

ΔΙΠΛΩΜΑΤΙΚΗ ΕΡΓΑΣΙΑ

Σπουδαστές:

Νάτσης Αινείας
Τρούλης Μιχάλης

Θέμα:

“Μελέτη και ανάπτυξη διαδικτυακής εφαρμογής, με βάση κειμενική πληροφορία και στόχο την εξαγωγή συναισθήματος από αυτό.”

Εισηγητής:

Παπαδάκης Νικόλαος

Περίληψη

Η παρούσα πτυχιακή εργασία πραγματεύεται το αντικείμενο της συναισθηματικής ανάλυσης σε ηλεκτρονικής μορφής έγγραφα, χρησιμοποιώντας κατηγοριοποιητές κειμένων είτε βασισμένους σε λεξικά αλλά και Naïve Bayes οι οποίοι βασίζονται στην εκπαίδευση απο άλλα κείμενα. Για να γίνει επαλήθευση των αποτελεσμάτων του αλγορίθμου που εξήγαγε τα καλύτερα αποτελέσματα θεωρήθηκε σκόπιμη η δημιουργία μιας διαδικτυακής εφαρμογής στην οποία οι χρήστες συνδέονται και αξιολογούν με την σειρά τους τα ίδια κείμενα. Η εφαρμογή συγκρίνει τα αποτελέσματα του χρήστη με του αλγορίθμου και βγάζει στο ποσοστό επιτυχίας του.

Abstract

The following thesis deals with the object of sentiment analysis of electronic textual information, with the use of text categorisation based on dictionaries and Naive Bayes algorithms, which are trained with other texts. To evaluate the best algorithm, was considered the creation of a web application in which the users sign in and evaluate the same textual information. The application compares the user's and the algorithm's results and results the algorithm's percentage of success.

Table of Contents

Περίληψη	3
Abstract	4
Κατάλογος εικόνων	8
Κατάλογος Πινάκων	9
1 ΚΕΦΑΛΑΙΟ: Εισαγωγή	10
1.1 Πρόλογος	10
1.2 Δομή	11
2 ΚΕΦΑΛΑΙΟ: Εξόρυξη Γνώσης- Γνώμης (opinion mining)	13
2.1 Άνάλυση ορισμών και βασικής ορολογίας	15
2.2 Μέθοδοι που χρησιμοποιήθηκαν	17
2.2.1 Μέθοδος χρήσης λεξικών	17
2.2.2 Πλεονεκτήματα – Μειονεκτήματα χρήσης λεξικών	19

2.2.3	Μέθοδος Μηχανικής Μάθησης	20
2.2.4	Είδη κατηγοριοποιήσεων κειμένου – Μέτρηση ακρίβειας	21
3	ΚΕΦΑΛΑΙΟ : Αποτελέσματα από την Πειραματική Διαδικασία	24
3.1	Ο αλγόριθμος “Sentimental” βασισμένος σε λεξικό	25
3.1.1	Τα ενσωματωμένα συναισθηματικά λεξικά	26
3.1.2	Αποτελέσματα πειραμάτων	28
3.2	Ο αλγόριθμος “Sad Panda” βασισμένος σε λεξικό	32
3.2.1	Τα ενσωματωμένα συναισθηματικά λεξικά	33
3.2.2	Αποτελέσματα πειραμάτων	35
3.3	Κατηγοριοποιητής “Classifier”	37
3.3.1	Data set - Συλλογές δεδομένων	39
3.3.2	Πειράματα - Αποτελέσματα	40
3.4	Προβλήματα - Συμπεράσματα	52
4	ΚΕΦΑΛΑΙΟ : Γενική Αρχιτεκτονική	54
4.1	Έρευνα Τεχνολογιών	54
4.1.1	Ruby [1]	54
4.1.2	Ruby On Rails [2]	55
4.1.3	MySQL	59
4.1.4	Git	61
4.2	Επιλεγμένες Τεχνολογίες	63
4.2.1	Βάση δεδομένων – PostgreSQL	63
4.2.2	Twitter Bootstrap	64
4.2.3	Highcharts	65
4.2.4	RVM [18]	66
4.2.5	Bundler	66
4.2.6	Λόγοι χρήσης επιμέρους τεχνολογιών	67
5	ΚΕΦΑΛΑΙΟ : Η Εφαρμογή αξιολόγησης Analymnt	70
5.1	Η δημιουργία του σκελετού της εφαρμογής	70
5.2	Τα αρχεία που απαρτίζεται η εφαρμογή	73
5.3	Πώς σχεδιάστηκε η εφαρμογή	76
5.4	Λειτουργίες χρηστών	77
5.4.1	Δημιουργία – σύνδεση χρήστη	78
5.4.2	Λειτουργίες διαχειριστών - απλών χρηστών	83
5.5	Προβλήματα - Συμπεράσματα και μελλοντική εξέλιξη της εφαρμογής	92
	Παράρτημα	95
	Κατάλογος κώδικα – Κώδικας	95
	Αναφορές	118
	Βιβλιογραφία	119

