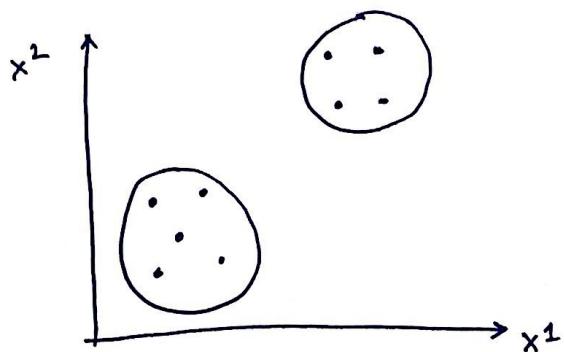


WEEK - 8

UNSUPERVISED LEARNING

* Clustering: don't have any labels.

Training set: $\{x^1, x^2, \dots, x^m\}$ without labels y^i



* K-means Algorithm

- ① • cluster Assignment }
 - ② • Move centroid }
- Iterative.

Inputs:

- K (number of clusters)
- Training set $\{x^1, x^2, \dots, x^m\}$
 $\hookrightarrow x^i \in \mathbb{R}^n$ {drop $x_0=1$ convention}

Algorithm :

- Randomly initialize K cluster centroids
 $u^1, u^2, \dots, u^K \in \mathbb{R}^n$

- Repeat

{

|| for $i = 1:m$

cluster assignment step

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

$$c^{(i)} = \min_k \|x^{(i)} - u_k\|$$

e.g. $\begin{cases} c^{(1)} = 3 \\ c^{(2)} = 3 \\ c^{(3)} = 5 \end{cases}$ } cluster 3
 $\begin{cases} c^{(4)} = 3 \\ c^{(5)} = 3 \\ c^{(6)} = 5 \end{cases}$ } cluster 5.

for $k = 1:K$ || $u_k := \text{average(mean)} \text{ of points assigned to cluster } k$

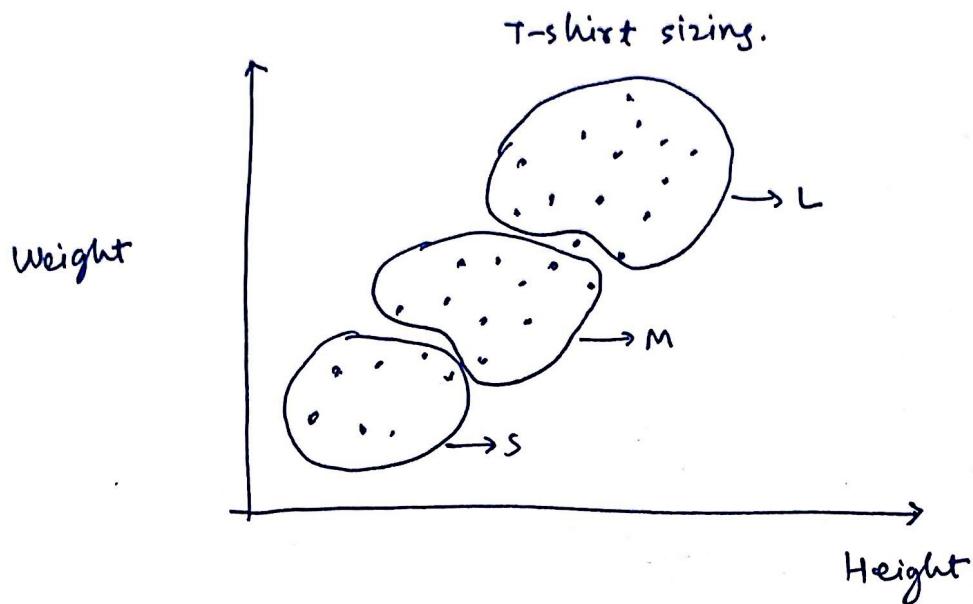
move centroid step

e.g. $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)} \Rightarrow \begin{cases} c^1 = 2 \\ c^5 = 2 \\ c^6 = 2 \\ c^{10} = 2 \end{cases}$

$$\therefore u_2 = \frac{1}{4} [x^1 + x^5 + x^6 + x^{10}] \in \mathbb{R}^n$$

}

④ K-means for non-separated clusters.



Optimization Objective : K-means

c^i : index of cluster to which x^i is assigned.

