

Lec. 18
04-10-24

"unsupervised Learning"

Goal of su learning \rightarrow Prediction

unsupervised \rightarrow finding the groups / patterns / similarities in the sample

2 methods:

- i) Principal components analysis
- ii) clustering

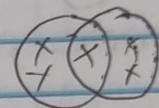
clustering:

finding subgroups / clusters

i) k-means clustering

$c_1 \cup c_2 =$ union \rightarrow

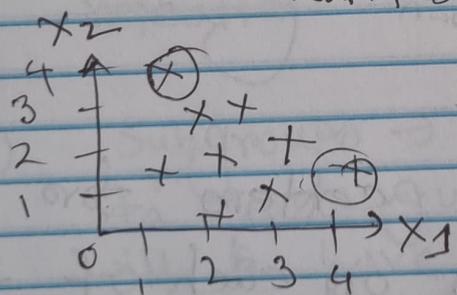
intersection $\rightarrow c_1 \cap c_2 = \emptyset$



shouldn't be any intersection

clusters should be distinct

samples



WCR

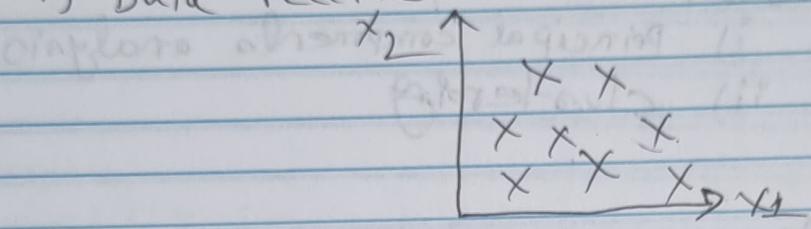
↓
within cluster variation

Euclidean distance:

"from no clusters to
↓
samples with clusters"

k-means clustering algorithm:

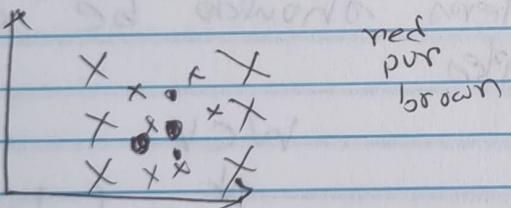
i) Data receives



$k = \text{num of clusters that we want to get}$
assume, we want to get 3 clusters

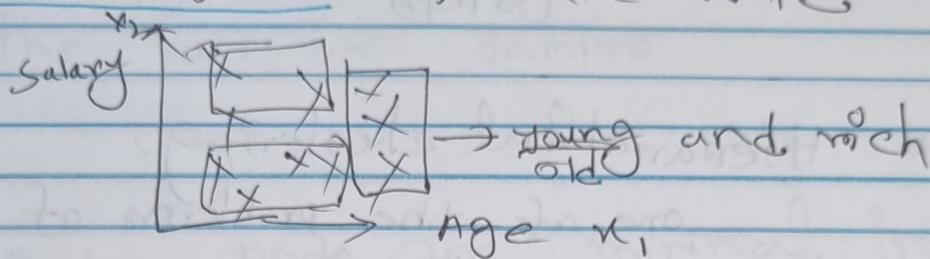
$$k = 3$$

1st step: assigning randomly samples



finding the centers for both
3 clusters or each of the
clusters

Reassigning each sample \rightarrow based on the nearest distance \rightarrow the sample



Then finding the new centroid based on the new clusters

2 problems / Disadvantages:

- i) randomly selecting samples
1st step usually doesn't give the consistent results

So need to run multiple times and finding out which one is repeating

- ii) Subgroups \rightarrow we specified 3/5 but it's not guaranteed that data will have this no of subgroups \rightarrow model is doing good but underfitting

$\uparrow \uparrow \rightarrow k \rightarrow$ more than sum \rightarrow overfitting
 $\downarrow \downarrow \rightarrow$ underfitting

steps summary

- no of clusters
- selecting k objects randomly
- then each observation to their closest centroid

overcoming: by doing cr in multiple steps.

