

2^η Εργασία κατανεμημένα λειτουργικά συστήματα

Νίκος Γεωργιάδης αεμ 2043

Θοδωρης Φασουλας αεμ 2096

Δημητρης Αγτζιδης αεμ 2040

ΠΡΟΒΛΗΜΑ ΤΑΞΙΝΟΜΗΣΗΣ

Στην εργασία αυτή καλούμαστε να υλοποιήσουμε με τη χρήση του Hadoop ένα δέντρο αποφάσεων και πιο συγκεκριμένα τον αλγόριθμο Id3 για να προσεγγίσουμε ένα πρόβλημα ταξινόμησης. Το σύνολο δεδομένων που θα χρησιμοποιήσουμε περιέχει αξιολογήσεις σε αυτοκίνητα με βάση ορισμένα χαρακτηριστικά.

ΥΛΟΠΟΙΗΣΗ

Το πρόγραμμα ξεκινάει μπαινοντας σε ένα βρόγχο με 10 επαναλήψεις όπου είναι και η υλοποίηση του ten-fold cross validation. Ο αλγόριθμος θα τρέξει 10 φορές κάθε φορά με διαφορετικό αρχείο εκπαίδευσης και διαφορετικό αρχείο αξιολόγησης. Στη συνέχεια μπαίνει σε ένα άλλο loop (while) όπου ελέγχεται με ένα flag αν όλοι οι mappers σε μια map/reduce φάση δεν έχουν επικοινωνήσει με τον reducer ή αν είναι η πρώτη φορά που ξεκινάει η map/reduce φάση για μια από τις 10 επαναλήψεις του αλγόριθμου.

Αφού φτιαχτώ κατάλληλα τις εισόδους και εξόδους για κάθε map/reduce φάση, τα στέλνω σαν παραμέτρους στην Start_Job() από όπου και ξεκινάει η διαδικασία του hadoop. Αυτό που θα τρέχει παράλληλα στο hadoop θα είναι οι κομβοί-χαρακτηριστικά που βρίσκονται κάθε φορά στο ίδιο υψός του δέντρου αποφάσεων. Δηλαδή θα τρέχουν διαφορετικά μονοπάτια που στην προηγούμενη map/reduce φάση είχαν το ίδιο μονοπάτι. Το key και value που θα εξαγονται από την mapper και την reducer κάθε φορά θα είναι τύπου text αφού το key στην mapper είναι το μονοπάτι μαζί με το αρχείο εκπαίδευσης και το value είναι μόνο το αρχείο εκπαίδευσης ενώ στην reducer το value θα είναι null για λόγους ευκολίας.

Map(key=δεν μας χρειάζεται,value=μονοπάτι_αρχείο εκπαίδευσης)

Εάν το count==0 δηλαδή εάν είναι η πρώτη map/reduce φάση μιας επαναλήψης του αλγόριθμου τότε το μονοπάτι θα είναι κενό, αλλιώς το value περιέχει μόνο το αρχείο εκπαίδευσης όπου και το παίρνουμε στο samples[[]]. Εάν είναι !=0 τότε σπάζουμε το value και στο path αποθηκεύουμε το μονοπάτι ενώ στο samples[[]] το αρχείο εκπαίδευσης. Στη συνέχεια εάν το path είναι κενό σημαίνει ότι μπήκαμε για πρώτη φορά στην map και για αυτό θα πάρουμε ένα κομμάτι από το samples[[]] όπου θα το χρησιμοποιήσουμε σαν αρχείο αξιολόγησης (samples _ [[]]). Παρακάτω εάν το μονοπάτι είναι κενό, βρίσκουμε το Gain όλων των χαρακτηριστικών με την Gain_Calculation() αλλιώς βρίσκουμε το Gain από τα χαρακτηριστικά τα οποία δεν υπάρχουν ήδη μέσα στο μονοπάτι. Βρίσκουμε το max Gain

μετα και αν το max Gain είναι ==0 αυτο σημαίνει είτε οτι εχουμε χρησιμοποιησει ολα τα χαρακτηριστικα μεσα στο μονοπατι είτε οτι το αρχειο εκπαιδευσης είναι μιας κλασης.Αυτο σημαίνει οτι ειμαστε σε φυλλο αρα θα γραψουμε το μονοπατι σε ενα αρχειο paths.txt. και δεν θα χρησιμοποιησουμε καθολου τον Reducer.Εαν είναι !=0 τοτε βρισκουμε ολες τις τιμες του χαρακτηριστικου που επιλεξαμε και δημιουργουμε τοσα μονοπατια τα οποια θα τα εκχωρησουμε ως key μαζί με το αντιστοιχο αρχειο εκπαιδευσης του και ως value θα εκχωρησουμε το αρχειο εκπαιδευσης. **(context.write(key=path+αρχειο εκπαιδευσης,value=αρχειο εκπαιδευσης)).**

Reducer(key,value)

Εδω ελεγχουμε εαν το value δηλαδη το αρχειο εκπαιδευσης είναι κενο δλδ αν δεν εχω στείλει κανενα αρχειο αρα ειμαι σε φυλλο.Αν δεν είναι κενο τοτε γραφω το key που εχω ως τωρα και value θα το εχω κενο αλλιως δεν κανω τιποτα.

