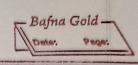
	(ab-5					
14/25	Bafna Gold		lassification	1		
210/	KNN	4	us) pear	and the state of t		
					a t. test	Anto
ony	Consider the f	ollowing	dataset	for K=	1 21 101	Owa
0111	(x, 35, (00) as	Person,	Age, Salary	16/30	IVE COSTE	29
011	Knn Classifier	& orca	get the to	orget.	Distante	Pont
1 27/17	Person	Dage	Salaryk	Carga	Distance 52.8	Jon
	A	18	50	N	46.6	-
	В	23	. 55	N	31.9	2
	6	24	70	Y	40.4	3
	P and	41	60	Y	31.1	-
	E+84 11.00	43	40	4.	60.1	divid
	F	38	40	9	1842-41P	
	X	35	100	nikimi ?	16 N 80/4	
	Nation Louis	merchant	They now	2+ Calary	42 - Salory	(1)2
~	Euclidean distan	ia: d=	1 (Ayer Ayer	ne moisse	ma) to	orkins.
3tep-14		MABIL	716	Losved	JA S	4
tons	A-> J(35-18) +	(100 - 30)	- 9 2.8		Sel phou	步.
	B > 46.6	ÀN.	Clare the	Rossilinas	12 Misc	
	C= 31.9	trende to	Thou or	Nuh!	molle	1
	MINISTER	2 ransiche	, KN G		1791	9
	E -> . 31.1				373	41
	F->60.1					1
0 1	Identify 3 nea	rist neig	hbours/		1/2/1	
Step:28	r. (31) u)					
12	E: (31.1, y) C: (31.9, N)					
2	D: (40, 4, 4)	1				
)	7. (, ,				4 100	
9tcp:34	Majority roting	j :			2 Mil	
	Yes				9	
->	x (35, 100) is	У,	. was arth		NO NUL	
			man bar	Alexander	9 1	_
						_
F54 5G						



On For ions dedaset

It flow to Choose the Knalue ? Denonstrate using

accuracy rate & correr rate. On: 2 For diabetes dataset.

What is the purpose of feature scaling ? How to

perform it? * Measure accuracy (higher is better) & error

rate (lower is better) The best K is where accuracy is highest & error rate is lowest (typically 5-10) features can dominate the smaller ones.

Standardization ensures all features contribute equally. * Scaling improves accuracy in prevents biased predictions.

axy F54 5G