Κατάλογος εικόνων

ΕΙΚΟΝΑ 2.2.1-1: ΚΕΙΜΕΝΙΚΗ ΕΞΟΡΥΞΗ ΓΝΩΣΗΣ	14
ΕΙΚΟΝΑ 2.2.1-1: ΥΠΟΚΕΙΜΕΝΙΚΗ ΚΑΙ ΑΝΤΙΚΕΙΜΕΝΙΚΗ ΤΑΞΙΝΟΜΙΣΗ	17
ΕΙΚΟΝΑ 2.2.1-1: Η ΚΑΤΗΓΟΡΙΟΠΟΙΗΣΗ ΛΕΞΙΚΩΝ.....	19
ΕΙΚΟΝΑ 2.2.4-1: ΟΡΙΣΜΟΣ ΤΩΝ PRECISION ΚΑΙ RECALL.....	23
ΕΙΚΟΝΑ 3.1.1-1: ΤΟ ΠΡΟΣΗΜΑΣΜΕΝΟ ΛΕΞΙΚΟ ΤΟΥ ΑΛΓΟΡΙΘΜΟΥ SENTIMENTAL	27
ΕΙΚΟΝΑ 3.1.1-2: ΤΑ ΕΜΟΤΙΧΩΝ ΤΟΥ ΑΛΓΟΡΙΘΜΟΥ SENTIMENTAL	28
ΕΙΚΟΝΑ 3.1.2-1: ΠΡΩΤΟ ΠΕΙΡΑΜΑ ΒΙΒΛΙΟΘΗΚΗΣ ΜΕ ΤΟ ΑΡΝΗΤΙΚΟ ΑΡΧΕΙΟ.....	31
ΕΙΚΟΝΑ 3.2.1-1: ΑΠΟΔΟΣΗ ΠΟΛΙΚΟΤΗΤΑΣ ΣΤΙΣ ΛΕΞΕΙΣ ΑΠΟ 1-10.....	34
ΕΙΚΟΝΑ 3.2.1-2: ΛΕΞΙΚΟ ΥΠΟΚΕΙΜΕΝΙΚΟΤΗΤΑΣ ΤΟΥ ΑΛΓΟΡΙΘΜΟΥ	34
ΕΙΚΟΝΑ 3.2.1-3: ΛΕΞΙΚΟ ΣΥΝΑΙΣΘΗΜΑΤΩΝ ΤΟΥ ΑΛΓΟΡΙΘΜΟΥ.....	35
ΕΙΚΟΝΑ 3.2.2-4: ΠΕΙΡΑΜΑ ΤΟΥ ΑΛΓΟΡΙΘΜΟΥ ΟΠΩΣ ΑΠΟΤΥΠΩΝΕΤΑΙ ΣΤΗΝ ΟΘΟΝΗ.....	36
ΕΙΚΟΝΑ 4.1.1-1: ΛΟΓΟΤΥΠΟ ΤΗΣ RUBY.....	54
ΕΙΚΟΝΑ 4.1.2-1: ΛΟΓΟΤΥΠΟ ROR.....	55
ΕΙΚΟΝΑ 4.1.2-2: ΑΝΑΠΑΡΑΣΤΑΣΗ MVC ΑΡΧΙΤΕΚΤΟΝΙΚΗΣ	56
ΕΙΚΟΝΑ 4.1.2-3: ΔΟΜΗ ΦΑΚΕΛΩΝ ΤΩΝ ΕΦΑΡΜΟΓΩΝ RUBY ON RAILS.....	59
ΕΙΚΟΝΑ 4.1.3-1: ΛΟΓΟΤΥΠΟ GIT.....	61
ΕΙΚΟΝΑ 4.1.4-3: ΕΝΗΜΕΡΩΣΗ ΤΟΠΙΚΟΥ ΚΑΙ ΑΠΟΜΑΚΡΥΣΜΕΝΟΥ REPOSITORY	62
ΕΙΚΟΝΑ 4.1.4-4: ΔΙΑΚΛΑΔΩΣΕΙΣ ΣΤΟ GIT	63
ΕΙΚΟΝΑ 4.2.1-1: ΛΟΓΟΤΥΠΟ ΤΗΣ POSTGRESQL	63
ΕΙΚΟΝΑ 4.2.4-1: ΛΟΓΟΤΥΠΟ BOOTSTRAP	64
ΕΙΚΟΝΑ 4.2.6-1: ΛΟΓΟΤΥΠΟ HIGHCHARTS.....	65
ΕΙΚΟΝΑ 4.2.7-1: : ΛΟΓΟΤΥΠΟ ΤΟΥ RVM	66
ΕΙΚΟΝΑ 4.2.9-1: ΤΟ ΑΡΧΙΚΟ SCHEMA ΒΑΣΗΣ ΔΕΔΟΜΕΝΩΝ	71
ΕΙΚΟΝΑ 4.2.9-2	71
ΕΙΚΟΝΑ 4.2.9-3: SCHEMA ΒΑΣΗΣ ΔΕΔΟΜΕΝΩΝ ΠΡΟΓΡΑΜΜΑΤΙΣΤΙΚΑ.....	72
ΕΙΚΟΝΑ 4.2.9-1: ΔΟΜΗ ΦΑΚΕΛΩΝ ΤΗΣ ΕΦΑΡΜΟΓΗΣ ANALYMENT	75
ΕΙΚΟΝΑ 4.2.9-1: ΔΟΜΗ ΚΑΙ ΣΥΝΔΕΣΗ ΤΩΝ CONTROLLERS.....	76
ΕΙΚΟΝΑ 4.2.9-1: ΑΡΧΙΚΗ ΣΕΛΙΔΑ ΓΙΑ ΟΛΟΥΣ ΤΟΥΣ ΧΡΗΣΤΕΣ.....	77
ΕΙΚΟΝΑ 5.4.1-1: ΦΟΡΜΑ ΣΥΝΔΕΣΗΣ ΧΡΗΣΤΗ.....	79
ΕΙΚΟΝΑ 5.4.1-2: ΦΟΡΜΑ ΔΗΜΙΟΥΡΓΙΑΣ ΧΡΗΣΤΗ	81
ΕΙΚΟΝΑ 5.4.1-3: ΛΕΙΤΟΥΡΓΙΑ ΑΝΑΚΤΗΣΗΣ ΚΩΔΙΚΟΥ	83
ΕΙΚΟΝΑ 5.4.2-1: ΑΡΧΙΚΗ ΣΕΛΙΔΑ ADMIN	84
ΕΙΚΟΝΑ 5.4.2-2: ΌΛΕΣ ΟΙ ΚΑΤΗΓΟΡΙΕΣ ΤΩΝ ΠΕΙΡΑΜΑΤΩΝ.....	85
ΕΙΚΟΝΑ 5.4.2-4: ΣΕΛΙΔΑ ΔΙΑΧΕΙΡΗΣΗΣ ΤΩΝ ΚΕΙΜΕΝΩΝ	86
ΕΙΚΟΝΑ 5.4.2-5: ΣΥΓΚΡΙΣΗ ΑΠΟΤΕΛΕΣΜΑΤΩΝ ΑΝΑ ΚΑΤΗΓΟΡΙΑ	87
ΕΙΚΟΝΑ 5.4.2-6: ΣΕΛΙΔΑ ΣΥΓΚΡΙΣΗΣ ΑΠΟΤΕΛΕΣΜΑΤΩΝ ΑΝΑ TEST OBJECT	88
ΕΙΚΟΝΑ 5.4.2-7: ΣΕΛΙΔΑ ΧΡΗΣΤΩΝ	88
ΕΙΚΟΝΑ 5.4.2-8: ΕΜΦΑΝΙΣΗ ΟΛΩΝ ΤΩΝ ΚΕΙΜΕΝΩΝ	89
ΕΙΚΟΝΑ 5.4.2-9: ΠΕΙΡΑΜΑΤΙΚΗ ΔΙΑΔΙΚΑΣΙΑ ΓΙΑ ΤΟΝ ΔΙΑΧΕΙΡΙΣΤΗ	90
ΕΙΚΟΝΑ 5.4.2-11: ΓΡΑΦΗΜΑ ΜΕ ΤΗΝ ΟΛΟΚΛΗΡΩΣΗ ΤΗΣ ΠΕΙΡΑΜΑΤΙΚΗΣ ΔΙΑΔΙΚΑΣΙΑΣ.....	91
ΕΙΚΟΝΑ 5.4.2-12: ΑΝΑΘΕΣΗ ΒΑΘΜΟΛΟΓΙΑΣ ΣΤΑ ΚΕΙΜΕΝΑ.....	92

Κατάλογος Πινάκων

ΠΙΝΑΚΑΣ 3.1.2-1: ΑΠΟΤΕΛΕΣΜΑΤΑ ΤΟΥ ΑΛΓΟΡΙΘΜΟΥ SENTIMENTAL	32
ΠΙΝΑΚΑΣ 3.2.2-1: ΤΑ ΑΠΟΤΕΛΕΣΜΑΤΑ ΤΟΥ ΑΛΓΟΡΙΘΜΟΥ	36
ΠΙΝΑΚΑΣ 3.3.2-1: ΠΡΩΤΟ ΠΕΙΡΑΜΑ.....	41
ΠΙΝΑΚΑΣ 3.3.2-2: ΔΕΥΤΕΡΟ ΠΕΙΡΑΜΑ ΜΕ ΙΔΙΟ TRAINING SET.....	41
ΠΙΝΑΚΑΣ 3.3.2-3: ΤΡΙΤΟ ΠΕΙΡΑΜΑ	41
ΠΙΝΑΚΑΣ 3.3.2-4: ΠΡΩΤΟ ΠΕΙΡΑΜΑ ΜΕ ΘΕΤΙΚΑ ΣΥΝΟΛΑ ΔΟΚΙΜΗΣ	42
ΠΙΝΑΚΑΣ 3.3.2-5: ΔΕΥΤΕΡΟ ΠΕΙΡΑΜΑ ΜΕ ΘΕΤΙΚΑ ΣΥΝΟΛΑ ΔΟΚΙΜΗΣ	42
ΠΙΝΑΚΑΣ 3.3.2-6: ΠΕΙΡΑΜΑ ΜΕ ΑΡΝΗΤΙΚΟ ΣΥΝΟΛΟ	43
ΠΙΝΑΚΑΣ 3.3.2-7: ΠΕΙΡΑΜΑ ΚΕΙΜΕΝΩΝ CORNELL ΩΣ TRAINING SET	44
ΠΙΝΑΚΑΣ 3.3.2-8: ΠΡΩΤΟ ΠΕΙΡΑΜΑ ΘΕΤΙΚΩΝ ΣΥΝΟΛΩΝ ΔΟΚΙΜΩΝ	45
ΠΙΝΑΚΑΣ 3.3.2-9: ΔΕΥΤΕΡΟ ΠΕΙΡΑΜΑ ΑΡΝΗΤΙΚΩΝ ΣΥΝΟΛΩΝ ΔΟΚΙΜΩΝ.....	45
ΠΙΝΑΚΑΣ 3.3.2-10: ΠΕΙΡΑΜΑ ΘΕΤΙΚΟΥ ΣΥΝΟΛΟΥ ΜΕ ΤΟ ΝΕΟ ΣΥΝΟΛΟ ΕΚΠΑΙΔΕΥΣΗΣ.....	46
ΠΙΝΑΚΑΣ 3.3.2-11: ΔΟΚΙΜΗ ΑΡΝΗΤΙΚΟΥ ΣΥΝΟΛΟΥ ΜΕ ΝΕΟ ΣΥΝΟΛΟ ΕΚΠΑΙΔΕΥΣΗΣ	46
ΠΙΝΑΚΑΣ 3.3.2-12: ΠΡΩΤΟ ΠΕΙΡΑΜΑ ΜΕ ΜΕΤΡΗΣΕΙΣ ΩΣ ΠΡΟΣ ΤΗΝ ΑΚΡΙΒΕΙΑ ΚΑΙ ΤΗΝ ΑΝΑΚΛΗΣΗ.....	48
ΠΙΝΑΚΑΣ 3.3.2-13: ΔΕΥΤΕΡΟ ΠΕΙΡΑΜΑ ΜΕ ΜΕΤΡΗΣΕΙΣ ΩΣ ΠΡΟΣ ΤΗΝ ΑΚΡΙΒΕΙΑ ΚΑΙ ΤΗΝ ΑΝΑΚΛΗΣΗ	49
ΠΙΝΑΚΑΣ 3.3.2-14: ΤΡΙΤΟ ΠΕΙΡΑΜΑ ΜΕ ΜΕΤΡΗΣΕΙΣ ΩΣ ΠΡΟΣ ΤΗΝ ΑΚΡΙΒΕΙΑ ΚΑΙ ΤΗΝ ΑΝΑΚΛΗΣΗ	49
ΠΙΝΑΚΑΣ 3.3.2-15: ΤΕΤΑΡΤΟ ΠΕΙΡΑΜΑ ΜΕ ΜΕΤΡΗΣΕΙΣ ΩΣ ΠΡΟΣ ΤΗΝ ΑΚΡΙΒΕΙΑ ΚΑΙ ΤΗΝ ΑΝΑΚΛΗΣΗ	50
ΠΙΝΑΚΑΣ 3.3.2-16: ΠΕΜΠΤΟ ΠΕΙΡΑΜΑ ΜΕ ΜΕΤΡΗΣΕΙΣ ΩΣ ΠΡΟΣ ΤΗΝ ΑΚΡΙΒΕΙΑ ΚΑΙ ΤΗΝ ΑΝΑΚΛΗΣΗ.....	51
ΠΙΝΑΚΑΣ 4.1.2-1: ΑΝΤΙΣΤΟΙΧΙΣΗ HTTP ΜΕΘΟΔΩΝ, ΔΙΑΔΡΟΜΩΝ (PATHS), ΛΕΙΤΟΥΡΓΙΩΝ ΒΑΣΗΣ	56
ΠΙΝΑΚΑΣ 4.2.1-1: ΟΡΙΑΚΕΣ ΤΙΜΕΣ ΕΝΕΡΓΩΝ ΕΓΚΑΤΑΣΤΑΣΕΩΝ POSTGRESQL.....	64

