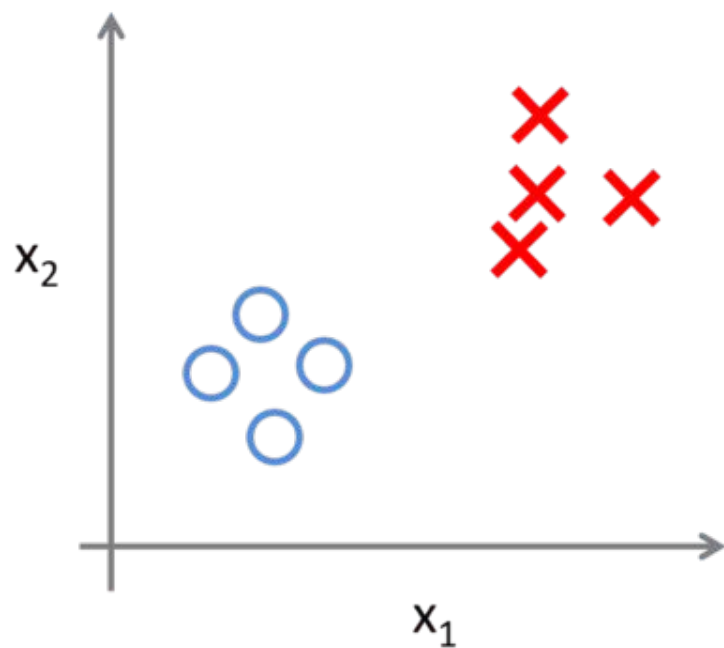


# SEMANA 8

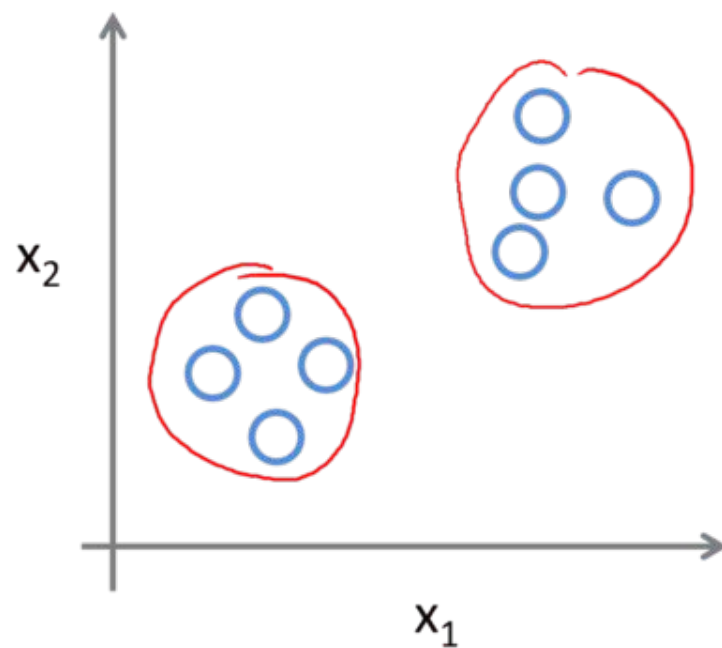
## UNSUPERVISED LEARNING

- ❏ Clustering
  - > K-means Algorithm
- ❏ Reducción de Dimensionalidad
  - > PCA Algorithm