(context.write(key=path+αρχειο εκπαιδευσης, value=null)

Gain_Calculation()

Εδω μεσα υπολογιζω το Κερδος Πληροφοριας του χαρακτηριστικου που εξεταζω καθε φορα.

$$\text{Total_Entropy} = p_1 * (\log(p_1) / \log 2) + p_2 * (\log(p_2) / \log 2) + p_3 * (\log(p_3) / \log 2) + p_4 * (\log(p_4) / \log 2)$$

$$S(u=\text{value}) = p_1(\text{value}) * (\log(p_1(\text{value})) / \log 2) + p_2(\text{value}) * (\log(p_2(\text{value})) / \log 2) + p_3(\text{value}) * (\log(p_3(\text{value})) / \log 2) + p_4(\text{value}) * (\log(p_4(\text{value})) / \log 2)$$

$S(u=\text{value}) \dots$

$S(u=\text{value})..$

$$\text{Gain} = \text{Total_Entropy} - (\text{new_c_all} / \text{c_all}) * S(u=\text{value}) + \dots$$

Το new_c_all είναι το ποσα δειγματα εχουν την συγκεκριμενη τιμη απο το χαρακτηριστικο που εξεταζουμε και το c_all είναι το ποσα δειγματα υπαρχουν συνολικα ανεξαρτητα απο την τιμη του χαρακτηριστικου.

ΠΙΣΩ ΣΤΗΝ MAIN() Evaluation Implementation

Αφου εχει τελειωσει η διαδικασια της εκπαιδευσης του ταξινομητη με μια επαναληψη του αλγοριθμου με την βοηθεια hadoop, το επομενο κομματι είναι η αξιολογηση του ταξινομητη. Διαβαζουμε το αρχειο με τα μονοπατια που δημιουργηθηκε και το αποθηκευω στο results[[]].Στην συνεχεια για να το εχω σε καλυτερη μορφη και να γινει ευκολοτερη η επεξεργασια το αποθηκευω στο new_results .Για καθε κλαση καλω την Evaluation()

Evaluation()

Στη συναρτηση αυτη με το αρχειο αξιολογησης που εχω (samples _ [[]]) εξεταζω για καθε δειγμα του σε ποιο μονοπατι αντιστοιχει.Βρισκω αυτο το μονοπατι και περνω την κλαση του.Εαν η κλαση είναι ιδια με αυτη του δειγματος τοτε True positive++ αλλιως False positive++ . Εαν εξεταζουμε αλλη κλαση και προβλεψουμε την παραπανω κλαση τοτε False positive ++

Οταν τελειωσω αυτη την διαδικασια για ολες τις κλασεις τοτε για καθε κλαση θα βρω το precision , recall , fmeasure και στην συνεχεια τον μεσο ορο γιαυτες τις μετρικες απο ολες της κλασεις.(average_precision,average_recall,average_fmeasure)

Οταν εχω πραγματοποιησει το ten fold cross validation δηλαδη εχω τρεξει τον αλγοριθμο 10 φορες, θα υπολογισω τον μεσο ορο απο ολους τους μεσους ορους που εχω βρει για καθε μετρικη.(final_precision,final_recall,final_fmeasure) .Τελος θα τα γραψω σε ενα αρχαιο με ονομα "results.txt".

Τα αποτελεσματα απο το αρχαιο ειναι : Final_precision 81% Final_recall 68%
Final_fmeasure 65%

Τελος προγραμματος

ΣΥΝΔΕΣΗ

ΜΕ ΤΗΝ ΒΟΗΘΕΙΑ 2 VIRTUAL BOX ΤΟ ΕΝΑ ΛΕΙΤΟΥΡΓΩΝΤΑΣ ΩΣ MASTER ΚΑΙ SLAVE ΚΑΙ ΤΟ ΑΛΛΟ ΩΣ SLAVE ΚΑΝΟΝΤΑΣ ΤΟ ΚΑΤΑΛΛΗΛΟ SETUP ΚΑΙ ΜΕ ΤΗΝ ΒΟΗΘΕΙΑ ΤΟΥ HADOOP ΠΡΑΓΜΑΤΟΠΟΙΗΘΗΚΑΝ ΟΙ MAP/REDUCE ΦΑΣΕΙΣ.