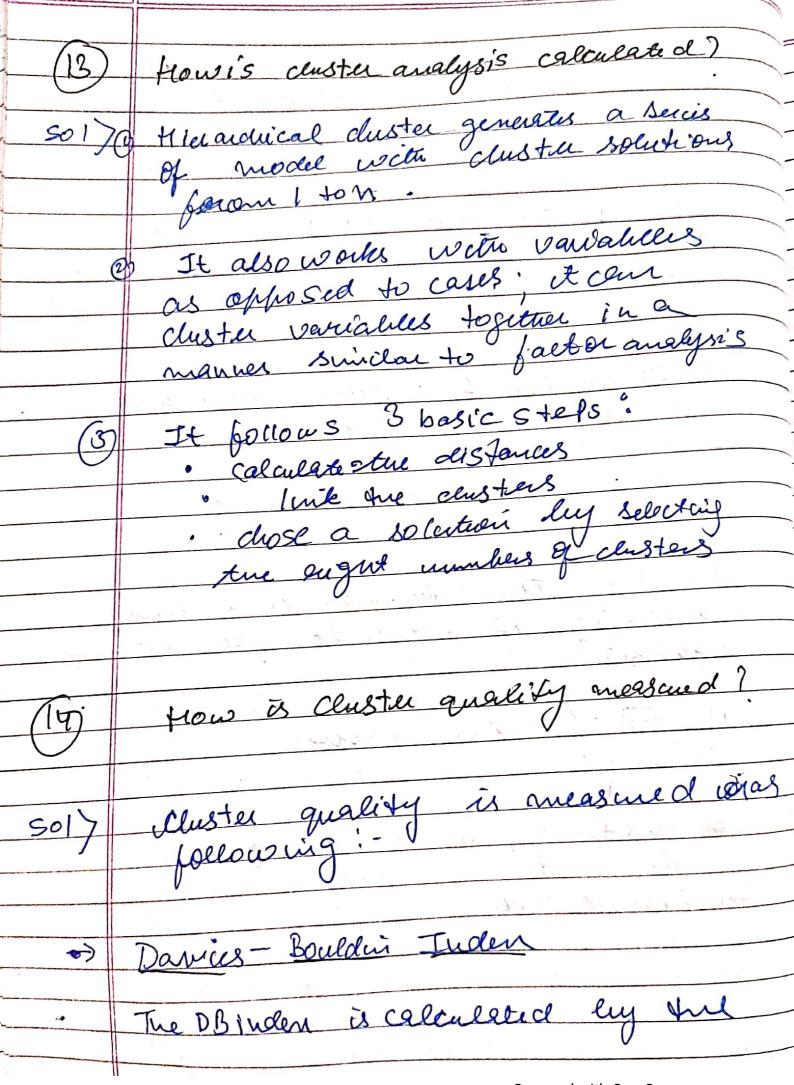
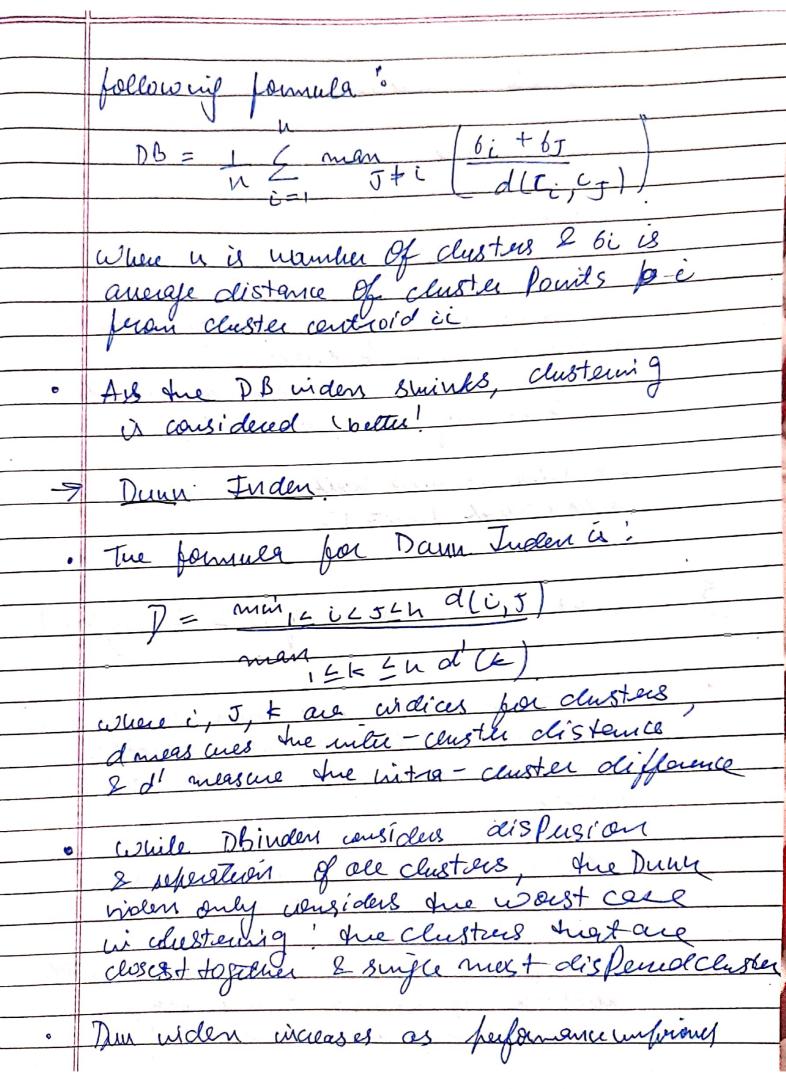
	Machine Learning - Worksheet - 1 By!
	Dy!
	Yatika Taneja
	V
0	for the data fonts represented by the following dendogram:
	for the data fourts represented by the
	following dendogram!
	The state of the s
Sol	(b) 4
/	100000000000000000000000000000000000000
(2)	In which of the following cases will K-means clustering fail to give good results?
-	K-means clustering fact to five for
-	Hesules.
Co.1\	12 2 24
501>	(d) 1,2 and 4.
	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
(8)	The most infod and fact of is selectory
	The most infodant fait of is selecting the variables on which clustering is based?
	0
Sol	(d) formulating the clustering problem
/	
(Ca)	The party of many of manager.
(Y)	of interior is the
	of similarity is the _ or its square
Sal	
201	(ai) enclidean distante

