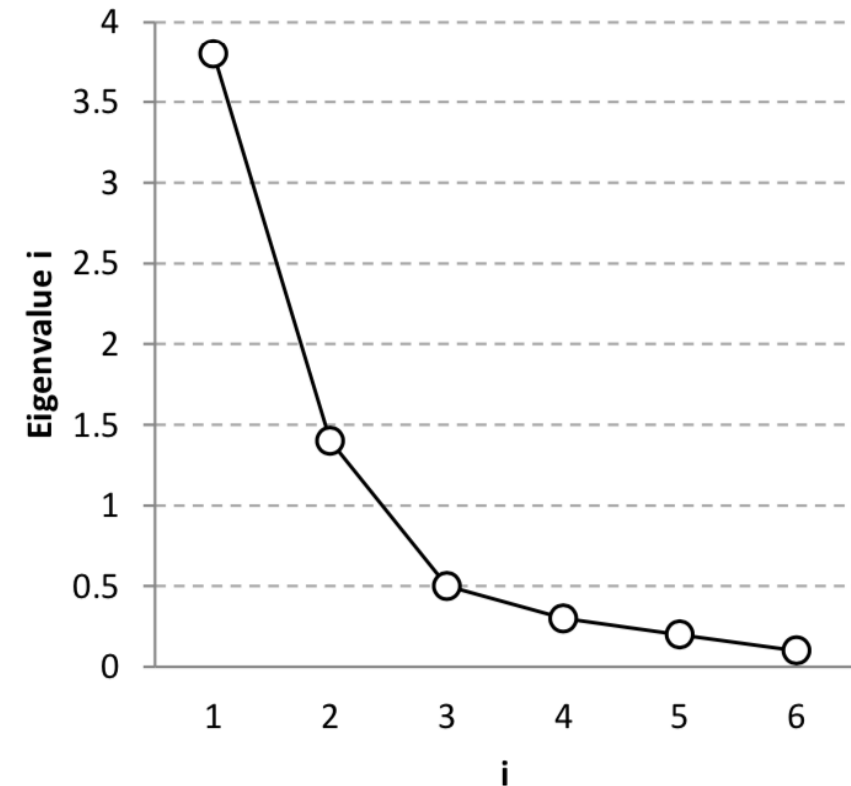
Example 1 - Number of PCs

- How many principal components?
 - There is no definitive answer
- Scree plot: a popular visual aid
 - Larger eigenvalues \Rightarrow more important eigenvectors
 - Small eigenvalues may be ignored without loss of important information
- Scree plot
 - Plot of λ_i versus i (sorted)
- To determine appropriate number of components (p), look for an elbow

Scree Plot



$p = 2$ is appropriate



$p = 2$ may be appropriate

Eigenvalues are a measure of the amount of variance accounted for by a factor

Example 2 - Iris Dataset

- Number of observations: 150
- Number of attributes: 4 numeric, predictive attributes and the class
- Attribute Information:
 - Sepal length in cm
 - Sepal width in cm
 - Petal length in cm
 - Petal width in cm
- Classes: Iris Setosa, Iris Versicolour, Iris Virginica



Iris Setosa



Iris Versicolour



Iris Virginica

Step 1: Calculate Covariance Matrix

- $A \in \mathbb{R}^{150 \times 4}$ (excluding class attribute)
- $\Sigma = \frac{1}{n} A^T A$ (assuming columns of A have zero mean)

Sepal.Length Sepal.Width Petal.Length Petal.Width

- $\Sigma = \begin{matrix} \text{Sepal.Length} \\ \text{Sepal.Width} \\ \text{Petal.Length} \\ \text{Petal.Width} \end{matrix} \begin{pmatrix} 0.6857 & -0.0424 & 1.2743 & 0.5163 \\ -0.0424 & 0.1899 & -0.3297 & -0.1216 \\ 1.2743 & -0.3297 & 3.1163 & 1.2956 \\ 0.5163 & -0.1216 & 1.2956 & 0.5810 \end{pmatrix}$

