# WEEK 8 - Unsupervised learning
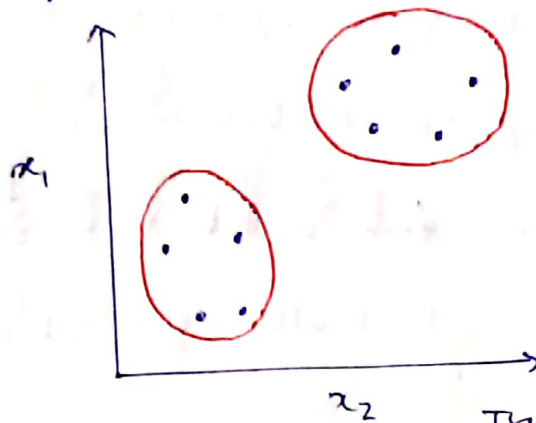
## Clustering algorithm



Training set: $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$
no labels $y^{(i)}$

## Application:

- Market segmentation
- Social network analysis
- Organize computing clusters
- Astronomical data analysis

## K-Means algorithm

- Randomly initialize 2 points (cluster centroids)
- Two steps in iteration:
  ○ Cluster assignment step
    Go through every example and see if its close to cluster centroid 1 or 2 and assign it to that group
  ○ Move centroid step
    Move the centroids to the new centroid of each group 1 & 2

- Here further iteration doesn't change the centroid positions

Input:

    - K (number of clusters)

    - Training set $\{x^{(1)}, x^{(2)}, \ldots, x^{(m)}\}$

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

k-means algorithm: ~~this called Lisfan~~)

Randomly initialize K cluster centroids $\mu_1, \mu_2, \ldots \mu_k \in$,

Repeat $\{$

**Cluster** **for** $i = 1$ to m

**Assignment**
**Step**
$$c^{(i)} := \text{index (from 1 to k) of cluster}$$
$$\text{centroid closest to } x^{(i)}$$
$$= \min_k \|x^{(i)} - \mu_k\|^2$$

**Move**
**Centroid**
**for** $k = 1$ to K

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

$\}$
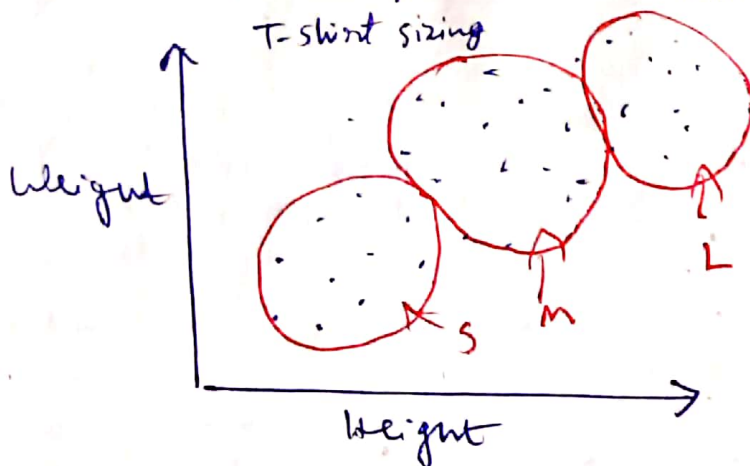$\text{fxa: } \mu_2 = \frac{1}{4}[x^{(1)} + x^{(5)} + x^{(4)} + x^{(10)}]$ if $c^{(1)}, c^{(5)}, c^{(6)}, c^{(10)} = 2$

$\uparrow$

Not a nested loop
    for 5
    5
    to 5
    3

K mean for non-seperated classes



T-shirt sizing
Weight
Weight

K - means optimisation objective

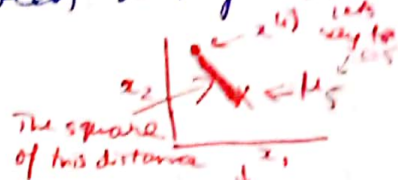$c^{(i)}$ = index of cluster $(1, 2, ..., k)$ to which example $x^{(i)}$ is currently assigned.