(3)	all poppets start out as ut one giant cluster. Clusters are formed by dividing the cluster with smaller b smaller dusting
	all payects start out as ut one joint
1 1	cluster. Clusters are formed by duriding
	fre cluster inte smallert smaller
0-11	
2017	(b) Divisure censter
6)	Wui ou of the following is required for k-means clustering?
	for k-means clustering?
	and the
<u>Sol</u>)	all answells are concet.
-	
$\widehat{\mathcal{T}}$	The goal of dusterers g is to -
Sol	divide que dota fonits ente
	grafis
(A)	01 tugia à a -
(8)	Clusteraig is a -
SOI	(b) renscipervised learning
3 (9)	wed ou of the following the browlend
	algoritums suffer prom the problems of convergence at local optimis?
(102	af Komens algorithm
	d) All of the above.
	Scanned with CamScanner

-	
(10)	anost sensitive de orothèrs!
	destan of chistian agora
	ands! sensure to grotter!
301	
201	(A) K-means clustum algoritum
-	
(It)	1,00
_(")	What of the following is a had characterstic of a dayaget for clustering analysis?
	Characters 4ic of a dayaget for consist
	analysis?
Soi>	(D) All of the above
. ,	
	and a first thought the same and the same an
(12)	Sor dustering, use do not require
	John Clus Terring, and
SOL	(A) Labeled Darka
/	
	No transfer to the second seco





