MET Bhujbal Knowledge City

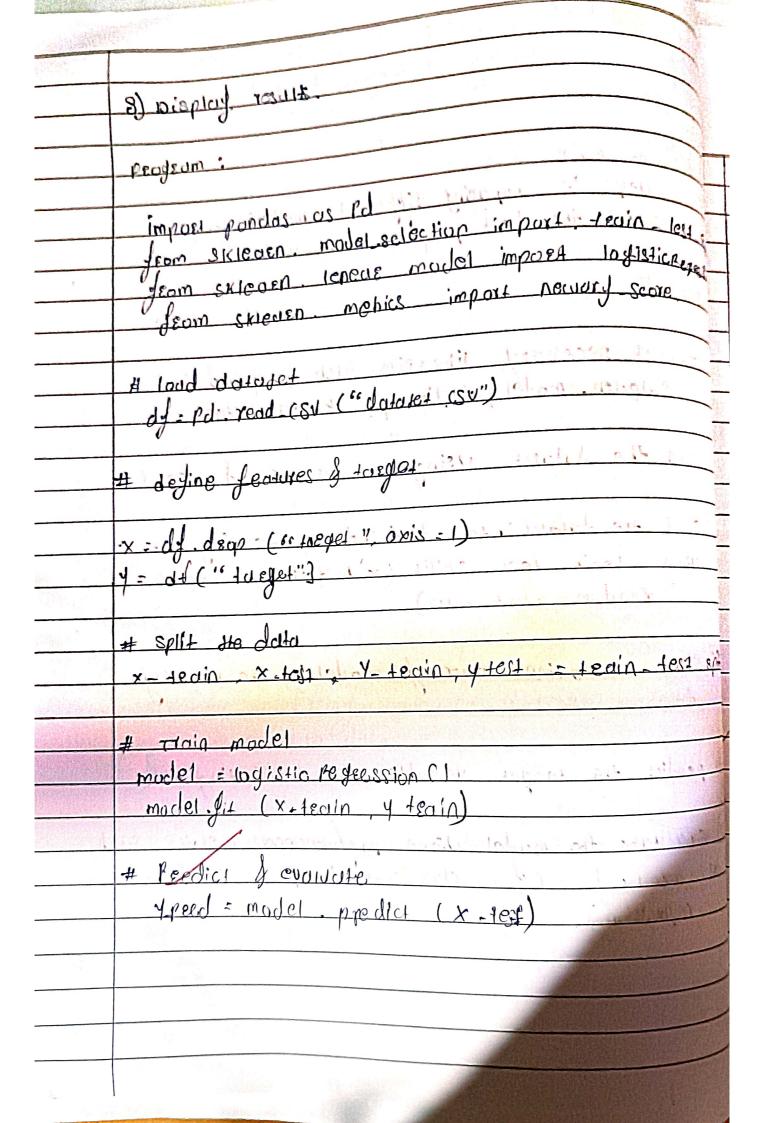
Institute of Technology- Polytechnic Adgaon, Nashik- 422 003.

Date: Peactical No: 3 AIM: To implement teain-test splip on a given duriset to evaluate the performance of a machine learning A daraset (eg CSV file) contains multiple features output: of the malk in the cheer he of mean squeed essos (depending on the model Used) Theory: In ml duta is esercity devided into two set I training set: When to team the model's allowing it to leave pottens from the data. 2) Testing set: used to evaluate the model's posto smore on unseen duto. The teain - test split is a techique where the datatest

is divided into teaning of testing subsets, commanded on do to se to 30 eating the sciki-lener of a split is the sciki-lener of split is the proces. Split is to facilate the proces. Need for teach test split: Need for teach test split: No proces of classes they have they have the seen before to ensure proper generalization.
• Need for train-test split: • MI models moust be ovaluated on date they have • MI models moust be ovaluated on date they have • MI models moust be ovaluated on date they have • MI models moust be ovaluated on date they have • MI models moust be ovaluated on date they have
Teaining a model on the entire deturet con lead to overfitting "where the model mome en the data instead of learning potterns."
• A sepecte test set allows ust to measure value model perforace ensuring its objility to hunde how inputs effectively.
· significance of touin - tosa split:
- prevents overfitting: gnsuses that the modelyend
- peavites unbiased evaluation: The test data act
-Helps in Hyperposeunctes: Spitting date allows fine tining of model peremetes perox And development

MET Bhujbal Knowledge City Institute of Technology- Polytechnic

	Adgaon, Nashik- 422 003. Date:
-efficient mode	1 selection:
Helps in	comparing different models to
select the	s helt performance one
- MARKETTANA	the state of the s
Algorithm:	ed in the state of the second of the
Impost nece	model solection; & skieden meters
Skieden, i	model salection & skieden meters
. 1 21 Y	
100d the da	taset using pandas read - CSU ()
4) Split the da	tatet into teatining & : testing sets.
Using teain	taket into tecining & :- testing sets : _ test _ spilit (x, y, test _ stree = 0.2, - 1.
Yandom	-State= 42).
	and the second s
s) Train a ma	dine leasning model on the training
fe↓.	σ
-6	tal
) peolice the	tager voiable on the tuting tel.
	Invict we chart it is talen.
Cralyane to	model ising a changing medice such
in and the	goe clussification) de RMSF (foe:
as a covey	(90 & C10) Sigication) 0 & R11 (5 (108)
regelssion)	The Mark of the Contract of th



MET Bhujbal Knowledge City

Institute of Technology- Polytechnic Adgaon, Nashik- 422 003.

Date:

occuracy = occuracy gross (4-test 4-red)
peint (" newacif:", Uccueacy)
candusion:
The teain test spit method is in essential stp
in machine reasoning to evaluate model perforcace
By keeping a purting of the date for testing,
we can coseen that the model generalized
well to unseen duta-treept reducing overefiting
& improve predictive dicieucit
Reference:
isciket leven documentation:
Action (and) laborations of the second
https:// scikit loven org/stuble/ modules/cos-
Varidation html. 101210
a) muching langhing
2) muchine leaching
hill and the second control
Allos: 11 www. conered. Ord.

