

Γλυκός Αθανάσιος Παναγιώτης p3190046

Γλυκού Χριστίνα Μαρία p3190047

Τεχνητή Νοημοσύνη: Εργασία 2

- αφελής ταξινομητής Bayes με πολυμεταβλητή μορφή Bernoulli

Το αρχείο μας ονομάζεται imdb

Αρχικά η πρώτη αλγόριθμος μάθησης που επιλέξαμε να υλοποιήσουμε είναι ο αφελής ταξινομητής Bayes με πολυμεταβλητή μορφή Bernoulli. Εισάγουμε στο πρόγραμμα μας τα δεδομένα από το imdb που αφορούν κάποιες κριτικές ταινιών και δημιουργούμε το λεξιλόγιο ανάλογα με τις παραμετρους μας .

Ύστερα δημιουργούμε τον πίνακα **x** ο οποίος περιέχει για κάθε review του **x_train** και για κάθε λέξη του λεξιλογίου αν περιέχεται στο κάθε review. Στην συνέχεια φτιάχνουμε έναν αντίστοιχο πίνακα **x2** αλλά για τα **x_test**.

Ύστερα με τη χρήση της μεθόδου **find_possibility** δημιουργούμε έναν πίνακα με πιθανότητες : **possibility** οποίος θα περιέχει 4 πιθανότητες, την πιθανότητα η λέξη να υπάρχει μέσα στο review και να κατατάσσεται στα θετικά reviews , η δεύτερη είναι να μην υπάρχει η λέξη και να κατατάσσεται στα θετικά reviews, η τρίτη να υπάρχει η λέξη και να κατατάσσεται στα αρνητικά reviews και τέλος να μην υπάρχει η λέξη και να κατατάσσεται στα αρνητικά reviews. Για να υπολογίσουμε τις πιθανότητες αυτές χρησιμοποιούμε τις μεταβλητές **yes_exists** , **yes_den_exists** , **no_exists** , **no_den_exists** και **sumyes** (πλήθος θετικών reviews) , **sumno** (πλήθος αρνητικών reviews). Χρησιμοποιούμε επίσης την εκτιμήτρια Laplace για να μην μηδενίζονται οι πιθανότητες.

Στη συνέχεια υπολογίζουμε τα ποσοστά των θετικών και των αρνητικών review που θα μας χρειαστούν αργότερα με τη μέθοδο **possibility_of_pos_neg**.

Έπειτα καλούμε τη μέθοδο **results** όπου υπολογίζουμε τις πιθανότητες **p_yes_y** και **p_no_y** που αντιστοιχούν στις πιθανότητες ενώ γνωρίζουμε

το review y , το review να είναι θετικό και αντίστοιχα ενώ γνωρίζουμε το review y , το review να είναι αρνητικό. Αυτό επιτυγχάνεται ως εξής : πχ για την πιθανότητα p_yes_y

Δημιουργούμε το γινόμενο

$$P(C=1|\vec{X}) = \frac{P(C=1) \cdot \prod_{i=1}^m P(X_i = x_i | C=1)}{P(\vec{X})} \quad (\text{ο παρονομαστής δεν μας απασχολεί})$$