1 ΚΕΦΑΛΑΙΟ: Εισαγωγή

1.1 Πρόλογος

Η έμφυτη ανάγκη του ανθρώπου να αποκτά αλλά και να συλλέγει την γνώση έχει οδηγήσει στην αύξηση εγγράφων κειμένων που δημοσιεύονται καθημερινά στα μέσα κοινωνικής δικτύωσης αλλά και σε όλες τις ιστοσελίδες. Αυτό όμως που έχει περισσότερο ενδιαφέρον για τον σύγχρονο άνθρωπο είναι η εύρεση της κοινής γνώμης για όλα τα κοινωνικά ζητήματα, αλλά και για τον ίδιο εφ'όσον πρόκειται για δημόσιο και γνωστό πρόσωπο.

Η αλματώδης αύξηση των εγγράφων καθιστά δύσκολη την συλλογή χρήσιμης πληροφορίας καθώς υπάρχει πάντα και η άχρηστη πληροφορία ή ο θόρυβος. Παρόλα αυτά όλος ο όγκος της πληροφορίας περιέχει επαρκή δεδομένα και μπορεί να χρησιμοποιηθεί για να διαχωρίσουμε την χρήσιμη για εμάς πληροφορία από τον θόρυβο.

Οι τεχνικές συναισθηματικής ανάλυσης ανήκουν σε τρεις κυρίως κλάδους, αυτόν της τεχνητής νοημοσύνης και μηχανικής μάθησης, της υπολογιστικής γλωσσολογίας (computational linguistics) και της εξόρυξης κειμένου (text mining). Μέσω του ευρύτερου κλάδου της τεχνητής νοημοσύνης στον οποίο συγκαταλέγονται οι τεχνικές της εξόρυξης γνώσης από δεδομένα (data mining), η εξόρυξη κειμένου (text mining) και η μηχανική μάθηση (machine learning), επιτυγχάνουμε την κατανόηση κειμένων όπως και την συλλογή της επιθυμητής πληροφορίας.

Η εργασία αυτή πραγματεύεται την ανάλυση κειμένων ως προς το συναίσθημα του περιεχομένου τους (sentiment analysis) και στην συνέχεια παρουσιάζει το τελικό εργαλείο που αξιολογεί (μέσω της ανθρώπινης παρουσίας) τις βιβλιοθήκες που χρησιμοποιήσαμε sentiment detection classifiers. Στην παρούσα πτυχιακή εργασία οι τεχνικές που χρησιμοποιήθηκαν είναι η εξόρυξη κειμένου και η μηχανική μάθηση. Για την ανάπτυξη του λογισμικού χρησιμοποιήθηκε η γλώσσα Ruby και για την εφαρμογή αξιολόγησης το framework Ruby On Rails 4.0.0 το οποίο χρησιμοποιείται για την ανάπτυξη ιστοσελίδων.

1.2 Δομή

Η εργασία αποτελείται από έξι κεφάλαια συνολικά. Αυτό είναι το πρώτο κεφάλαιο στο οποίο γίνεται αναφορά στο γενικό περιεχόμενο της, αλλά και στο πώς είναι διαρθρωμένα τα κεφάλαια.

Στο δεύτερο κεφάλαιο εξηγούνται οι ορολογίες των τεχνικών μηχανικής μάθησης που εφαρμόστηκαν για να ακολουθήσουν στην συνέχεια όλες οι βιβλιοθήκες που χρησιμοποιήθηκαν καθώς και τα αποτελέσματα που εξήγαγαν. Αυτό αποτελεί και το τρίτο κεφάλαιο της εργασίας στο οποίο επιχειρείται η καταγραφή των πειραμάτων που διεξήχθησαν από τις βιβλιοθήκες που χρησιμοποιήθηκαν με σκοπό την σύγκριση των αποδόσεων τους και την επιλογή της βιβλιοθήκης με τα πιο ικανοποιητικά αποτελέσματα. Για την αξιολόγηση των βιβλιοθηκών έπρεπε να συλλεχθούν κείμενα τα οποία παρουσιάζονται στο ίδιο κεφάλαιο όπως και τα προβλήματα που

αντιμετωπίστηκαν και τις λύσεις που δώθηκαν σε αυτά.

Στο τέταρτο κεφάλαιο καταγράφονται τα μέρη από τα οποία απαρτίζεται η διαδικτυακή εφαρμογή, οι αρχιτεκτονικές και τα είδη των τεχνολογιών που χρησιμοποιήθηκαν κατά την διάρκεια της κατασκευής της.

Στο πέμπτο και τελευταίο κεφάλαιο εξηγούνται όλες οι λειτουργίες της εφαρμογής με την αναλυτική παρουσίαση εικόνων και πώς συνεργάζονται για να παρέχουν στον χρήστη το τελικό περιβάλλον. Ακολουθεί και μια αναφορά σε όλα τα σχεδιαστικά προβλήματα και διλήμματα που αντιμετωπίστηκαν και την τελική λύση αυτών, καθώς και τις μελλοντικές αλλαγές και προσθήκες μπορούν να γίνουν στην εργασία.

