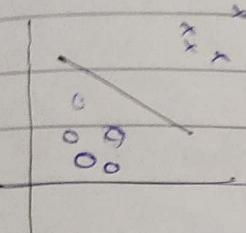


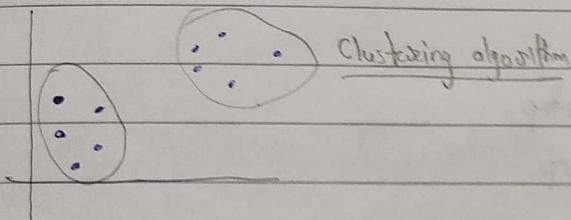
① Unsupervised learning Introduction

• Supervised learning



Training set $\rightarrow \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots (x^{(n)}, y^{(n)})\}$

• Unsupervised learning

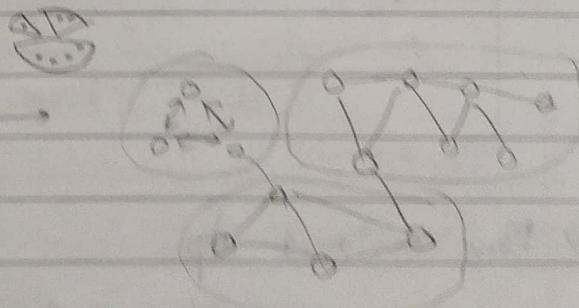


Training set $\rightarrow \{x^{(1)}, x^{(2)}, x^{(3)}, \dots x^{(n)}\}$

In unsupervised learning, what we do is we give this set of unlabeled training set to an algorithm and we just ask the algorithm to find some structure for us.

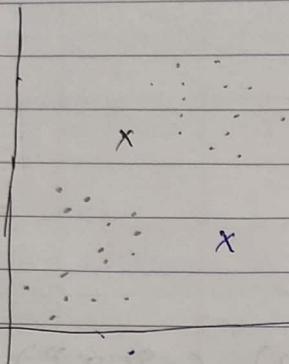
• Applications of clustering

- market segmentation
- social network analysis.
- organise computing clusters.
- astronomical data analysis.



--① K-Means Algorithm

A clustering algo.^{iterative}

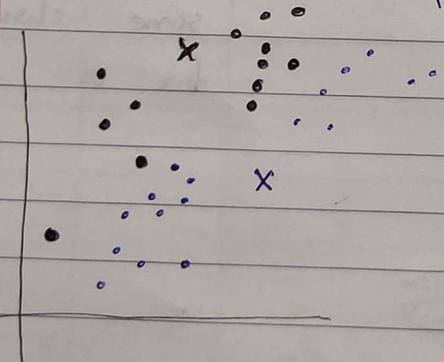


s1) Randomly initialize two centroids. centroid.

2 steps →

- cluster assignment step.
- Move centroid step.

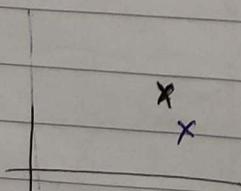
s2) Cluster assignment step.



color each depending its closer
to blue or black cluster
centroid.

s3) Move centroid step.

Move centroid to average of
points coloured the same color.



s4) (Repeat s2, s3) ^{my companion} ^{Repeat}

cluster
assignm
st

move
cent
sta

What if cluster center is assigned with zero points.

- Eliminate their cluster centroid
- $k-1$ clusters

- Randomly reinitialize the cluster centroids
- k clusters.

127

After many iterations,

K-Means converged.

centroids will not change,

colours of points will not change.

K-means algorithm.

Input :

- K (no of clusters)
 - Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

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

~~x_{0+10}~~ ← we will not use this
so (loop $x_0=1$ convention)

Randomly initialize k cluster centroids $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$

$\times \mu_1$ $\times \mu_2$

Repeat {

for i = 1 to m

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

cluster

assignment

step

Same

(but we use this) \rightarrow

[Q12] Hence pick $c^{(i)} = \text{cluster centroid with smallest squared distance to training e.g. } x^{(i)}$

$$\min_k \|x^{(i)} - u_k\|^2$$

mc ve

centooid

stop

my companion

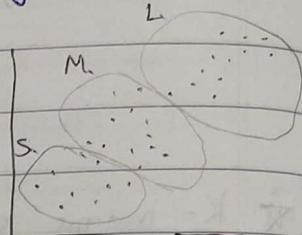
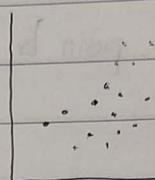
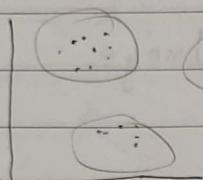
$\mu_{lk} := \text{average of points assigned to cluster } k$

$$M_2 = \frac{1}{4} \left[x^{(1)} + x^{(2)} + x^{(3)} + x^{(4)} \right] \text{EIR}$$

$C^{(1)} = 2, C^{(2)} = 2, C^{(3)} = 2, C^{(4)} = 2$
 This e.g. were assigned
 to clusters
 centroid 2.

