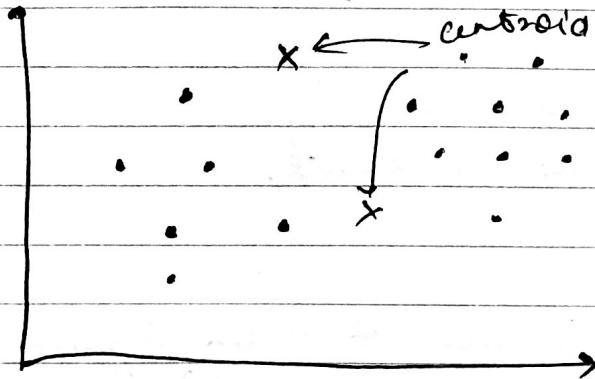


K-means Algorithm

- Randomly assign 2 points (clustering centroids)



Step 1 :- cluster assignment

Step 2 :- Move centroid.

(move centroid is mean of data)

i/p - No. of clusters (k)

Training set. $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

Algo.

- Randomly initialize K cluster centroids ($\mu_1, \mu_2, \dots, \mu_K$)

- Repeat {

minimizes $\sum_{i=1}^m \|x^{(i)} - \mu_j\|^2$ for $j = 1 \dots k$

$c^{(i)} := \text{index } (1 \dots k) \text{ of cluster centroid}$
closest of $x^{(i)}$ ($\min \|x^{(i)} - \mu_k\|^2 = c^{(i)}$)

minimizes $\sum_{j=1}^k \|x^{(i)} - \mu_j\|^2$ for $i = 1 \dots m$

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

if a cluster has no points, eliminate it.

What if there are not well separated clusters.

Optimization objective.

$c^{(1)}$ → index of clusters

μ_k → cluster centroid k .

$\mu_c(i)$ = mean centroid of cluster which example $x^{(i)}$ has been assigned

$x^{(i)}$ → s

$c^{(1)}$ → s

$\mu_c(i) = \mu_s$; that means $x^{(i)}$ is assigned of cluster

Optimization objective:

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

Find c & μ for min J .

Random Initialization

How do randomly initialize cluster centroids. μ_1, \dots, μ_K

→ should have $K \leq m$.

→ randomly pick K training examples

→ set μ_1, \dots, μ_K equal to these K examples

K-means can end at local minima,

If you don't want it, try multiple initialization.

~~eg:-~~ eg:-

Running k-means 100 times.

for $i = 1 \text{ to } 100 \}$

randomly initialize k means

Run k means

Get $c^{(i)}$ & $\mu^{(i)}$

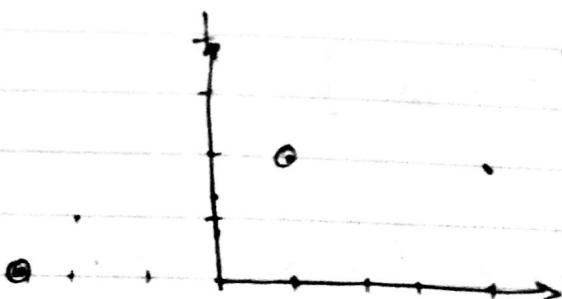
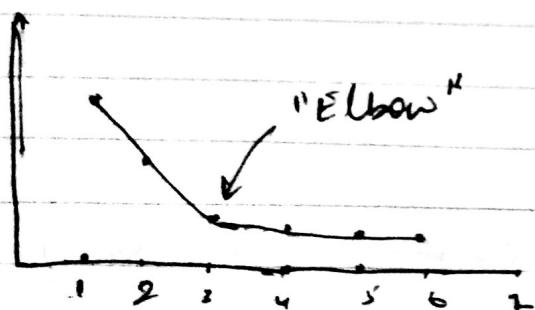
$J(c^{(i)}, \dots, c^{(m)}, \mu, \dots, \mu_m)$

Pick min cost $J(c^{(i)}, \dots, c^{(m)}, \mu, \dots, \mu_m)$

If k is small, random initialization will work definitely.
(2-10)

Choosing the No. of clusters.

Elbow method.



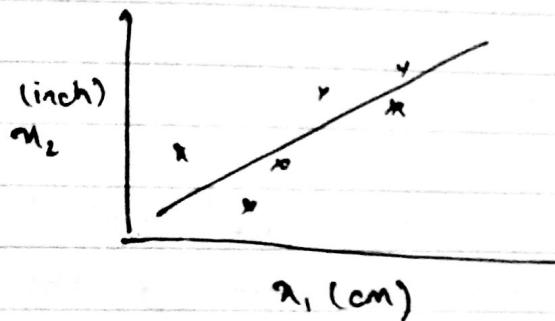
Not good method.

To choose no. of clusters - human insights.

Dimensionality Reduction

* Data compression

→ Maybe if you have a lot of features, and if you have a feature in length (cm) and a feature in inches. Then your dimensionality is reduced.



reduce dimension from 2D to 1D.

This method algorithm runs quickly.

FOR 3D, if we project it on to a plane. And from n, z_1, z_2 , you get z_1, z_2 .

