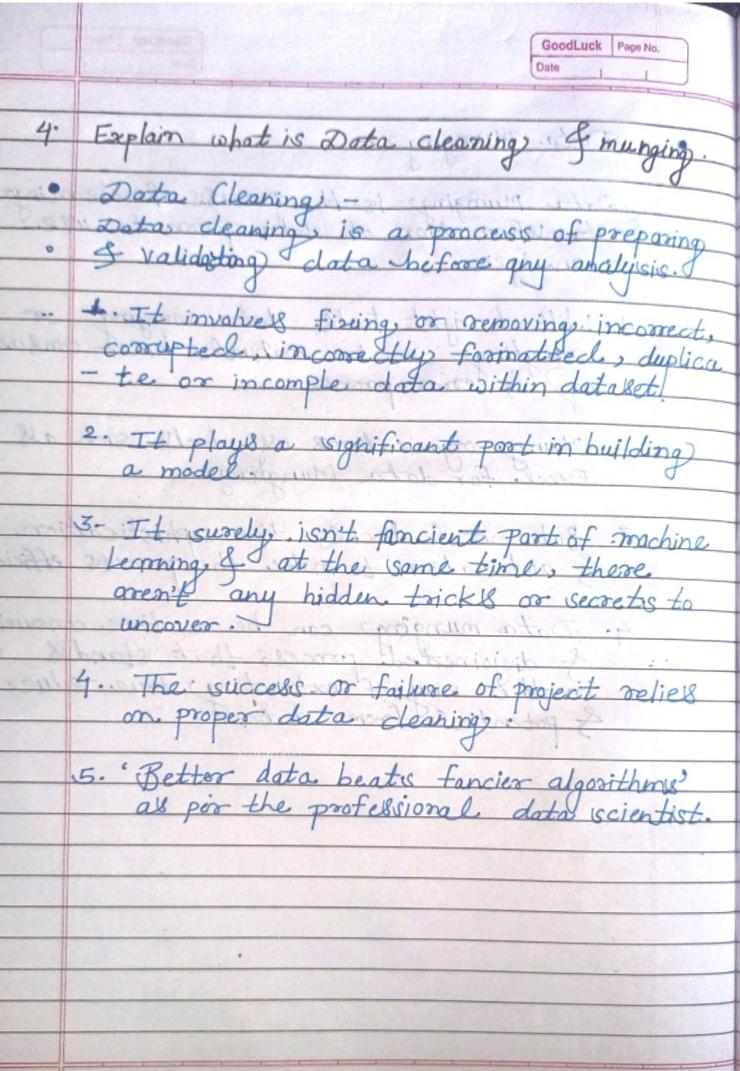
	Asya N Kumbhar GoodLuck Page No. Date 26 1 11 12
	Asya N Kumbhar Goodlyck Page No
	Aya N Kumbhar GoodLuck Page No. A47 (2122000375) Date 26 11 23.
	* Deta Science *
13	Explain : 4 that is mean his toring flows to the
	Explain What is mean by train /text isplit
•	Train test split -
A NET	these technique in used to estimate
	the performance of machine, learning
1221531	algorithms, which ore used to mit
	production on data not used to the
	the model.
	- Chamber of State of the state
30-50	Splitting dataset is essential for an unbiased evaluation of prediction performance.
richard	unbiased evaluation of prediction performance
parade 9	The train of test split produces aprocedure
Therese	is appropriate, when you have a port love
	charasely a costly model to train or
	require a good Estimate model performance.
9/11/10 4	It can be used for classification or
V. mit	regression problems of can be used for
	any supervised lemning algorithm.
2000	V. Provide Company Aspaire From R. T.
	Train Data set - used to fit the machine
	learning model.
•	Test. Dotalet - used to evaluate the
	fit machine learning model.
	Wil. D. Las tregoation.
The state of the s	· siii. Interespectabilitis - fasily converte - to
ager	madifyed to little among. I

Proces	s for train/te	no market	
Proces	s for train/te	101 -111	
		St SRUES -	
	· Taking a dinto stop s	atalet &	dividing it
	into otwo is	ubsetis.	
ii	First subset	is used to	fit the model.
	I reffered a	& training	dataset.
li i-	The second si	ebiset is pr	S compared.
	input element t	the mode	I then
	predictions 98	e made	& compared
	to the expec	ted values	?
24 0	The last to be a		
iv.	That one is test dataset	reffered.	as the
	test dataset	all Some and	9
	A floor of the same		

3,	GoodLuck Page No. Date 11111
	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2.4	White advantages of Numpy of mateplotlib.
	Numerical python (Numpy) is a powerful library in python that provides support for large, multi-demensional growys & matric along with mathematical functions.
	Advantagel of Numpy - i. Efficiency: Optimized for numerical computations, with of fortain
	ii. Array operation: Supports efficient operation - ion on multidimensional arrays.
	iii. Broadcasting
	iv. Mathematical functions: Extensive collection of math functions.
	V. Provides comprehensive functions of Linear algebra. operations.
	vi. Random number generation.
	vii- Integration.
	viii. Interoperability - Faxily converts to & fro from other python data (strectures.
	in Large Community support of comprehensive documentation.

	GoodLuck Page No. Date
IAI	Platploblib + A A A Mary of Cantion
	Itis a widely - used data visualization
-	library in Python.
- 1	and the state of t
Aut	Advantagel of Matplotlib -
	The state of the s
141	i. Matplotlih provided a wide variety of
- Control	Plot typell. This versatility make it.
Trans.	Shitable too vanous ages visquesting
	ii. Publication - guality Plots.
-	e is if the coins in Health come
	111. customization - offers extensive customiza
deta	-tions options for plot.
a Free	iv. blide Apotional.
	and the same of th
1 11	Notebook: Notebook:
	Aires Land September State of
	vi. It supportes 3D plots.
41	is. To assente data malus of using.
	vii. Platplotlib is a large ecayaten.
15	
C,	viii. Active development.
1.71	me on had with an interference of the second
	in. Matplotlib is cheing, open - says
Marie Const	-ce, It was Freely available.
	James La gradieties madel



	GoodLuck Page No.
1	Data Plunging
	Lota Mungings is the process of cleaning
1	Data Mungings is the process of cleaning - & transformation of data prime to use
	ing & error-prone.
	Execute for data Mungling.
	treet for data trungings.
	3. But excel lacks, the sophistication. S automation to make the process efficient
	4. Data munging can be a time consuming
324	S disjoineded process that stands in the way of extracting true value & potential from data
	3 potential from data:
1	Setting data beater favoir alonithe
A	at pie the professions to datal cein