μ_k : cluster centroid

μ_{c^i} = cluster centroid of cluster i to which x^i has been assigned

$$x^i \geq 5$$

$$c^i = 5.$$

$$\mu_5 \in \mu_{c^i}$$

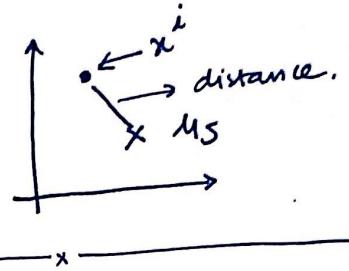
Optimization Objective : Distortion fn.

$$\left[J(c^1, c^2, \dots, c^m, u_1, u_2, \dots, u_K) = \frac{1}{m} \cdot \sum_{i=1}^m \|x^i - u_{c_i}\|^2 \right]$$

$$\min_{c^1, c^2, \dots, c^m} J(c^1, c^2, \dots, c^m, u_1, \dots, u_K)$$

$$c^1, c^2, \dots, c^m$$

$$u_1, u_2, \dots, u_K$$



↑ cluster assignment step: minimize $J(\dots)$ with resp. to.
 c^1, c^2, \dots, c^m holding (u_1, u_2, \dots, u_K)
fixed

↑ Move centroid step: min. $J(c\dots)$ wr.t u_1, u_2, \dots, u_K
holding c^1, c^2, \dots, c^m fixed

Random Initialization

① should have $K < m$

② Randomly pick K training examples

③ set u_1, \dots, u_K equal to these K examples

for $j = 1 : 100$

{

① Randomly initialize k-means.

② Run k-means, Get c^j, u_k for all

③ Compute $J(\dots)_{c, u} \rightarrow$ cost/distortion fn.

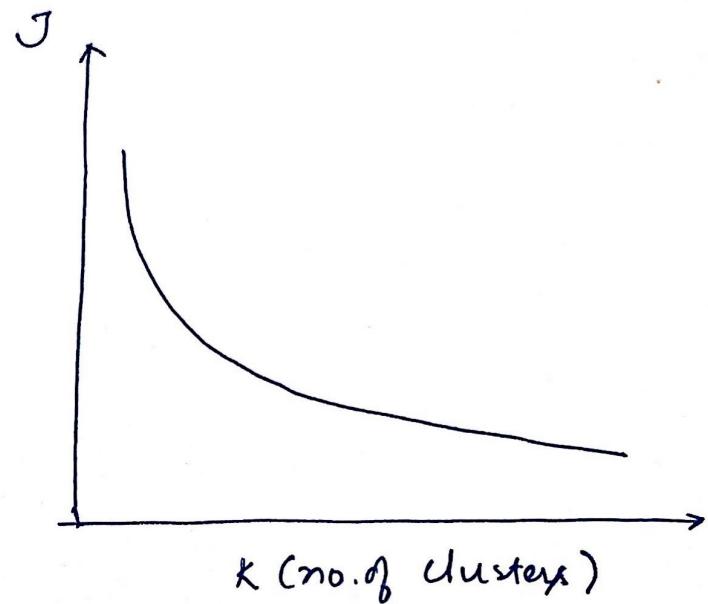
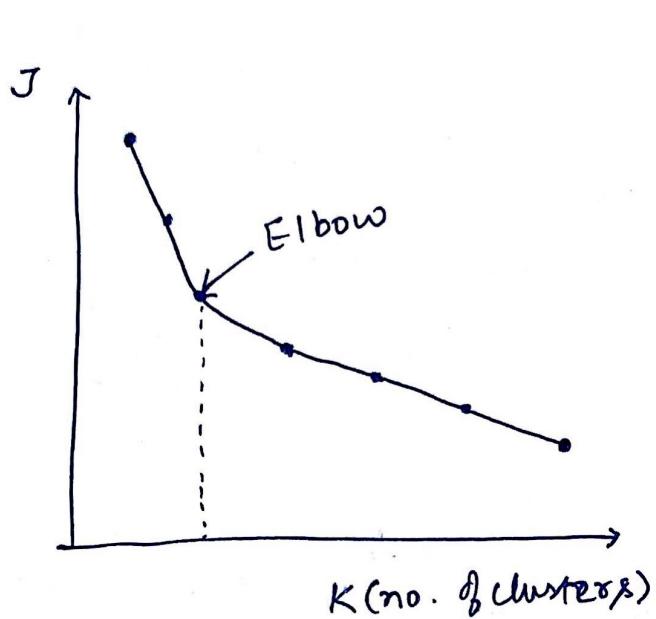
}

• pick clustering that gave lowest cost $J(c^1, c^2 \dots c^m, u_1, \dots u_K)$

(99)

Choosing the Number of clusters

Elbow Method.