## Supervised Learning



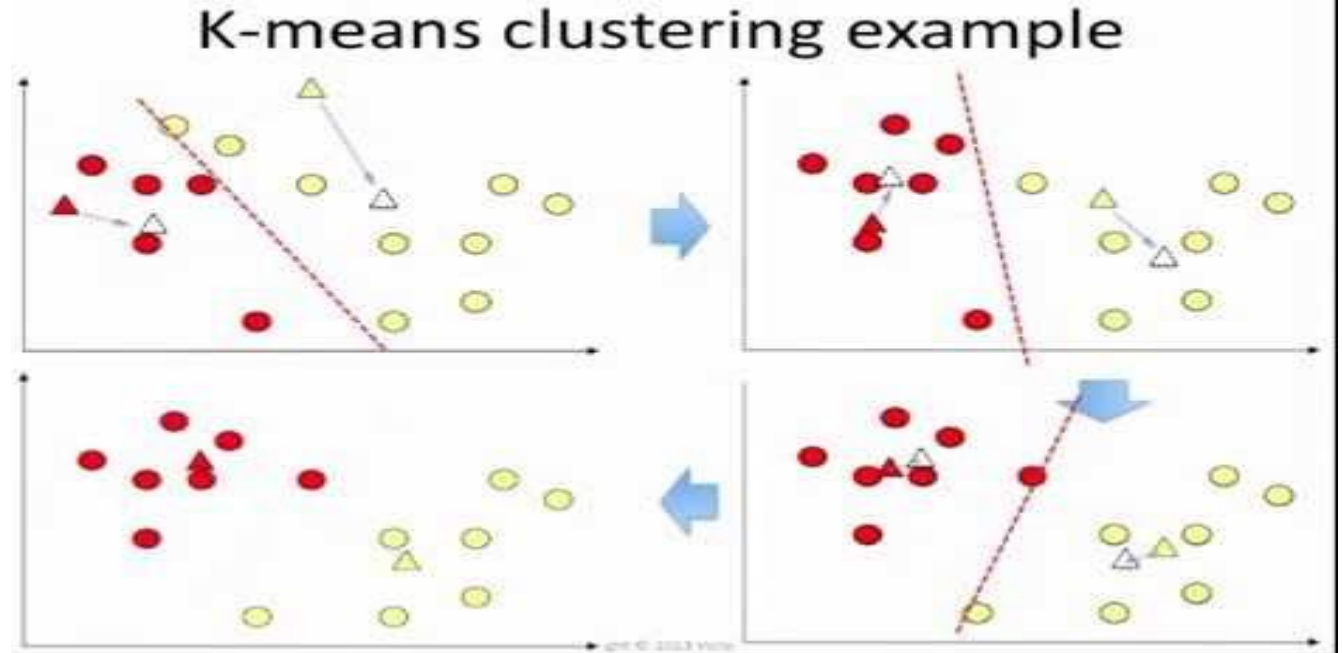
## Unsupervised Learning



# K-means Algorithm

Tiene 2 fases:

- F. Asignación
- F. Movimiento



# Objetivo De Optimización

## K-means optimization objective

- $c^{(i)}$  = index of cluster  $(1, 2, \dots, K)$  to which example  $x^{(i)}$  is currently assigned

$\mu_k$  = cluster centroid  $k$  ( $\mu_k \in \mathbb{R}^n$ )

$\mu_{c^{(i)}}$  = cluster centroid of cluster to which example  $x^{(i)}$  has been assigned

$K$   
 $k \in \{1, 2, \dots, K\}$   
 $x^{(i)} \rightarrow 5$      $c^{(i)} = 5$      $\mu_{c^{(i)}} = \mu_5$

Optimization objective:

$$\rightarrow J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \boxed{\|x^{(i)} - \mu_{c^{(i)}}\|^2} \leftarrow$$