Step 2: Calculate Eigenvalues and Eigenvectors

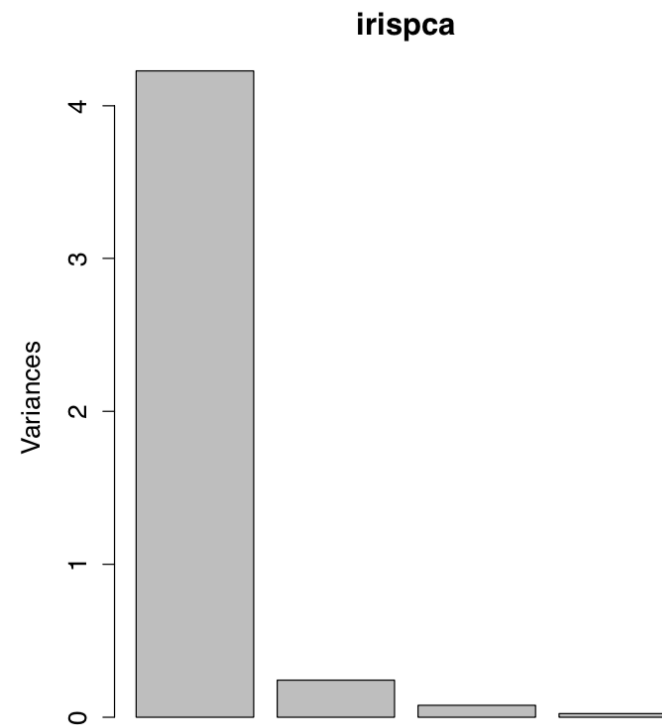
- Eigenvalues = (4.2282, 0.2427, 0.0782, 0.0238)

- $\Sigma = \begin{matrix} & \begin{matrix} PC1 & PC2 & PC3 & PC4 \end{matrix} \\ \begin{matrix} Sepal.Length \\ Sepal.Width \\ Petal.Length \\ Petal.Width \end{matrix} & \begin{pmatrix} 0.3614 & -0.6566 & -0.5820 & 0.3155 \\ -0.0845 & -0.7302 & 0.5979 & -0.3197 \\ 0.85671 & 0.1734 & 0.0762 & -0.4798 \\ 0.3583 & 0.0755 & 0.5458 & 0.7537 \end{pmatrix} \end{matrix}$

- The first principal component is the most important (largest eigenvalue), while others are not very significant.

Step 3: Scree plot

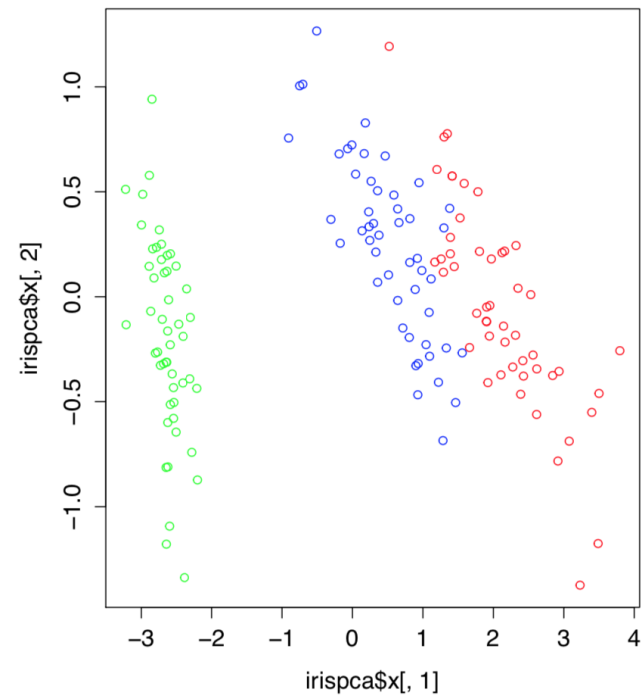
- Eigenvalues = (4.2282,0.2427,0.0782,0.0238)



- Confirm that using only one or two components is enough!

Step 4: Projection to Smaller Dimension

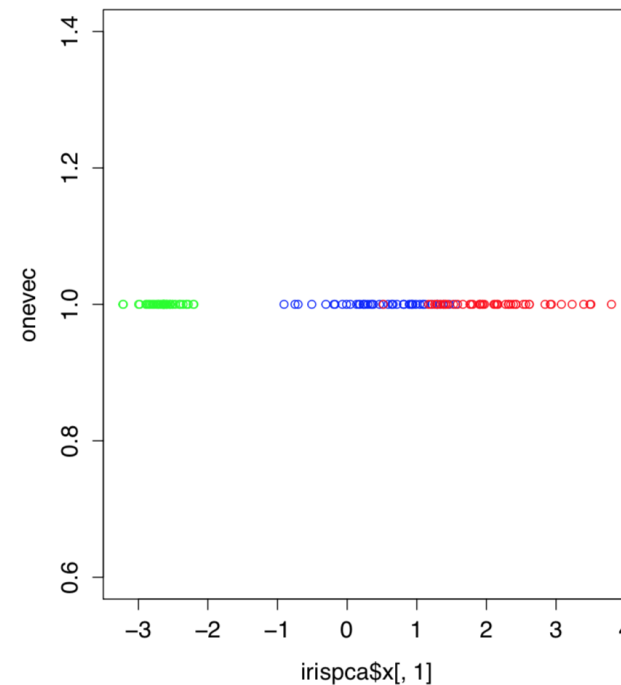
- With $p = 2$, projected data is obtained by $Y = AX \in \mathbb{R}^{150 \times 2}$
- Visualize Y with class attribute (different class in different colors)



