

HW5

1a. Multinomial Naive Bayes

Υπολογίζω τις πιθανότητες των κλάσεων TV|yes και TV|no

$$P(\text{TV}|\text{yes}) = (3/5) = 0,6$$

$$P(\text{TV}|\text{no}) = (2/5) = 0,4$$

Υπολογίζω τις πιθανότητες των λέξεων στα έγγραφα στις κλάσεις TV|yes και TV|no

$$P(\text{πρόγραμμα}|\text{yes}) = (3+1)/(14+8) = 0,182$$

$$P(\text{πρόγραμμα}|\text{no}) = (0+1)/(6+8) = 0,071$$

$$P(\text{επεισόδιο}|\text{yes}) = (5+1)/(14+8) = 0,273$$

$$P(\text{επεισόδιο}|\text{no}) = (1+1)/(6+8) = 0,143$$

$$P(\text{σειρά}|\text{yes}) = (1+1)/(14+8) = 0,091$$

$$P(\text{σειρά}|\text{no}) = (0+1)/(6+8) = 0,071$$

$$P(\text{κανάλι}|\text{yes}) = (3+1)/(14+8) = 0,182$$

$$P(\text{κανάλι}|\text{no}) = (0+1)/(6+8) = 0,071$$

$$P(\text{ταινία}|\text{yes}) = (1+1)/(14+8) = 0,091$$

$$P(\text{ταινία}|\text{no}) = (0+1)/(6+8) = 0,071$$

$$P(\text{ειδήσεις}|\text{yes}) = (1+1)/(14+8) = 0,091$$

$$P(\text{ειδήσεις}|\text{no}) = (1+1)/(6+8) = 0,143$$

$$P(\text{γήπεδο}|\text{yes}) = (0+1)/(14+8) = 0,045$$

$$P(\text{γήπεδο}|\text{no}) = (2+1)/(6+8) = 0,214$$

$$P(\text{ομάδα}|\text{yes}) = (0+1)/(14+8) = 0,045$$

$$P(\text{ομάδα}|\text{no}) = (2+1)/(6+8) = 0,214$$

Υπολογίζω τις πιθανότητες του κατηγοριοποιητή

$$\mathbf{P(d6|Yes)} = P(\text{TV}|\text{yes}) * P(\text{επεισόδιο}|\text{yes})^2 * P(\text{γήπεδο}|\text{yes})^3 * P(\text{ειδήσεις}|\text{yes}) = 0,6 * (0,273)^2 * (0,045)^2 * 0,091 = \mathbf{0,0000003810199028}$$

$$\mathbf{P(d6|No)} = P(\text{TV}|\text{no}) * P(\text{επεισόδιο}|\text{no})^2 * P(\text{γήπεδο}|\text{no})^3 * P(\text{ειδήσεις}|\text{no}) = 0,4 * (0,143)^2 * (0,214)^3 * 0,143 = \mathbf{0,00001147481067}$$

Άρα το έγγραφο d6 θα κατηγοριοποιηθεί στην κατηγορία **P(TV|no)** με την μεγαλύτερη πιθανότητα.

1β. Bernulli Naive Bayes

Υπολογίζω τις πιθανότητες των κλάσεων TV|yes και TV|no

$$P(\text{TV}|\text{yes}) = (3/5) = 0,6$$

$$P(\text{TV}|\text{no}) = (2/5) = 0,4$$

Υπολογίζω τις πιθανότητες της μοναδικής παρουσίας των λέξεων στα έγγραφα στις κλάσεις TV|yes και TV|no