* K-means for clustering.

T-shirt sizing



That's how k-means works
(Market segmentation like.)

"The number of clusters should be decided first."

Ex. 4 X

3 X

Distance should be (of all points) minimum

or of from

Math - Explan

Explain

date

-@ Optimization Objective - - -

upto now each algo had optimization objective/cost fn which we were trying to minimize.

k-means also has cost fn.

cost fn Useful for →

i) To debug learning algo. (Learning right/wrong)

ii) How k-means avoid local minima by finding better cost fn.

Notations

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

$x^{(i)}$ is currently assigned.

μ_k = cluster centroid k ($\mu_k \in \mathbb{R}^n$) $k \in \{1, 2, \dots, K\}$

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

$x^{(i)} \rightarrow 5$ $x^{(i)}$ is assigned to cluster no. 5.

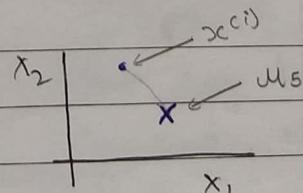
$c^{(i)} = 5$

$\mu_{c(i)} = \mu_5$

Optimization objective.

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

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



Distortion cost function/

Distortion of k-means algo.

★ K-means algo.

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

Repeat {

for $i=1$ to m

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

- Cluster assignment step
- minimize $J(\cdot)$ w.r.t $c^{(1)}, c^{(2)}, \dots, c^{(m)}$
(holding $\mu_1, \mu_2, \dots, \mu_K$ fixed)

for $k=1$ to K • Move centroid.

$\mu_k :=$ average (mean) of points assigned to cluster k .

μ_1, \dots, μ_K

creates 12 sets $c^{(1)}, \dots, c^{(12)}$ & μ_1, \dots, μ_{12}

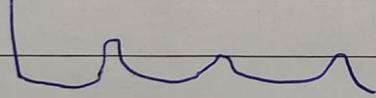
balanced good and bad

$\mu_1 = 0.16$

$\mu_2 = 0.33$

$\mu_3 = 0.50$

$J(\cdot)$



No of iterations.

→ It is not possible for the cost fun to sometimes increase. These must be bug in the code.

--- @ K-means Random Initialization ---

- K-means initialization (how)

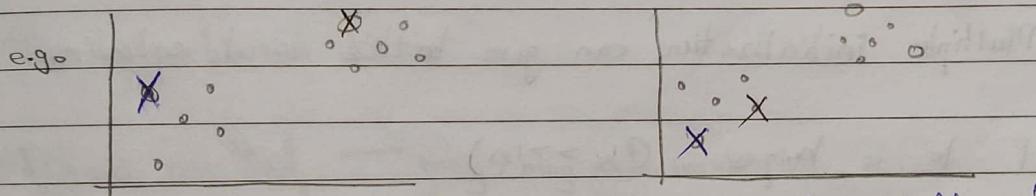
- How to we avoid local optimum.

* Random initialization.

- Should have $K = m$.

- Randomly pick K training examples.

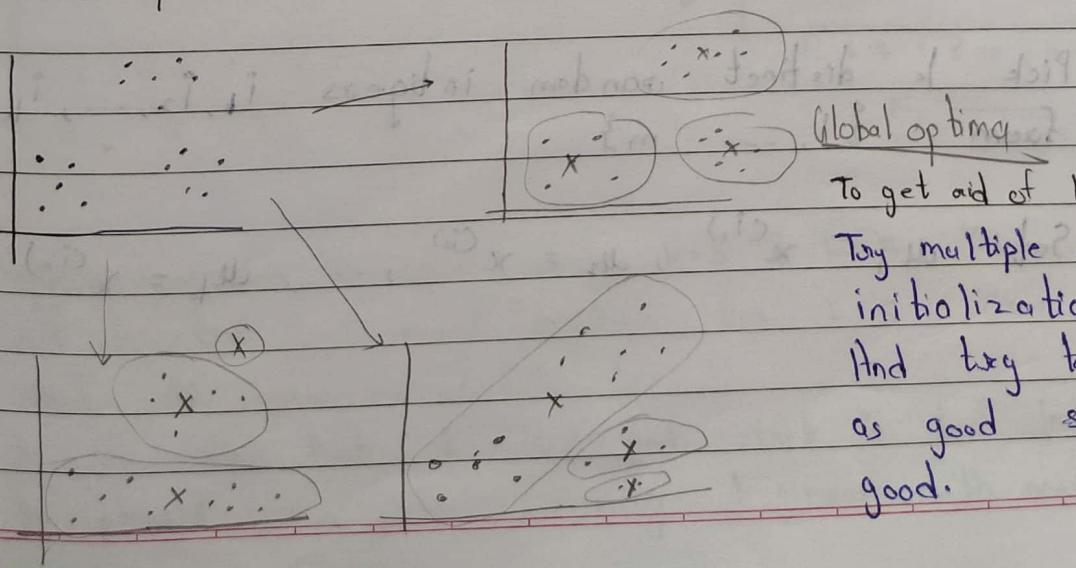
- Set μ_1, \dots, μ_K equal to these K examples.



