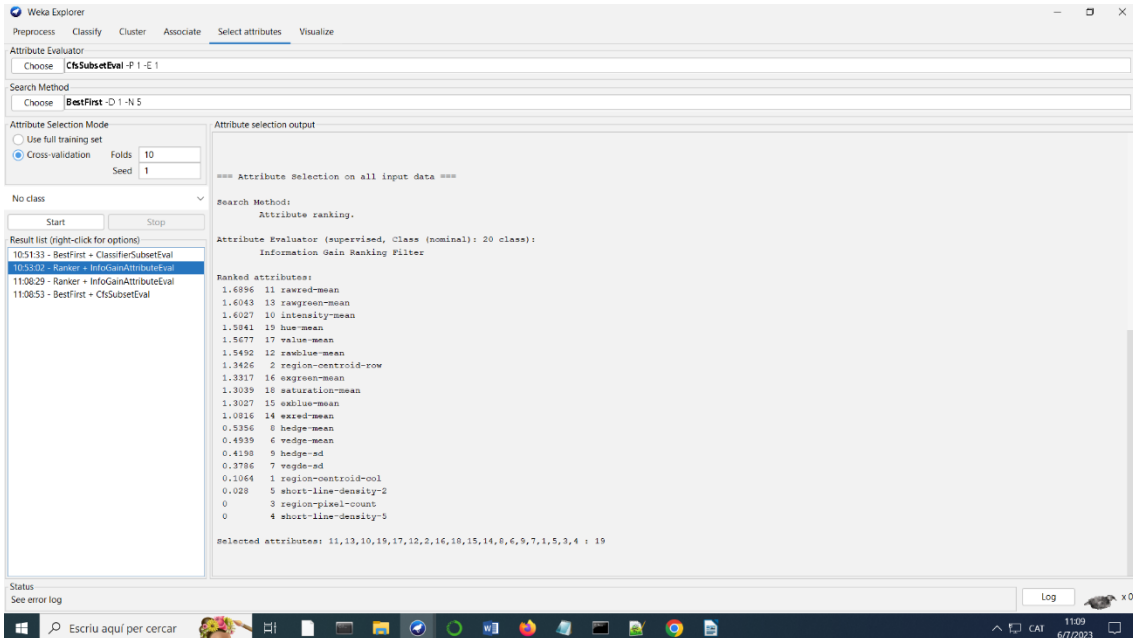


Pràctica final. Weka

Primer valoro quins atributs son menys rellevants per poder eliminar-los i que la meva predicció s'acosti millor al valor real.



Amb aquesta

Amb aquesta comprovació decideixo eliminar les dues darreres columnes que ens indica en la imatge de dalt que tenen menys incidència.

En la següent imatge també comprovem que la influencia es poca i puntual.

ARFF-Viewer - C:\Users\Alumne_mati\Desktop\segment-challenge.arff

segment-challenge.arff

Relations: segment

Id	1: region-centroid-col	2: region-centroid-row	3: region-pixel-count	4: short-line-density-2	5: short-line-density-5	6: wedge-mean	7: wedge-ad	8: hedge-mean	9: hedge-ad
1	38.0	189.0	9.0	0.0	0.0	1.0	0.22222	6.22222	33.3185
2	25.0	199.0	9.0	0.0	0.0	1.11111	0.607407	1.05556	0.462963
3	49.0	139.0	9.0	0.0	0.0	0.166667	0.0777778	0.0777778	0.0777778
4	63.0	220.0	9.0	0.0	0.0	3.05556	1.5	0.055556	0.136083
5	161.0	135.0	9.0	0.0	0.0	0.055556	0.136083	0.111111	0.172133
6	235.0	88.0	9.0	0.0	0.0	0.611111	0.240741	0.944445	0.32963
7	67.0	32.0	9.0	0.0	0.0	0.944444	1.06394	1.77778	1.31092
8	188.0	182.0	9.0	0.0	0.0	1.61111	0.742898	4.16667	2.12655
9	217.0	245.0	9.0	0.111111	0.111111	3.16667	3.01662	2.16667	1.24276
10	9.0	171.0	9.0	0.0	0.0	1.5	1.00554	2.77778	1.64204
11	149.0	117.0	9.0	0.222222	0.0	0.833333	0.255556	1.0	0.444445
12	136.0	139.0	9.0	0.0	0.0	0.666667	0.177778	0.777778	0.207407
13	214.0	79.0	9.0	0.0	0.0	1.27778	0.862963	1.33333	0.533333
14	118.0	125.0	9.0	0.0	0.0	0.333333	0.298142	0.888889	0.344265
15	83.0	53.0	9.0	0.0	0.0	0.555556	0.455419	0.722222	0.490652
16	70.0	38.0	9.0	0.0	0.0	1.22222	0.272166	1.0	0.760116
17	96.0	84.0	9.0	0.0	0.0	1.5	1.27778	1.61111	2.28519
18	151.0	129.0	9.0	0.0	0.0	8.05556	8.14157	0.722222	0.827758
19	98.0	144.0	9.0	0.0	0.0	0.222222	0.02963	0.333333	0.088889
20	146.0	140.0	9.0	0.0	0.0	1.05556	0.462963	1.0	0.577778
21	162.0	237.0	9.0	0.111111	0.0	2.27778	1.14342	2.66667	2.56472
22	155.0	27.0	9.0	0.0	0.0	1.05556	0.240737	1.44445	0.25185
23	83.0	198.0	9.0	0.0	0.0	2.38889	1.49691	4.83333	3.00925
24	94.0	104.0	9.0	0.0	0.0	6.61111	1.42075	1.0	0.516399
25	220.0	111.0	9.0	0.0	0.0	2.55556	1.65552	2.27778	1.59745
26	157.0	85.0	9.0	0.0	0.0	1.22222	1.24127	0.222222	0.172133
27	18.0	145.0	9.0	0.0	0.0	0.388889	0.018519	0.611111	0.374074
28	95.0	191.0	9.0	0.0	0.111111	3.16667	2.78687	6.0	5.18545
29	112.0	38.0	9.0	0.0	0.0	11.0556	11.0743	1.11111	1.06805
30	206.0	133.0	9.0	0.0	0.0	0.055556	0.136083	0.055556	0.136083
31	222.0	145.0	9.0	0.0	0.0	0.0	0.0	0.0	0.0
32	217.0	208.0	9.0	0.0	0.0	7.16667	7.66667	7.44445	7.44445

