

My
Wing
know
our score.

Unsupervised Learning Intro

Clustering - unsupervised learning algo

↳ learn from unlabeled data

Training set:

$$\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$$

In unsupervised learning

we give unlabelled
training set to an

algorithm & ask the
algo to find some structure

within the data. The structure found in the
dataset ~~are~~^{is} shown as clusters.

Clustering algo - an algo that finds clusters

Applications of clustering

- Market segmentation
- social network analysis

~~etc~~

- organise computing clusters
and etc

which computers
communicate between
each other the most
& hence change
network according
to net

K-Means Algorithm

↳ most ~~pop~~ popular, widely used clustering algo (an iterative algo)

Steps of k-means :

- 1) randomly initialize ~~#~~ points called cluster centroids - the # of these points depend on the # of clusters you wish to group your data into
- 2) Iterate over a, b
 - a) cluster assignment - each data point is assigned to a cluster centroid depending on how close the data point is from the cluster centroid (data point assigned to closest cluster centroid)
 - b) move centroid - the cluster centroid location is changed to the average of all of the points assigned to that cluster centroid
 - c) Iterate until no data points ~~switch~~ from one cluster to the other in assignment & ^{new} average position of cluster centroid does not change.

↳ Algorithm converges

$k \in \{1, 2, \dots, K\}$

Date: _____

Input to algorithm:

- k (number of clusters)
- training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

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

This convention is only present in supervised learning, not in unsupervised learning

K-means algorithm

capital "K" - total # of clusters
lowercase "k" - indexing into # of clusters, count variable

Randomly initialize K cluster centroids

$\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

$c^{(i)}$ holds closest cluster index of $x^{(i)}$

Repeat {

cluster assignment step [for $i=1$ to m
 $c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$]

$$c^{(i)} = \min_k \|x^{(i)} - \mu_k\|^2$$

adding square is simply convention
index into all possible clusters

cycle continued...

Date: _____

Move
Centroid

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

?

{ suppose $x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$ assigned to
cluster 2 then:
 $c^{(1)} = 2, \cancel{c^{(5)}} \cancel{c^{(6)}} c^{(10)} = 2, c^{(6)} = 2, c^{(10)} = 2$
so the move centroid step is then:
 $\mu_2 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}] \in \mathbb{R}^n$

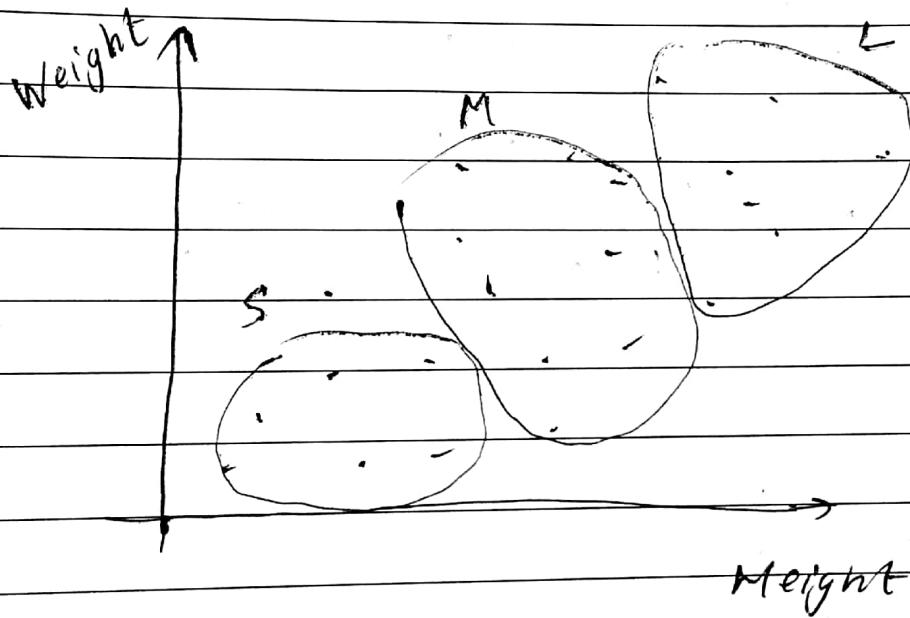
If a centroid has 0 points assigned to it
then how is the move centroid step done?

- 1) oftenly the ~~curr.~~ cluster centroid is dropped
- 2) seldomly (if ~~K~~ clusters are really needed)
then the cluster centroid is randomly initialised again

Date:

K means for non-separated clusters

Application scenario: a T-shirt manufacturer has gathered the weight & height of people within a certain population where he/she plans on selling T-shirts. Find the size of S, M, L that he can use to better fit the population.



Apply the clustering algo on this data set with 3 centroids

A sort of an example of market segmentation

Optimization Objective

k-means also has an optimization objective - a cost function it is trying to minimize

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

μ_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

keep in mind that:
 $k \in \{1, 2, \dots, K\}$

so if $x^{(i)}$ assigned to cluster 5
 then $c^{(i)} = 5 \rightarrow \mu_{c^{(i)}} = \mu_5$

Optimisation Objective :

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

Sometimes called distortion
 cost-function or
 distortion of k-means algo

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

wrt = with respect to

Date:

Cluster assignment step - it carries out this:

$\min J(\dots)$ wrt $c^{(1)}, c^{(2)}, \dots, c^{(m)}$
(holding μ_1, \dots, μ_K fixed)

Move centroid step carries this out:

$\min J(\dots)$ wrt μ_1, \dots, μ_K

~~This~~ ^{It} can be proved mathematically that
this is happening

This can be seen as a 2-step process
where initially $J(c, \dots)$ is minimised ~~with~~
wrt c & then it minimises it
wrt μ

~~J(...)~~ ^{can} ~~be used to~~

$J(\dots)$ can be used to find better
clusters & help k-means avoid local
optima.