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

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

optimisation objective

The square of this distance

$$J(c^{(1)}, ..., c^{(m)}, \mu_1, ..., \mu_k) = \frac{1}{m} \sum_{i=1}^{m} \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

distortion function

$$\min_{c^{(1)}, ..., c^{(m)}, \mu_1, ..., \mu_k} J(c^{(1)}, ..., c^{(m)}, \mu_1, ..., \mu_k)$$

Random initialization

• Should have $k < m$
• Randomly pick $k$ <u>training examples</u>
• Set $\mu_1, ..., \mu_k$ equal to these $k$ examples

Problem: Sometimes when we randomly initialise cluster centroids, it can result in local optima

Solution: Run the algorithm on different values of randomly initialised cluster centroids and pick the best

For $i=1$ to $100\ \{$ ← Running it 100 times

→ Randomly initialise k-means

→ Run k-means, get $c^{(1)},\ldots, c^{(m)}, \mu_1,\ldots, \mu_K$

→ Compute cost function (distortion)
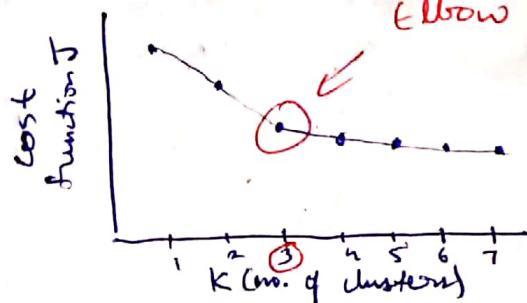$$J(c^{(1)},\ldots, c^{(m)}, \mu_1,\ldots, \mu_K)$$

$\}$

Pick clustering that gave lowest cost
$$J(c^{(1)},\ldots, c^{(m)}, \mu_1,\ldots, \mu_K)$$

Note, this is only good for problems where $k = 2-10$. For stuff like $k=100$, the first trial itself will give a good result

## Choosing number of clusters

- Common way is to visualise the data and manually decide

- Elbow method - called elbow cus the graph looks like a hand & we choose the elbow.



Elbow

cost function J

K (no. of clusters)

∴ K = 3

Note : Most times we can't make out the elbow so we don't use this much
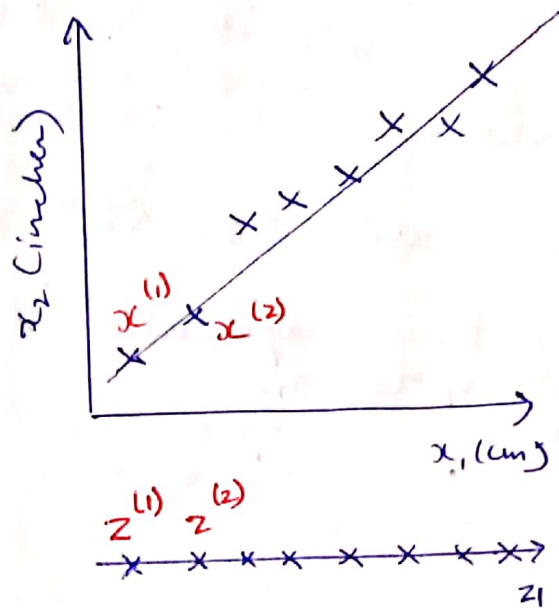
- Depending on problem.
  Eg: If we need to categorise t-shirt sizes into S, M, L. Then we have $k=3$

# Data compression

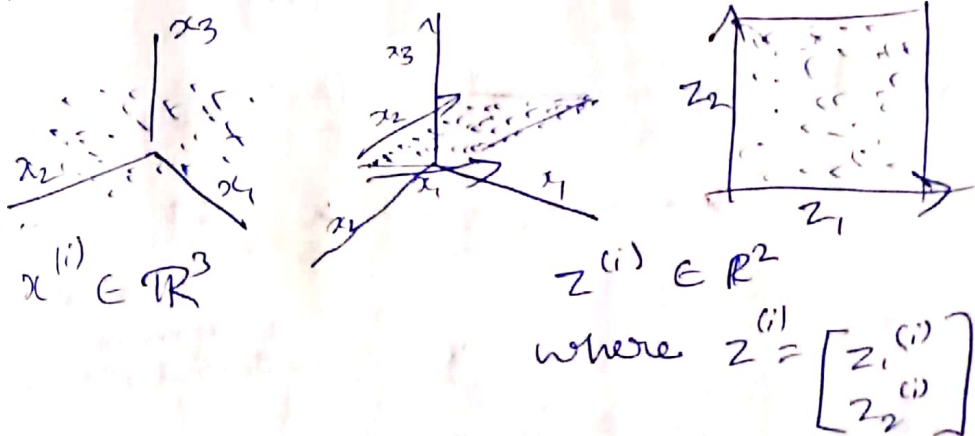Sometimes we have redundant features that we can compress into 1 feature.