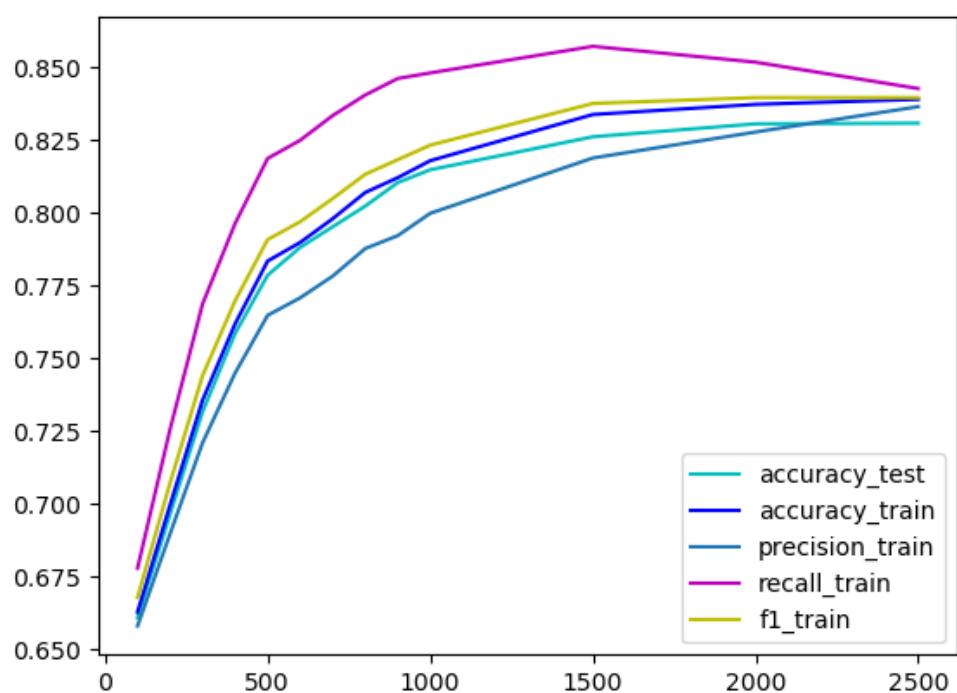
Όπου έχουμε ήδη υπολογίζει την πιθανότητα των θετικών review ,οπότε υπολογίζουμε το γινόμενο των $PP(X=x_i|C=1)$ το οποίο το βρίσκουμε από τον πίνακα possibility που δημιουργήσαμε πριν. Ομοίως και για το p_no_y .

Αν το p_yes_y είναι μεγαλύτερο από το p_no_y τότε κατηγοριοποιούμε το review y ως θετικό αλλιώς ως αρνητικό, και αποθηκεύουμε το αποτέλεσμα σε έναν πίνακα.

Στη συνέχεια συγκρίνουμε τον πίνακα των αποτελεσμάτων που υπολογίσαμε με τον πίνακα των κανονικών αποτελεσμάτων και υπολογίζουμε τα ποσοστά

Αυτή την διαδικασία την εκτελούμε για τα x_train και y_train και στο τέλος εμφανίζουμε τα accuracy κτλ.

λέξεις	100	200	300	400	500	600	700	800	900	1000	1500	2000	2500
Accuracy test	0.66104	0.6966	0.73176	0.7586	0.77856	0.78812	0.7952	0.80228	0.81036	0.81476	0.82608	0.83052	0.83076
Accuracy train	0.66296	0.69968	0.73568	0.76188	0.78344	0.7898	0.79796	0.80696	0.81204	0.81788	0.83372	0.8372	0.83892
precision	0.65820130 47530288	0.689809885931 5589	0.72107 158937 41558	0.74500 411645 8349	0.76479 820627 80269	0.77080 062794 34851	0.77825 924542 39821	0.78771 745650 86983	0.79215 040071 90473	0.79981 890892 62808	0.81880 015284 67711	0.82770 953195 4595	0.83641 705709 52116
recall	0.678	0.72568	0.76872	0.79632	0.81864	0.82488	0.83336	0.8404	0.84608	0.848	0.85712	0.85168	0.84264
F1	0.66795397 22572508	0.707290448343 0799	0.74413 381863 2386	0.76980 781872 31738	0.79080 370942 81299	0.79692 390926 30521	0.80486 768398 68648	0.81320 637869 63926	0.81822 753471 81928	0.82320 506348 77489	0.83752 198553 84014	0.83952 369686 9332	0.83951 699677 20082



- ID3, με πριόνισμα

Το αρχείο μας ονομάζεται imdb2

Αρχικά, διαβάζουμε τα δεδομένα μας με τον ίδιο τρόπο που τα διαβάσαμε και στον bayes και δημιουργούμε ξανά δυο πίνακες, τον x που αφορά τα train και τον x_2 που αφορά τα test, όπου για κάθε feature του λεξιλογίου έχουν τιμή 1 αν το περιέχουν και τιμή 0 αν όχι.

