

Ανάπτυξη Λογισμικού για Πληροφοριακά Συστήματα Εργασία:2η

Μέλη της ομάδας:

Βουρτζούμη Ουρανία A.M:115201600024

Κοσμάς Αλέξανδρος A.M:1115201700299

Βαρώνος Διονύσης A.M:1115201600017

UniTesting:

Για το uni testing χρησιμοποιούμε την βιβλιοθήκη acutest.h

Με την εντολή `make test` δημιουργείται το εκτελέσιμο `test` το οποίο εκτελείται με την εντολή `make test-run`. Με την εντολή `make test-clean` διαγράφονται διαγράφεται το εκτελέσιμο `test` και όσα .o αρχεία δημιουργήθηκαν. Το εκτελέσιμο δημιουργείται μέσα στον φάκελο `test` αλλά όλες οι παραπάνω εντολές δίνονται στον αρχικό φάκελο.

Μεταγλώττιση κώδικα:

Η υλοποίησή μας για το δημιουργείται να εκτελέσιμο `sigmod` το οποίο δημιουργείται με την εντολή `make`. Με την εντολή `make clean` διαγράφεται το εκτελέσιμο όλα τα .o αρχεία

και όλα τα .csv αρχεία τα οποία έχουν δημιουργηθεί κατά την εκτέλεση του προγράμματος.

Εντολή εκτέλεσης προγράμματος :

Το εκτελέσιμο εκτελείται με την εντολή :

```
make run ARGS="- d Directory -w datacsv"
```

ή απλά :

```
./sigmod - d Directory -w datacsv
```

Όπου:

Directory:

ο φάκελος που έχει τους φακέλους με τα .json αρχεία δηλαδή ο φάκελος 2013_camera_specs στην συγκεκριμένη άσκηση (ο φάκελος δεν υπάρχει στον φάκελό μας στο github αλλά θα πρέπει να υπάρχει στο directory που τρέχει το εκτελέσιμο sigmod)

ΣΗΜΑΝΤΙΚΟ : ως directory θα πρέπει να δοθεί το όνομα του φακέλου και όχι ο φάκελος σαν path . Δηλαδή το σωστό είναι 2013_camera_specs και όχι 2013_camera_specs/

datacsv:

το όνομα του csv αρχείου που περιέχει τα ταιριασμένα των κάμερων(πχ sigmod_medium_labelled_dataset.csv ή sigmod_large_labelled_dataset.csv)

Επίσης στην εντολή εκτέλεσης μπορεί να προστεθεί η προαιρετική σημαία -o.

Εάν υπάρχει η συγκεκριμένη σημαία το πρόγραμμα δοκιμάζει το μοντέλο για όλα τα στοιχεία του datasetX που δεν έχουν δοκιμαστεί μεταξύ τους.

Εάν δεν υπάρχει η σημαία το πρόγραμμα τερματίζει αφού εκτυπώσει τα αποτελέσματα του testing.

Ο λόγος που επιλέξαμε αυτή την υλοποίηση είναι επειδή ο έλεγχος όλων των στοιχείων απαιτεί πολύ χρόνο (τουλάχιστον 30 λεπτά) και θεωρήσαμε σημαντικό να μπορεί ο χρήστης να τρέξει το πρόγραμμα χωρίς αυτόν τον έλεγχο ώστε να διαπιστώσει ότι όλες οι υπόλοιπες λειτουργίες του προγράμματος εκτελούνται ομαλά.

Output προγράμματος:

Το πρόγραμμα κατά την εκτέλεση του δημιουργεί ένα νέο αρχείο με όνομα Same.csv που περιέχει όλες τις θετικές συσχετίσεις που δημιουργήθηκαν από τις κλικες μέσω του 60% των δεδομένων του dataset W και ένα αρχείο με όνομα Different.csv με τις αντίστοιχες αρνητικές συσχετίσεις.

Στην συνέχεια εκτυπώνει το Success rate του μοντέλου κατά την διαδικασία του testing το οποίο γίνεται με το επόμενο 20% του datasetw.

Τέλος εκτυπώνει τις συσχετίσεις που έχουν όριο πιθανότητας 0.001(δηλαδή τις αρνητικές με πιθανότητα τριασματος μικρότερο του 0.001 και τις θετικές με πιθανότητα τριασματος μεγαλύτερη του 0.999) οι προέρχονται από τα

αποτελέσματα που μας δίνει το μοντέλο όταν το εκτελούμε για κάθε στοιχείο (δηλαδή κάθε μερά) του datasetX με όλα τα υπόλοιπα. (εάν ο χρήστης δώσει την σημαία -o στην εντολή εκτέλεσης)

Κατά την διάρκεια του προγράμματος εκτυπώνονται κάποιες προτάσεις που μας ενημερώνουν πως το πρόγραμμα μόλις τελείωσε κάποια συγκεκριμένη διαδικασία (πχ όταν διαβάσει όλα τα δεδομένα του datasetX)

Ροή προγράμματος :

Το πρόγραμμα δημιουργεί δημιουργεί μια δομή Hash στην οποία αποθηκεύονται τα στοιχεία του κάθε json αρχείου από το datasetX και μια δομή LHash η οποία αντιπροσωπεύει το vocabulary όλων των json.