Random Initialization

↳ how to initialise k-means & how to avoid local optima

How do you randomly initialise the cluster centroids?

- you should have $K \leq m$.
- randomly pick K training examples in your ~~data set~~ training data set
- set μ_1, \dots, μ_K equal to these K examples then:

$$\mu_1 = x^{(i)}$$

$$\mu_2 = x^{(j)} \quad (\text{for some random } i, j, \dots, l)$$

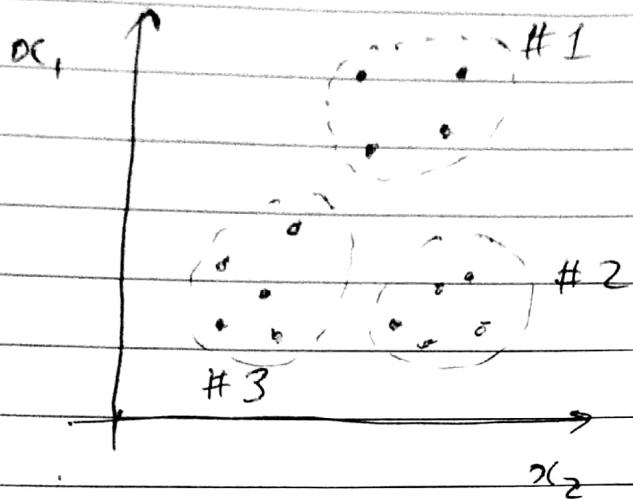
$$\vdots$$

$$\mu_K = x^{(m)}$$

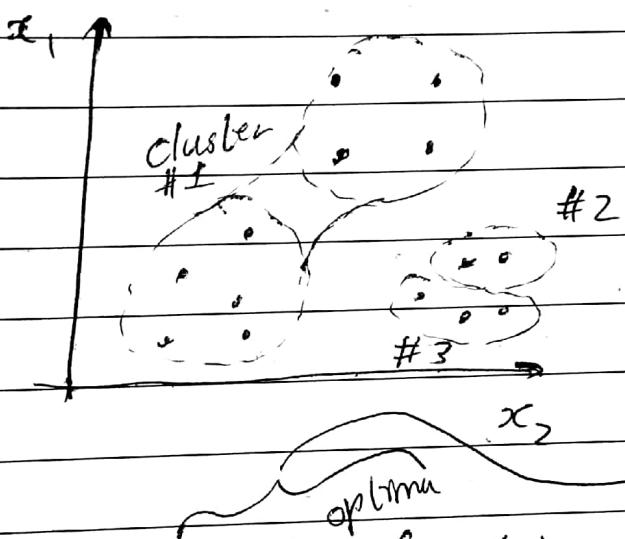
K-means may end up converging to different solutions depending on how the clusters were initialised

↳ K-means can end up with different solutions & can end up at a local optima

Suppose the data set:



Global optima of this data set



Local optima of this data set

The global ~~function~~ minimises function

$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$ optimally

while the ~~local~~ local optima does not

To avoid the local optima run ~~the~~ k-means multiple times with random initialization & find the best solution possible this way

→ 50 - 100

Date: _____

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

Randomly initialize k-means

Run k-means. Get $c^{(1)}, \dots, c^{(2)}, \mu_1, \dots, \mu_K$

Compute cost function (distortion)

$J(c^{(1)}, \dots, c^{(2)}, \mu_1, \dots, \mu_K)$

Run k-means multiple of
times & then pick clustering
that gives the lowest cost

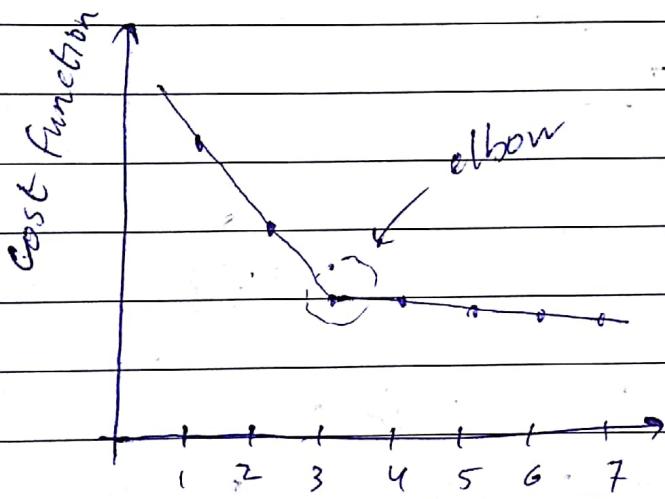
If there are few K clusters (till up till 10)
then doing random initializations can help
in choosing a better optima. But if
 K is large (greater than 10) then
doing multiple random initialisation is less
likely to make a great difference.

Choosing the # of clusters.