Green = setosa, blue = versicolour, red = virginica

Step 4: Projection to Smaller Dimension

- With $p = 1$, projected data is obtained by $Y = AX \in \mathbb{R}^{150 \times 1}$
- Visualize Y with class attribute (different class in different colors)



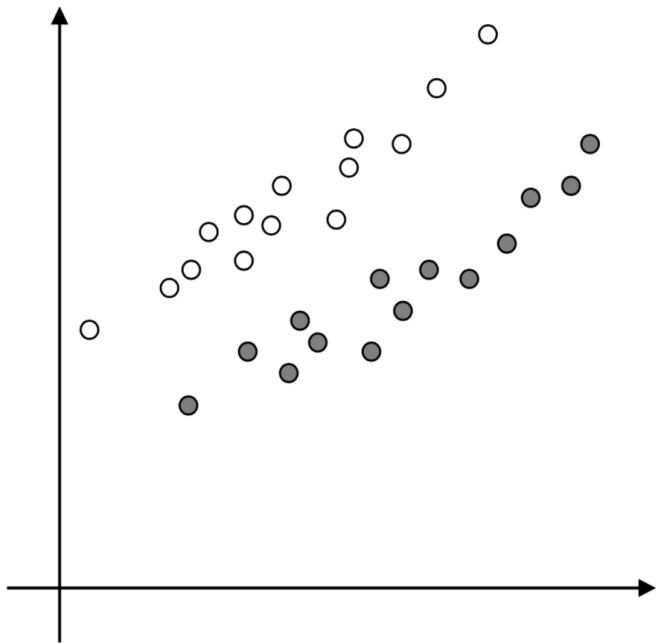
Green = setosa, blue = versicolour, red = virginica

PCA: Limitations

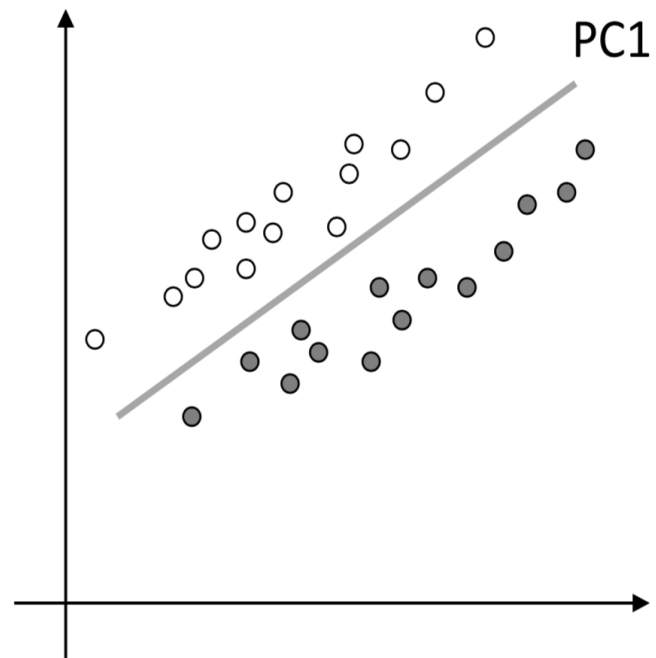
- PCA assumes a linear relationship between features i.e. it cannot capture non-linear structure in the data (as in many real-world applications).
- PCA assumes a correlation between features.
- PCA is sensitive to the scale of the features
- PCA is not robust against outliers
- Low interpretability of principal components.
- There is a trade-off between information loss and dimensionality reduction
- Technical implementations often assume no missing values

Issues: Data with Labels

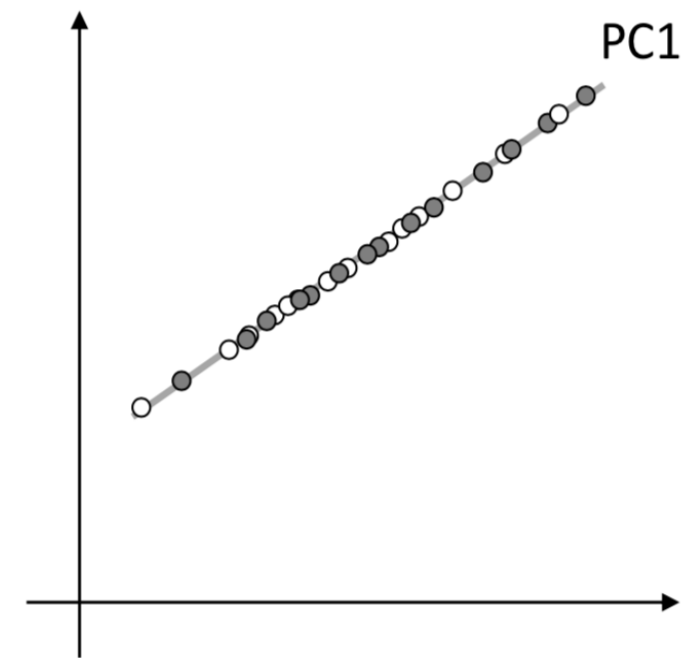
- A problematic example



Data with labels (white and gray)