Για κάθε αρχείο ο json αποθηκεύεται το id της κάθε μεριάς στην δομή Hash και στην συνέχεια διαβάζουμε το json λέξη-λέξη ό που κάθε λέξη αποθηκεύεται στο λεξιλόγιο και σε μια δομή WHash που αποθηκεύεται στην αντίστοιχη θέση που αποθηκευτηκε το id του json μέσα στην δομή Hash.

Για κάθε νέα λέξη που βρίσκει το λεξιλόγιο την αποθηκεύει και για κάθε λέξη που έχει ήδη αυξάνει κατά ένα την μεταβλητή wordperj που αντιστοιχεί στην συγκεκριμένη λέξη και αντιπροσωπεύει τον αριθμό των json στα οποία βρέθηκε αυτή η λέξη.

Αντίστοιχα η δομή WHash αποθηκεύει κάθε νέα λέξη του json και για κάθε ήδη υπάρχουσα αυξάνει κατά ένα την μεταβλητή τάδε που αντιστοιχεί στο πόσες φορές βρέθηκε η λέξη αυτή σε αυτό το json.

Μόλις διαβαστούν όλες οι λέξεις του json το WHash υπολογίζει το tf της κάθε λέξης του και στην συνέχεια για κάθε μία από τις λέξεις του δίνει το tf στο λεξιλόγιο το οποίο ο το προσθέτει στην μεταβλητή tfcount της αντίστοιχης λέξης και με αυτόν τον τρόπο κρατάει το άθροισμα των tf της κάθε λέξης για όλα τα json.

Αφού τελειώσει η παραπάνω διαδικασία για όλα τα αρχεία του datasetX η δομή LHash υπολογίζει το μέσο tf-idf της κάθε λέξης και το αποθηκεύει στην μεταβλητή isf και αποθηκεύει το idf της λέξης στην μεταβλητή tfcount.

Στην συνέχεια το λεξιλόγιο ταξινομεί τις λέξεις με βάση το μεγαλύτερο μέσο tf-idf και κρατάει τις 1000 πρώτες.

Στην συνέχεια για κάθε στοιχείο της δομής Hash μια δομή Hvector η οποία είναι ένας πίνακας κατακερματισμού που αντιπροσωπεύει ένα sparse array.

Δηλαδή για κάθε μια λέξη από τις 1000 σημαντικές του λεξιλογίου εάν υπάρχει στην κάθε μερα αποθηκεύεται η θέση της και η τιμή $tf*idf$ της αντίστοιχης λέξης για την συγκεκριμένη μερα.

Με αυτόν τον τρόπο δεν αποθηκεύουμε στο vector μας τις τιμές που ισούνται με 0 και είναι περιττή πληροφορία.

Μολις δημιουργείται το κάθε vector διαγράφεται η δομή WHash αυτής της θέσης.

Στην συνέχεια δημιουργούνται οι κλικες βάσει το 60% των θετικών και το 60% των αρνητικών συσχετίσεων του dataW και οι αρνητικές συσχετίσεις μεταξύ των κλικών που δεν τερματίζουν.

Επίσης δημιουργούνται τα αρχεία Testing.csv και Validation.csv που το κάθε ένα περιέχει 20% αρνητικών και θετικών συσχετίσεων και θα χρειαστούν στην συνέχεια για να πάρουμε τα δεδομένα στις αντίστοιχες διαδικασίες.

Αφού δημιουργηθούν οι κλικες και οι αρνητικές συσχετίσεις το πρόγραμμα δημιουργεί τα αρχεία Same.csv και Different.csv που σε αυτά αποθηκεύονται όλες οι θετικές

συσχετισείς βάρη των κλικών της δομής και όλες οι αρνητικές συσχετισείς αντίστοιχα.

Στην συνέχεια μέσω της συνάρτησης Training δημιουργού με μια δομή model η οποία αντιπροσωπεύει το μοντέλο μας. Το μοντέλο μας εκπαιδεύεται με τα δεδομένα των αρχείων Same.csv και Different.csv και πιο συγκεκριμένα για τις διπλάσιες αρνητικές συσχετισείς από τις θετικές.

Η διαδικασία της εκπαίδευσης είναι ως εξής :

Για κάθε ζευγάρι των αρχείων Same.csv και Different.csv υπολογίζουμε το concatenation των vectors των αντίστοιχων κωμερών και την τιμή της πρόβλεψης του μοντέλου μείον την σωστή τιμή της συσχέτισης.

Στην συνέχεια για κάθε θέση του concatenation υπολογίζουμε το γινόμενο της τιμής της συγκεκριμένης θέσης επί της παραπάνω τιμής.