$$\rightarrow \min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

$\rightarrow \mu_1, \dots, \mu_K$

Distortion

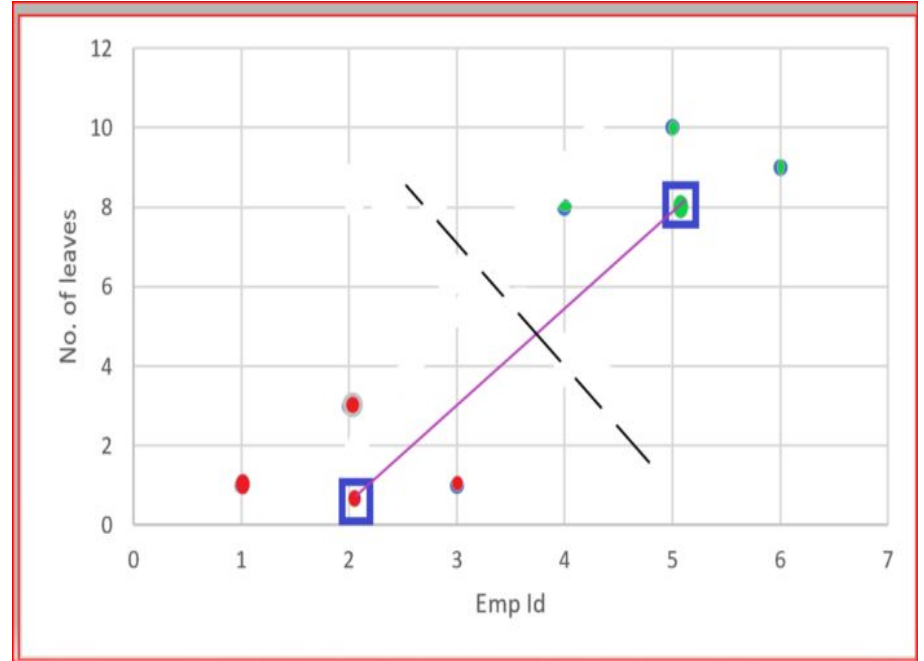


# Random Initialization

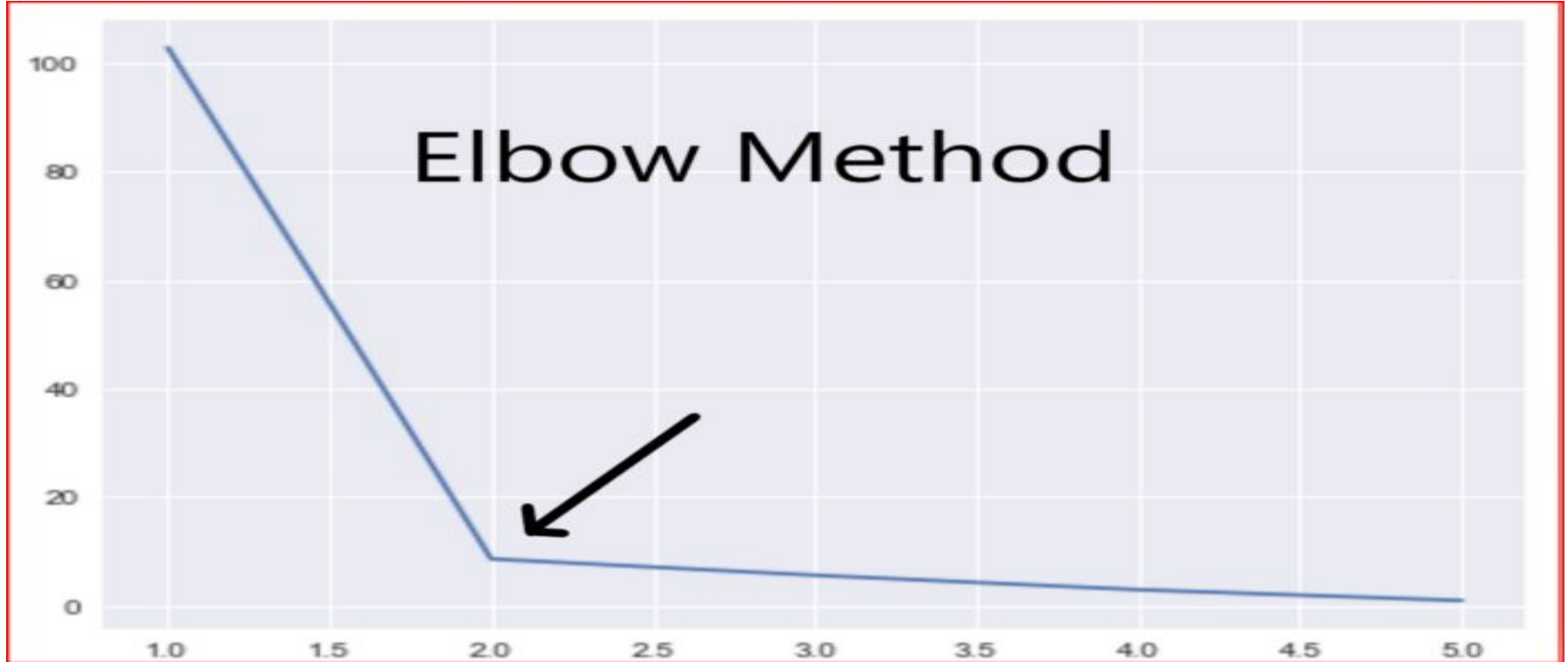
Paso 1: elegimos no. de grupos K

Paso 2: elegimos aleatoriamente centroides

Paso 3: Realizar las 2 fases del algoritmo explicados anteriormente.



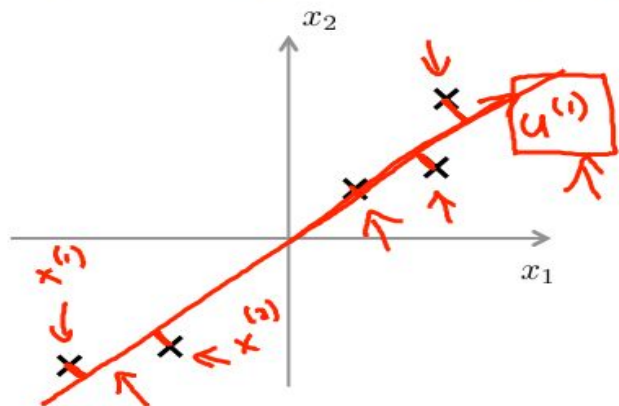
# ¿Cómo elegir el numero de clusters?



# PCA Algorithm

- Se obtiene rapidez en la ejecución del algoritmo.
- Los datos se pueden visualizar más fácilmente

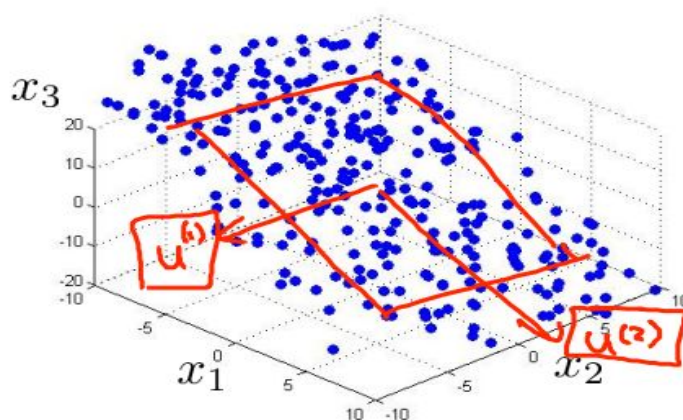
## Principal Component Analysis (PCA) algorithm



Reduce data from 2D to 1D

$$x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}$$

$z = [z_1]$



Reduce data from 3D to 2D

$$x^{(i)} \in \mathbb{R}^3 \rightarrow z^{(i)} \in \mathbb{R}^2$$

$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$

# APLICACIÓN DE PCA

1. Etapa de Pre-Procesamiento de datos.
2. Calculamos la Matriz de Covarianzas.
3. Con la función **svd** hallamos la matriz **U**
4. De este matriz **U** tomamos nuestras primeras vectores **K**.
5. Finalmente hallamos nuestra matriz reducida de variables **Z**.