2 ΚΕΦΑΛΑΙΟ: Εξόρυξη Γνώσης- Γνώμης (opinion mining)

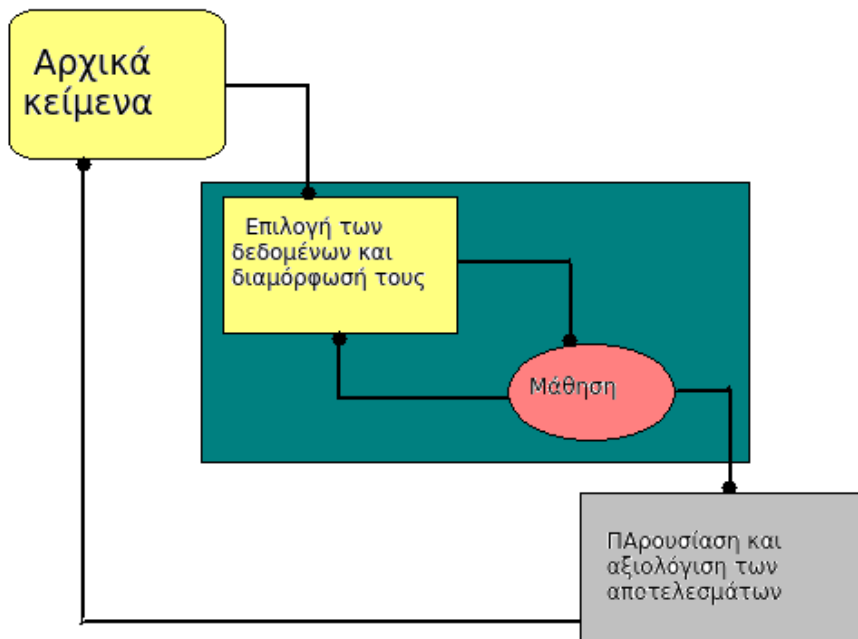
Η εξαγωγή συναισθημάτων από κείμενα διαφορετικής θεματολογίας είναι ένα σημαντικό παιδί έρευνας και τείνει να γίνει ακόμα μεγαλύτερο με την εισχώρηση της έννοιας big data στην επιστήμη των υπολογιστών. Ο γενικότερος όρος που χρησιμοποιείται και έχει γίνει γενικά αποδεκτός είναι η Εξόρυξη Γνώμης (Opinion Mining), ο οποίος δεν περικλείει την συναισθηματική ανάλυση, αλλά αποτελεί υποκατηγορία της Εξόρυξης Δεδομένων (Data mining).

Μπορούμε να ορίσουμε την κειμενική εξόρυξη γνώσης (text mining) ως μια αυτοματοποιημένη διαδικασία ανακάλυψης άγνωστων πληροφοριών που προέρχονται από διαφορετικές κειμενικές πηγές. Για να καταλήξουμε στην διαδικασία της εξόρυξης πρέπει αρχικά να έχει προεπεξεργαστεί το κείμενο (εξάλειψη τυχών κειμενικών ιδιαιτεροτήτων, αναγραμματισμών και αποθήκευση των τελικών δεδομένων σε μια βάση δεδομένων) να εφαρμοστεί το μοντέλο επεξεργασίας όπως η μηχανική μάθηση, και τελικά να αξιολογηθούν τα τελικά αποτελέσματα.

Αναφέρουμε μερικές από τις κατηγορίες των μεθόδων που χειρίζεται η εξόρυξη κειμένου:

- Γλωσσικός προσδιορισμός (Language Identification)
- Εξαγωγή χαρακτηριστικών γνωρισμάτων (Feature Extraction)
- Περιληπτική Παρουσίαση της Πληροφορίας (Summarization)
- Κατηγοριοποίηση, κατάταξη με επίβλεψη (Categorization, Supervised Classification)
- Ομαδοποίηση, μη επιβλεπόμενη κατάταξη (Clustering, Unsupervised Classification)

Δεν πρέπει να συγχέεται ο όρος text mining με την απλή αναζήτηση στον παγκόσμιο ιστό (google). Με την απλή αναζήτηση ο χρήστης ψάχνει να βρεί μια πληροφορία η οποία υπάρχει σιγουρα στον ιστό, σε αντίθεση με εξόρυξη γνώσης η οποία έχει ως στόχο να ανακαλύψει άγνωστη πληροφορία χρησιμοποιώντας την ήδη υπάρχον όγκο δεδομένων. Παραδείγματα απο εξόρυξη γνώσης μπορεί να είναι η κατηγοριοποίηση κειμένων, η ομαδοποίησή τους καθώς και η εξαγωγή και ανάλυση συναισθήματος. Στην παρακάτω εικόνα παρουσιάζονται σχηματικά οι παραπάνω έννοιες:



Εικόνα 2.2.1-1: Κειμενική εξόρυξη γνώσης.

Οι πηγές κειμένου που βρίσκονται στην διαθεσή μας στο διαδίκτυο είτε από τα δωμάτια συζητήσεων (chat rooms) είτε απο τις ηλεκτρονικές βιβλιοθήκες καθώς και απο το ηλεκτρονικό μας ταχυδρομείο μπορούν να αποτελέσουν χρήσιμο εργαλείο σε πολλές εμπορικές εφαρμογές. Όλος αυτός ο όγκος πληροφορίας που αποθηκεύεται ως αδόμητο η ημιδομημένο κείμενο μπορεί να

βοηθήσει στην εξόρυξη πληροφορίας από κείμενα ο οποίος είναι ένας τομέας που αφορά την επιστήμη της ανάκτησης πληροφοριών, της μηχανικής μάθησης, της στατιστικής, της εξόρυξης δεδομένων, καθώς και της υπολογιστικής γλωσσολογίας.

Ο όρος ο οποίος χρησιμοποιείται περισσότερο στην εργασία είναι αυτός της Συναισθηματικής Ανάλυσης (Sentiment Analysis). Μια τεχνική στην εξόρυξη κειμένου είναι η εύρεση ομοιότητας των αρχείων κειμένου με κάποια θεματική περιοχή (στην κατηγοριοποίηση κειμένων) είτε και μεταξύ των αρχείων (στην ομαδοποίηση αρχείων κειμένου).

Καταλήγουμε ότι συνδυάζοντας τις δύο μεθόδους (Εξόρυξη Γνώμης & Ανάλυση Συναισθήματος) μπορούν να προσδιοριστούν τόσο το συναίσθημα είτε είναι θετικό είτε αρνητικό, η υποκειμενικότητα-αντικειμενικότητα αλλά και άλλου είδους συναισθήματα των χρηστών που δημιουργούν τις αναρτήσεις(θυμός, χαρά κλπ).

Όπως αναφέρθηκε και στην αρχή της εργασίας η εξόρυξη συναισθήματος βρίσκει εφαρμογή σε πολλούς τομείς όπως τον πολιτικό, τον κοινωνικό, τον οικονομικό, τον καλλιτεχνικό και όλους τους τομείς που κατα κάποιο τρόπο εξαρτώνται από την κοινή γνώμη.

