# BITS Pilani

**BITS** Pilani
Hyderabad Campus

Prof.Aruna Malapati
Professor
Department of CSIS

# Dimension Reduction using PCA
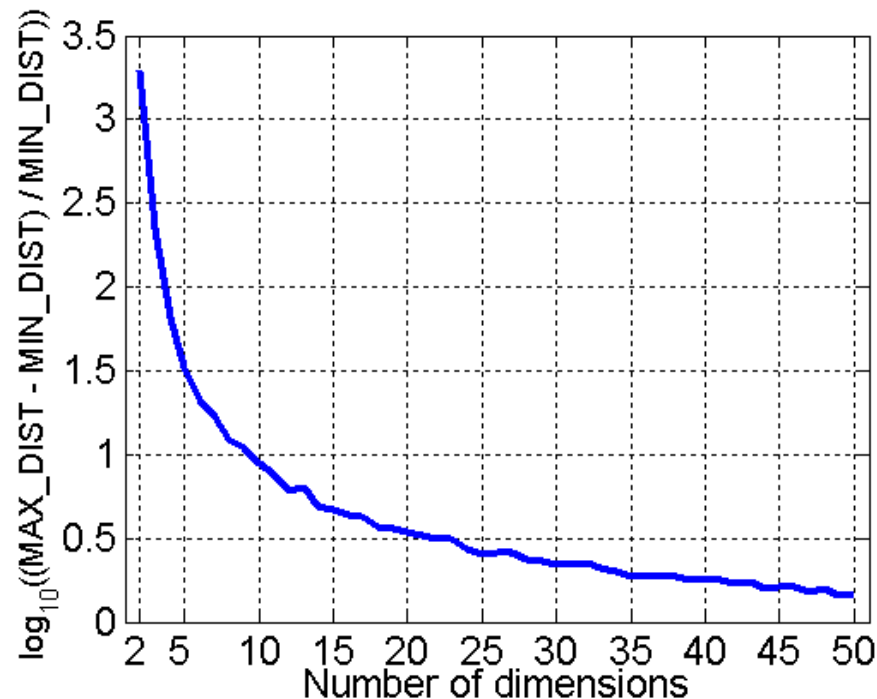
# Today's Agenda

- Curse of Dimensionality

- Introduction to Dimension Reduction

- Motivation for PCA

- PCA

# Curse of Dimensionality

- When dimensionality increases, data becomes increasingly sparse in the space that it occupies

- Definitions of density and distance between points, which is critical for clustering and outlier detection, become less meaningful



- **Randomly generate 500 points**
- **Compute difference between max and min distance between any pair of points**
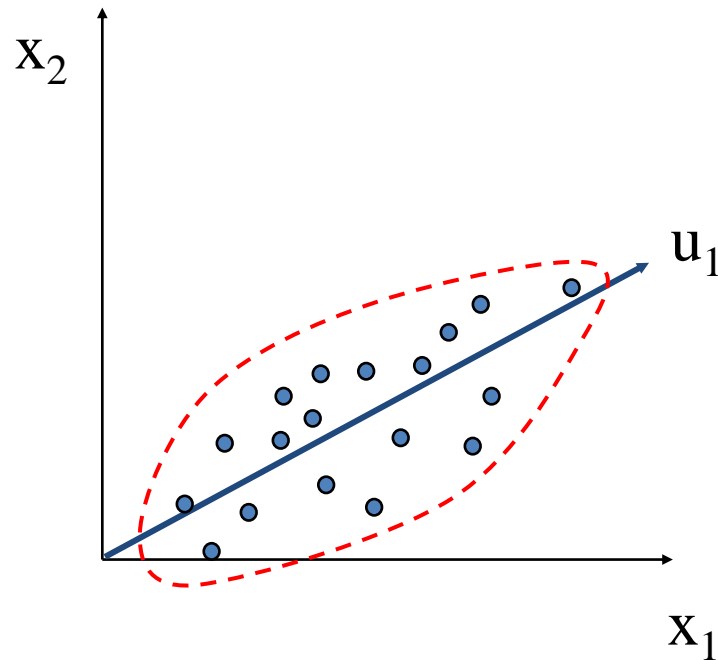
# Dimensionality Reduction