Hierarchical clustering:

one of the problems of in k-means can be solved.

finding all the no of clusters within the data

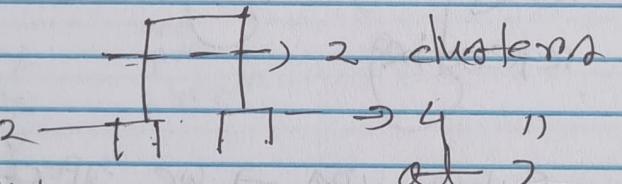
↓
bottom up approach

↓
individual app sample

Dendrogram → will guide us

to see how the clusters look like like

cutting



if we cut at 2

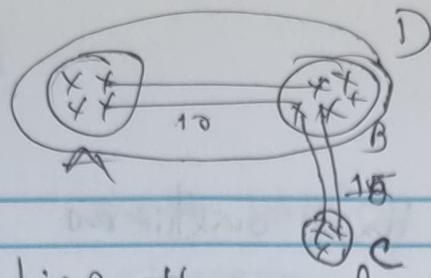
* Types of linkage:

linkage method

complete linkage method:

clusters are closer → similar

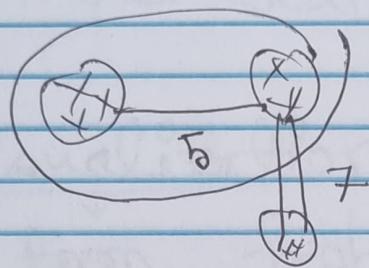
find the maximal intercluster dissimilarity



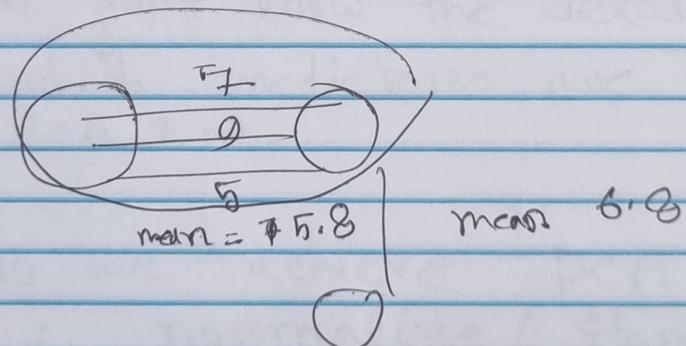
→ now 2 clusters
D and C

Finding the largest distance between the pairwise samples

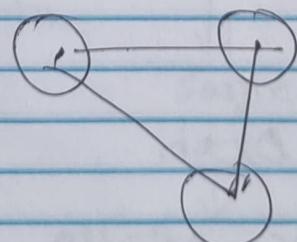
ii) single: finds the minimal interclass dissimilarities



iii) Avg: mean inter-cluster dissimilarity



iv) centroid:
also sensitive
to the outliers



complete and single → are sensitive to outliers



If data has outliers

then ↓

centroid

and avg point

scaling metrics: mean / normalized

Lec: 19
04-17-24

Deep-learning

PCA:

Reducing the dimension ~~frommer~~ for example from 6000 to only 2.

Principal Component analysis: Advantage:

- i) for reducing the dimension correlated features are pulling together

ii) visualization

from high dimension to 2-dimension

Disadvantage :

- i) we don't know the details which predictors are making PCA 1 →

How do we achieve PCA :

Step 1: normalize / standardizing the data

↓
same scale of the data (not vast difference)

↓
all the values will fall between 0 and 1.