↳ done mostly by hand

The reason why at times choosing the # of clusters is difficult is because it's ambiguous how many clusters there even are in the dataset

Elbow method

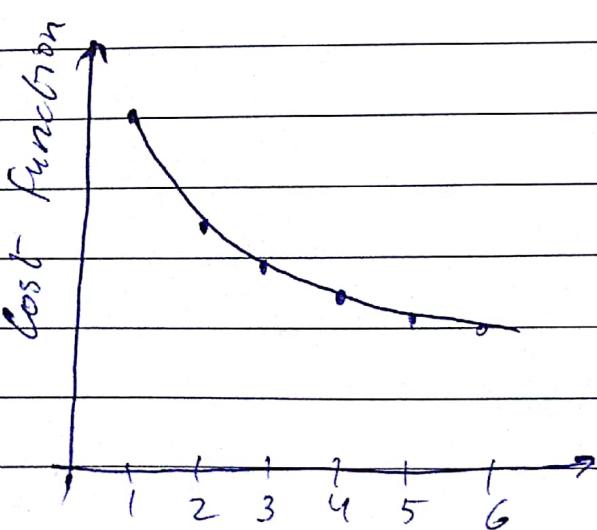


It's ~~clear~~ clear that 3 is a good candidate for the # of clusters here as it is the elbow

However the elbow method can be ambiguous at times as well:

No elbow hence ambiguous.

The elbow method can be used but chances are it may not work & won't be able to give an answer



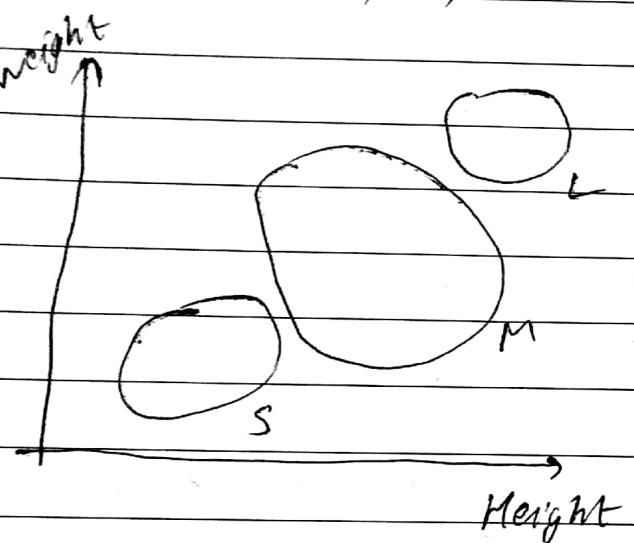
Victory

Page No.

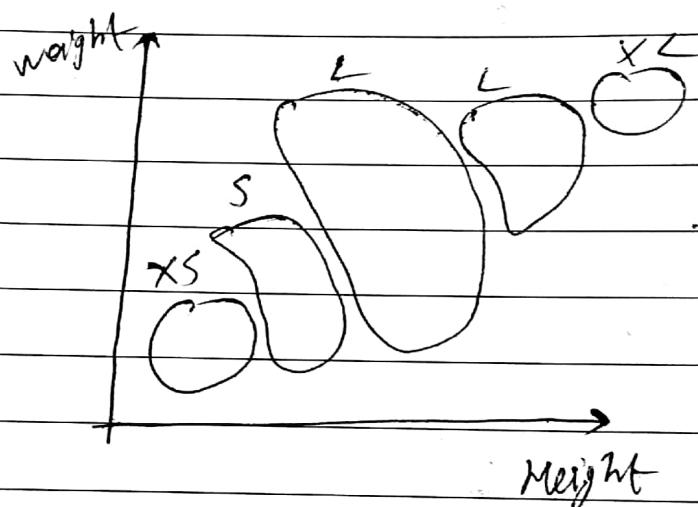
Date: _____

Sometimes, you're running k-means to get clusters ~~for~~ to use for some later/downstream purpose. Evaluate k-means based on a metric for how well it performs for that later purpose. Looking back at the t-shirt manufacturer example:

$$K=3 : S, M, L$$



$$K=5 : XS, S, M, L, XL$$



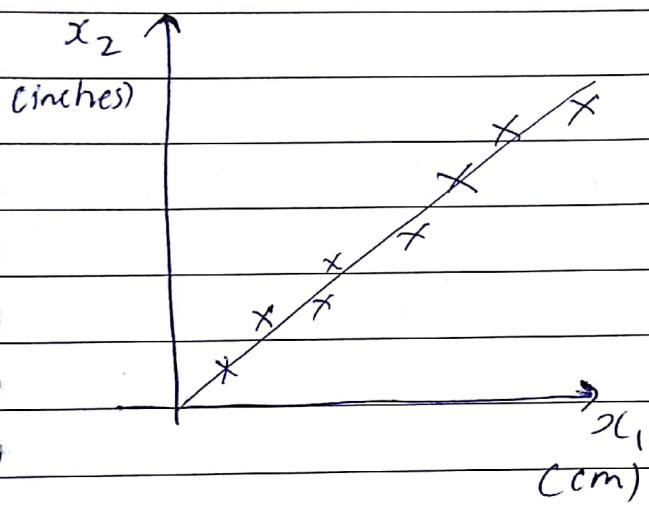
The T-shirt business may have their own performance metrics as to whether 3 or 5 sizes will fit them better.

Data compression

Dimensionality reduction - type of unsupervised algo

↳ can help achieve data compression

↳ allows to speed up learning algorithms



Suppose that x_1 & x_2 both measure the length of some object in cm & inches respectively. The data points don't lie perfectly on the line because the features could have been rounded off to the nearest-inch or cm.