2.1 Άναλυση ορισμών και βασικής ορολογίας

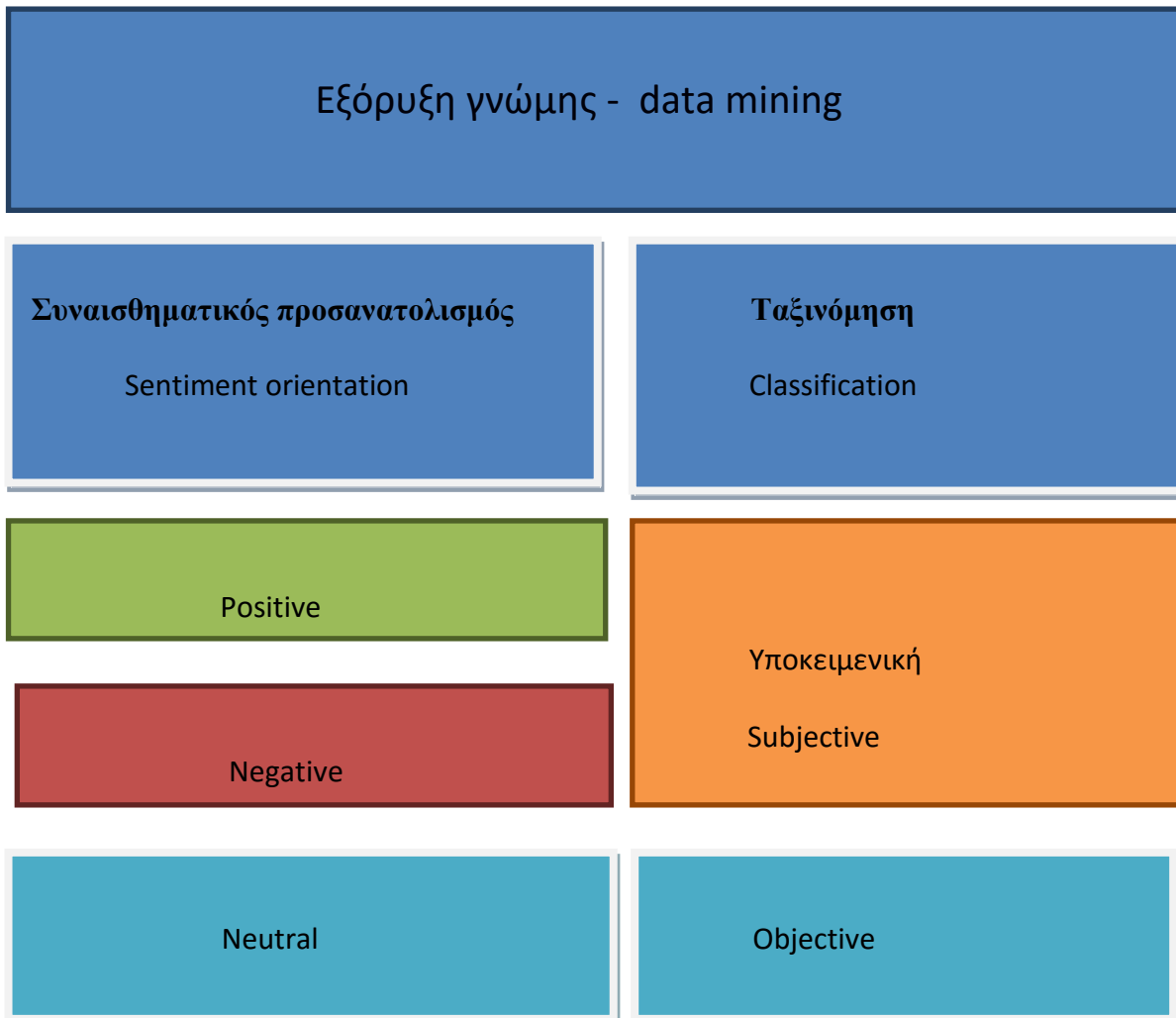
Το κεφάλαιο αυτό αφιερώνεται στην ανάλυση και ερμηνία των βασικότερων όρων που θα χρησιμοποιηθούν, για να γίνουν πιο κατανοητοί στον αναγνώστη.

Η σημασιολογική ανάλυση (text mining) εκτός των άλλων στοχεύει και στον προσδιορισμό μιας άποψης που εκφράζεται σε ένα κείμενο, το οποίο περιλαμβάνει και την διαδικασία της διάκρισης μεταξύ υποκειμενικής και αντικειμενικής άποψης [30].

Όταν αναφερόμαστε στην κατηγοριοποίηση κειμένου βάση της υποκειμενικότητας της, εννοούμε

την ταξινόμηση των προτάσεων σε δύο κλάσεις, υποκειμενικές ή αντικειμενικές. Αντικειμενικές θεωρούνται οι αμερόληπτες προτάσεις, ενώ υποκειμενικές όσες μεταφέρουν πληροφορίες με προσωπικές απόψεις. Η τεχνική αναγνώρισης και ταξινόμησης των προτάσεων με άποψη ονομάζεται «Υποκειμενική Ταξινόμηση» (Subjectivity Classification)

Υπάρχουν δυο είδη κατηγοριοποίησης, η μία είναι σε όλο το έγγραφο (Document-level), και η άλλη γίνεται σε κάθε πρόταση ξεχωριστά (Sentence level) κατηγοριοποιώντας μία μία τις προτάσεις σαν μία υποκειμενικές ή αντικειμενικές, και στην συνέχεια οι υποκειμενικές κατηγοριοποιούνται σε θετικές ή αρνητικές. Για να συμβεί αυτό, εφαρμόζονται τεχνικές επιβλεπόμενης μάθησης όπως ακριβώς και στην ταξινόμηση συναισθήματος εγγράφων καθώς και τεχνικές οι οποίες βασίζονται σε λεξικά συναισθήματος. Στη βιβλιογραφία συναντούμε συχνά τους συνώνυμους όρους «Συναισθηματική Κατηγοριοποίηση» (Sentiment Classification) και «Εξαγωγή Γνώμης» (Opinion Extraction).



Εικόνα 2.2.1-1: Υποκειμενική και αντικειμενική ταξινόμηση

Ορισμός : *Αντικειμενική - Υποκειμενική πρόταση:*

Η αντικειμενική πρόταση μεταφέρει αυτούσια μια πληροφορία, ενώ μια υποκειμενική εκφράζει συναίσθημα, απόψεις η και σκέψεις.

Ορισμός : *Πρόταση με άποψη:* είναι η πρόταση που εκφράζει άμεση ή έμμεση, θετική ή αρνητική άποψη. Μια τέτοια πρόταση μπορεί να κατηγοριοποιηθεί υποκειμενική είτε αντικειμενική.

ning) και την μέθοδο χρήσης λεξικών (lexicon) για να καταλήξουμε στη μέθοδο με τα καλύτερα αποτελέσματα για την διαδικτυακή εφαρμογή.

2.2.1 Μέθοδος χρήσης λεξικών

Η μέθοδος χρήσης λεξικών βασίζεται στη χρήση ειδικών λεξικών τα οποία περιέχουν προσημασμένες λέξεις, στις οποίες αποδίδονται βάρη ανάλογα με το νοήμα που έχουν και τί ένταση έχει η καθεμία στην χρήση της.

Η χρήση των λεξικών κατηγοριοποιείται ανάλογα με το άν

- Η χρήση λεξικών πραγματοποιείται ολιστικά, δηλαδή χρησιμοποιούνται στο σύνολο του κειμένου και υπολογίζεται με τη βοήθεια αλγορίθμου η θετική ή αρνητική στάση απέναντι στη θεματολογία (document level analysis),
- Ή η χρήση τους γίνεται τμηματικά, τα οποία διαχωρίζονται κατά την ανάλυση τους το οποίο

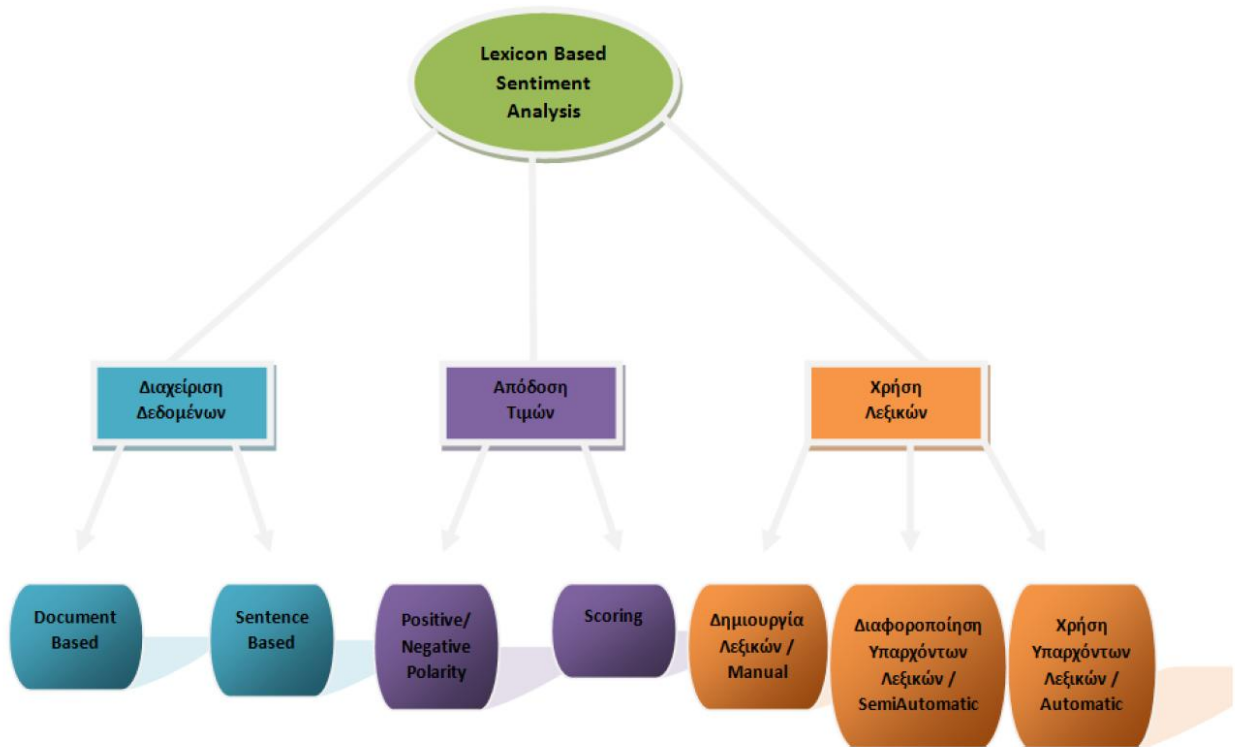
αποδίδει καλύτερα αν θέλουμε να εστιάσουμε σε συγκεκριμένη θεματολογία (sentence level analysis).