J goes ~~out~~ low rapidly
still Elbow.

Dimensionality Reduction : Data Compression



$$x^i \in \mathbb{R}^2$$

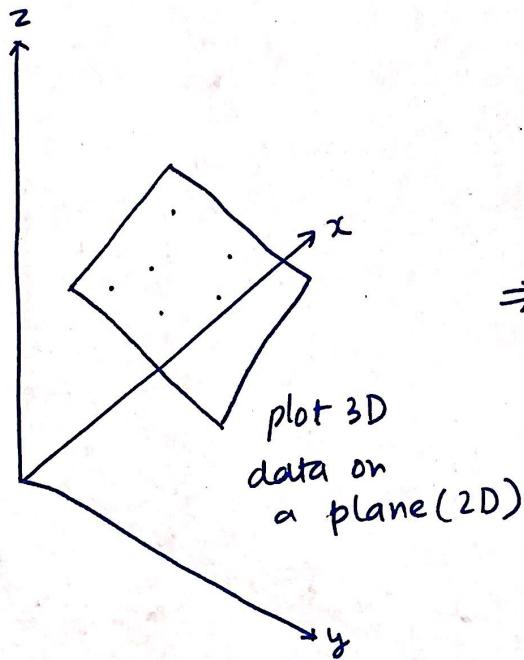
Reduce data from 2D to 1D



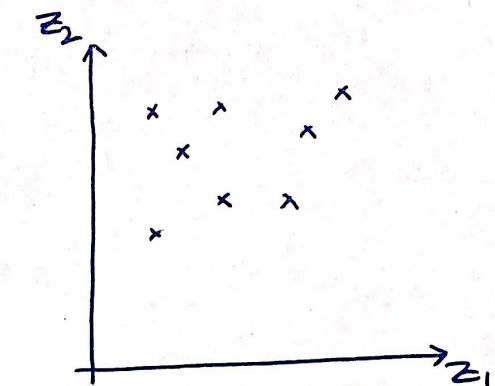
$$z^i \in \mathbb{R}^1$$

* approximation of the training set

Reduce data from 3D to 2D



$$x^i \in \mathbb{R}^3$$



$$z^i \in \mathbb{R}^2$$

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

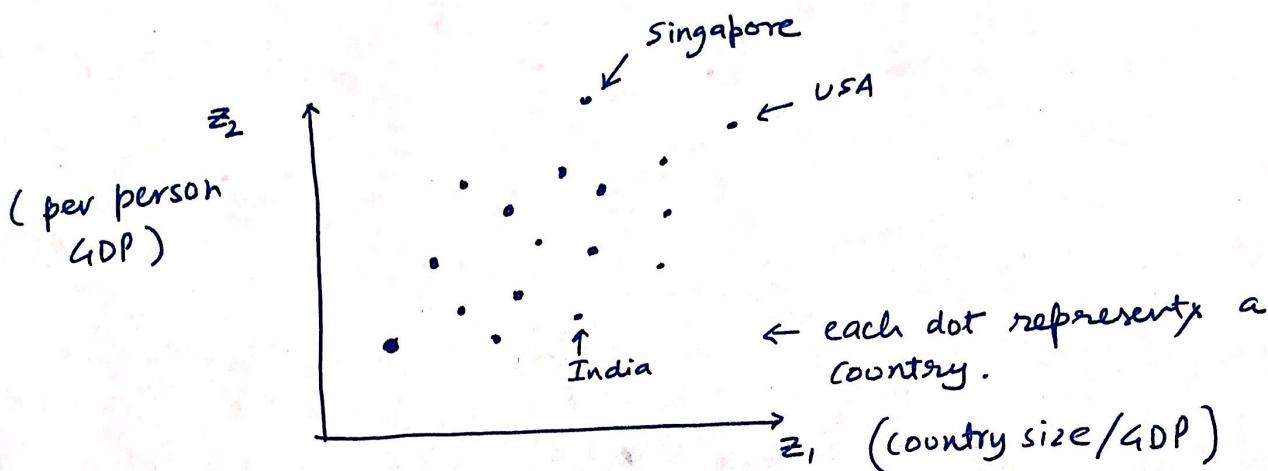
Dimensionality Reduction : Data Visualization