K-means can end up converging to different solutions depending on exactly how clusters were initialized,
(so random initialization)

K-mean can actually end up at local optima.

- * local optima. local optima vectors to local optima of distortion fun.



To get rid of local optima

Try multiple random initialization

And try to get as good solution as good.

Random initialization.

for $i=1$ to $100 \rightarrow 50-100$ Compute for 50-100 times.

{ Randomly initialize centers

Run k-means. Get $c^{(1)}, \dots, c^{(k)}$ initialize multiple times.

(compute cost J^*)

$J(c^{(1)}, \dots, c^{(k)})$

}

Run code for multiple /

initialize multiple times.

Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(k)}, u_1, \dots, u_k)$

If k is small ($k=2-10$)

Multiple initializations can give better local optimum.

If k is large ($k > 10$)

Multiple initializations will not make huge difference.

First random initialization will also be close to global optima.

Q1 Recommended way to initialize k-means.

→

Pick k distinct random integers i_1, i_2, \dots, i_k from $\{1, \dots, m\}$

Set $u_1 = x^{(i_1)}, u_2 = x^{(i_2)}, \dots, u_k = x^{(i_k)}$

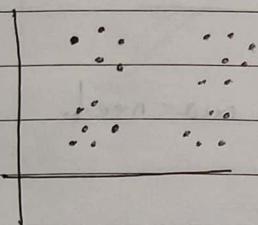
-- @ Choosing the Number of clusters --

No good method to choose automatically.

We have to choose by visualization / exp. of clustering algo manually.

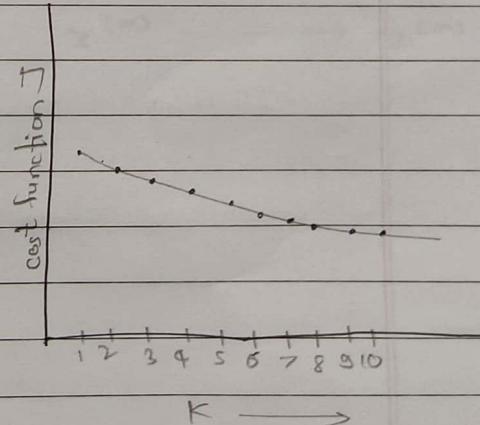
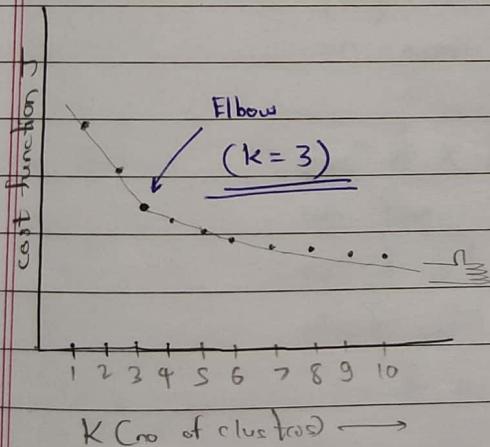
By hand (manually)

- What is value of k -- here it might be 2/4.



- Choosing value of K .

#Elbow method → may / may not work



- Q1 Suppose you run k-means using $k=3$ & $k=5$. You find that the cost function J is much higher for $k=5$ than for $k=3$. What can you conclude?

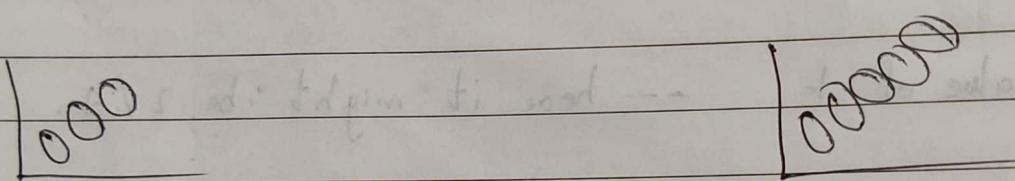
→ In the run $k=5$, k-means got stuck in bad local minima. You should ~~re~~ try running k-means with multiple times.