$$P(\text{πρόγραμμα}|\text{yes}) = (2+1)/(3+2) = 0,6$$

$$P(\text{πρόγραμμα}|\text{no}) = (0+1)/(2+2) = 0,25$$

$$P(\text{επεισόδιο}|\text{yes}) = (3+1)/(3+2) = 0,8$$

$$P(\text{επεισόδιο}|\text{no}) = (1+1)/(2+2) = 0,5$$

$$P(\text{σειρά}|\text{yes}) = (1+1)/(3+2) = 0,4$$

$$P(\text{σειρά}|\text{no}) = (0+1)/(2+2) = 0,25$$

$$P(\text{κανάλι}|\text{yes}) = (2+1)/(3+2) = 0,6$$

$$P(\text{κανάλι}|\text{no}) = (0+1)/(2+2) = 0,25$$

$$P(\text{ταινία}|\text{yes}) = (1+1)/(3+2) = 0,4$$

$$P(\text{ταινία}|\text{no}) = (0+1)/(2+2) = 0,25$$

$$P(\text{ειδήσεις}|\text{yes}) = (1+1)/(3+2) = 0,4$$

$$P(\text{ειδήσεις}|\text{no}) = (1+1)/(2+2) = 0,5$$

$$P(\text{γήπεδο}|\text{yes}) = (0+1)/(3+2) = 0,2$$

$$P(\text{γήπεδο}|\text{no}) = (2+1)/(2+2) = 0,75$$

$$P(\text{ομάδα}|\text{yes}) = (0+1)/(3+2) = 0,2$$

$$P(\text{ομάδα}|\text{no}) = (2+1)/(2+2) = 0,75$$

Υπολογίζω τις πιθανότητες του κατηγοριοποιητή

$$\begin{aligned} P(\mathbf{d6}|\mathbf{Yes}) &= P(\text{TV}|\text{yes}) * P(\text{επεισόδιο}|\text{yes}) * P(\text{ειδήσεις}|\text{yes}) * P(\text{γήπεδο}|\text{yes}) * (1- \\ &P(\text{πρόγραμμα}|\text{yes})) * (1- P(\text{σειρά}|\text{yes})) * (1- P(\text{κανάλι}|\text{yes})) * (1- P(\text{ταινία}|\text{yes})) * (1- P(\text{ομάδα}|\text{yes})) = \\ &0,6*0,8*0,4*0,2*(1-0,6)*(1-0,4)*(1-0,6)*(1-0,4)*(1-0,2) = \mathbf{0,001769472} \end{aligned}$$

$$\begin{aligned} P(\mathbf{d6}|\mathbf{No}) &= P(\text{TV}|\text{no}) * P(\text{επεισόδιο}|\text{no}) * P(\text{γήπεδο}|\text{no}) * P(\text{ειδήσεις}|\text{no}) * (1- P(\text{πρόγραμμα}|\text{no})) * \\ &(1- P(\text{σειρά}|\text{no})) * (1- P(\text{κανάλι}|\text{no})) * (1- P(\text{ταινία}|\text{no})) * (1- P(\text{ομάδα}|\text{no})) = \\ &0,4*0,5*0,75*0,5*(1-0,25)*(1-0,25)*(1-0,25)*(1-0,25)*(1-0,75) = \mathbf{0,005932617188} \end{aligned}$$

Άρα το έγγραφο d6 θα κατηγοριοποιηθεί στην κατηγορία **P(TV|no)** με την μεγαλύτερη πιθανότητα.

2. Κατηγοριοποίηση κειμένου με Rapid Miner

2.1 Εισαγωγή

Έχουμε ένα σύνολο δεδομένων από τη συλλογή με Newsgroups άρθρα.

Έχουμε 9 κατηγορίες:

Pc Hardware

Mac Hardware

Sport Baseball

Sport Hockey

Science Space

Misc ForSale

Politics

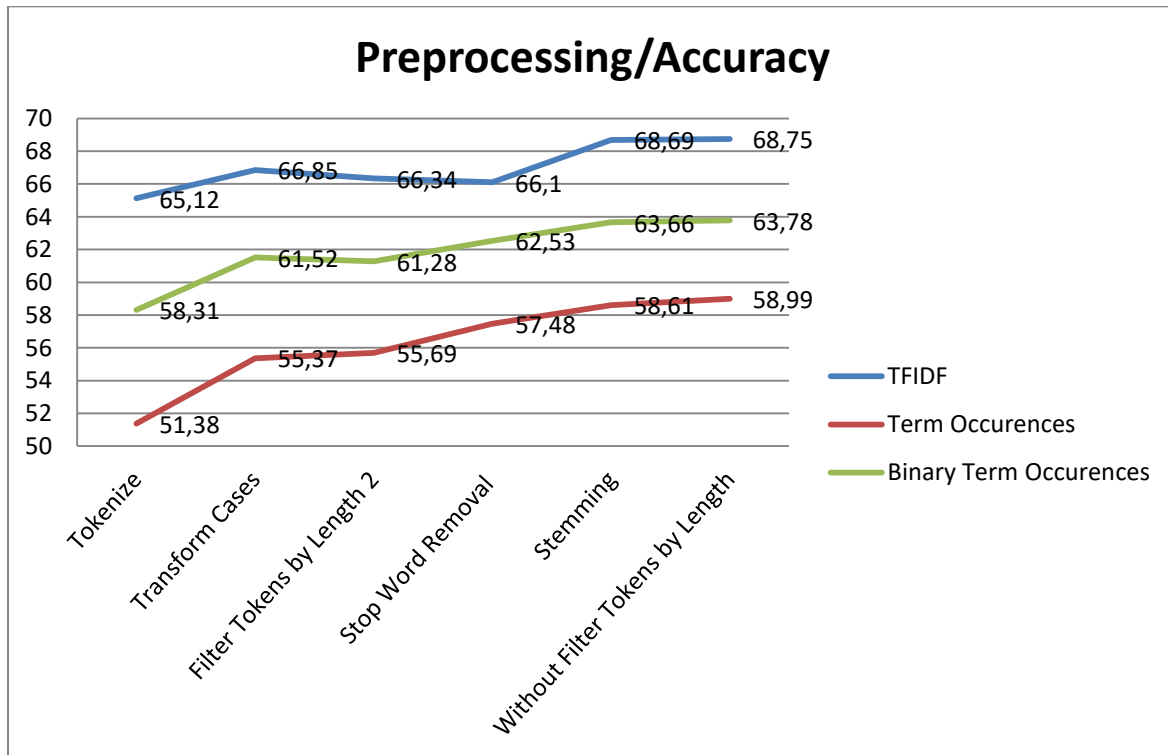
Religion

Atheism

2.2 Preprocessing

Σχεδιάσα την διαδικασία της προεπεξεργασίας ως εξής:

Δοκίμασα τις τεχνικές αυξητικά με κριτήριο το accuracy και με τα 3 διανύσματα, TFIDF, Term Occurrences και Binary Term Occurrences. Επίσης το pruning το έθεσα στο 3%/97% για να έχω γρήγορους χρόνους επεξεργασίας (~10sec). Αποτέλεσμα το παρακάτω γράφημα:



Γράφημα 1. Preprocessing

Μετά από το βασικό Tokenize, εφάρμοσα το Transform Cases, το οποίο βελτίωσε το accuracy και στα 3 διανύσματα, στην συνέχεια φιλτράρισα τις λέξεις με 1 στοιχείο αλλά εδώ δεν είχα βελτίωση, στο TFIDF και στο Binary Term Occurrences υπήρξε ελαφριά μείωση του accuracy. Στη συνέχεια εφάρμοσα τα Stop Word Removal και Stemming(Porter), τα οποία βελτίωσαν την ακρίβεια αρκετά.

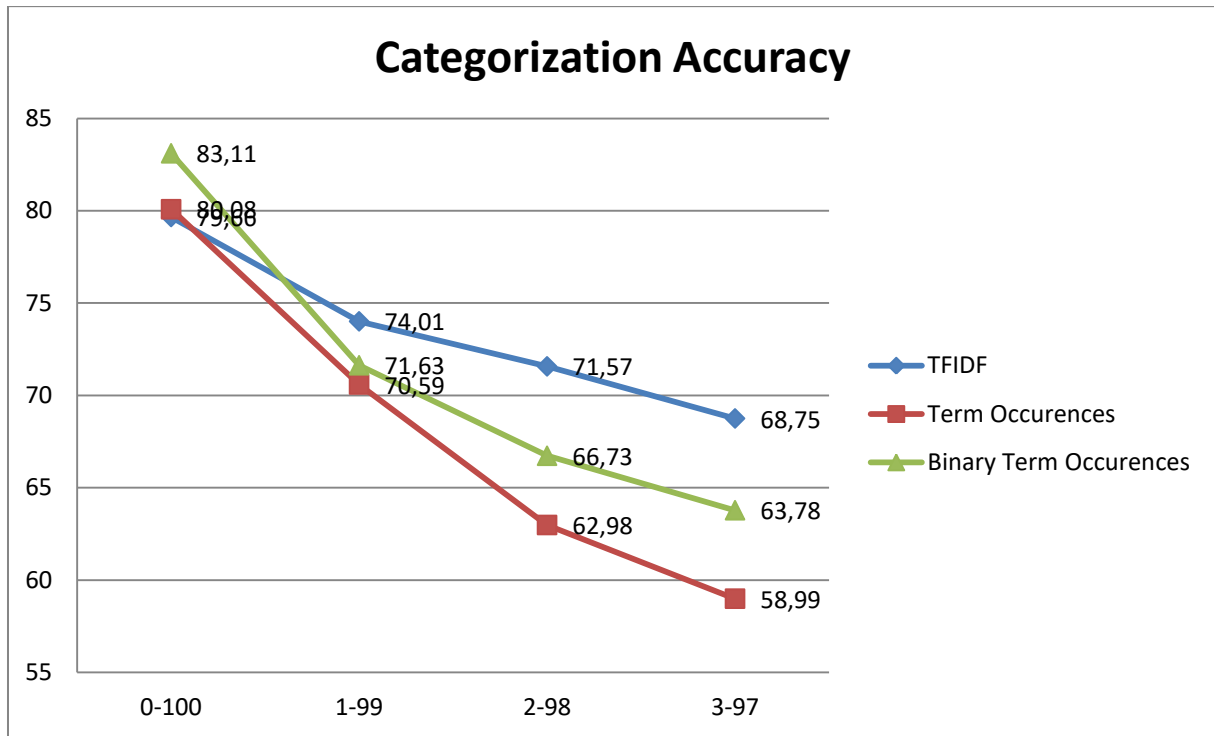
Τέλος δοκίμασα να αφαιρέσω το φίλτρο Filter Tokens by Length και είδα βελτίωση στο accuracy οπότε αποφάσισα να μην το χρησιμοποιήσω για την επεξεργασία.