$$x_i = \{ \text{GDP, per capita, HDI, life exp., Gini coeff, ...} \}$$

~~#~~ $x^i \in \mathbb{R}^{50}$

↓ Reduce from 50D to 2D

$$z^i = \{ z_1, z_2 \} \quad z^i \in \mathbb{R}^2$$

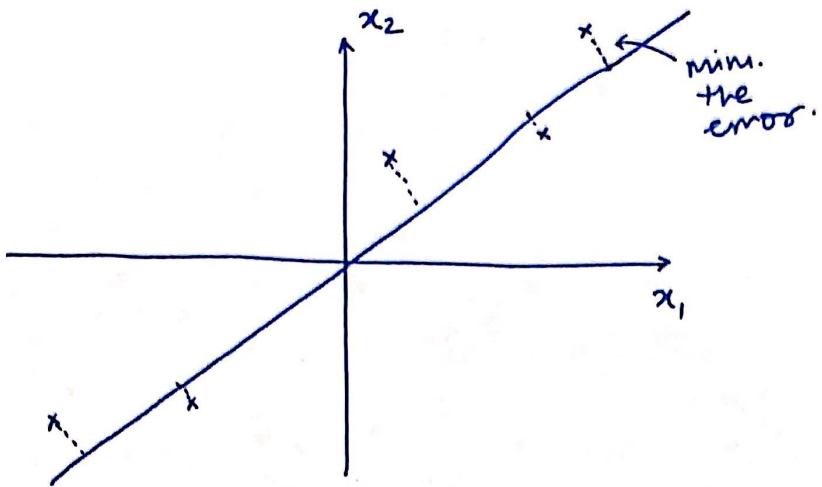


- # Suppose, I have a data $x^i \in \mathbb{R}^n$. In order to visualize it, we apply dimensionality reduction & get $z^i \in \mathbb{R}^k$ where $z^i \in \mathbb{R}^k$.

① $K \leq n$

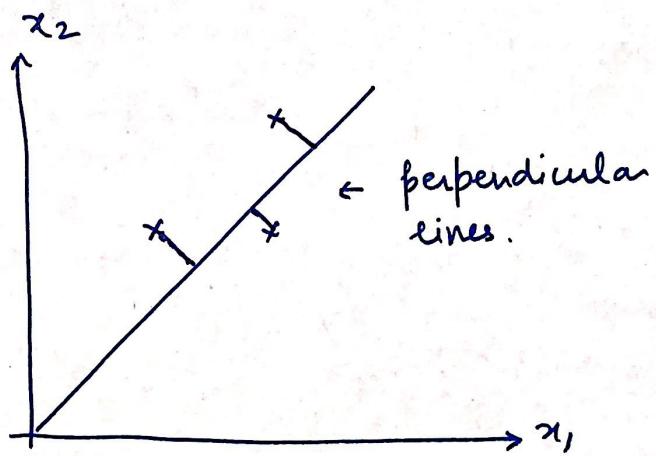
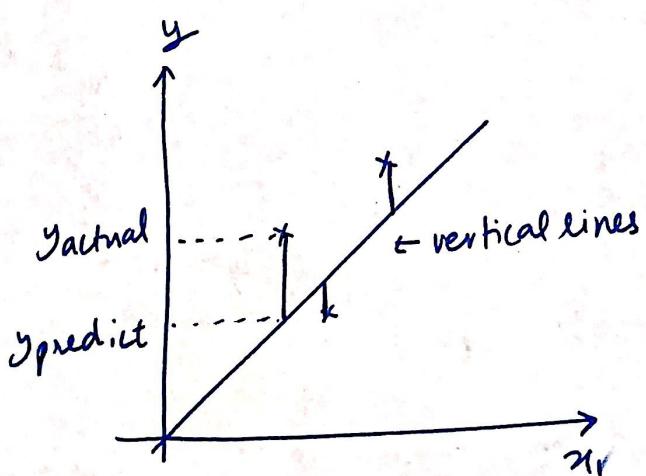
- ② $K=2$ or 3 , because don't have ways to visualize higher dimension data.