Standardization

$$[0, 1]$$

Min Max



SD → 1
upper/lower
value can
vary

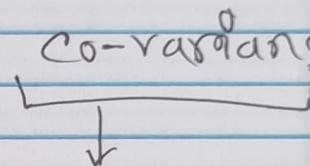


Bringing the values

around the mean
of the numerical
values

If the data's are
sensitive to min or max,
then we need to normalize

Step 2: calculating the covariance matrix

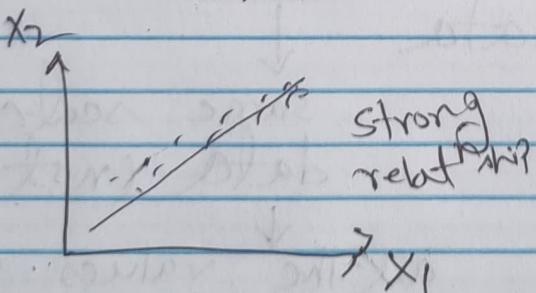


how different
the predictors are
relating to each other etc

$$\begin{matrix} & x_1 & x_2 & x_3 \\ x_1 & 1 & 0.9 & 0.1 \\ x_2 & & & \\ x_3 & & & \end{matrix}$$

→ we wanna measure
how strongly x_1
measures x_1

3x3
PXP matrix



Data : $X_{n \times p}$
 normalizing
 finding co-variance
 $\text{cov}_{p \times p}$

Step 3 :

Taking cov matrix

calculating Eigenvalues and Eigenvectors
from covariance matrix

Eigenvalues

↓
Represent the principal
component

Eigenvectors

↓
represent
the variance
explained by
the principal
components

Shape: $P \times P$

$x_1 \quad x_2 \quad x_3$

$$\begin{bmatrix} PC_1 & PC_2 & PC_3 \\ PC_1 & - & - \\ PC_2 & - & - \\ PC_3 & - & - \end{bmatrix}$$

$1 \times P$

$$\begin{bmatrix} PC_1 & PC_2 & PC_3 \end{bmatrix}$$

captured
information
from PC_1

$if, PC_1 > PC_2$

then PC_1 will
contain the more information

from the data

Step 4: Project the principal component onto the normalized or standardized data.

eigenvectors \circledcirc scaled data $n \times p^T$

matrix transpose:

$$A = \begin{bmatrix} 2 & 1 \\ 2 & 4 \\ 3 & 1 \end{bmatrix}$$

$P \times P \circledcirc P \times n = p \times n$

Scaled data eigenvectors

$$A^T = 2 \times 3 \times 2$$
$$\begin{matrix} 3 \times 2 \\ 2 \times 3 \end{matrix} \quad \begin{bmatrix} 3 & 2 & 2 \\ 1 & 4 & 1 \end{bmatrix}$$

$$P \times P \quad P \times P$$
$$= n \times p$$

[no need to matrix transpose
as internal $P \cdot P$ are same]

Summary:

- normalizing (only numerical features
not categorical ")
 - covariance matrix \rightarrow eigenvectors \rightarrow PC.
 \downarrow eigenvalues \rightarrow amount of info
that each PC carries
 - projecting PC onto the standardized data
- Then selecting the principal component, M,
where $M \leq P$
 $\Rightarrow n \times M$

Selecting M based on eigenvalues.

Select the principle component, M ,
where $M \leq P$

$$\Rightarrow n \times M$$

selecting M based on eigenvalues.

youtube:

Eigen vectors and eigen values:

(i) covariance matrix between features;

(ii) Eigen vectors and eigen values
will be found out from this
covariance matrix

(iii) Eigen vector \rightarrow eigen value \rightarrow

magnitude of the eigen vector
captures the main variance.

Linear transformation of matrix

$$A v = \lambda v$$

Steps to calculate eigen value and vectors:

① covariance of features:

$$\begin{bmatrix} x, y \end{bmatrix} \rightarrow X^T$$
$$\text{cov}(x, y) = \frac{w}{N-1} \sum_{i=1}^w (x_i - \bar{x})(y_i - \bar{y})$$

→ Standardizedⁿ

→ covariance matrix of x and y

→ find out eigen vectors and values

$$A v = \lambda v$$

$$\begin{bmatrix} \lambda_1 & \lambda_2 \end{bmatrix} \rightarrow \text{eigen values}$$

\downarrow \downarrow
PC1 PC2

Gaussian distribution
outliers

Standardizedⁿ → to transform the data
to have a mean of zero
and SD of 1

mean of 0
and SD of 1
any upper/lower
values or Z-score normalizationⁿ

Normalizedⁿ → to make feature values
range from 0 to 1.