Τελικό setup για την περαιτέρω επεξεργασία:

Tokenize, Transform Cases, Stop Word Removal, Stemming

2.2 Κατηγοριοποίηση

Με δεδομένο το Preprocessing Pipeline και την εφαρμογή του Naive Bayes classifier πήρα τα παρακάτω αποτελέσματα για διαφορετικές τιμές Prunning και τα 3 διαφορετικά διανύσματα.

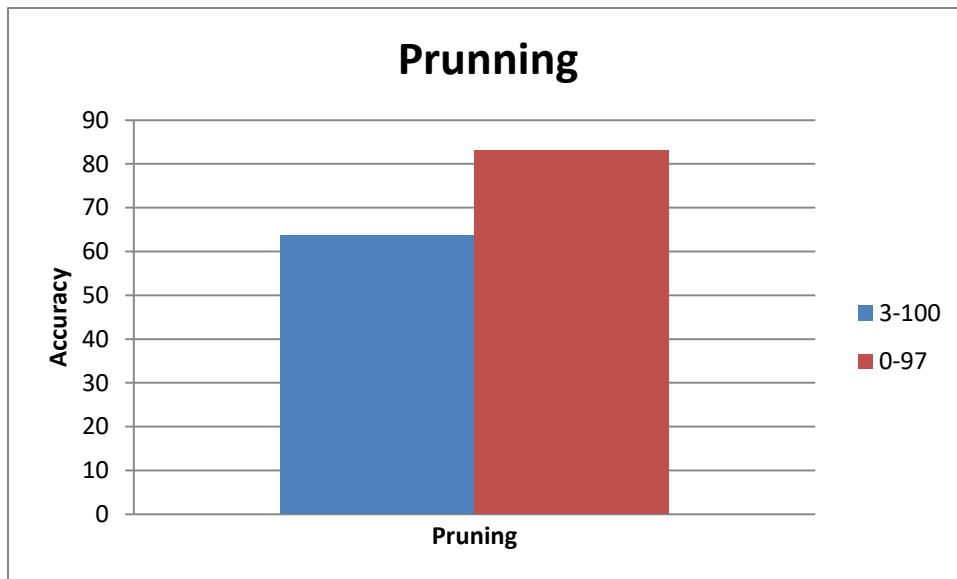


Γράφημα 2. Prunning levels

Όπως φαίνεται κι απ το γράφημα, το Prunning έχει αρνητικά αποτελέσματα στην ακρίβεια του μοντέλου. Το Binary Term Occurences διάνυσμα πετυχαίνει τη μεγαλύτερη ακρίβεια με απουσία Prunning.

Θέλω να δω όμως και σε ποια πλευρά του κλαδέματος οφείλεται η απώλεια κουρέματος, στις λέξεις που εμφανίζονται σπάνια ή στις λέξεις που εμφανίζονται συχνά;

Οπότε με δοκιμή στο Binary Term Occurences με την υψηλότερη ακρίβεια και κλάδεμα διαδοχικά 3% και 97% έχουμε:



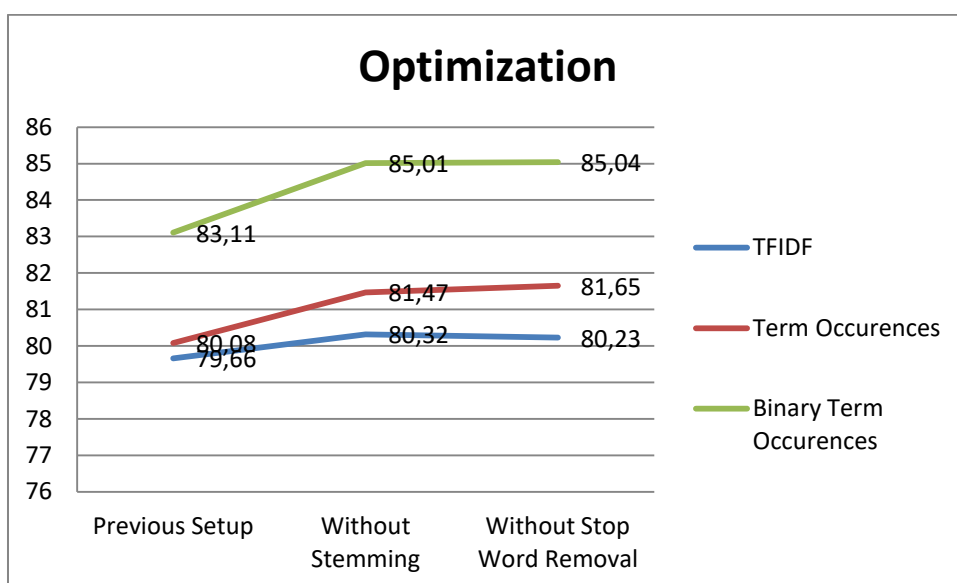
Γράφημα 3. Pruning Sides

Βλέπουμε λοιπόν ότι με κλάδεμα στις σπάνιες λέξεις στο 3% του κειμένου, η ακρίβεια πέφτει στο 63% περίπου ενώ με κλάδεμα στο 97% των συχνών λέξεων η ακρίβεια μένει σχεδόν ανεπηρέαστη.

Αυτό οφείλεται στο γεγονός κάποιες από τις κατηγορίες που έχουμε είναι κοντινές μεταξύ τους, π.χ. οι pc hardware/mac hardware, atheism/religion, baseball/hockey συνεπώς χρειαζόμαστε σπάνιες λέξεις στα έγγραφα μας για να πετύχουμε καλύτερη κατηγοριοποίηση.

2.3 Βελτιστοποίηση Κατηγοριοποίησης

Θέλω να δώ αν είναι δυνατό ένα καλύτερο ποσοστό ακρίβειας οπότε με δεδομένο το κλάδεμα στο 0/100, θα γυρίσω στις παραμέτρους της προεπεξεργασίας και θα αφαιρέσω το Stemming και το Stop Word Removal διαδοχικά.



Γράφημα 4. Βελτιστοποίηση

Η αφαίρεση του Stemming έχει θετική επιρροή στην ακρίβεια όλων των διανυσμάτων αλλά το η αφαίρεση του Stop Word Removal δεν έφερε μεγάλη βελτίωση στην ακρίβεια.

Συνεπώς το τελικό Pipeline και μοντέλο που πετυχαίνει την μεγαλύτερη ακρίβεια είναι 85,01%:

Preprocessing: Tokenize, Transform Cases, Filter Stopwords

Categorization: Naive Bayes, No Pruning, Binary Term Occurrences

2.4 Recall Precision

accuracy	true hardware	true sportBaseball	true hardwareMac	true sportHockey	true sciSpace	true miscForSale	true politics	true religion	true atheism	class precision
85,01%										
pred. hardware										
rePC	303	0	45	1	4	38	0	0	4	76.71%
pred. sportBaseball										
seball	1	349	1	6	1	6	0	0	0	95.88%
pred. hardwareMac										
reMac	58	3	301	2	1	40	1	0	0	74.14%
pred. sportHockey										
ockey	1	22	3	379	3	8	3	0	0	90.45%
pred. sciSpace										
e	11	10	23	2	366	21	13	0	5	81.15%
pred. miscForSale										
rsale	15	4	10	3	3	267	0	0	0	88.41%
pred. politics										
politics	1	7	2	5	13	6	248	0	10	84.93%
pred. religion										
religion	2	2	0	0	1	2	31	375	29	84.84%
pred. atheism										
m	0	0	0	1	2	2	14	2	271	92.81%
class recall	77.30%	87.91%	78.18%	94.99%	92.89%	68.46%	80.00%	99.47%	84.95%	

Πίνακας 1. Recall/Precision

Ο πίνακας 1 μας δείχνει τις τιμές precision-recall για κάθε κατηγορία του μοντέλου με την μεγαλύτερη ακρίβεια(85,01%).

Το μοντέλο μας έχει πολύ καλές επιδόσεις στο precision σε κάποιες κατηγορίες, baseball, hockey, misc ForSale, atheism, politics και religion ενώ στα pc hardware και mac hardware δεν τα πήγε και τόσο καλά. Εντυπωσιακό είναι το ποσοστό του recall για την κατηγορία religion με 99,47%.

Στο recall πάλι έχουμε καλές επιδόσεις στις κατηγορίες baseball, hockey, science space, religion και atheism ενώ στις κατηγορίες hardware pc, hardware mac, misc ForSale και politics είχε μέτρια αποτελέσματα.