- Purpose:
  - Avoid curse of dimensionality
  - Reduce amount of time and memory required by data mining algorithms
  - Allow data to be more easily visualized
  - May help to eliminate irrelevant features or reduce noise

- Techniques
  - Principle Component Analysis
  - Singular Value Decomposition
  - Others: supervised and non-linear techniques
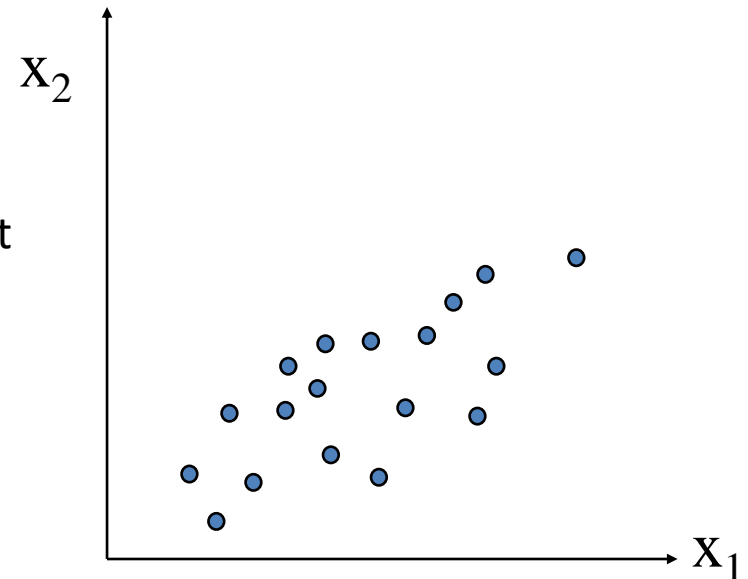
# Dimensionality Reduction: PCA

- Goal is to find a projection that captures the largest amount of variation in data

# PCA

- Reduce High Dimensional data into something that can be explained in fewer dimensions.

- We need PCA since we suspect that in our interesting data set not all measures are independent i.e there exist correlations.

Assume the data set represent height and weight of people in a region.
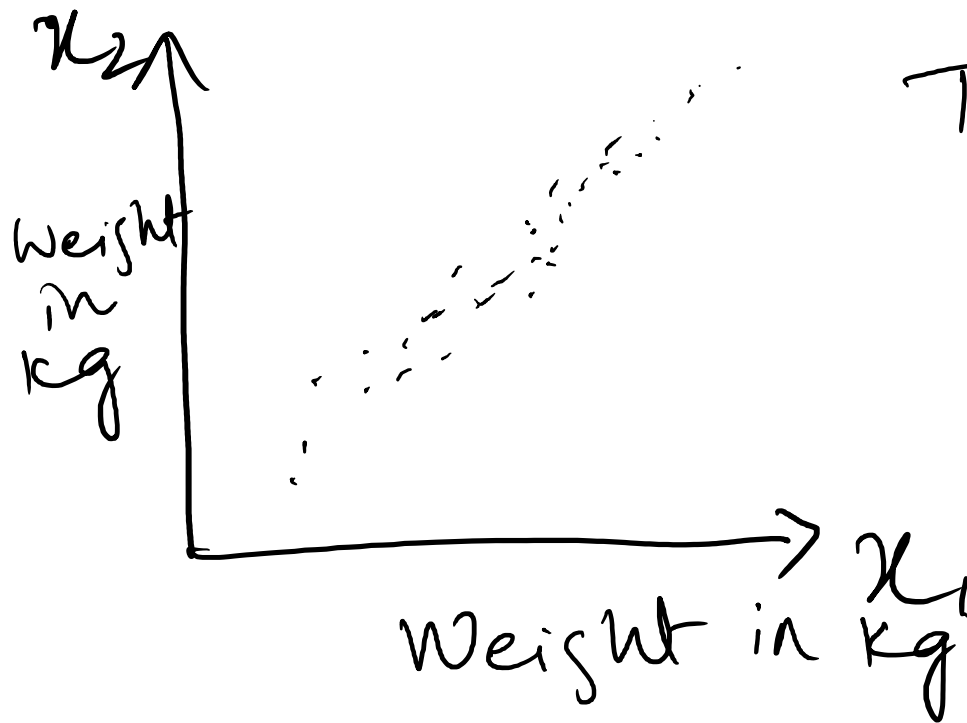
$x_2$

$x_1$

# Principal Component Analysis (PCA)

- Reduce higher dimension data into something that can be explained in fewer dimensions and gain an understanding of the data.