Ex:

$x_2$ (inches)

$x^{(1)}$  $x^{(2)}$

$x_1$ (cm)

$z^{(1)}$  $z^{(2)}$

$z_1$

Reduce data from 2D to 1D

$$x^{(1)} \in \mathbb{R}^2 \;\Rightarrow\; z^{(1)} \in \mathbb{R}$$
$$x^{(2)} \in \mathbb{R}^2 \;\rightarrow\; z^{(2)} \in \mathbb{R}$$
$$\vdots$$
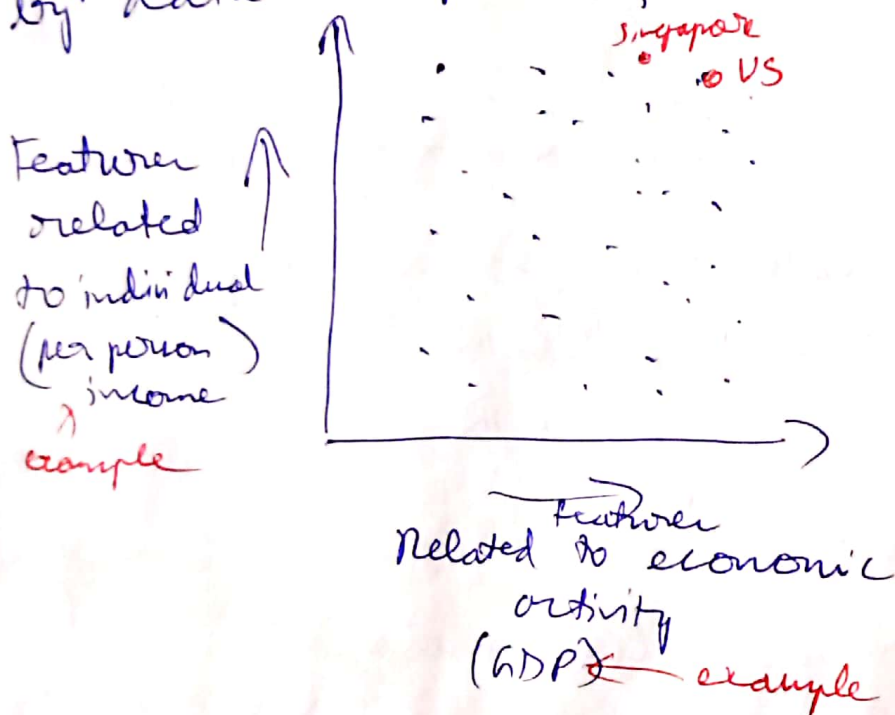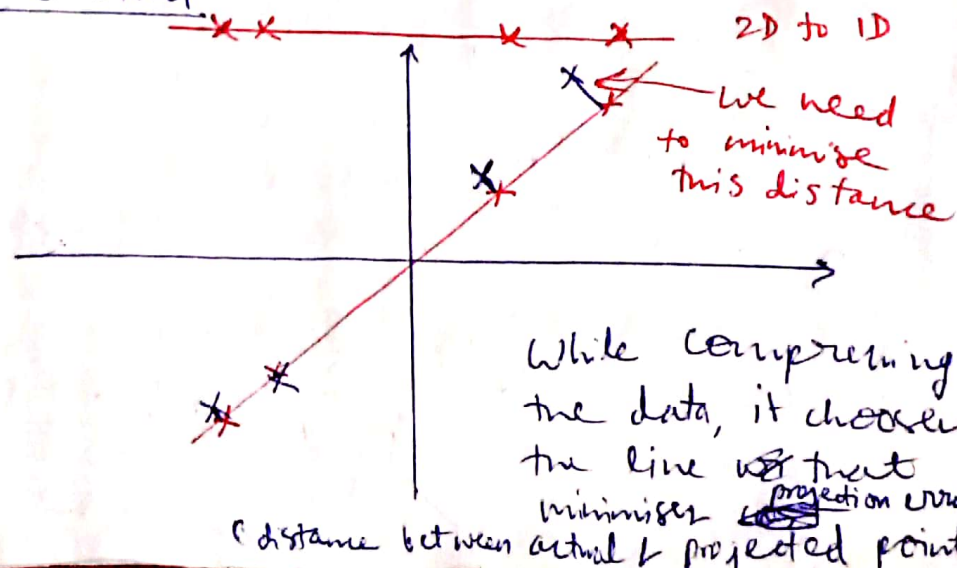$$x^{(m)} \in \mathbb{R}^2 \;\rightarrow\; z^{(m)} \in \mathbb{R}$$

In 3D to 2D also similar

$x_3$  $x_2$  $x_1$

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

$x_3$  $x_2$  $x_1$

$z_2$  $z_1$

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

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

# Data Visualisation

Consider a dataset of ~~information~~ 50 features (GDP, per person income, poverty, etc) of countries of the world. So how do you visualise for $x \in \mathbb{R}^{50}$.