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T$$

→  $[U, S, V] = \text{svd}(\text{Sigma});$

→  $\text{Ureduce} = U(:, 1:k);$

→  $z = \text{Ureduce}' * x;$

↑

↑

$x \in \mathbb{R}^n$

~~$x_0 = 1$~~

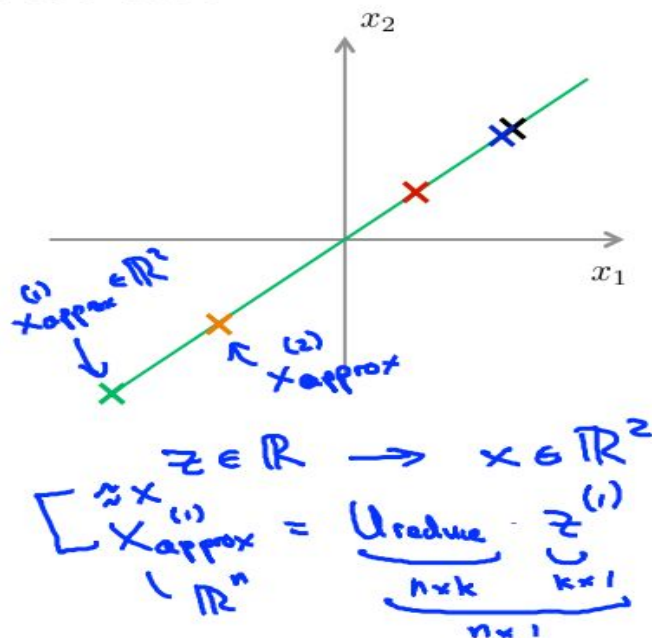
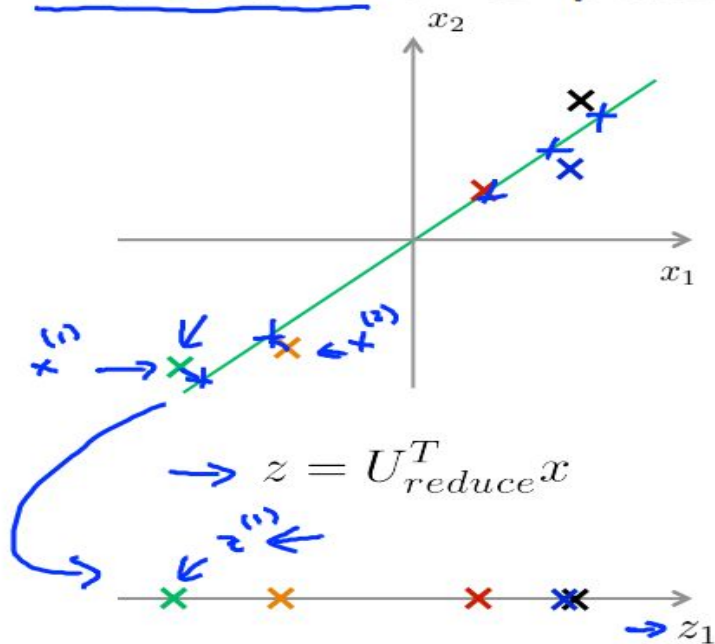
$$X = \begin{bmatrix} - & x^{(1)} & - \\ & \vdots & \\ - & x^{(m)} & - \end{bmatrix}$$

→  $\text{Sigma} = (1/m) * X' * X;$



# Reconstrucción a partir de la Representación Comprimida

## Reconstruction from compressed representation



# Elegir el número de componentes principales

→  $[U, S, V] = \text{svd}(\text{Sigma})$

Pick smallest value of  $k$  for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

$k=100$

(99% of variance retained)

# TIPS

- Debe definirse ejecutando PCA solo en el set de entrenamiento.
- No es recomendado usar PCA para prevenir Overfitting.

**PCA is not linear regression**