- We need PCA since we suspect that in our data set not all measures are independent and there exist correlations or structures or patterns.
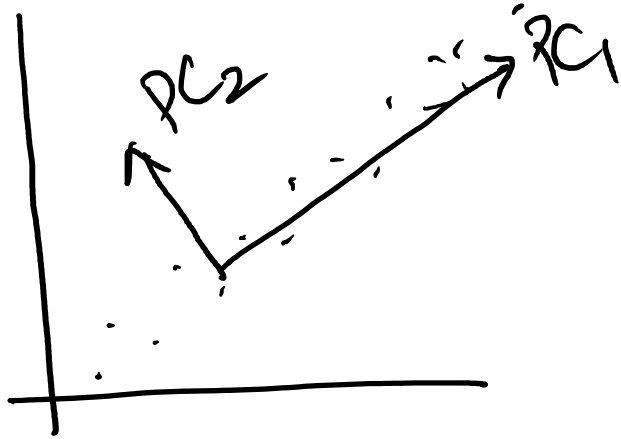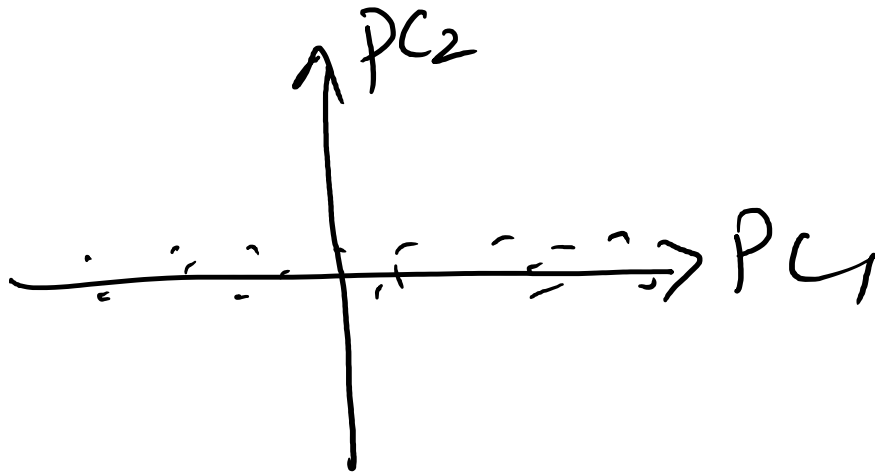
The data is not random, as $x_1$ increases $x_2$ also increases and hence +vely Co-related

$x_2$ — Weight in kg (vertical axis)

Weight in kg — $x_1$ (horizontal axis)