* Choosing the value of k.

Sometimes, you are running k-means to get clusters to use for some later/downstream purpose. Evaluate the k-means based on a metric for how well it performs for that later purpose.

$k=3$

$k=5$



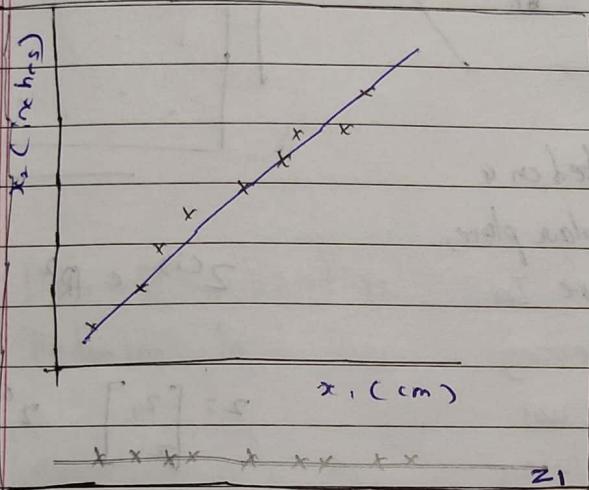
Take value of k, according to our needs.

optimization
to train faster.

Motivation I: Data Compression

Dimensionality reduction

Out of 100, let say x_1, x_2 are features (coincident features) & we want to reduce features to train faster.



Reduce dat from 2D to 1D

Some how using only one feature, which has features of x_1, x_2 both.

$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}$$

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

a line (best fit)

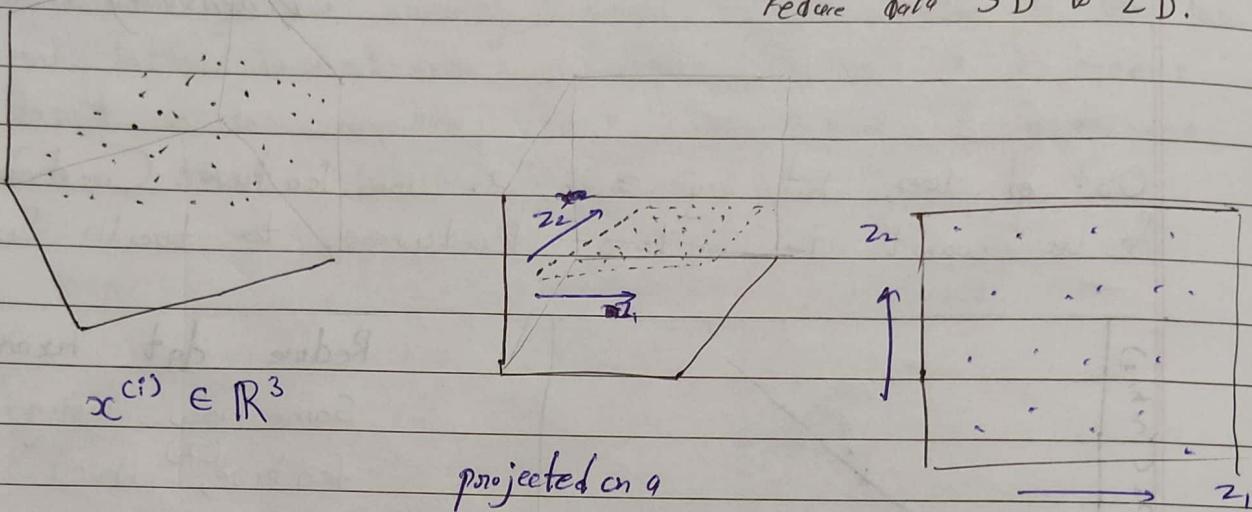
projection of x
on green line

$$x^{(m)} \rightarrow z^{(m)}$$

$z_{1,2,..} \rightarrow$ pos' of x from
on line,
so, it is feature
which will replace
 $x_1 \& x_2$

• Data compression.

Reduce data 3D to 2D.