Η διαφορά στον τρόπο χειρισμού των λεξικών έχει να κάνει με τον τρόπο που μπορούμε να χειριστούμε τα δεδομένα. Μπορούμε να διαχωρίσουμε τα λεξικά και από τον τρόπο που το κάθε λεξικό χειρίζεται τα βάρη-τιμές που θα αποδοθούν στις λέξεις του και υπάρχουν οι εξής δύο κατηγορίες:

1. Η πρώτη και η πιο απλή κατηγορία είναι αυτή της κατηγοριοποίησης των λέξεων σε αρνητικές και θετικές. Η λέξη “super” θεωρείται θετική, ενώ το “awful” είναι αρνητική. Έτσι οι λέξεις αποκτούν πολικότητα, και αποδίδεται σε όλο το κείμενο θετική ή αρνητική υποκειμενική άποψη.

2. Η δεύτερη μέθοδος αντί να ταξινομεί τις λέξεις τους αποδίδει τόσο αρνητικό όσο θετικό βάρος, ανάλογα με την σημασιολογία τους στις προτάσεις. Για παράδειγμα η λέξη “mad” είναι 0,4 positive και 0,6 negative, ενώ η λέξη “angry” είναι 1,0 negative και 0 positive. Ο αναμενόμενος χειρισμός σε αυτήν την περίπτωση μπορεί να πραγματοποιηθεί με δύο τουλάχιστον τρόπους:

Ακολουθεί ένα αντιπροσωπευτικό σχήμα όπου παρουσιάζονται οι κατηγοριοποιήσεις των λεξικών για τη σημασιολογική – συναισθηματική προσέγγιση.



Εικόνα 2.2.1-1: Η κατηγοριοποίηση λεξικών

2.2.2 Πλεονεκτήματα – Μειονεκτήματα χρήσης λεξικών

Το κυριότερο πλεονέκτημα της χρήσης των λεξικών είναι αρχικά η χαμηλή υπολογιστική ισχύς που χρειάζονται για να εξάγουν αποτελέσματα αλλά και το γεγονός ότι δεν χρειάζονται κανένα κείμενο για να εκπαιδευτούν, το οποίο γλυτώνει χρόνο και την ανθρώπινη εργασία. Η τεχνική αυτή χαρακτηρίζεται ως μη επιβλεπόμενη τεχνική μάθησης χρησιμοποιείται όταν δεν έχουν χτιστεί σύνολα εκπαίδευσης.

Για να εξάγει ένα λεξικό καλά αποτελέσματα πρέπει να επιλεχθούν με πολύ προσοχή οι λέξεις του. Πρέπει να προσεχθεί να μην χρησιμοποιηθούν λέξεις μόνο από συγκεκριμένη θεματολογία, διότι το λεξικό ειδικεύεται μόνο σε μία θεματολογία και θα χάνει τις υπόλοιπες. Για να είναι ένα λεξικό ολοκ-

ληρωμένο και να εξάγει αξιόλογα αποτελέσματα πρέπει να περιέχει λέξεις από πολύ μεγάλες συλλογές δεδομένων.

Ο συναισθηματικός προσδιορισμός που προέρχεται λεξικά έχει να κάνει με πιθανότητες (probabilistic results). Η σημασιολογία μιας πρότασης δεν εξαρτάται μόνο από την κατηγοριοποίηση των λέξεων ή/και των συνωνύμων τους σε θετικές ή αρνητικές. Η πρόταση «you are damn right» δεν μπορεί να ταξινομηθεί ως θετική, επειδή το επίθετο «damn» έχει κατηγοριοποιηθεί στο λεξικό ως αρνητική λέξη, η σημασιολογία των προτάσεων εξαρτάται και από άλλους παράγοντες, όπως οι σύνδεσμοι και οι λέξεις. Οι τεχνικές και οι μέθοδοι ποικίλουν και ενδεχομένως, να πρέπει να εμπλουτιστούν ακόμα περισσότερο αλλά εν κατακλείδι είναι μια τεχνική που ενισχύει σημαντικά τις διαδικασίες της σημασιολογική ανάλυσης.

2.2.3 Μέθοδος Μηχανικής Μάθησης

Η μηχανική μάθηση (machine learning) ανήκει στο πεδίο της τεχνητής νοημοσύνης η οποία περιλαμβάνει διάφορους αλγορίθμους όπως και μεθόδους που δίνουν την δυνατότητα στους υπολογιστές να «μαθαίνουν». Η μηχανική μάθηση καθιστά εφικτή τη κατασκευή *προσαρμόσιμων* προγραμμάτων τα οποία λειτουργούν με βάση την αυτοματοποιημένη ανάλυση συνόλων δεδομένων και όχι τη διαίσθηση των μηχανικών που τα προγραμμάτισαν. Η μηχανική μάθηση μοιάζει σε μεγάλο βαθμό με τη στατιστική, αφού και τα δύο πεδία μελετούν την ανάλυση δεδομένων. Η χρήση των μεθόδων της μηχανικής μάθησης στην ανάλυση συναισθήματος και τη σημασιολογική ανάλυση κειμένων αποσκοπεί στον εντοπισμό και στην χρήση του κατάλληλου αλγόριθμου για την εξαγωγή αποτελεσμάτων.

Όπως είναι φυσικό όμως η εξαγωγή επιθυμητών αποτελεσμάτων απαιτεί πολλούς πειραματισμούς

από τους ερευνητές του παιδιού με διαφορετικού τύπου αλγόριθμους, τους οποίους πρέπει να εκπαιδεύσουν σε πολλά και διαφορετικά σύνολα δεδομένων έτσι ώστε να κατηγοριοποιήσουν τις αγνώστες περιπτώσεις.

2.2.4 Είδη κατηγοριοποιήσεων κειμένου – Μέτρηση ακρίβειας

Ένα από τα προβλήματα της εξόρυξης κειμένων είναι η εκτίμηση ομοιότητας μεταξύ εγγράφων διαφορετικού περιεχομένου. Αυτό σημαίνει είτε διαχωρισμό των εγγράφων σε προκαθορισμένες κατηγορίες είτε ομαδοποίηση εγγράφων σε φυσικές ομάδες.

Το προαπαιτούμενο βήμα για να καταλήξουμε στην επιλογή κατάλληλου αλγόριθμου μηχανικής μάθησης, είναι η εξόρυξη κειμένου και η κατηγοριοποίηση του. Ένα σημαντικό πρόβλημα στην εξόρυξη κειμένου είναι να βρεθεί η ομοιότητα μεταξύ διαφορετικών εγγράφων.