Στη συνέχεια δημιουργούμε τις εξής μεθόδους:

- `twoCEntropy`: η μέθοδος αυτή επιστρέφει την εντροπία 2 μεταβλητών σύμφωνα με γνωστό μαθηματικό τύπο.
- `calculatelG`: η συνάρτηση αυτή υπολογίζει έναν πίνακα με τα information gain των ιδιοτήτων μας χρησιμοποιώντας την εντροπία (`twoCEntropy`) και τις πιθανότητες που υπολογίζουμε.
- `most_informative_feature`: αρχικά υπολογίζουμε τα IG με τη χρήση της συνάρτησης `calculatelG`, βρίσκουμε το max μεταξύ αυτών και παράλληλα ελέγχουμε αν αυτό το feature έχει ξαναχρησιμοποιηθεί(χρησιμοποιώντας τον πίνακα `xrisimoroihmenes_theseis`)
- `make_sub_tree`: Αρχικά υπολογίζουμε το πλήθος των reviews που περιέχουν το συγκεκριμένο feature (`count_1`) και το πλήθος αυτών που δεν το έχουν (`count_0`) . Επίσης ξεχωρίζουμε σε 2 πίνακες τα reviews που έχουν το feature και τα reviews που δεν το έχουν . Για κάθε περίπτωση υπολογίζουμε πόσα reviews είναι θετικά και πόσα είναι αρνητικά . Αναλυτικότερα, έχουμε δηλαδή 4 περιπτώσεις οι οποίες είναι :
`class0_count1` : όσα στην στήλη feature =1 το $y_{train}=1$
`class0_count0` : όσα στην στήλη feature =1 το $y_{train}=0$
`class1_count1` : όσα στην στήλη feature =0 το $y_{train}=1$
`class0_count1` : όσα στην στήλη feature =0 το $y_{train}=0$

Ύστερα σε έναν πίνακα `most_common` κρατάμε για το συγκεκριμένο feature αν τα περισσότερα reviews είναι θετικά ή αρνητικά το οποίο το χρειαζόμαστε για την συνέχεια. Δηλαδή

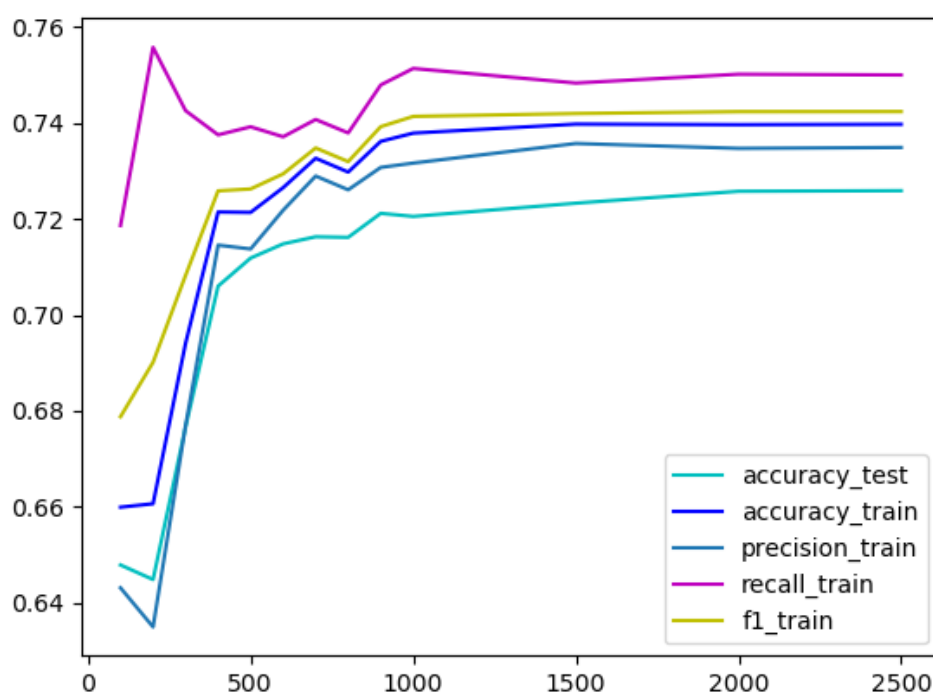
στην περίπτωση που μας τελειώσουν είτε οι ιδιότητες είτε τα παραδείγματα και ενώ ο κόμβος έπρεπε να επεκταθεί δεν μπορούμε να τον επεκτείνουμε περεταίρω.