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

projected on a

particular plane,

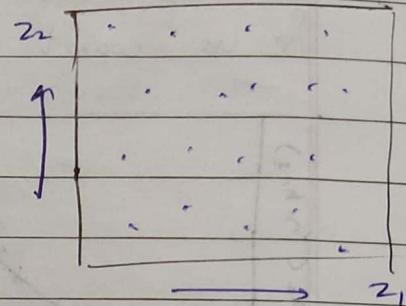
A plane In

which every

example was

almost lying

just little
bit distortion



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

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

vectorized

form.

ith example

or After applying dimensionality reduction on $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

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

we will get

lower dimension dataset $\{z^{(1)}, \dots, z^{(m)}\}$ where $z^{(i)} \in \mathbb{R}^k$
($k \leq n$)

Motivation II : Visualization

• Dataset.

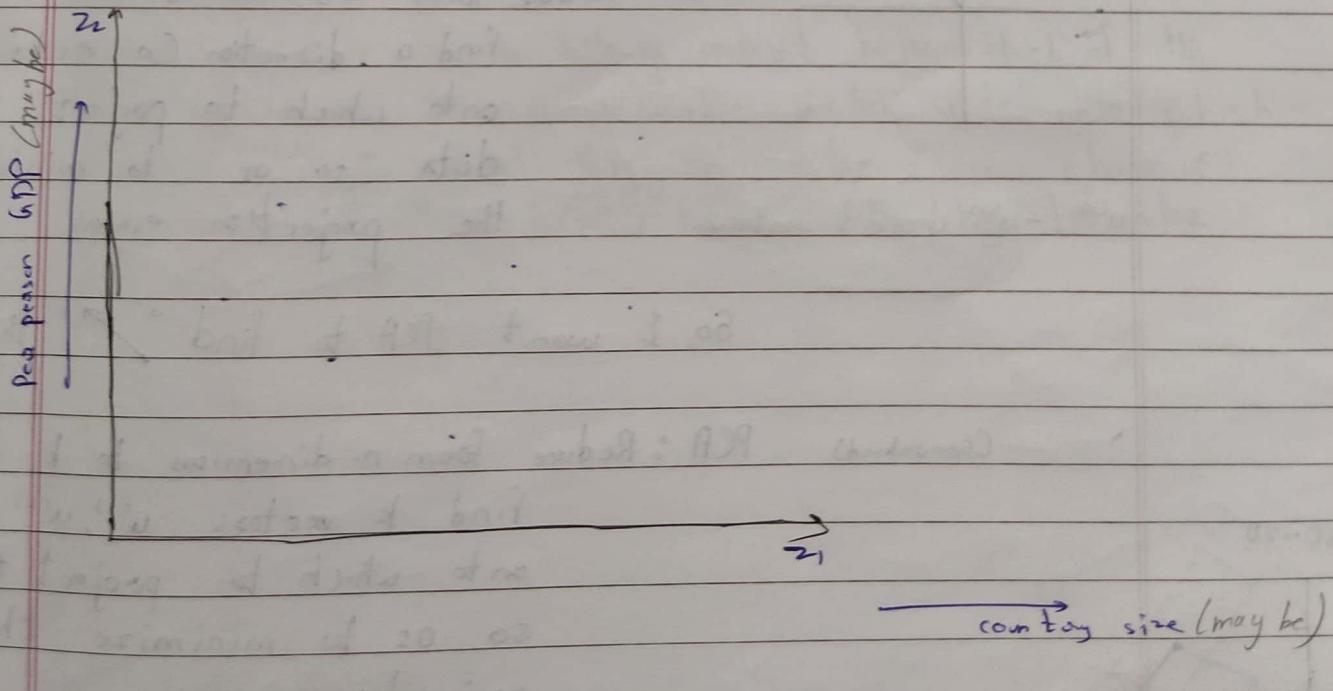
country x_1 GDP x_2 Per capita GDP x_3 life expectancy x_4 x_5 \dots $x_{C1} \in \mathbb{R}^{S0}$

• Data Visualization.

→ Reducing to

country z_1 z_2 $z_{C1} \in \mathbb{R}^2$

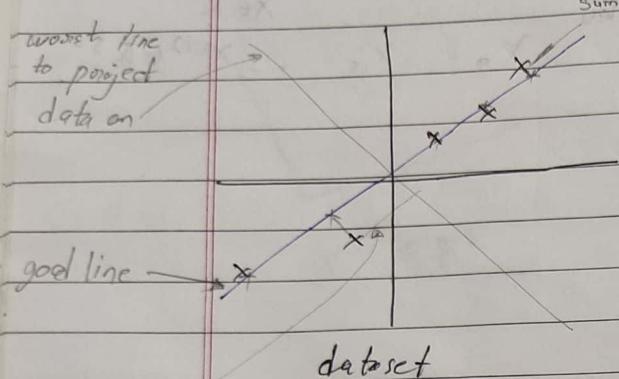
Reduce data from 50 D to 2 D.