Η κατηγοριοποίηση κειμένων είναι η διαδικασία κατά την οποία τα ηλεκτρονικά έγγραφα που είναι γραμμένα σε φυσική γλώσσα, ανατίθενται σε μια ή περισσότερες κατηγορίες, ανάλογα με το περιεχόμενό τους. Κατατάσσουμε την κατηγοριοποίηση κειμένων στις εξής τρεις κατηγορίες:

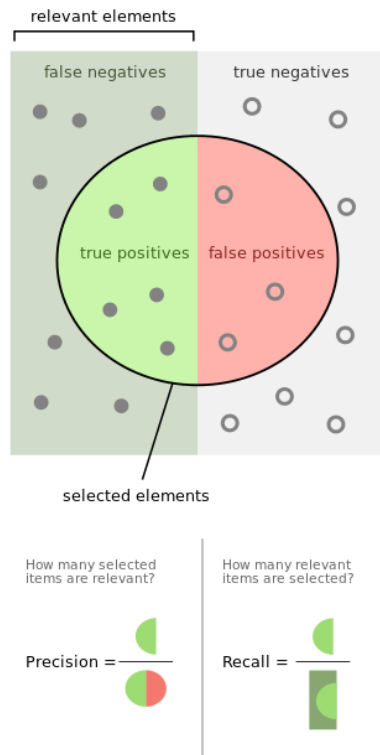
- Κατηγοριοποίηση με επίβλεψη (supervised classification), όπου κάποιος εξωτερικός μηχανισμός όπως ο άνθρωπος συμμετέχει στην σωστή κατηγοριοποίηση.
- Κατηγοριοποίηση με ημι-επίβλεψη (semi-supervised classification), όπου τμήματα των κειμένων έχουν ανατεθεί με ετικέτες από εξωτερικό μηχανισμό.
- Κατηγοριοποίηση χωρίς επίβλεψη (unsupervised classification), όπου η κατηγοριοποίηση γίνεται χωρίς καμία ανθρώπινη επίβλεψη.

Λόγω της αμφισιμότητας της γλώσσας και την δυσκολία να βρεθούν τα συναισθήματα η μέθοδος κατηγοριοποίησης που χρησιμοποιήθηκε στην παρούσα εργασία είναι αυτή με επίβλεψη. Το σύνολο κειμένων για να γίνουν τα πειράματα χωρίζεται σε δύο τμήματα το στο σύνολο εκπαίδευσης (training set) και το σύνολο δοκιμής (test set).

Το training set πρέπει να είναι πάντα μεγάλο για αυτο αποτελεί συνήθως το 80% της συλλογής, ενώ το test set το έχουμε για να δοκιμάσουμε πόσο καλά λειτουργεί ο αλγόριθμος για αυτό και αποτελεί το υπόλοιπο 20%. Το σύνολο δοκιμής ονομάζεται και ground truth διότι απο την πληροφορία που παρέχει εξαρτάται η ακρίβεια της κατηγοριοποίησης με επίβλεψη.

Τα έγγραφα του σύνολο εκπαίδευσης είναι γνωστό απο τον αλγόριθμο σε ποια κατηγορία ανήκουν αφού αυτά είναι τα δεδομένα με τα οποία λειτουργεί. Αρα είναι λογικό ότι όσο μεγαλύτερο είναι το σύνολο εκπαίδευσης τόσο πιο πολλά θα γνωρίζει ο αλγόριθμος και κατα συνέπεια θα αποδώσει καλύτερα. Ενώ για τα έγγραφα δοκιμής η πληροφορία της κατηγορίας υπάρχει αλλά δεν δίνεται στον αλγόριθμο. Ο αλγόριθμος αφού εκπαιδευτεί καλείται να τα κατατάξει στις σωστές κατηγορίες και στη συνέχεια τα αποτελέσματά του συγκρίνονται με τα πραγματικά δεδομένα.

Γνωστές τεχνικές κατηγοριοποίησης κειμένων είναι: κατηγοριοποιητής naïve Bayes, Latent Semantic Indexing (LSI), Support Vector Machines (SVM), tf-idf, τεχνητά νευρωνικά δίκτυα, αλγόριθμος k-κοντινότερων γειτόνων (kNN), concept mining και άλλοι. Οι κατηγοριοποιητές συνήθως δίνουν υψηλά ποσοστά ακρίβειας (high precision) αλλά χαμηλά ποσοστά ολοκλήρωσης (low recall). Οι όροι ακρίβεια (precision) και recall (ανάκληση) ορίζονται ως εξής:



Εικόνα 2.2.4-1: Ορισμός των *precision* και *recall*

Στην περίπτωση ενός ταξινομήτη που παίρνει την απόφαση για το αν ένα μήνυμα είναι θετικό ή αρνητικό ισχύουν τα εξής:

True positive: Το μήνυμα να ήταν θετικό και να προβλέφθηκε ως θετικό.

False positive: Το μήνυμα να ήταν θετικό και να μην προβλέφθηκε ως θετικό

False negative: Το μήνυμα να ήταν αρνητικό και να προβλέφθηκε ως θετικό

True negative: Το μήνυμα να ήταν αρνητικό και να προβλέφθηκε ως αρνητικό

Οι δύο όροι ανάκληση και ακρίβεια είναι αντιστρόφως ανάλογοι, οπότε συνήθως υπολογίζεται η ακρίβεια σε διάφορα επίπεδα ανάκλησης. Το μέτρο F είναι ενδεικτικό της ακρίβειας του συστήματος, ως συνάρτησης των recall και Precision λαμβάνοντας τιμές από 0 έως 1 και υπολογίζεται ως εξής:

$$F = 2 * \frac{\textit{Precision} * \textit{Recall}}{\textit{Precision} + \textit{Recall}}$$

Μια μέθοδος που έχει αξιόλογα αποτελέσματα σε ένα πεδίο μπορεί να μην είναι αξιόπιστη σε κάποιο άλλο. Το γεγονός αυτό οφείλεται στις ιδιαιτερότητες της φυσικής γλώσσας όπου η ερμηνεία πολλών λέξεων είναι αμφίσημη.

3 ΚΕΦΑΛΑΙΟ : Αποτελέσματα από την Πειραματική Διαδικασία

Οι τεχνικές λεξικών που χρησιμοποιήσαμε έδειξαν ότι η χρήση ενός σημασιολογικά βαθμολογημένου λεξικού ενισχύει σημαντικά τις μεθόδους της υποκειμενικής ταξινόμησης. Αν και το λεξικό που

χρησιμοποιήσαμε επιδέχεται βελτιώσεις, παρόλα αυτά συνέβαλε σημαντικά στο υψηλό ποσοστό των κειμένων που ταξινομήθηκαν και επιπλέον δεν εξαρτάται από ένα σημασιολογικό πεδίο γεγονός που ενισχύει την χρηστικότητά του.

Όπως ανφέραμε και σε προηγούμενο κεφάλαιο οι τεχνικές βασισμένες σε λεξικά τό μόνο που απαιτούν είναι σύνολα δοκιμών και όχι εκπαίδευσης. Το σύνολο δοκιμής που θα χρησιμοποιηθεί σε αυτά τα πειράματα είναι 5 csv αρχεία με 30 αρχεία το καθένα.

Η συνέχεια του κεφαλαίου αποτελείται από την πειραματική διαδικασία που ακολουθήθηκε για να εξεταστούν οι δύο βιβλιοθήκες που επιλέξαμε να δοκιμάσουμε.

3.1 Ο αλγόριθμος “Sentimental” βασισμένος σε λεξικό

Πρόκειται για μια απλή βιβλιοθήκη που κάνει συναισθηματική ανάλυση, είναι γραμμένη εξ’ ολοκλήρου στη γλώσσα προγραμματισμού Ruby. Μπορεί να εφαρμοστεί είτε σε επίπεδο πρότασης είτε σε επίπεδο κειμένου. Γίνεται διάσπαση στις προτάσεις και υπολογίζει την θετική ή την αρνητική πολικότητα, χαρακτηρίζει κάθε πρόταση ξεχωριστά και στην συνέχεια ολόκληρο το κείμενο.

Η προεπιλεγμένη πολικότητα είναι το 0.0. Αν μια πρόταση βγάλει άθροισμα 0.0 τότε θεωρείται ότι δεν έχει συναίσθημα και χαρακτηρίζεται ως ουδέτερη. Μεγαλύτερο άθροισμα από το 0.0 θεωρεί θετικό το αποτέλεσμα ενώ μικρότερο αρνητικό.