Reduce this data from 2D to 1D

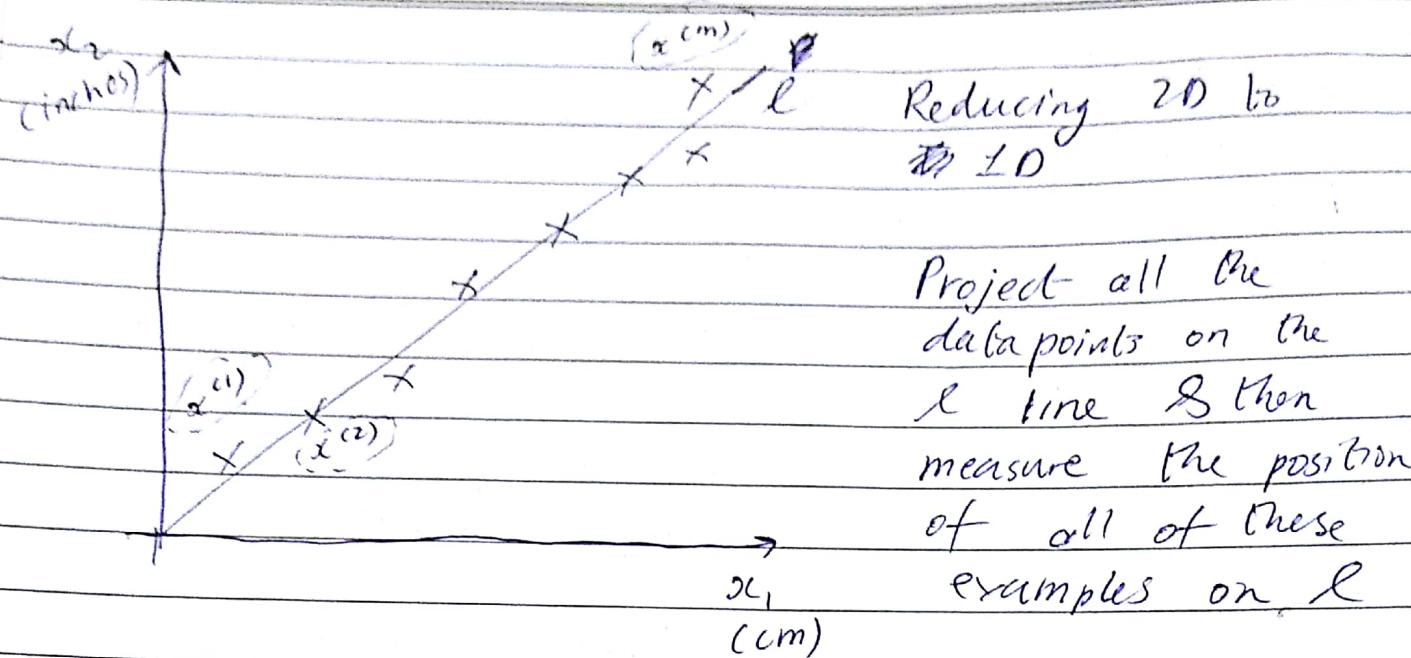
↳ reduces redundancy

If you have highly correlated feature then you might want to loose a dimension

(example: pilot enjoyment & pillow skill will be highly correlated but not the same thing)

Previously shown graph:

Date:



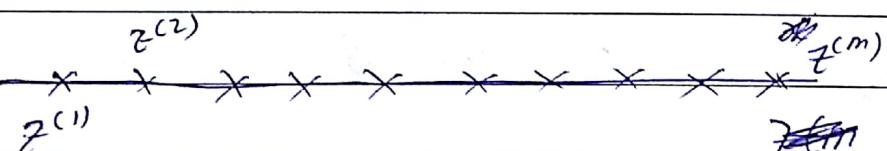
Create new feature z which is the position of data point on line l

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

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

:

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



Basically this method effectively replaces the need to have 2 dimensions. Now you can have just 1 (~~this is an approximation to the original~~)

Date: _____

the new data set is an approximation to the original data set

Consequence of doing this:

- halves memory requirement
- makes learning algorithm run quickly

It is possible to reduce 10000 to 1000 or 3D to 2D if needed

→ project all points in 3D space to a plane

$$x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \end{bmatrix} \in \mathbb{R}^3 \rightarrow z^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix} \in \mathbb{R}^2$$

Data Visualisation

→ second application of dimensionality reduction

You can take a very large dataset suppose that has 50 columns & reduce each data point into a 3D or 2D vector which you can plot to understand the data better.

Note: After dimensionality reduction, the new features are **Victory** not ascribed any meaning, the scientist must find that himself

Principal Component Analysis - algo that allows us to do dimensionality reduction

Date: _____

Principle component

Principal Component Analysis's Problem Formulation

PCA (Principal Component Analysis) - most widely used algo of dimensionality reduction

