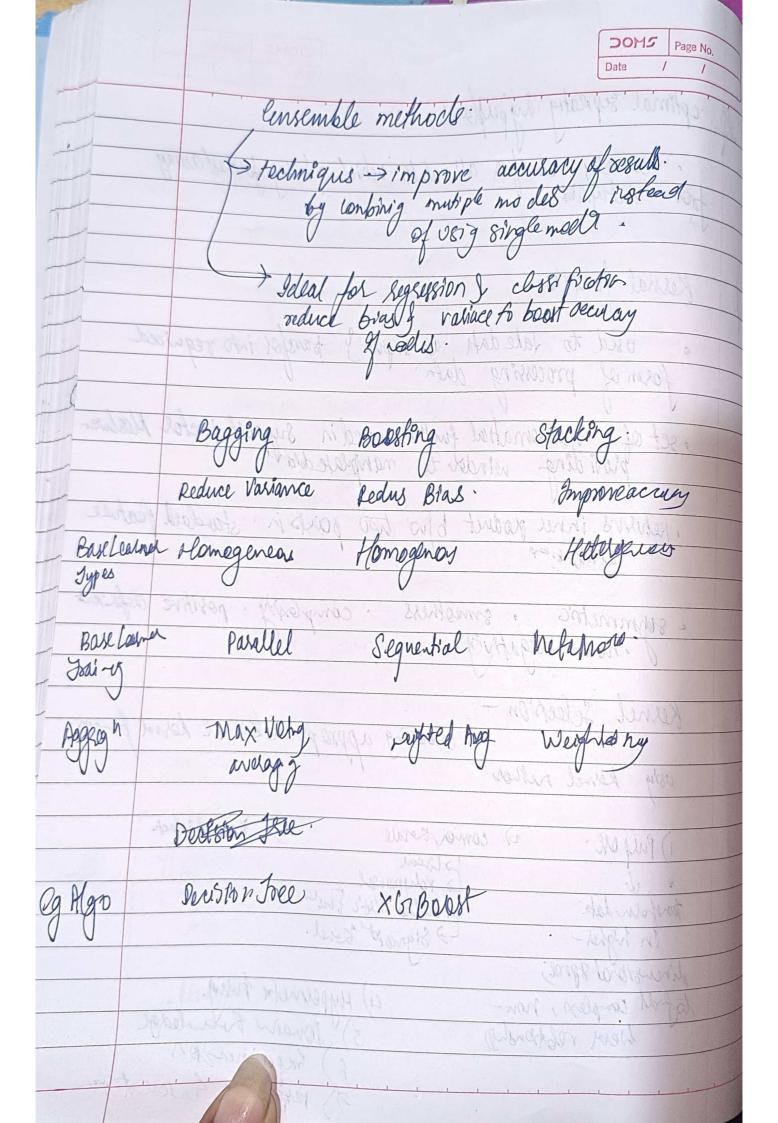
DS&ML DOM5 0-2 classification Methods Support Vector Machine: a

supervoused M. olgo

binary classification: > goal is to septate data points into · Hyperplane -> fonds it that best seprets. · Support vector > datapoint : closest to hyperplan : Keenel Joick - sum shands nonlineded by using beenel fuch. Multi Class Classifiel. " Sensatiti'll to scal · Strongths -> effethe in light directions spark. robust againt -> over fitty -> when hard-· Demerits -> , expensive for large dataself.
· model surpretation -> chology. Appl >) Image classiff, test categorium de Hyperplace in many my many into two or more regions. , lirectly proportional to no of feature

	DOM5 Page No.	
	Date / /	
D	aptimal separating hyperplaies-	
	from dutypeint cloudies all dute which being forthest away	y
	from dutypens	
	- Morastilla kion ta na na manda	The Air
la		
X	Kernal functions -	
1	Lernal functions - o used to take data as input I transfer into require form of processing data	
	o used to take data as input y transfer into require	ed
	form of processing data.	
	e set of mathematical functions used in Support Vector	Machie
	set of mathematical fuctions used in support vector providing window to manipulate data	
	REGINE VINTANCE REGINS BIRS. PARS.	^ .
	refulls inner product blo two pours in standard of	carre
	returns inner product b/w two pours in standard of	SANC LANGE
		778 ST.
	· summetoic · smeathers · complainty · positive a	lefn ki
	· summetric · smeathers · complainty · positive c	BOSELER
		At LOS
4	Kismed Solven on -	1
*	Kernel Selection - choosing appropriatelevel torrel volg kernel netters.	fran.
V	under have I so all one	UP
	Using Review Meroda	
	1) Pull Off. 2) comer Kerale 3) Kernelselsch	1000
	() God	
		wall.
	parfoln data Seadla Baric Fru	JAN 1
	in higher - Signord Konel.	7
	chrev-storal spral	
	to fell complex, non- 4) Hyperneta trung- breve rebetonship 5) Donair Reacledge	
	to fell complex, non- west relationship 5) Donary to a ledge 6) Engainer 100 7 Perforace Great	
	6) sagainer de	Year I
	of Perforace greater	

עומגעע דעטורוויד. אז



0	Date / /
-	Random toeest without sunt making
	· consists of multiple random declsson tree-
	· Two types of Jandomnesses are built into trees-
	, I each tole is built on randdom samplefrom original date.
	2) each tree is built on randdom samplefrom original date. 2) each tree rode - subserted feature are varidantly selected
	to generate best # split-
	· greater no. of trees leads to higher accurage prevents everfty
	· takes less training times when carpaced to others
	predicts output -> high accuracy > even for large details.
	· maintains acculary > when large portion of data is missy
	· Performs both classification & Rogers's task.
	can handle large dataset win high dinessorbit
^	vsed for both classif of legistron is not sustable for legislature. Randon Feople Selection
Deme	wife - Ramalon Feature Jevering
	16 model tuning > polimize model patolinee
	Hyperparameter tuning -> optimize model perfolioce
	Growing of Random Forest > process of creating muliple decision trees
4-	Consina of Random Forest
	Strawing of national to process of creating muliple decision trees
	with ensemble -
	o combines decision trees - improve accurage reduce overfity
	of all 21 mg preparation 2) Bookstrapping 3) Random France Sedection
	Steps - 1 Decision Tole Building 5) Parallelization
	Greenble Greation 7) Prediction Aggs of on
	8) Butol Bag who shorter 9) Featre Importable.
	Steps -> 1) Data preparation 2) Bookstrapping 3) Random Franke Stelection (4) Decision Tree Building 5) Parallelization (6) Ensemble Greation 7) Prediction Aggregation (8) But of Bog Error Strathon 9) Feature Importance. (10) Hypergal anether turky