Hence we can reduce from $2D \rightarrow 1D$

# Projection

PC₂ → PC₁

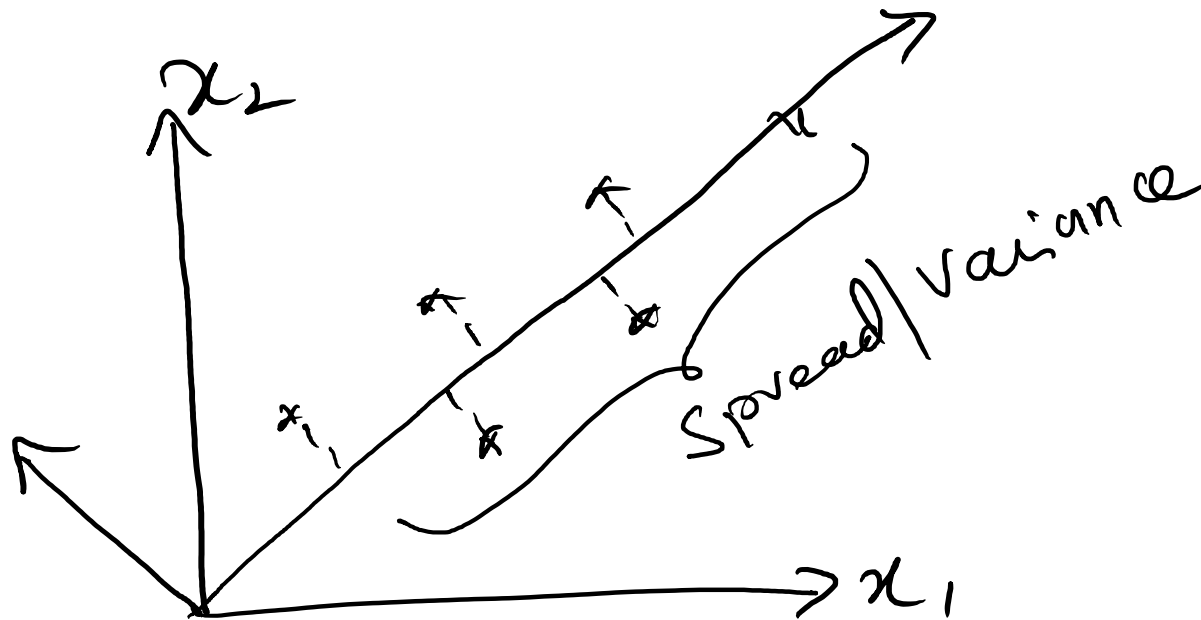PC$_1$ is the direction of the data, so that if we transform the data by projection.

PC₂

PC₁

Spread/variance

If data is projected on $x_1$ or $x_2$ we loose one dimension.

PCA helps us in identifying the best projection.

The goal is to find a lower dimensional surface on which to project such that the sum of squares errors are minimal.

We can project data on multiple vectors $(v_1 \ldots v_\alpha)$

But which is the best suitable vector for projection?

# Salient features of PCA

- Directions are in the order of % of variance explained.

- Every PC is orthogonal.


- PCA can be solved using

- Maximum Variance

- Minimum Error

$$\text{Data Set} = \{x_i, y_i\}_{i=1}^{S} \qquad S = \text{no of Samples}$$

$$\text{Each } x_i \in \mathbb{R}^N \qquad N = \text{no of dimensions}$$

Goal: Project data onto a space having $M \dim < D$ while maximizing the variance of projected data.

1. Mean center the data

2. Compute covariance matrix

3. Compute the eigen value decomposition of covariance matrix.

4. Find the best eigen vector by using eigen values.

Suppose you have a 2D matrix

Compute $X^T X$ Covariane matrix 2x2

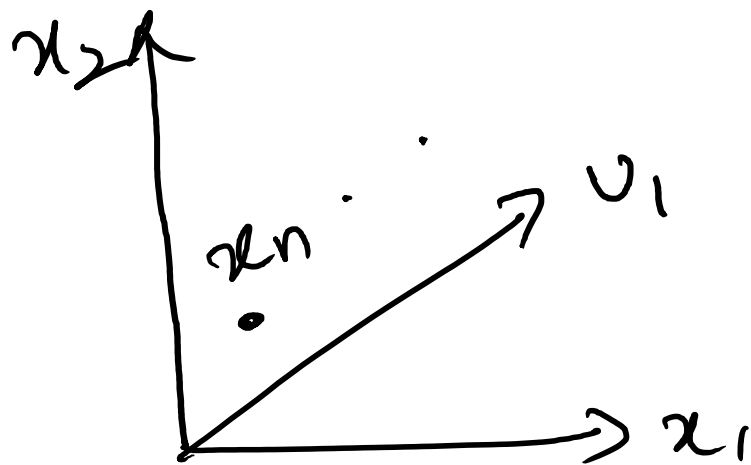$\Downarrow$

eigen decomposition

$\Downarrow$

2 eigen values

2 eigen vectors

Suppose our $M=1$ and let $v_1$ be a unit vector so that $v^T v = 1$ So that we are intrested in direction of $v_1$ not magnitude.