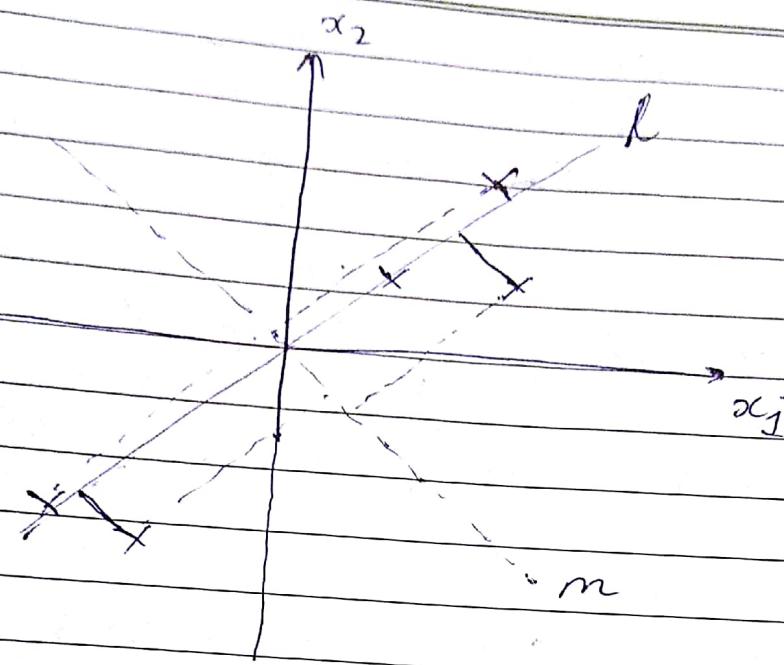
~~Problem~~: PCA finds a lower dimensional surface onto to which to project the data so that the sum of squares of the distances^A to the ~~line~~^B is minimized from data points surface

Projection error - ~~the~~ distance of ~~line~~^C from data point

Before applying PCA, it's standard practice to first perform mean normalisation so that all features have 0 mean & perform feature scaling (features have comparable ranges)

2D \rightarrow 1D

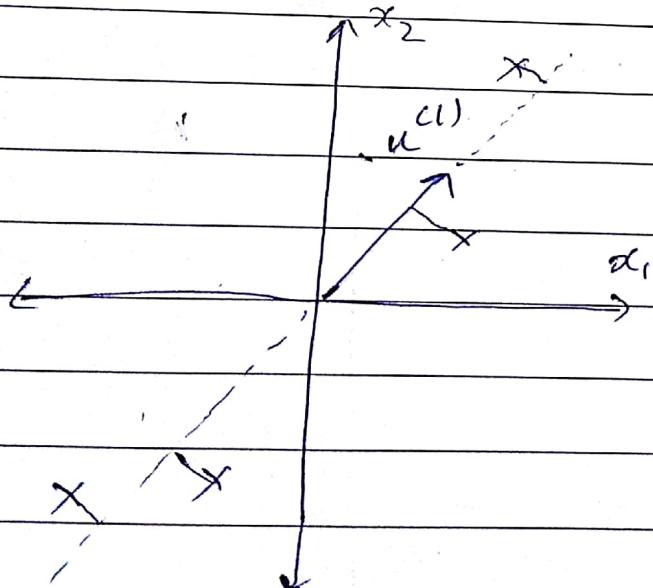
Date: _____



lower dimensional
surface here
is a line.
PCA chooses
such a line
so that the
distance is
minimized
between data
point & line

~~l~~ will be chosen as opposed to m
because distance is lower with l .

PCA problem formulation



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

Does not matter if PCA outputs
 $-u^{(1)}$ as the direction is still the same.

The projection is usually at 90° / ~~orthogonality~~
~~ortho~~ Orthogonal projection distance is
considered the error

Date:

Reduce from n -dimension 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

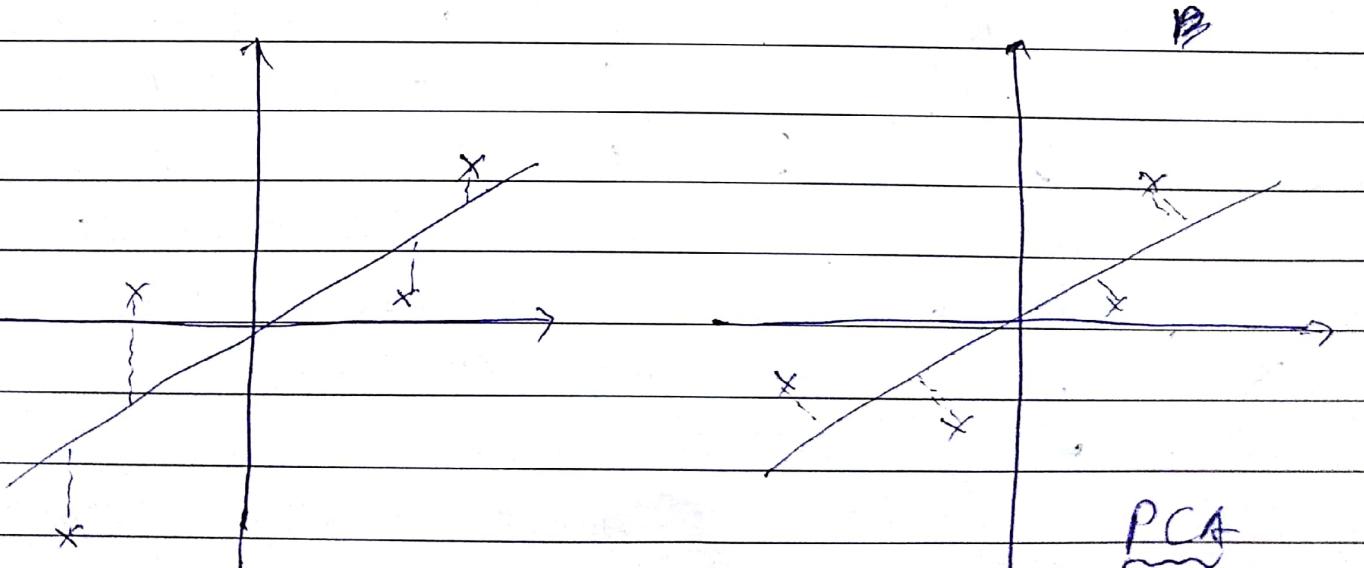
In the case of $3D \rightarrow 2D : k=2$