27	Silhouette Coefficient
	J. J. C.
,	It is mousure as!
	S(i) = b(0) - a(0)
	man [aev), b(i))
	when acil is arreage alstance of i perom
	Smallest average d'osserver of to all
-	Smallest average d'ésource of to all Points in any otrue cluster
0	It toppe us well-ansigned each
	It tells us how heell-assigned each inderied half fourt is
	January and January Danier Danier de
(15)	what is duster analysis & its typus!
	Cluster analysis is a class of techniques
	that are used to classify sheets
	Cluster analysis is a class of techniques that are used to classify slipeets of cases wito relative groups called
	clusters
6	Eluster analysis idencelo cirolores formulais a publicar, selecturg a distance
	a publicus selecting a destance
	measure Selecting a Clustening
	dusters interlegeting the runther of
1,4,1,5,1	clusters, interfereting que Profile
	Scanned with CamScanner

	valionity of clustering
	8
	Typus of cluster analysis are as follows:
	o 1
]=	Merarchical clustering
2.0	Partition clustering
	· · · · · · · · · · · · · · · · · · ·
30	Enclusine clustering
4.	queilapping clustering
Si	Luzzy clustury
	So f
B -	Complete clustreing

	Machine Learning - Worksheet -2
	Machine Learning-Worksheet-2
$\langle \hat{C} \rangle$	Movie recommendation septems are an example of!
	16 of
	are an example of
	and the second s
601	20 d 3 /D)
-3017	2 and 3 (D)
<u>@</u>	Sat it analysis is a enample
	Sentunent analysis is a enample
	1
Cal	(e) 12 and 4
<u>SOI</u>	(e) 1,2 and 4
(F)	Can decision teres be used to
(6)	and decision of the state of 7
	for perform dusturia.
9	
Cal	[A] True-
_	
(4)	Which of the following is the
	and indudant atheroficale strelegy
	1 and a los mais full one Dullonning
	for data cleaning before Performing
	clustering andaysis, given less
	fran desirable number of data
	louis'
	(1)
2017	A) louly
	0
	Scanned with CamScanner

(3)	so what is the minimum no. Of
	variables / peatures brequired to Perform
	clustering?
Sol	(B) 1
(B)	For 2 suns of K- mean Clustering
	For 2 suns of K-mean Clustering is it enpeoled to get same clustering sustering
	Iselutis.
	the state of the s
SOI	(B) NO
/	
	ia ment of
(0)	Is it possible that assignment of Object observations to cluster does not
	Objet Objetvations to cluster does hot change between successive teration on Emeans
	change het veeler succession
	ai Emeans
5017	(A) Yes
	The state of the s
(m)	which of the following Can act as a possible termination conditions?
8	a possible temination condivious
	k-means?
5017	(D) All of the above

as a possible temenation conditions in k-means. Same questrio 80 Culiais as 8 8 Sol word of the following and is a sensitive to outless? (10) (A) K-means clustering afforithm How can Clustering be ased to imperore the accuracy of linear regression model (11) SOI) (f) AN of the above what could be the possible reasons for penducing two different dendroquent using agglomerative des clustering algorithms for the same dataset SOI) (E) All of the above

	means
(3)	Is & sensitive to outliers?
SOI	Tes k means all clustering is senselvie to outliers as it uses que mean of cluster data Pouits to find the cluster
	to outliers as it uses one mean of
	clustree data Points to bijd the cluster
	Ceretel.
(4)	When is to seems botter?
0	Why is k means better?
co 1\	· Polatorial , Sociable to mindelement
501)	· Relateriely scriple to inflement
	· Scalar de Para deta sots
	· Scales to large data sets
	· Charanteel Convergence.
	· Guarantees convergence.
	· Can warm start the positions of
	centroid
	Centroid
	Cail Malte Land engules
	· Easily adapts to new enamples
	V
	· Generalizes to consten of différent Shapes and sizes, such as elliptical clusters
	and siges, such as ellerited
	(Kusters

B	Is & means a deterministic
	algoritum
501) .	K-mean is a non-deterministic
	algorithan.
٥	The mon-determentistic nature of
The second second	E-means is due to its random
17	Selection of data fourts as without
	indial controvels
-	