always range from
0 to 1

in case of neural
network

Project ++

Bean classes are : ✓

SEKER BAR BOM CALI THOT
DERMAZON, SIRIA

① Data Import + Descriptⁿ

② " Preprocessing :
→ Data cleaning

- Dropping the duplicated items
- Null data
- Data split + normalizatⁿ
- Data - visualizatⁿ

Train → LR (penalty = 'l2', verbose=1,
n-jobs = -1)

Penalty l2 = L2 regularization
Ridge ↓
"

This adds a penalty term
to the loss functⁿ

which penalizes the
↑↑ co-efficient values

so that overfitting can
be minimized

verbose=1 ; Detailed output
during the training
process
('Iteratⁿ" + convergence)

$n_jobs = -1$, the number of CPU
cores
all available cores

These parameters help improve
the efficiency and effectiveness
of the LR model for classificatⁿ.

model.get_params().keys() → returning
a list of all the parameter
names

solver → LR model → optimizatⁿ
algorithm

lbfgs → Limited memory Broyden-
Fletcher-Goldfarb-Shanno
↓
optimizatⁿ algorithm

model.selectⁿ, Grid Search CV;
Dimensionality Reductⁿ
Linear_model.LR → machine model
classifier

null → mean + median of specific column

dup → drop row

outli → proposed classifier doesn't deviate from its prediction

normalized → to standardize the scale of features

data was divided to a training set,
and a test set 80%

Pearson correlation matrix,
analyzing the inter-relationships
between the variables

Analys: 16 fea
correlat + Pearson plot

04-22-24

Deep Learning

Traditional / classical ML models:

*** SVM, Decision tree
KNN, Hierarchical

came from statistical

1990 → Deep learning:

at that time computational / hardware
was scarce



so classical deep-learning

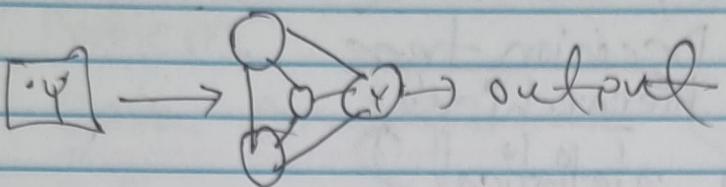
GPU → can compute heavy ~~matrices~~ and
" " run

then based on this → deep learning
models (large dataset)

↓
can start / experimenting

GANS → can generate images from
text
single layer NN

→ let the computer learn from data



each neuron → responsible for particular task

Artificial neuron that will learn from the input

$$y = mx + b$$

$$= w_n + b$$

↓ ↓ ↓
 input weights/parameters bias

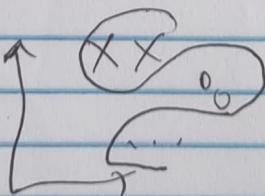
example:

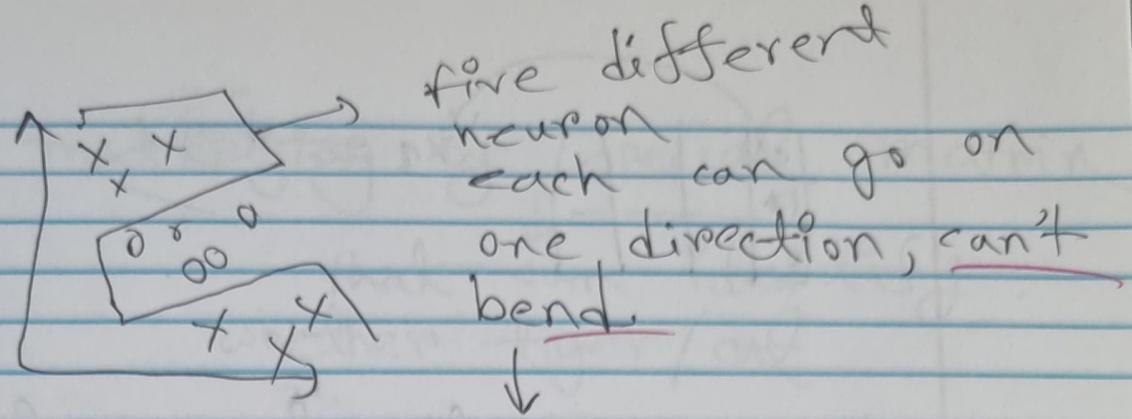
$$A_1 = w_1 x_1 + w_2 x_2 + \dots + w_4 x_4 + b$$

those parameters values need to checked to the true values

* Adding nonlinear behaviors from linear functn:

Activation:



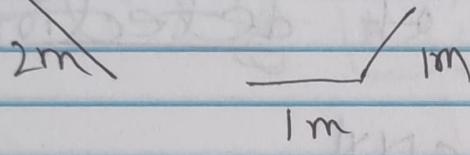


So it's difficult
to capture non-linearity

Rectified
(ReLU) Learninear :

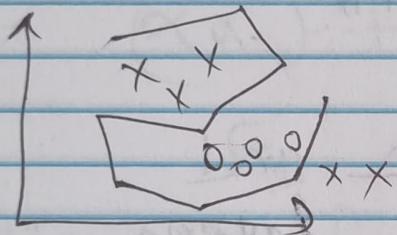
Rectified linear unit:

How works?



$$\text{Relu} = \max(A_i^0, 0)$$

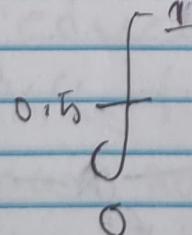
if the line is not
good \rightarrow then 0



Sigmoid : shaped

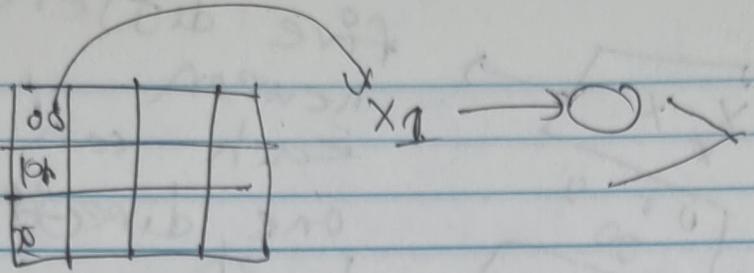
How it works?

between 0 and 1



don't learn
from this

If not 0,
then A_i^0
good learn



convolutional neural network

↓ CNN

for image classification

any kind of detection

↓
from CNN

2 main oper:

→ convolution and
→ Pooling layers

convolution filters

image → matrices → elements → number
pixels value
↑
Grey scale

white
↓
darker → 0

0 - 255

darker → ↑ pixels value

multiplying by special matrix

that can identify
whether tiger/cat

convolutn filters = $\begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix}$

$$\begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \cdot \begin{bmatrix} [a][b] & c \\ [d][e] & f \\ g & h \end{bmatrix} = \begin{bmatrix} F_1 & F_2 \\ F_3 & F_4 \end{bmatrix}$$

$$(\alpha \cdot a) + (\beta \cdot b) + (\gamma \cdot d) + (\delta \cdot e)$$

F_1

Hierarchy of filters

cat pic: $11 \rightarrow$ vertical \rightarrow lower pixel \rightarrow white
is more visible

Pooling

reducing the dimension from
 4×4 to 2×2 .

max pooling/ avg pooling

convolve \rightarrow pool \rightarrow convolve \rightarrow pool \rightarrow

flatten

like basic
neural network
output

AI model \rightarrow review

↓
differentiating whether
+ / -

featurized: Bag-of-Words

\rightarrow Dictionary of frequently used words

	x_1	x_2	x_3	$\rightarrow P$ Prediction
good	0	0	0	
bad	1	4	0	
happy	2	2	0	
⋮	⋮	⋮	⋮	⋮

\rightarrow Then learning from the sentences
point of the word in

a sentence (meaningful)

Recurrent neural Network Archi
↓
mothers of chatgpt