

Clustering:

Unsupervised Learning → Unlabeled training set

Good for:

- Market segmentation
- Social Network analysis
- Organizing computer clusters
- Astronomical data analysis

k-means Algorithm is the most popular algo.

Pseudo Code:

Randomly initialize K cluster centroids μ_1, \dots, μ_k

Repeat {

 for $i=1$ to m :

$c(i) = \text{index } (1 \text{ to } K) \text{ of cluster centroid}$
 closest to $x(i)$

 for $k=1$ to K :

$\mu_k = \text{average (mean) of points assigned to}$
 cluster k

}

step1 → first for loop $\rightarrow C(i) = \operatorname{argmin} \|x^{(i)} - \mu_k\|^2$

step2 → Second for loop $\rightarrow \mu_k = \frac{1}{n} [x^{(k_1)} + x^{(k_2)} + \dots + x^{(k_n)}]$

In two steps, cost function is minimized.

Sometimes K-means can get stuck in local optima. To get rid of it we need random initialization

for i = 1 to 100:

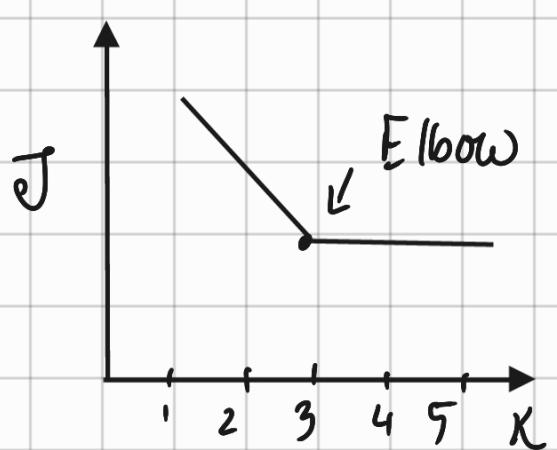
randomly initialize K-means

run K-means to get 'c' and ' μ '

compute the cost function $J(c, \mu)$

pick the clustering that give us the lowest cost

We can choose the number of cluster using "Elbow" method



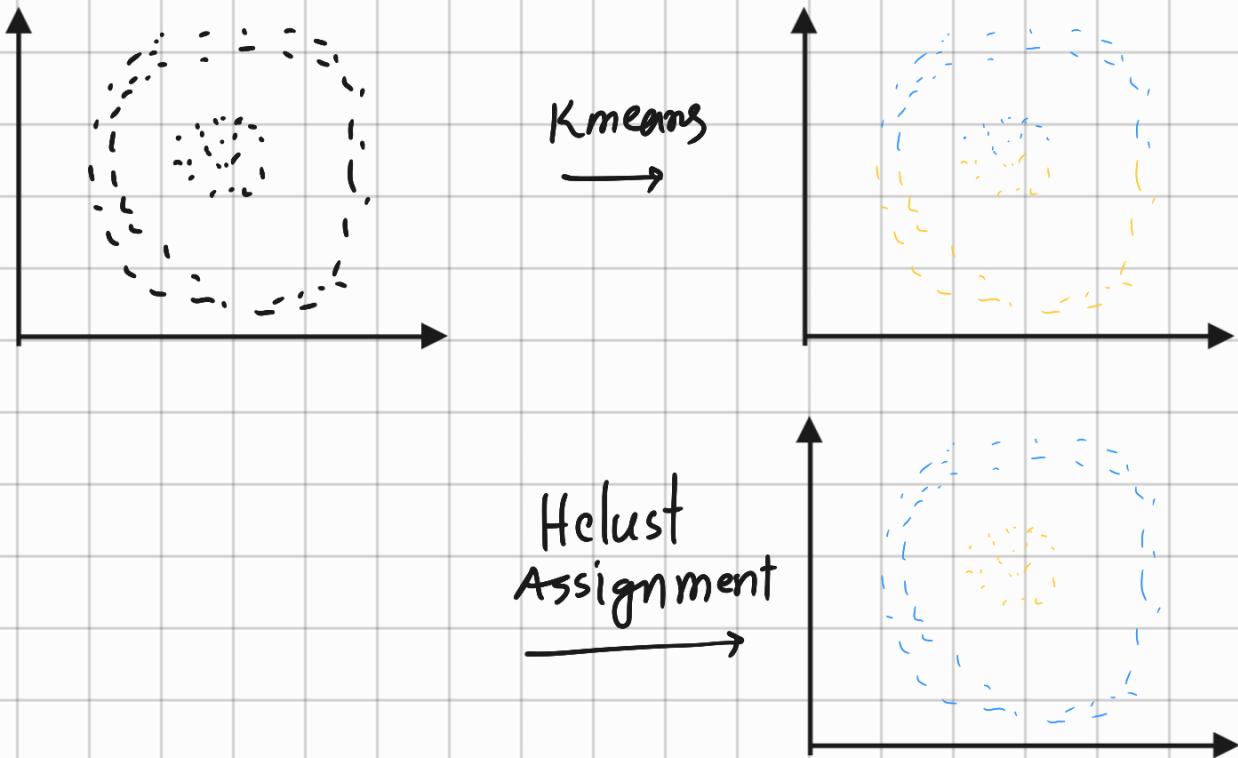
But sometimes curve has no elbow, then we'll find the later purpose like Tshirt sizing $\rightarrow K \rightarrow 3$ or 5