Principal Component Analysis Problem formulation.

→ used for dimension reduction

* PCA problem formulation.



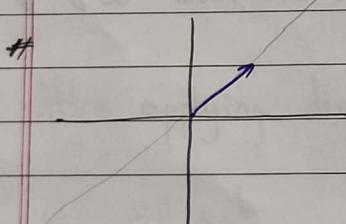
Sum of squares of this line must be minimized
(the dataset & the i.e.
projection length)
 $x \in \mathbb{R}^2$ (projection error)

We want to reduce dimension, so,
we will find projection of dataset
on a line.

Projection
error is
high.

before applying PCA, do first mean normalization at feature scaling. on x_1, x_2 .

PCA will choose " / " this line



PCA: Reduce from 2D to 1D :

find a direction (a vector $u^{(0)} \in \mathbb{R}^2$)
onto which to project the
data so as to minimize
the projection error.

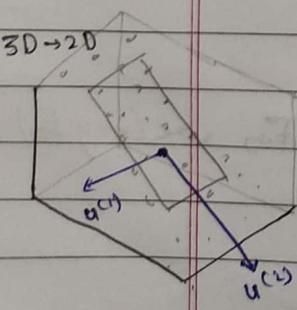
$u^{(0)}/-u^{(0)}$ doesn't matter.

So, I want PCA to find " ↗ " this vector

(Generalized)

PCA : Reduce from n-dimensions to k-dimensions

Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$
onto which to project the data,
so as to minimize the
projection error.

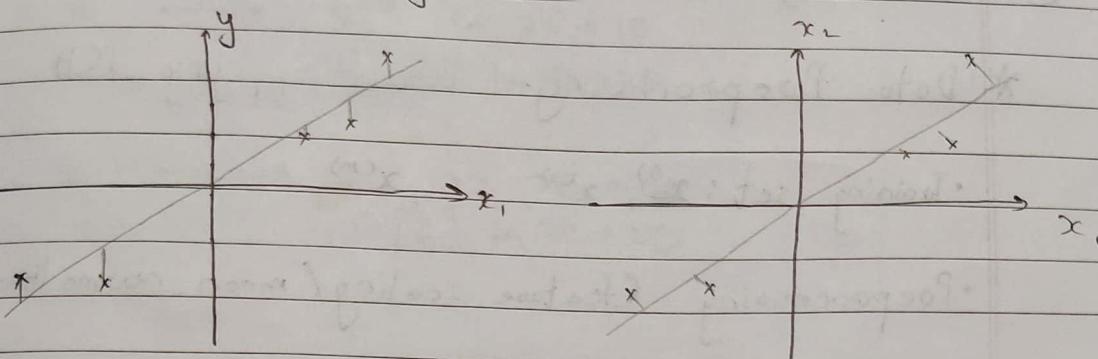


OR

Project data onto the linear subspace
spanned by set of k vectors.

Finding k directions (instead of 1) onto which to

PCA is not linear regression.



LR

PCA

- Vertical dist

$$(x_1^{(1)}, x_2^{(1)}, \dots, y^{(1)}), \dots$$

- \perp^a dist.

$$x_1, x_2, \dots, x_n$$

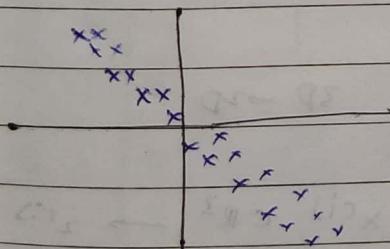
- To predict.

y is special variable.

- All variable are treated equally.

[Even in 3D we are treating all feature sim'ls]

Suppose you run PCA on following dataset below. Which of the following would be a reasonable vector $u^{(1)}$ onto which to project the data? (By convention, we choose $u^{(1)}$ so that $\|u^{(1)}\| = \sqrt{(u_1^{(1)})^2 + (u_2^{(1)})^2}$, the length of vector $u^{(1)}$, equals 1.).



✓ $u^{(1)} = \begin{bmatrix} -1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$

① PCA algorithm

* Data Preprocessing: Before applying PCA

- Training set: $x^{(0)}, x^{(1)}, \dots, x^{(m)}$

- Preprocessing (feature scaling / mean normalization):

$$\bar{u}_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each $x_j^{(i)}$ with $x_j - \bar{u}_j$

If different features on different scales

(e.g. x_1 = size of house

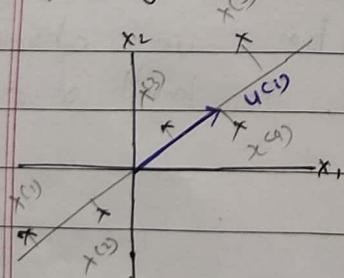
x_2 = no. of bedrooms)

Scale features to have comparable range of values.