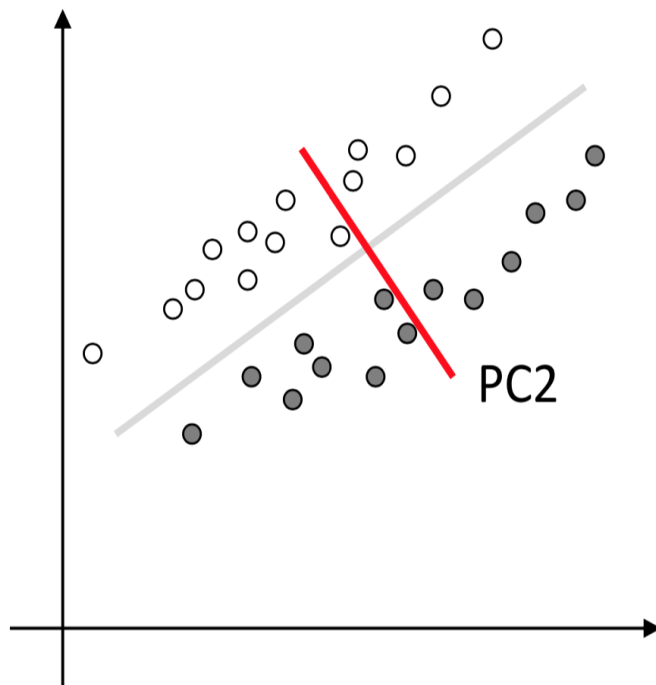
The first PC that explains most of the variance



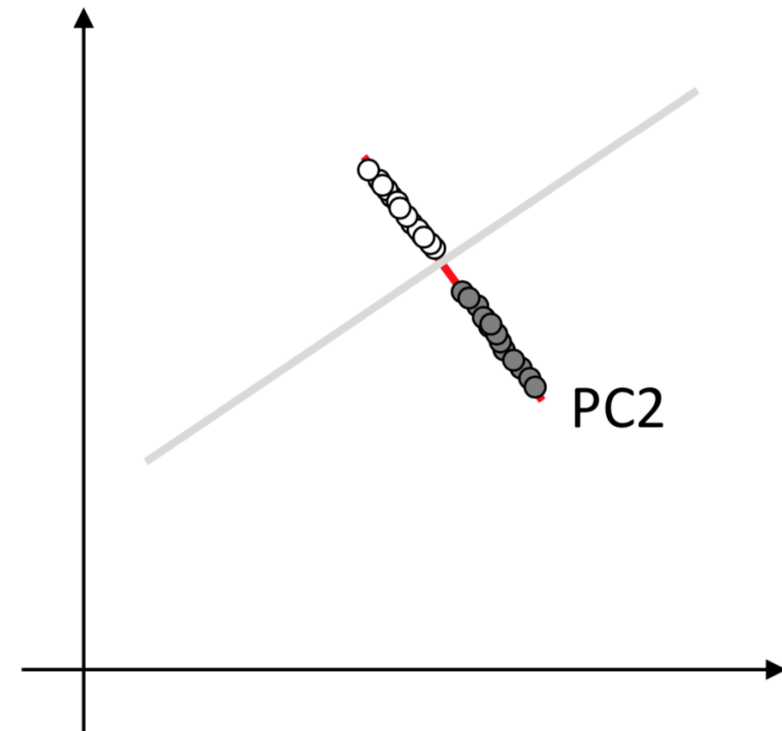
After projection, the result is useless

Issues: Data with Labels

- A problematic example



In fact, if we use the second PC,
we obtain a better result



In fact, if we use the second PC,
we obtain a better result

LINEAR DISCRIMINANT ANALYSIS (LDA)

Linear Discriminant Analysis (LDA)

- LDA is a supervised dimensionality reduction technique that also achieves classification of the data simultaneously.
- PCA focuses on capturing the direction of maximum variation in the data set.
- LDA focuses on finding a feature subspace that maximizes the separability between the groups.

LDA Solution Approach

- LDA projects the data points onto new axes such that these new components maximize the separability among categories while keeping the variation within each of the categories at a minimum value.
- As with PCA, LDA assumes that your data is centered around the origin and that your features are uncorrelated with one another
- The LD1 the first new axes created by LDA will account for capturing most variation between the groups or categories and then comes LD2 and so on.