Each $x_n$ is Projected on $v_1$ by taking dot product $v_1^T \cdot x_n$

Mean of projected data will be

$$U_1^T \cdot \bar{x} \quad \text{where } \bar{x} \text{ is the set mean}$$

$$\bar{x} = \frac{1}{N} \sum_{n=1}^{N} x_n$$

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

$$\text{Variance } \sigma^2 = \frac{1}{N} \sum_{n=1}^{N} \left( U_1^T x_n - U_1^T \bar{x} \right)^2$$

$$\sigma^2 = \frac{1}{N} \sum_{n=1}^{N} \left( U_1^T \left( x_n - \bar{x} \right) \right)^2$$

using linear algebra

$$\left( A \cdot B \right)^2 = \left( A \cdot B \right) \left( A \cdot B \right)^T$$

$$\frac{1}{N} = \left( A \cdot B \right) B^T \cdot A^T$$

$$\sigma^2 = \frac{1}{N}\sum_{n=1}^{N}\left(v_1^T(x_n-\bar{x})\right)\left(v_1^T(x_n-\bar{x})\right)^T$$

$$= \frac{1}{N}\sum_{n=1}^{N}\left(v_1^T(x_n-\bar{x})\left(v_1(x_n-\bar{x})^T\right)\right.$$

Since $v_1$ & $v_1^T$ has nothing to do with data points

$$= v_1^T\left(\frac{1}{N}\sum_{n=1}^{N}(x_n-\bar{x})(x_n-\bar{x})^T\right)v_1$$

$$= U_1^T S U_1 \quad --- \text{①}$$

where $S$ is the Covariance Matrix

$$S = \frac{1}{N} \sum_{n=1}^{N} (x_n - \bar{x})(x_n - \bar{x})^T --- \text{②}$$

$$x = \begin{bmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nD} \end{bmatrix} \quad \text{and} \quad \bar{x} = \begin{bmatrix} \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_D \end{bmatrix}$$

mean of 1st feature

$$x_n - \bar{x} = \begin{bmatrix} x_{n1} - \bar{x}_1 \\ x_{n2} - \bar{x}_2 \\ \vdots \\ x_{nD} - \bar{x}_D \end{bmatrix}$$

$$(x_n - \bar{x})^T = \begin{bmatrix} x_{n1} - \bar{x}_1, & x_{n2} - \bar{x}_2 \cdots & x_{nD} - \bar{x}_D \end{bmatrix}$$

# How to derive S?

$$(x_n - \bar{x})(x_n - \bar{x})^T =$$

$$(x_{n1} - \bar{x_1})(x_{n1} - \bar{x_1}) \quad (x_{n1} - \bar{x_1})(x_{n2} - \bar{x_2}) \cdots (x_{n1} - \bar{x_1})(x_{nD} - \bar{x_D})$$

$$(x_{n2} - \bar{x_2})(x_{n1} - \bar{x_1}) \quad (x_{n2} - \bar{x_2})(x_{n2} - \bar{x_2}) \cdots (x_{n2} - \bar{x_2})(x_{nD} - \bar{x_D})$$

$$= \begin{bmatrix} (x_{n1} - \bar{x_1})^2 & (x_{n2} - \bar{x_2})^2 & - - - - - \\ & - & - - - \end{bmatrix}$$

Substitute this matrix in S eqn ②

$$= \frac{1}{N} \sum_{n=1}^{N} \begin{bmatrix} (x_{n1} - \bar{x}_1)^2 & & \\ & (x_{n2} - \bar{x}_2)^2 & \\ & & \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{N} \sum_{n=1}^{N} (x_{n1} - \bar{x}_1)^2 & \frac{1}{N} \sum_{n=1}^{N} (x_{n1} - \bar{x}_1)(x_{n2} - \bar{x}_2) \cdots \\ & \end{bmatrix}$$

# How to derive S?

The result is a Co-variance matrix of the form.

$$\begin{pmatrix} \text{Var } x_{n_1} & \text{Co-var}(x_{n_1}, x_{n_2}) & \cdots & \text{Cov}(x_{n_1}, x_{nD}) \\ \text{Cov}(x_{n_1}, x_{n_2}) & & & \end{pmatrix}$$