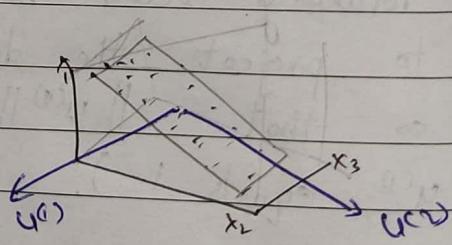
$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \bar{u}_j}{s_j}$$

max - min / standard deviation

PCA algo

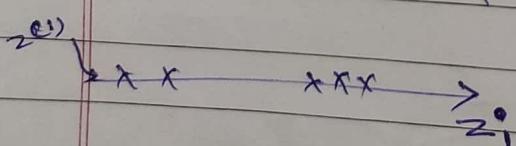


2D \rightarrow 1D



3D \rightarrow 2D

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



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

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

We have to find :

$$u \longrightarrow u^{(1)} \text{ in } 2D \rightarrow 1D$$

$$u^{(1)}, u^{(2)} \text{ in } 3D \rightarrow 2D$$

$$z \longrightarrow z_1 \text{ in } 2D \rightarrow 1D$$

$$z = [z_1 \ z_2] \text{ in } 3D \rightarrow 2D$$

* PCA algo

before this • Reduce data from n-dimensions to k-dimensions

do

compute

Compute "covariance matrix" :

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

store this variable called Sigma.

Compute "eigen vectors" of matrix Σ :

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

/ eig(Sigma)

we need only U

single value decomposition.

Covariance matrix always satisfies a mathematical property
symmetric positive definite

Sigma $\rightarrow n \times n$ matrix

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

k

To make k Dimensions pick first k vectors. $u^{(1)}, \dots, u^{(k)}$

PCA algo.

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

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

$\underbrace{\quad\quad\quad}_{K}$

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

$$Z^{(i)} = \begin{bmatrix} | & & | \\ u^{(1)} & \dots & u^{(k)} \\ | & & | \end{bmatrix}^T X^{(i)} = \begin{bmatrix} -(u^{(1)})^T & - \\ -(u^{(k)})^T & - \end{bmatrix} X^{(i)}$$

$\underbrace{\quad\quad\quad}_{n \times k}$ $\underbrace{\quad\quad\quad}_{k \times n}$ $\underbrace{\quad\quad\quad}_{n \times 1}$

U_{reduce}

$$z \in \mathbb{R}^k$$

★ PCA algo summary.

After mean normalization (ensure every feature has zero mean) and optionally feature scaling.

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

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

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

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

$$U_{\text{reduce}} = U[:, 1:k];$$

$$z = U_{\text{reduce}}^T * x;$$

$$x \in \mathbb{R}^n \quad x \neq 1$$

All In PCA we obtain $z \in \mathbb{R}^k$ from $x \in \mathbb{R}^n$ as follows.

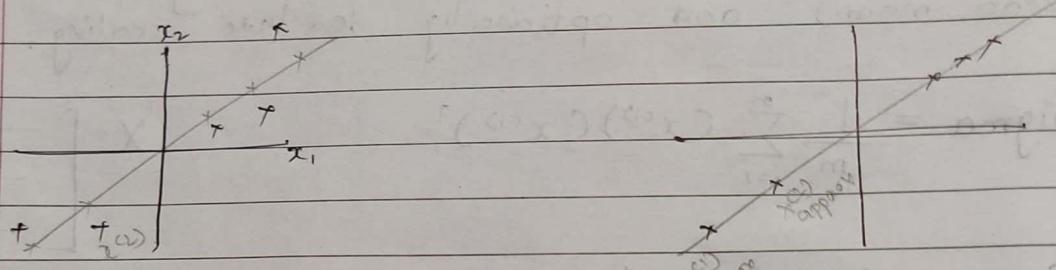
$$z = \begin{bmatrix} 1 & 1 & 1 \\ u^{(1)} & u^{(2)} & \cdots u^{(k)} \\ 1 & 1 & 1 \end{bmatrix}^T x = \begin{bmatrix} (u^{(1)})^T \\ \vdots \\ (u^{(k)})^T \end{bmatrix} x$$

$$\checkmark z_j = (u^{(j)})^T x$$

- We haven't proved mathematically PCA of getting minimized squared distance.