You reduce data from 50D to 2D by selecting two important features by data compression

Feature related to individual (per person) income → example



Feature related to economic activity (GDP) → example

## Principal Component Analysis (PCA) problem formulation



2D to 1D

we need to minimise this distance

While compressing the data, it chooses the line that minimises projection error (distance between actual & projected points)

Reduce from 2D to 1D: Find a direction (a vector $u^{(i)} \in \mathbb{R}^n$) onto which to project the data so as to minimise the projection error..

Reduce from n-dimension to k-dimension: Find k vectors $u^{(1)}, u^{(2)}, \ldots, u^{(k)}$ onto which to project the data, So as to minimise the projection error
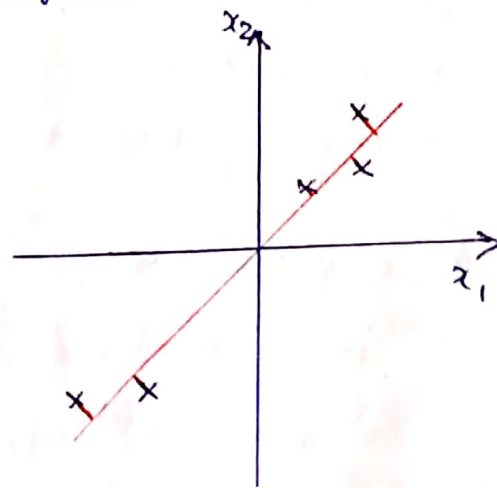
## PCA is not linear regression



Linear regression
The lines are vertical drops from the point to the line

We are predicting y

PCA
The lines are shortest distance between the point and the line

We are reducing dimensionality

# Data preprocessing

Training set: $x^{(1)}, x^{(2)}, \ldots, x^{(m)}$

Preprocessing (feature scaling & mean normalisation)

- $\mu_j = \frac{1}{m} \sum\limits_{i=1}^{m} x_j^{(i)}$  — Mean normalisation

  *Average of all $x_j$* (mean)

  Replace each $x_j^{(i)}$ with $x_j - \mu_j$
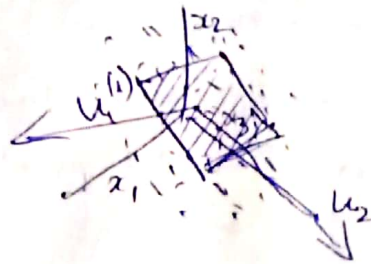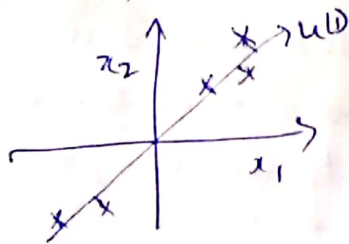
- If different features on different scales
  (e.g. $x_1 =$ size of house, $x_2 =$ no. of bedrooms)
  scale features to have comparable
  range of values.  — feature scaling

# PCA algorithm
$(2D \rightarrow 1D)$

We need to find a vector $u^{(1)}$ (1D)
or $u^{(1)}$ and $u^{(2)}$ $(3D \rightarrow 2D)$ that fits
the graph well.



Reduce data from $n$-dimensions to $k$-dimensions

Compute "covariance matrix"

sigma $\rightarrow \Sigma = \frac{1}{m} \sum\limits_{i=1}^{n} \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^{T}}_{1 \times n}$  ← $n \times n$ matrix

compute eigenvectors of matrix $\Sigma$:

$[U, S, V] = \text{svd}(\text{Sigma})$;  ← singular value decomposition ← can even use eig(Sigma)

$\underset{\text{n×n matrix}}{U} = \begin{bmatrix} u^{(1)} & u^{(2)} & u^{(3)} & \ldots & u^{(n)} \\ | & | & | & & | \end{bmatrix}$

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

Ureduce $\downarrow^{k}$

$$\begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \cdots & u^{(k)} \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times k}$$