Principal Component Analysis (PCA)



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

PCA is not linear regression



PCA Algorithm :

- ① Data preprocessing (feature scaling / mean normalization)

Training set: $x^1, x^2 \dots x^m$

Step 0:

- ① compute mean, μ_j for each feature.
- ∴ $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^i$
- ② Replace x_j^i with $x_j^i - \mu_j$

- ② PCA: Reduce data from n-dimension to k-dimension

Step 1: Covariance Matrix

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

sigma

Step 2: Compute "eigenvectors" of matrix Σ

$$[U, S, V] = \text{svd}(\Sigma) \quad \begin{matrix} \text{\% svd:} \\ \text{\% singular} \\ \text{\% value} \\ \text{\% decomposition} \end{matrix}$$

octave
command

$$U = \left[\underbrace{\begin{array}{c|c|c|c} 1 & & & \\ \hline u^1 & u^2 & \dots & u^m \end{array}}_K \right]$$

$$U \in \mathbb{R}^{n \times n}$$

↪ Take k vectors.

(104)

$$x \in \mathbb{R}^n \longrightarrow z \in \mathbb{R}^k$$

step 3:

$$z^k = \begin{bmatrix} | & | & | \\ u^1 & u^2 & \dots & u^k \\ | & | & | \end{bmatrix}^T * x$$

$\underbrace{\quad\quad\quad}_{U_{\text{reduce}}}$

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

$\hookrightarrow \underbrace{(k \times n)}_{k \times 1} \otimes \underbrace{(n \times 1)}_{k \times 1}$

$$\therefore z^k \rightarrow \begin{bmatrix} z_1 \\ z_2 \\ \vdots \end{bmatrix}_{k \times 1}$$

Reconstruction from compressed representation

$$z = (U_{\text{reduce}})' * x$$

: compression

$$x_{\text{approx}} = \underbrace{U_{\text{reduce}}}_{n \times k} * \underbrace{z}_{k \times 1}$$

: De-compression

choosing the number of principal components.

k for PCA

→ Avg. squared projection error : $\frac{1}{m} \sum_{i=1}^m \|x^i - x_{\text{approx}}^i\|^2$

→ Total Variation in the data : $\frac{1}{m} \sum_{i=1}^m \|x^i\|^2$

→ Typically, choose k to be smallest value so that

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

"99% variance is retained"

(ratio ≤ 0.05 , means 95% is retained)

Algorithm: (inefficient)

step 1: Try PCA with $K=1$

step 2: Compute & reduce, z^i , x_{approx}^i

step 3: compute & check if

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

iterate until step 3 is satisfied.

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

$$S = \begin{bmatrix} S_{11} & & & \\ & S_{22} & \dots & \\ & & \ddots & \\ & 0 & & S_{nn} \end{bmatrix}_{n \times n}$$

square matrix
of $n \times n$ size

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^i - x_{\text{approx}}^i\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^i\|^2} \Rightarrow 1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}}$$

for a given k .

$$\left. \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99 \right.$$

} iterate

check

\Downarrow
pick smallest value of k . **

Advice for applying PCA

→ Supervised learning speedup

$$(x^1, y^1) (x^2, y^2) \dots (x^m, y^m)$$

$$x^i \in \mathbb{R}^n \quad [\text{where } n \text{ is large}]$$

eg: computer vision

100x100 pixel
image

$$n = 10,000$$

- Extract inputs:

unlabeled inputs dataset:

$$\underbrace{x^1, x^2, \dots, x^m}_{\downarrow \text{PCA}} \in \mathbb{R}^{10,000}$$

$$\downarrow$$

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

$$\begin{array}{c} x^{\text{test}} \\ \downarrow \\ z^{\text{test}} \\ \downarrow \end{array}$$

- New Training set:

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

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T z}}$$

Note: Mapping $x^i \rightarrow z^i$ should be defined by running PCA only on the training set (not on x_{cv}^i or x_{test}^i)

Application of PCA• compression

- Reduce memory/disk needed to store data
 - Speed up learning algorithm. ** (good!)
- choose K by % of variance retain

• Visualization

choose $K=2$ or $K=3$

Bad use of PCA: To prevent overfitting

✓ use regularization instead

✗ less features, less likely to overfit (! Bad)
 { does not look at y^i } thus throwing away
 some useful information

Note: Before implementing PCA, first run whatever you want to do with the original/raw data (x^i). Only if that doesn't do what you want, then implement PCA & consider using (z^i)

Octave application

K-means Algorithm :

if $x \in [m \times n]$
~~and~~ then,
 $K \in [k \times n]$

} if x^i has n -dimensions
then u_k will also have
 n -dimensions.

$\|x - u\| = \sqrt{\sum \{(\text{power}(x - u), 2)\}}$;