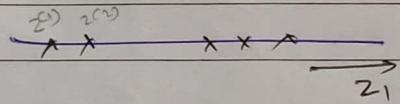
--① Reconstruction from Compressed Representation -

3D to 2D now 2D \rightarrow 3D going back in this video.

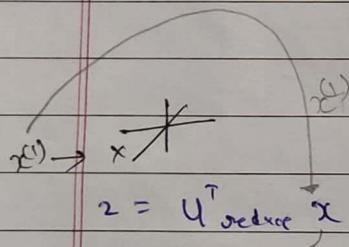


$$z = U_{\text{reduce}}^T x$$

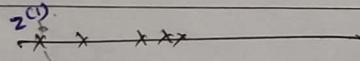
$$z \in \mathbb{R} \rightarrow x \in \mathbb{R}^2$$



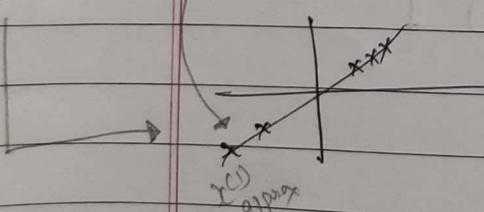
$$x_{\text{approx}} = U_{\text{reduce}} z$$



$$z = U_{\text{reduce}}^T x$$



$$x_{\text{approx}} = U_{\text{reduce}} z^{(1)}$$



$$z^{(1)}_{\text{approx}} = U_{\text{reduce}} z^{(1)}$$

Suppose we run PCA with $k=n$ so that the dimension of the data is not reduced at all. (This is not useful in practice but is a good thought exercise.)

Recall that the percent / fraction of variance retained is given by : $\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}}$. Which of the following will be true?

✓ U_{reduce} will be an $n \times n$ matrix

✓ $x_{\text{approx}} = x$ for every example x

✓ The percent of variance retained will be 100%.

• We have that $\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 1$

— (c) Choosing value of Principal components.

Choosing k

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

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

→ On average how far away my vectors from origin (training ex.)

Typically, choose k to be smallest value so that.

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

99% variance is retained.

95%.

90%.

85 above is good but try 95% 100%.

* Choosing k (number of principal components)

Algorithm:

• Try PCA with $k=1$

• Compute $U_{\text{reduce}}, z^{(1)}, \dots, z^{(m)}$, $x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$.

• Check if

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

How we deal with algorithm.

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

$$\Sigma = \begin{bmatrix} S_{11} & 0 & & \\ 0 & S_{22} & 0 & 0 \\ & 0 & S_{33} & 0 \\ & & 0 & \ddots & S_{nn} \end{bmatrix}$$

For given k ,

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01 \quad \leftarrow \text{check this cond.}$$

find minimum k

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

which satisfies

Q Previously we said that PCA chooses a direction $u^{(1)}$ (for k directions $u^{(1)}, \dots, u^{(k)}$) onto which to project the data so as to minimize the squared projection error. Another way to say is PCA tries to minimize

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

① Advice for applying PCA: - - -

★ Supervised learning speed up

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

$$\text{suppose } x^{(1)} \in \mathbb{R}^{16,000}$$



Extract inputs.

$$\text{unlabelled dataset: } z^{(1)}, z^{(2)}, \dots, z^{(n)} \in \mathbb{R}^{16,000}$$

↓ PCA

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

New training set:

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

Usecase

$$\text{accordingly hypothesis } h_0(z) = \frac{1}{1 + e^{-\alpha^T z}}$$

Note:- Mapping $z^{(1)} \rightarrow z^{(1)}$ should be defined by running PCA only on training set. This mapping can be applied as well as to the examples $z_{cv}^{(1)}$ & z_{test} in the cross validation set.

Applications of PCA.

• Compression.

- Reduce memory/disk needed to store data.
 - Speed up learning algo
- choose k by % of variance retain.

• Visualization. $\rightarrow k=2 / k=3$

Bad use of PCA:

E Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the number of features, less likely to overfit

~~Bad application of PCA~~

This might work OK, but isn't a good way to address overfitting. Use regularization method instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

In regularization knows about label y , PCA don't know about y , PCA may throw important information.

Note: E

Before implementing PCA, first try running whatever you want to do with the original/slow data $x^{(i)}$. Only if that doesn't do what you want then implement PCA & consider $z^{(i)}$.

- Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$.
- Get Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$.
- Train model on $\{(z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots, (z^{(m)}, y^{(m)})\}$
- Test on test set: Map $x_{\text{test}}^{(i)}$ to $z_{\text{test}}^{(i)}$. Run $h_{\theta}(z)$ on $\{(z_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}), \dots, (z_{\text{test}}^{(m)}, y_{\text{test}}^{(m)})\}$