Hence find $u^{(1)} \& u^{(2)}$ that creates
a plane onto which data points can be
projected

Linear algebra jargon: find set of vectors
 $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ & project the data onto
the linear subspace spanned by this set of
 k vectors

PCA is not linear regression

↳ There are some similarities but not
exactly the same different algos



linear regression:

Victory non-orthogonal distances minimized

Page No.

In PCA the orthogonal distances are minimised

linear regression has a special variable y that we are trying to predict, PCA has no such variable

By convention, we choose $u^{(1)}$ so that
 $\|u^{(1)}\| = \sqrt{[u_1^{(1)}]^2 + [u_2^{(1)}]^2} = 1$

PCA Algorithm

Data preprocessing executed

Before any PCA is done the following must be done:

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

Preprocessing (feature scaling, mean normalization):

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

mean normalization of individual feature

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

If different features on different scales, scale features to have comparable range of values

$$x_j^{(i)} \leftarrow x_j^{(i)} - \mu_j$$

$$s_j \leftarrow \text{standard deviation of feature } j$$

The PCA calculates 2 things:

- $u^{(1)}, u^{(2)}, \dots, u^{(k)}$

- the values that ~~are~~ higher dimensional points are mapped on to in the lower dimension

example: 3D \rightarrow 2D

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

\hookrightarrow find this value

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

The mathematical proof of how to find these values is extremely complicated

Procedure of how to compute the wanted values:

PCA algorithm

Reduce data from n-dimensions to k-dimensions:

Compute "covariance matrix":

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

greek letter sigma $n \times 1$ $1 \times n$
 $n \times n$

Date: _____

Compute "eigenvectors" of matrix Σ :

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

↳ singular value decomposition

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

take the ^{first}
k dimensions

required to create lower
dimensional surface

reduced
version of
 U

↳ $u^{(1)}, \dots, u^{(k)}$

$$\text{Define } U_{\text{reduce}} = \begin{bmatrix} 1 & 1 & 1 \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & | \end{bmatrix}$$

To find out Z in : $x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$
do this :

$$z^{(i)} = U_{\text{reduce}}^T x^{(i)} = \begin{bmatrix} u^{(1)} \\ u^{(2)} \\ \vdots \\ u^{(k)} \end{bmatrix} x^{(i)}$$

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

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

$k \times 1$

done when features
take on very different
ranges of values

Date:

~~Practical~~ PCA algo summary

After mean normalization (ensure every feature has 0 mean) & optionally feature scaling :

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

$$[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\text{Sigma}); \quad \begin{matrix} \text{Vectorised} \\ \text{implementation} \end{matrix}$$
$$\mathbf{U}_{\text{reduce}} = \mathbf{U}(:, 1:k); \quad \text{where}$$
$$\mathbf{Z} = \mathbf{U}_{\text{reduce}}^\top * \mathbf{x};$$
$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}^{(1)} \\ \vdots \\ -\mathbf{x}^{(m)} \end{bmatrix}$$