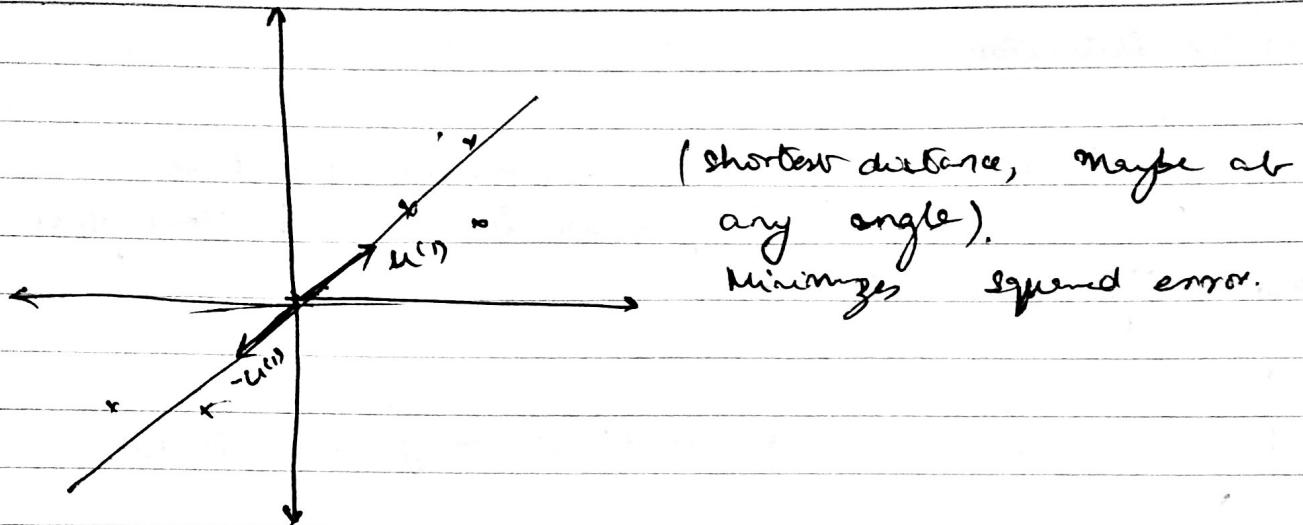
* Data Visualization

. How to plot so feature (20-D) You could take our 2 most important features

PRINCIPAL COMPONENT ENVIRONMENT (for Dimensionality Reduction)

before pca, perform mean normalization & feature scaling, so the feature have comparable range.

$$u^{(1)} = \sqrt{\left(\frac{1}{\sqrt{2}}\right)^2 + \left(\frac{1}{\sqrt{2}}\right)^2}$$



- Reduce 2D to 1D : find direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.
- Reduce n-dim to k-dim : find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$

Data processing.

* Training set

* Preprocess (feature scaling / mean normalization)

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

replace $x_j^{(i)} = x_j^{(i)} - \mu_j$

- # Reduce data from n-dimensions to k-dimensions
- ↳ Compute "covariance matrix":

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

↳ compute "eigenvectors" of matrix Σ

$$\rightarrow [U, S V] = \text{svd}(\Sigma) \dots \text{(singular value decomposition)}$$

you can use `eig(sigmas)`

Sigma will be $n \times n$ matrix.

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

$$U = \begin{bmatrix} | & | & | & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & | & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

(row)

Take first K dimensions.

You will get Ureduce [] with k-dim

$$Z = Ureduce^T X$$

$$= \begin{bmatrix} -u^{(1)} \\ -u^{(2)} \\ -u^{(3)} \end{bmatrix} X$$

$$Ureduce = U(:, 1:k);$$

$$Z = Ureduce^T \cdot x;$$

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

$$= \left(\frac{1}{m} \right)^{\frac{1}{2}} X^T \cdot X \quad \dots \quad X = \begin{bmatrix} -x^{(1)} \\ \vdots \\ -x^{(m)} \end{bmatrix}$$

Reconstruction from compressed representation

$$X \text{approx} = \underbrace{Ureduce}_{n \times k}^T \underbrace{Z}_{k \times 1} \underbrace{\cdot}_{m \times 1}$$

k = No. of principle components.

→ PCA minimizes avg. squared projection error.

$$\text{Squared projection error} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$$

$$\text{Total variation} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

↑ distance from origin.

we choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

↓ choose k .

"99% of variance is retained"

Algo.

Try PCA; $k=1$

Compute U vector, $z^{(1)}, z^{(2)}, \dots, z^{(m)}$
 ~~$x^{(1)}_{\text{approx}}, \dots, x^{(m)}_{\text{approx}}$~~

check if

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

This is less efficient

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

$$S = \begin{bmatrix} s_1 & & & \\ & s_2 & & \\ & & s_3 & \\ & & & \ddots \end{bmatrix}$$

for given k

$$s_i = \sqrt{s_{ii}}$$

$$1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^m s_{ii}} \leq 0.01$$

$$\frac{\sum_{i=1}^k \zeta_{ii}}{\sum_{i=1}^n \zeta_{ii}} \geq 0.99$$

Supervised learning speedup

$$\rightarrow (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)}).$$

one important use is to use it in computer vision.

for 100x100 pixels

with PCA, our algo will run smoothly.

1). Extract inputs

Unlabeled datasets: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{10000}$

\downarrow PCA

$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{1000}$

2) New training set

$(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})$

~~for some examples~~

~~Mapping should be done only on training set.
Run PCA only on training set. Apply mapping to
cross validation & test set.~~

Application

- compression
- visualization.

* Problem: It shouldn't be use to reduce overfitting.
Use regularization to reduce overfitting.