J will always decrease as K is increased unless it got stuck at bad local optima.

However, K means fails in some cases.

- It always makes all cluster spherical
- All axes have the same distribution and thus variance

One example:



Data Compression:

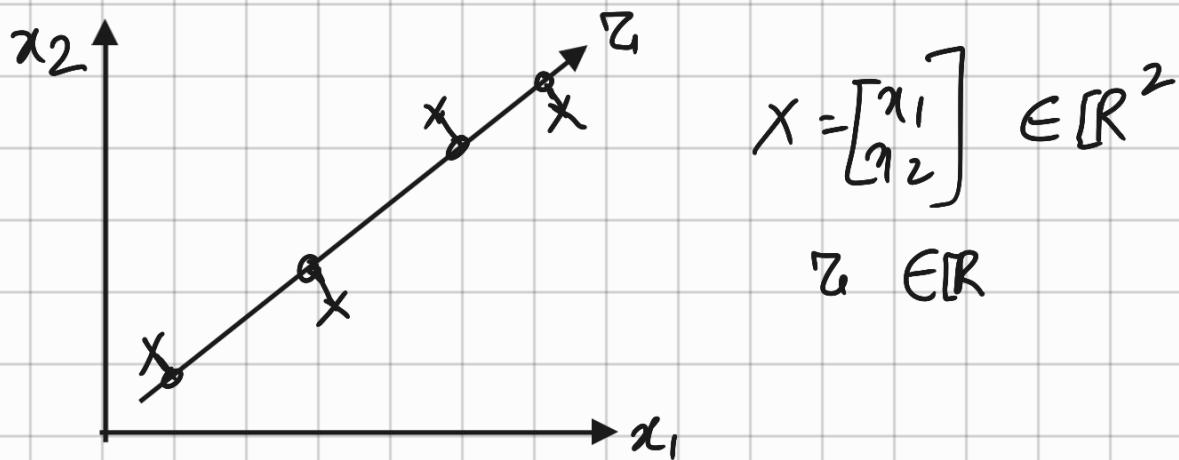
- Reduce dimension of our features [3d to 2d]
- Easy to visualize
- Summarize all other features

$$500 \rightarrow 2D$$

Principle Component Analysis (PCA)

The goal of PCA is to reduce the average of all the distances of every feature to the projection line.

In PCA, we are minimizing the shortest distance or shortest orthogonal distances, to our data point



Algorithm:

Data Preprocessing :

→ Given dataset x_1, x_2, \dots, x_m

→ Mean normalization $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$

→ Replace $x_j^{(i)}$ with $x_j^{(i)} - \mu_j$

→ Feature scaling, $x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{s_j}$

→ Reduce 2d to 1

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

↓

Covariance Matrix

$$\text{Sigma} = \frac{1}{m} (X^T * X);$$

→ eigenvector of covar matrix

$$[U, S, V] = \text{svd}(\text{sigma}), \quad [\text{eig}(\text{sigma})]$$

→ Take first k column of U matrix

$$U_{\text{reduce}} = U(:, 1:k);$$

→ Finally, $Z = X * U_{\text{reduce}}$

Reconstruction from compressed representation

$$X_{\text{approx}} = Z * U^T \quad [\text{as } U^T = U^{-1}]$$

Choosing # of PCA:

Choose k to the smallest value such that,

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \rightarrow \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}}$$

99% of the variance is retained.

Algo:

$S = \text{diagonal mat}$

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

iterate $k = 1 : \infty$,

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

Application of PCA:

- speed up supervised learning
- Data compression [mostly used in \downarrow]
- Visualization of data
 100×100 pixel image
 $= \mathbf{x} \in \mathbb{R}^{10000}$

Bad use of PCA

- prevent overfitting due of huge # of features.
 - ↳ Using regularization is more effective in that case

"First use Full Machine Learning Algorithm without PCA, if needed then use PCA."