Στην συνέχεια αποφασίζουμε αν ο κόμβος στον οποίο βρισκόμαστε θα επεκταθεί ή όχι χρησιμοποιώντας την πιθανότητα για τις περιπτώσεις το feature να είναι 1 και αν είναι 0 .Αν το ποσοστό των θετικών ή των αρνητικών reviews είναι ίσο ή ξεπερνα το 70% τότε δεν τον επεκτείνουμε άλλο και του κατατάσσουμε το review ως θετικό ή αρνητικό αντίστοιχα. Αλλιώς το μαρκάρουμε με « ? » που σημαίνει ότι μπορεί να επεκταθεί παραπάνω.

- `generate_tree`: μέσω αναδρομικών κλήσεων όσο δεν μας έχουν τελειώσει τα reviews και οι ιδιότητες που μπορούμε να χρησιμοποιήσουμε, βρίσκουμε το most informative feature και υπολογίζουμε το δέντρο και το νέο πίνακα με reviews και αποτελέσματα, βρίσκουμε το νέο κόμβο που αντιστοιχεί στο δέντρο μας και αν είναι δυνατή η επέκτασή του συνεχίζουμε αναδρομικά της κλήσεις μας.
- `Id3`: σε αυτή τη συνάρτηση απλά καλούμε την `generate_tree` και επιστρέφουμε το δέντρο μας.
- `decision`: μας δίνεται ως είσοδος το δέντρο που δημιουργήσαμε πριν και ένα review. Προσπελαύνουμε τους κόμβους του δέντρου ελέγχοντας την αντιστοιχία με το review ώστε να καταλήξουμε σε φύλλο του δέντρου που θα είναι το αποτέλεσμα του υπολογισμού κατηγοριοποίησης του review.
- `Evaluate`: Τέλος έχουμε την μέθοδο `evaluate` όπου για κάθε review ενός set υπολογίζουμε τα ποσοστά επιτυχίας κτλ.(Αν ένα φύλλο είχε την τιμή ? χρησιμοποιούμε τον πίνακα most common και αποφασίζουμε αν το review θα κατηγοριοποιηθεί ως θετικό ή αρνητικό.

Οι διαδικασίες εκτελείται για τα test και τα train και εμφανίζουμε τις μετρήσεις που μας ζητούνται.

λέξεις	100	200	300	400	500	600	700	800	900	1000	1500	2000	2500
Accuracy test	0.64788	0.64484	0.67664	0.706	0.71184	0.7148	0.71628	0.71616	0.72116	0.72052	0.72328	0.72576	0.72588
Accuracy train	0.65992	0.66064	0.69408	0.72144	0.72136	0.72652	0.73264	0.72976	0.73616	0.73788	0.73976	0.73964	0.73972
precision	0.64311282 93241696	0.634964376932 3834	0.67692 532088 68145	0.71454 038133 6227	0.71373 397188 32072	0.72181 746964 35566	0.72894 032435 83688	0.72607 052896 72544	0.73073 315616 69532	0.73163 511723 92304	0.73572 439830 10854	0.73470 187260 04858	0.73489 064827 15372
recall	0.71864	0.75576	0.74256	0.73752	0.7392	0.73712	0.74072	0.73792	0.74792	0.75136	0.74832	0.75016	0.75
F1	0.67878192 53438114	0.690116151654 6132	0.70822 524034 79322	0.72584 835839 69767	0.72624 381042 20702	0.72938 848208 98476	0.73478 295373 38306	0.73194 730995 08016	0.73922 669407 76469	0.74136 638118 1671	0.74196 874752 12183	0.74235 047302 37899	0.74236 845231 02507



Τα προγράμματα τα τρέχουμε **δίνοντας 100,...,2500 ως number of words και 0 ως skip top words**. Δηλαδή τα διαγράμματα αφορούν αυτές τις περιπτώσεις.