Iniciem el primer model predictiu:

NaiveBayes:

Scheme: weka.classifiers.bayes.NaiveBayes

Relation: segment-weka.filters.unsupervised.attribute.Remove-R4-5

Instances: 1500

Test mode: split 80.0% train, remainder test

Time taken to build model: 0 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances	250	83.3333 %
--------------------------------	-----	-----------

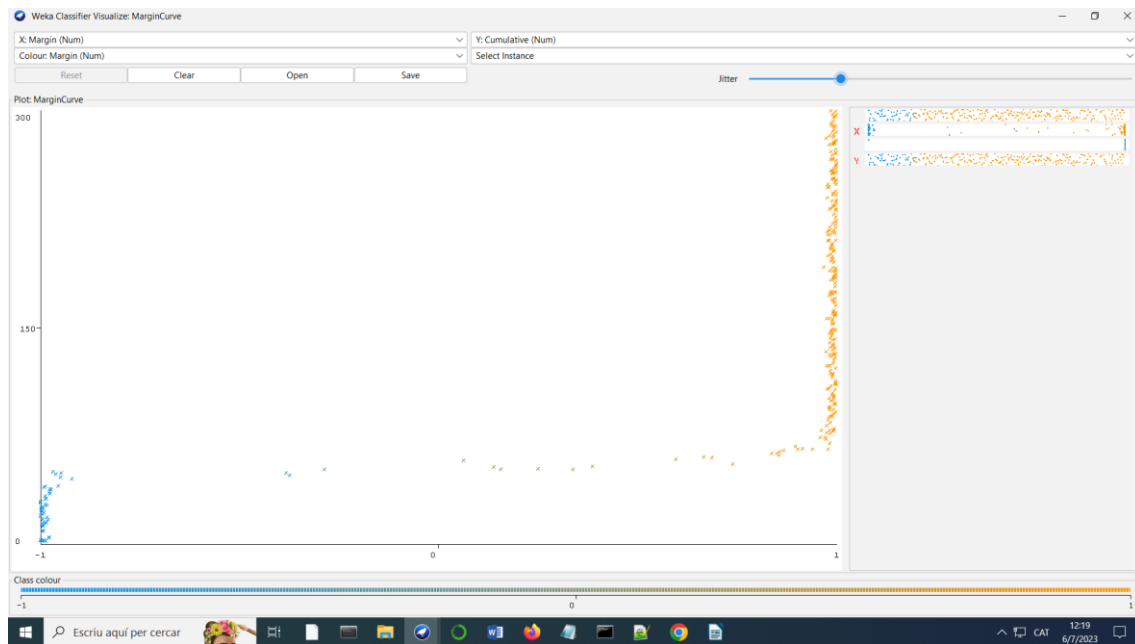
Incorrectly Classified Instances	50	16.6667 %
----------------------------------	----	-----------

=== Detailed Accuracy By Class ===

Weighted Avg.	0,833	0,023	0,847	0,833	0,816	0,807	0,977	0,897
---------------	-------	-------	-------	-------	-------	-------	-------	-------

Decideixo fer servir percentatge Split 80% perquè el procés és més ràpid.