We need to $\quad x \in \mathbb{R}^{n} \rightarrow z \in \mathbb{R}^{k}$

$$\therefore z = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \cdots & u^{(k)} \\ | & | & & | \end{bmatrix}^{T} x = \underbrace{\begin{bmatrix} -(u^{(1)})^{T}- \\ \vdots \\ -(u^{(k)})^{T}- \end{bmatrix}}_{k \times n} \underset{\underset{n \times 1}{\uparrow}}{x}$$

$$\underbrace{\qquad}_{k \times 1}$$

$\therefore$ Algorithm summary

$\rightarrow$ Preprocessing (mean normalization & feature scaling)

$\rightarrow$ sigma $= \dfrac{1}{m} \sum\limits_{i=1}^{m} (x^{(i)})(x^{(i)})^{T} \in$

In vector form

Sigma $= (1/m) * X' * X$;

given $X = \begin{bmatrix} - x^{(1)T} - \\ - x^{(2)T} - \\ - x^{(m)T} - \end{bmatrix}$

$\rightarrow [U, S, V] = svd(sigma)$;

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

$\rightarrow z =$ Ureduce' $* x$;

Reconstruction from compressed representation (1D to 2D)

"If we have $\qquad \times \quad \times \qquad \times \quad \times\times \longrightarrow z_1$ (1D)

how to get back 2D?

$\boxed{X_{approx} = U_{reduce} \cdot z}$

Since $z = U_{reduce}^{T} x$



notice that the point will be on the line in the approximation unlike

# Choosing k (no. of principal components)

Avg. squared projection error: $\frac{1}{m} \sum\limits_{i=1}^{m} \|x^{(i)} - x^{(i)}_{approx}\|^2$

Total variation in the data: $\frac{1}{m} \sum\limits_{i=1}^{m} \|x^{(i)}\|^2$

Typically choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum\limits_{i=1}^{m} \|x^{(i)} - x^{(i)}_{approx}\|^2}{\frac{1}{m} \sum\limits_{i=1}^{m} \|x^{(i)}\|^2} \leq 0.01$$

→ For this example we say
"99% of variance is retained"

For 0.05 → 95%.
For 0.1 → 90%. and so on

## Algorithm

→ try PCA with k=1    $\overleftarrow{\quad}$   k=2, k₁=3, k=4...    till it is ≤0.01

→ compute Vreduce, $z^{(1)}, z^{(2)}, \dots z^{(m)}, x^{(1)}_{approx}, \dots, x^{(m)}_{approx}$

→ check if $\dfrac{\frac{1}{m} \sum\limits_{i=1}^{m} \|x^{(i)} - x^{(i)}_{approx}\|^2}{\frac{1}{m} \sum\limits_{i=1}^{m} \|x^{(i)}\|^2} \leq 0.01$

→ Repeat for another k value

However the simpler method is to use

$[U, S, V] = $ svd (sigma)

Here $S = \begin{bmatrix} S_{11} & & & \cdots & 0 \\ & S_{22} & & & \\ 0 & & S_{33} & & \\ & & & \ddots & \\ & & & & S_{nn} \end{bmatrix}_{n \times n}$

$\dfrac{\{\text{selected } S_{ii} \text{ values}}{\{\text{All } S_{ii} \text{ Values}}$

then → Erea for k=3

For given k,

$$1 - \frac{\sum\limits_{i=1}^{k} S_{ii}}{\sum\limits_{i=1}^{n} S_{ii}} \leq 0.01 \quad \text{or} \quad \frac{\sum\limits_{i=1}^{k} S_{ii}}{\sum\limits_{i=1}^{n} S_{ii}} \geq 0.99$$

## Applying PCA

We can speed up supervised learning by reducing the no. of features without affecting the accuracy of the classification.

Note: We should only define mapping (learning Vsreduce) $x^{(i)} \to z^{(i)}$ by running PCA only on the training set. Later this mapping can be applied to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation & test sets.

Application:

- Compression:

Choose k by % of variance retained [- Reduce memory/disk needed to store data
- Speed up learning algorithm]

- Visualisation:
  - reduce dimensionality to $k=2$ or $k=3$ so it can be visualized

Bad use of PCA: To prevent overfitting

Since reducing features from $n$ to $k$ ($k < n$) results in overfitting happening less likely (since of less features), ppl use it. But don't use it, use regularization instead. (uz it throws away some information because it doesn't consider y labels)

Note: While designing an ML system, ppl first apply PCA on the data and then train logistic regression and test it. Don't do this. Run it on the raw data, later if it doesn't do what you want, then try implementing PCA and consider using $z^{(i)}$ instead of $x^{(i)}$