And as it is with

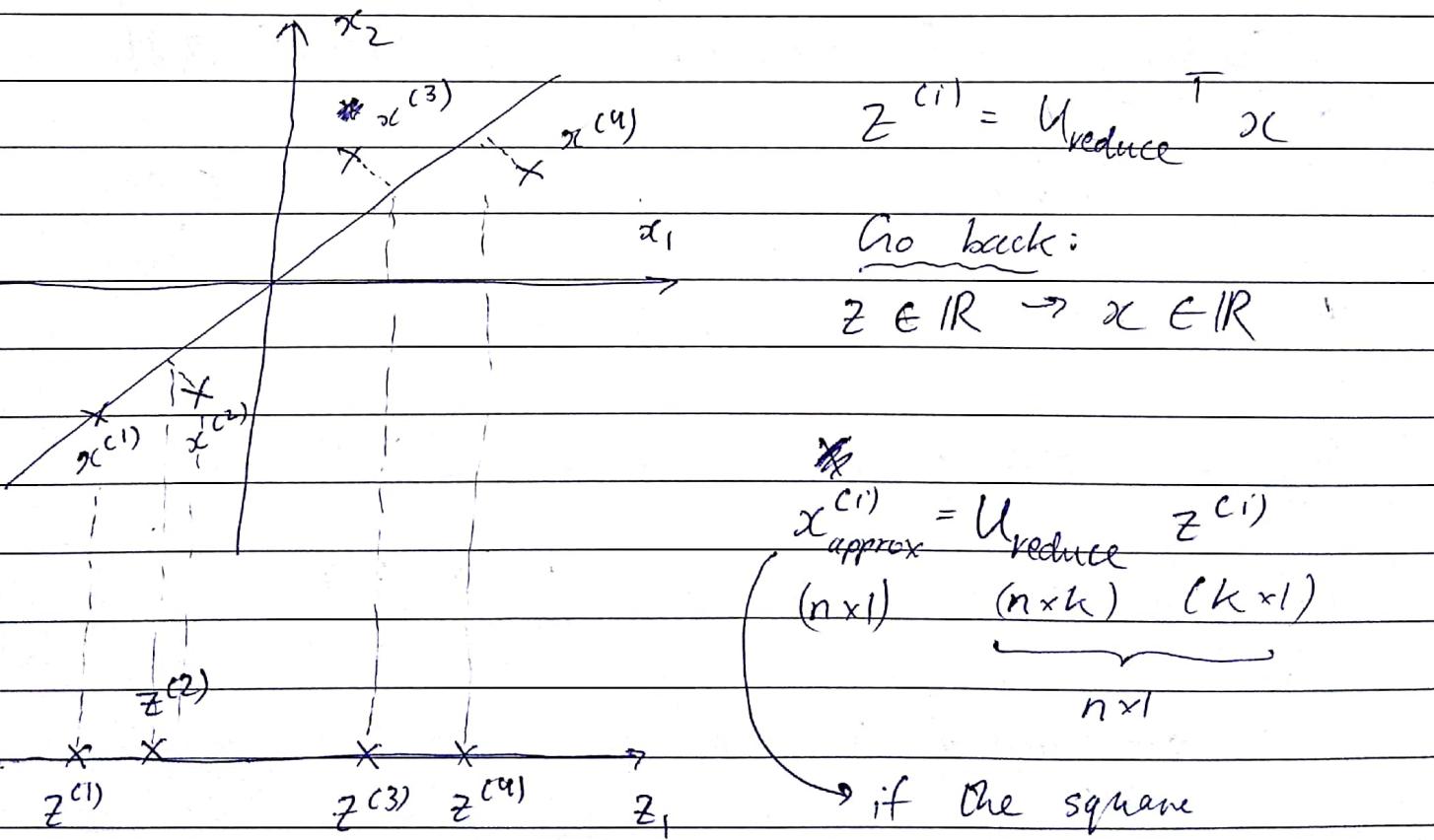
unsupervised learning $\mathbf{x} \in \mathbb{R}^n$

$\hookrightarrow \mathbf{x}_0 = 1$ is not used or
does not exist

Reconstruction from compressed representation

PCA is a compression algorithm where a 1000D data is compressed into a 100D feature vector

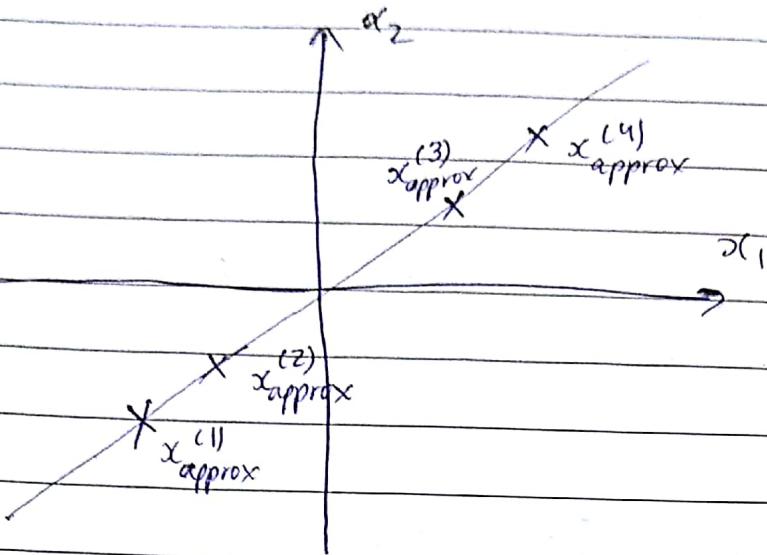
How do you go back from the compressed representation to the original representation (or an approximation of it)



if the square projection error is small then $x_{\text{approx}}^{(i)} \approx x^{(i)}$

square projection error
when PCA was applied

Reconstructed values of x :



The values
reconstructed
will be on the
line

Choosing the Number of Principal Components

PCA
minimizes
this

k = # of principal components

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

$$\text{Total variation in the data: } \frac{1}{m} \sum_{i=1}^m \|x_i\|^2$$

↳ how much the data
varies on average

Typically, choose k to be the smallest value so that :

$$\frac{\text{Avg squared projection error}}{\text{Variation in data}} \leq 0.01 \rightarrow 1\%$$

"99% variance is retained"

If σ^2 value changed then variance retained changes as well :

$$\frac{\text{Avg squared projection error}}{\text{Variation in data}} \leq 0.05$$

"95% variance retained"

$$\frac{\text{Avg squared projection error}}{\text{Variation in data}}$$

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

Variance retained
lends to be chosen to be 95% - 99% usually

This eqn involves k , where by reducing k the variance retained may drop so we must find such a value of k where variance retained remains in the 95% - 99% range

Even for this value, k is significantly reduced as most real life features tend to be highly correlated

Date: _____

Implementation of algo:

Inefficient version:

~~Try PCA~~

try PCA with $k=1 \rightarrow 8$ onwards until (1) becomes true

Compute U_{reduce} , $z^{(1)}, z^{(2)}, \dots, z^{(m)}$,
 $x^{(1)}_{\text{approx}}, \dots, x^{(m)}_{\text{approx}}$

Check if :

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

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

chooses the smallest value of
K where condition (1) becomes true

Efficient version:

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

find the smallest
k that satisfies (2)

For given k:

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

mathematically
eqv. to (1)

Victory

Date: _____

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$$

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

$$\rightarrow S = \begin{bmatrix} S_{11} & & & \\ & S_{22} & & \\ & & S_{33} & \\ & 0 & & \\ & & & S_{nn} \end{bmatrix}$$

Only diagonal
in this matrix
has non-zero
values. All other
values are zeros

values used from here

In the inefficient implementation:

- 1) U_{reduce} needs to be defined iteratively
- 2) $Z^{(1)}, \dots, Z^{(m)}$ must be calculated iteratively
- 3) $X_{\text{approx}}^{(1)}, \dots, X_{\text{approx}}^{(m)}$ must be calculated iteratively
- 4) Formula (1) evaluated

↳ all this creates a large overhead

Date:

If you simply run:

$$\frac{\sum_{i=1}^k S_{ii}}{m}$$

$$\frac{\sum_{i=1}^k S_{ii}}{m}$$

This will give the % of variance that will be retained

↳ tells how well the new dataset will approximate the original dataset

Advice for Applying PCA

Supervised learning speed up with PCA

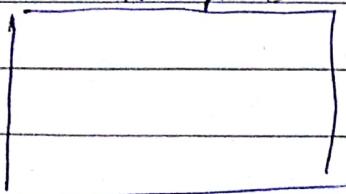
$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$$

Suppose $x^{(i)}$ is in high dimension

↳ such as in computer vision

100 pixels

100 pixels



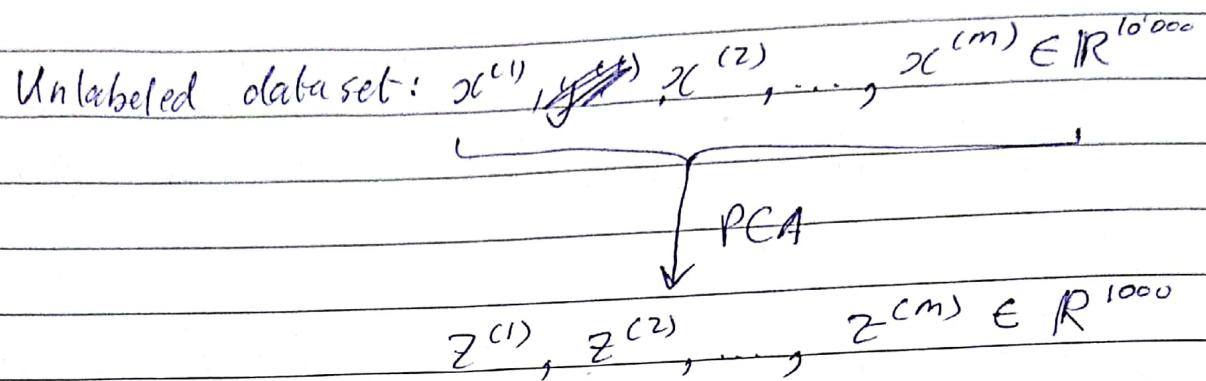
= 10'000 pixels/features of input

image

Having so many features can cause algo to run slowly

Date: _____

Extract the inputs:



New training set:

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

Now feed this new training set into a learning algorithm such as logistic regression:

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

\uparrow
PCA
 $z \leftarrow x$

Mapping $x^{(i)} \rightarrow z^{(i)}$ on training set creates some parameters that apply only to the training set, not any other data set:

{ feature normalization (making mean = 0),

{ feature scaling

- U reduce

"not sure about these"

none of these values obtained from training set can be applied to CV or test set

Date: _____

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)} \& x_{test}^{(i)}$ in the ~~CV~~ test sets

Learning algo can run much faster now due to reduced dimensionality.

Application of PCA

- compression:
 - { - reduce memory / disk needed to store data
 - { - speed up the learning algo & usually 99% of variance is retained by users for this purpose
 - choose k by % of variance retained

- Visualisation

- ↳ k is chosen to be 2 or 3 as only 2D or 3D can be plotted

% of variance retained - measure of how much info is thrown away when finding lower dimensional representation.

Date: _____

Bad use of PCA : To prevent overfitting

Use $z^{(i)}$ instead of $x^{(i)}$ to reduce the # of features to $k < n$

(suppose $k = 1000, n = 10^6$)

Thus, fewer features, less likely to overfit

Bad!

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

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

Reason for not choosing PCA for overfitting :

PCA

When finding lower-dimensional approximation, PCA throws away some info without knowing the values of y

↳ okay if 99% variance is retained but not okay if some valuable info is thrown away

Regularization will work better to prevent overfitting because minimization problem knows values of y hence less likely to throw away valuable info.

PCA is sometimes used where it shouldn't be

Design of ML system:

- Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- Run PCA to reduce $x^{(i)}$ in dimension to get $z^{(i)}$
- Train logistic regression on $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
- Test on test set: Map $x_{\text{test}}^{(i)}$ to $z_{\text{test}}^{(i)}$. Run $h_0(z)$ on:
 $\{(z_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}), \dots, (z_{\text{test}}^{(m)}, y_{\text{test}}^{(m)})\}$
- How about doing the whole thing without using PCA?
- Before implementing PCA, first try running whatever you want to do with the original / raw data $x^{(i)}$. Only if that doesn't work do what you want, then implement PCA & consider using $z^{(i)}$

Date: _____

Run PCA if your learning algo ends up running too slowly, or requirement for memory is too large.