Once Covariance matrix is derived

$$X^T X$$

$$\Downarrow$$

$$eign\left(X^T X\right)$$

$$\Downarrow$$

$$\left(\lambda_1 \quad \lambda_2 \; - \; - \; \lambda_D\right) \left(W_1 \quad W_2 - - W_D\right)$$

eigen values        'gen vectors
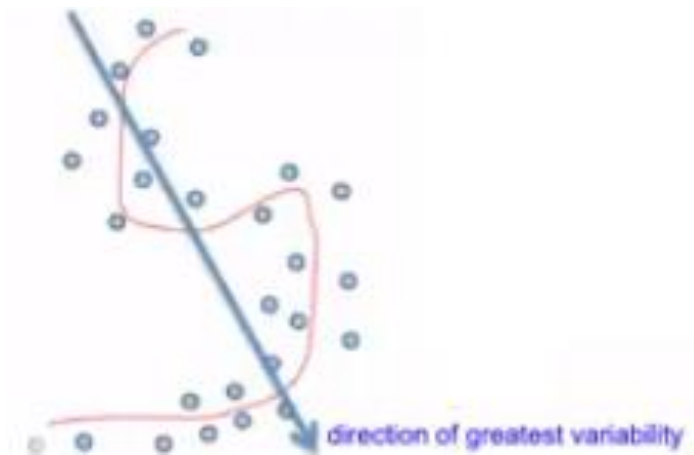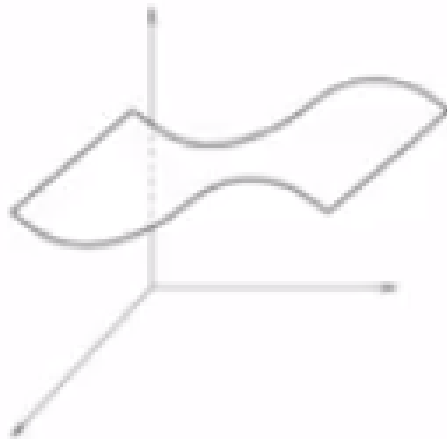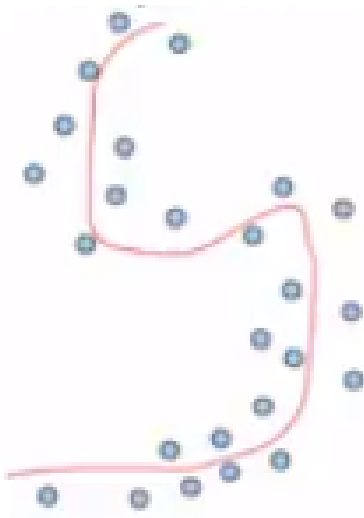
# PCA Limitations

- Covariance is extremely sensitive to large values

  - Multiply some dimension by 1000

  - Dominates covariance

  - Becomes principal component

- Normalize each dimension to zero mean and unit variance.
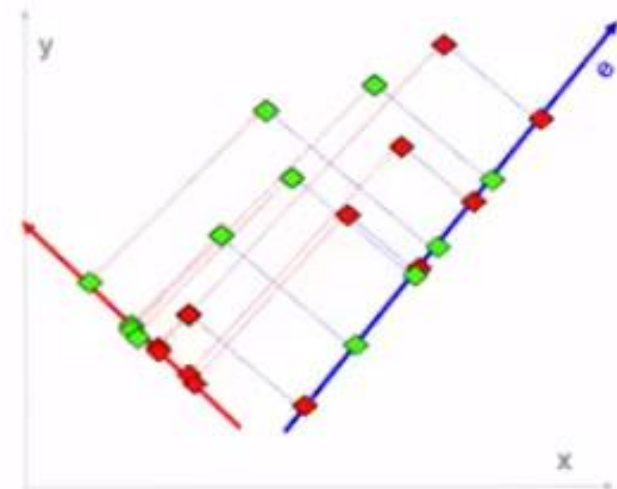
    X'=(X-mean)/standard-deviation

# PCA Limitations

- PCA assumes underlying subspace is linear.
- 1D –line
- 2D - plane



direction of greatest variability

# PCA and classification

- PCA is unsupervised

- Maximize overall variance of the data along a small set of directions

- Does not known anything about class labels

- Can pick direction that makes it hard to separate classes

# Take home message

- As the number of dimensions increases, the complexity and computational power required to build the model also increases.

- Dimension reduction methods are employed to find the best representation of data.

- PCA finds the best vectors on which the maximum variance in the data can be preserved.