Αυτό το γινόμενο αφαιρούμε κάθε φορά από το βάρος της αντίστοιχης θέσης

Αυτή η διαδικασία επαναλαμβάνεται 3 φορές και μόλις τελειώσει η συνάρτηση επιστρέφει το μοντέλο

Στην συνέχεια μέσω της συνάρτησης Testing για κάθε ζευγάρι του αρχείου Testing.csv υπολογίζουμε την πρόβλεψη

του εκπαιδευμένου μοντέλου και αφού αυτό γίνεται για όλα τα ζευγάρια το πρόγραμμα εκτυπώνει το success rate του μοντέλου δηλαδή το ποσοστό των σωστων προβλέψεων

Στην συνέχεια εάν ο χρήστης δώσει την σημαία -o στην εντολή εκτέλεσης το πρόγραμμα εκτυπώνει τις συσχετίσεις που έχουν όριο πιθανότητας 0.001(δηλαδή τις αρνητικές με πιθανότητα τεριασματος μικρότερο του 0.001 και τις θετικές με πιθανότητα τεριασματος μεγαλύτερη του 0.999)οι οποίες προέρχονται από τα αποτελέσματα που μας δίνει το μοντέλο όταν το εκτελούμε για κάθε στοιχείο (δηλαδή κάθε μερά) του datasetX με όλα τα υπόλοιπα.

Τέλος αποδεδειγνται όλες οι δομές και το πρόγραμμα τερματίζει.

Δομές :

Η δομή Hash είναι η δομή που αποθηκευονται τα δεδομένα του κάθε json αρχείου.

Είναι ένας δυναμικός πίνακας κατακερματισμού με bucket-list ό που ως κλειδί χρησιμοποιεί το id κάθε json (πχ www.ebay.com//567)και κάνει rehash κάθε φορά που φτάνει 80% πλήροτητα.

Για την υλοποίηση του bucket-list χρησιμοποιείται η δομή NList που σε κάθε κόμβο της αποθηκεύεται:

το id του json στην μεταβλητή camera τύ που char*

οι λέξεις του json αρχείου στην μεταβλητή spear που είναι τύ που WHash*

το αντίστοιχο vector του json που είναι τύ που HVector*

ένας δείκτης σε δομή CList που αντιστοιχεί στην κλικα την οποία ανήκει η camera.

Η δομή CList είναι μια συνδεδεμένη λίστα που για την υλοποίηση των κλικων.

Σε κάθε θέση της αποθηκεύει:

το όνομα της κάμερας

έναν δείκτη σε NList που αντιστοιχεί με τον κόμβο NList που είναι αποθηκευμένα τα στοιχεία της κάμερας.

Ο πρώτος κόμβος της κάθε CList δεν αποθηκεύει δεδομένα μιας κάμερας αλλά μια λίστα TList(συνδεδεμένη λίστα που αποθηκεύει δείκτες σε CList) που αποθηκεύει τις κλικες με τις οποίες δεν τερματίζει η κλικα.

Η δομή WHash είναι μια δομή Πίνακα κατακερματισμού χωρίς bucket-list που αποφεύγει τα collision πηγαίνοντας στην επόμενη διαθέσιμη κενή θέση και κάνει rehash όταν έχει 80% πληρότητα.

Η δομή αυτή χρησιμοποιείται για να αποθηκεύει τις λέξεις του κάθε json αρχείου με κλειδί την κάθε λέξη.

Σε κάθε bucket αποθηκεύει μια λέξη και το tf αυτής της λέξης για το συγκεκριμένο json.

Η δομή Hvector είναι μια δομή πίνακα κατακερματισμού χωρίς bucket-list που αποφεύγει τα collision πηγαίνοντας στην επόμενη διαθέσιμη κενή θέση και κάνει rehash όταν έχει 80% πληρότητα.

Η δομή αυτή απαριστά το vector της κάθε μέρας

Σε κάθε bucket αποθηκεύει την θέση και την tf*idf της αντίστοιχης θέσης για κάθε λέξη από τις 1000 πιο σημαντικές του λεξιλογίου ένα υπάρχει στην κάθε μέρα.

Η δομή LHash είναι μια δομή Πίνακα κατακερματισμού χωρίς bucket-list που αποφεύγει τα collision πηγαίνοντας στην επόμενη διαθέσιμη κενή θέση και κάνει rehash όταν έχει 80% πληρότητα.

Η δομή αυτή χρησιμοποιείται για την υλοποίηση του λεξιλογίου όλων των json με κλειδί την κάθε λέξη.

Σε κάθε bucket αποθηκεύεται η κάθε λέξη το idf και το μέσω tf-idf της κάθε λέξης

Η δήλωση όλων των δομών λιστών της δομής Hvector και της δομής WHash και τα πρότυπα των συναρτήσεων για την διαχείριση τους βρίσκονται στο αρχείο list.h και οι υλοποιήσεις των συναρτήσεων στο αρχείο list.c

Η δήλωση όλων των δομών Hash και WHash και τα πρότυπα των συναρτήσεων για την διαχείριση τους βρίσκονται στο

αρχεί ο hash.h και οι υλοποιή σεις των συναρτή σεων στο
αρχεί ο hash.c

Η δή λωση της δομή ς Model και τα πρό τυπα των συναρτή σεων
για την διαχεί ριση της βρί σκονται στο αρχεί ο logistic.h και οι
υλοποιή σεις των συναρτή σεων στο αρχεί ο logistic.c