Aquest model és el pitjor dels tres escollits perquè en el Summary els valors escollits (a dalt) son més baixos que en els altres dos, a més Accuracy també dona un resultat més baix.



Logistic:

Scheme: weka.classifiers.functions.Logistic -R 1.0E-8 -M -1 -num-decimal-places 4

Relation: segment-weka.filters.unsupervised.attribute.Remove-R4-5

Instances: 1500

Test mode: split 80.0% train, remainder test

Time taken to build model: 0.58 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0 seconds

=== Summary ===

Correctly Classified Instances	286	95.3333 %
--------------------------------	-----	-----------

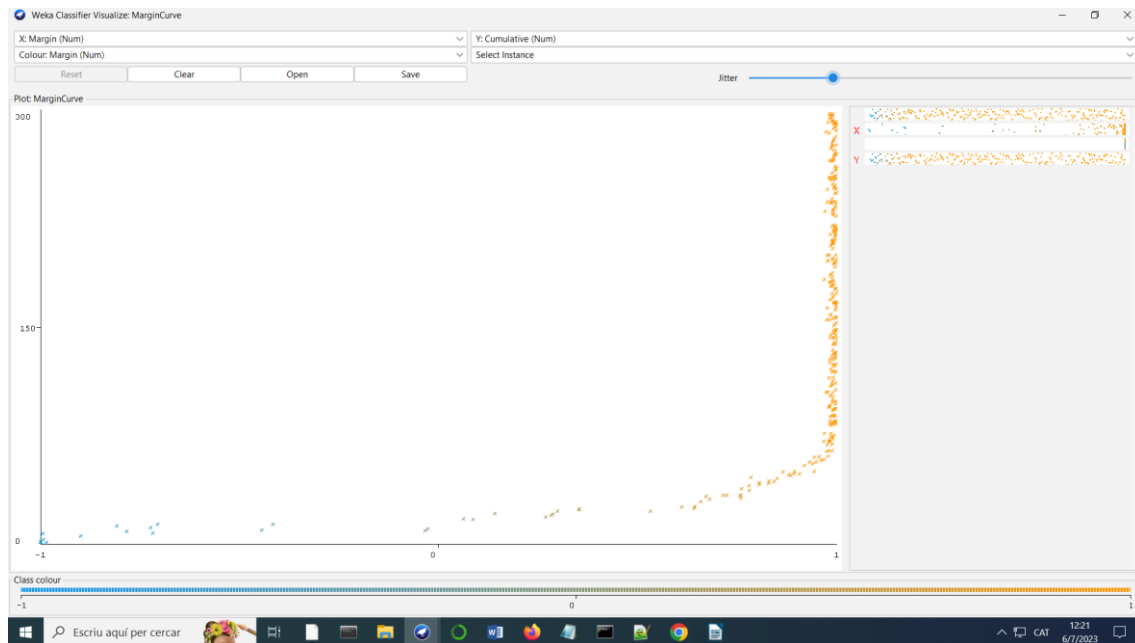
Incorrectly Classified Instances	14	4.6667 %
----------------------------------	----	----------

=== Detailed Accuracy By Class ===

Weighted Avg.	0,953	0,007	0,956	0,953	0,954	0,947	0,996	0,981
---------------	-------	-------	-------	-------	-------	-------	-------	-------

Escullo en aquest cas el percentatge Split per la mateixa raó que abans, el temps del procés del model és més curt.

Aquest model s'acosta més a unes dades predictives més valides amb valors de les instàncies correctes i incorrectes més fiables. També noto millora en la precisió. Ens acostem a un model millor d'acord amb les dades.



RandomForest:

Scheme: weka.classifiers.trees.RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1

Relation: segment-weka.filters.unsupervised.attribute.Remove-R4-5

Instances: 1500

Test mode: split 80.0% train, remainder test

=== Evaluation on test split ===

Time taken to test model on test split: 0.01 seconds

=== Summary ===

Correctly Classified Instances	298	99.3333 %
Incorrectly Classified Instances	2	0.6667 %

=== Detailed Accuracy By Class ===

Weighted Avg. 0,993 0,001 0,994 0,993 0,993 0,992 1,000 0,997

Amb la mateixa raó que els altres dos models escullo la modalitat del test de 80%/20% perquè és més ràpida en la execució.

Comprovant les dades determino que aquest model RandomForest és el que millors dades predictives ens donarà. Podem comprovar que els valors són més alts en tots els valors presentats en aquest document.

He escollit sempre els mateixos tipus de valors per poder fer una comparació més fiable, tot hi que cada model dona més importància i informació sobre valors diferents (com pot ser el cas de les mitjanes en el model NaiveBayes)

Finalment he posat les imatges de la MarginCurbe on es veu que la NiveBayes és la més fluixa i com els altres dos models s'acosten més a una millor predicció.