Υπάρχει η δυνατότητα αλλαγής της προεπιλεγμένης πολικότητας σε μη μηδενική τιμή όπως 0.5, άρα τα αποτελέσματα ορίζονται ως εξής:

- Θετικό άθροισμα > 0.5
- Ουδέτερα άθροισμα $-0.5 - 0.5$
- Αρνητικά άθροισμα < -0.5

3.1.1 Τα ενσωματωμένα συναισθηματικά λεξικά

Στην παρακάτω εικόνα εμφανίζονται πως είναι προσημασμένες οι λέξεις στην βιβλιοθήκη Sentimental. Συνολικά έχουν χρησιμοποιηθεί 18550 λέξεις όπου τους έχει αποδοθεί πολικότητα [-1, 1]

```

2 -1.0 )-:
3 -1.0 ):
4 -1.0 )o:
5 -1.0 8-0
6 -1.0 8/
7 -1.0 8\
8 -1.0 8c
9 -1.0 :'(
10 -1.0 :!-(
11 -1.0 :(
12 -1.0 :*(
13 -1.0 :,(
14 -1.0 :-(
15 -1.0 :-/
16 -1.0 :-S
17 -1.0 :-\
18 -0.50 :-|
19 -0.50 :/
20 -0.25 :0
21 -0.25 :5
22 -0.25 :\
23 -0.25 :|
24 -1.0 =(
25 -1.0 >:(
26 -1.0 D:
27 -1.0 sux
28 1.0 (o;
29 1.0 8-)
30 1.0 ;)
31 1.0 ;o)
32 1.0 %-)
33 1.0 (-:
34 1.0 :-)
35 1.0 (:
36 1.0 (o:
37 1.0 8)
38 1.0 :)
39 1.0 :-D
40 1.0 :-P
41 1.0 :D
42 1.0 :P
43 1.0 :P
44 1.0 :|
45 1.0 :o)
46 1.0 :p
47 1.0 ;^)
48 1.0 <3
49 1.0 &lt;3
50 1.0 =)

```

```

13658 -0.41666666667 grumbling
13659 -0.41666666667 fragile
13660 -0.41666666667 endure
13661 -0.41666666667 earthy
13662 -0.41666666667 dilute
13663 -0.41666666667 delinquent
13664 -0.41666666667 defenseless
13665 -0.41666666667 crude
13666 -0.41666666667 crooked
13667 -0.41666666667 bald
13668 -0.41666666667 annoyance
13669 -0.41666666667 angry
13670 -0.41666666667 alien
13671 -0.425 unresponsive
13672 -0.425 unfriendly
13673 -0.425 slack
13674 -0.425 mindless
13675 -0.4375 yahoo
13676 -0.4375 wistful
13677 -0.4375 whammy
13678 -0.4375 wasteful
13679 -0.4375 waive
13680 -0.375 war
13681 -0.4375 vulgar
13682 -0.4375 vanity
13683 -0.4375 unwell
13684 -0.4375 untried
13685 -0.4375 untested
13686 -0.4375 unshaded
13687 -0.4375 unquestioning
13688 -0.4375 unpredictable
13689 -0.4375 unofficially
13690 -0.4375 uninvited
13691 -0.4375 uninspired
13692 -0.4375 unhappily
13693 -0.4375 unfledged
13694 -0.4375 undertow
13695 -0.4375 underestimate
13696 -0.4375 underage
13697 -0.4375 uncovered
13698 -0.4375 unconvincing
13699 -0.4375 uncontrollable
13700 -0.4375 unconditional
13701 -0.4375 unclassified
13702 -0.4375 uncharitable
13703 -0.4375 unceremonious
13704 -0.4375 uncalled-for
13705 -0.4375 unawakened
13706 -0.4375 unalterability
13707 -0.4375 unaided

```

Εικόνα 3.1.1-1: Το προσημασμένο λεξικό του αλγορίθμου Sentimental

Εικόνα 3.1.1-2: Τα emoticons του αλγορίθμου Sentimental

Επειδή οι ιστοσελίδες και τα κοινωνικά δίκτυα όπως το twitter χρησιμοποιούν ειδικούς χαρακτήρες δηλαδή “φατσούλες” για να εκφράσουν συναίσθημα πολλά λεξικά χρησιμοποιούν τα ειδικά emoticons, η συγκεκριμένη χρησιμοποιεί 56 τέτοιους ειδικούς χαρακτήρες όπως φαίνεται και παραπάνω. Αυτή η ανάγκη προκύπτει από το γεγονός ότι στα κοινωνικά δίκτυα, που χρησιμοποιούν το Micro – Blogging, οι επιτρεπόμενοι χαρακτήρες ανάρτησης είναι περιορισμένοι. Έτσι λοιπόν, προκειμένου να γίνει πιο περιεκτική η επικοινωνία ανάμεσα στους χρήστες, αναπτύχθηκε και η χρήση των emoticons τα οποία μπορούν να εκφράζουν τα συναισθήματα των χρηστών με τη χρήση των ειδικών χαρακτήρων.

3.1.2 Αποτελέσματα πειραμάτων

Το πρώτο σύνολο πειραμάτων θα πραγματοποιηθεί με την Sentimental βιβλιοθήκη.

Η εξαγωγή του συναισθήματος θα γίνει μέσω της σύγκρισης των δεδομένων που έχουμε συλλέξει με τη χρήση του λεξικού που έχει η Sentimental. Όπως αναφέρθηκε και πιο πάνω το λεξικό έχει «Θετι-

κές» και «Αρνητικές» σημασιολογικά λέξεις. Υπάρχουν περιπτώσεις στις οποίες μπορεί να προκύψουν και μηδενικά σκόρ, αν πέσουμε στην περίπτωση να έχουμε τον ίδιο αριθμό θετικών και αρνητικών λέξεων.

Μετά τις προσθαφαιρέσεις προκύτουν οι τρεις παρακάτω κατηγορίες:

- *Θετική – Positivee*, αν το αποτέλεσμα της προσθαφαίρεσης είναι μεγαλύτερο από το 0.
- *Αρνητική – Negative*, αν το αποτέλεσμα της προσθαφαίρεσης είναι μικρότερο από το 0.
- *Ουδέτερη – Neutral*, αν το αποτέλεσμα της προσθαφαίρεσης είναι ίσο με το 0.

Παρακάτω παρουσιάζονται όλα τα πειράματα για όλα τα σύνολα δεδομένων, με τις εντολές που έτρεξαν και τους πίνακες με τα αποτελέσματα τους.

Το αρχείο “Gemfile” πού περιέχει τις βιβλιοθήκες που για να στήσουμε και να τρέξουμε το περιβάλλον του πρώτου συνόλου πειραμάτων με την βιβλιοθήκη Sentimental.

```
1 source 'http://rubygems.org'
2 gem 'sentimental' //Η βιβλιοθήκη
3 gem 'pry' //Εργαλείο για debugging
4 gem 'pry-rescue'
5 gem 'pry-stack_explorer'
```

Το εκτελέσιμο .rb που προγραμματίσαμε για να τρέξουμε τα πειράματα:

```
1 require 'csv'
2 require 'sentimental'
3 require 'open-uri'

4 csv_filename = ARGV[0]

5 def load_class_table(csv_filename)
6 csv_arxeio = File.new(csv_filename, "r")
7 # onoma arxeiou , class
8 table = CSV.read(csv_file, headers: true, header_converters: :symbol ,
9 col_sep: " , ")
9 csv_file.close
10 return pinakas
11 end
12 def sentiment_classify(pinakas , classif)
13 pragmatika = 0
```

```

14 pinakas.each do |rep|
15 dokimi_object = open(rep[:filename]).read
16 dokimi_klasi = classif.get_senti(dokimi_object)
17 alithini_klasi = rep[:class].to_sym
18 if alithini_klasi == dokimi_klasi
19   a. pragmatika += 1
20 end
21 akriveia = 0.0
22 akriveia = pragmatika / Float(pinakas.count)
23 return akriveia , pragmatika
24 end

25 def setup_classifier(dataset_filename = nil)
26 Sentimental.load_defaults
27 Sentimental.load_senti_file(dataset_filename) if dataset_filename
28 classif = Sentimental.new
29 return classif
30 end

31 dikimi_objects = load_class_table(csv_filename)

32 dokimi_objects.each do |rep|
33 puts "#{rep[:filename]} #{rep[:class]}"
34 end

35 classif = setup_classifier()
36 acc, tp = sentiment_classify(dokimi_objects,classif)
37 puts "Sunoliki akrivia: #{acc} kai #{tp} antikeimena katigoriopoihthikan
swsta"

```

Για την διευκόλυνσή μας χρησιμοποιήσαμε διαφορετικό σύνολο βιβλιοθηκών (gems) που λέγεται gemspec, για να μην παρουσιαστεί κάποιο σφάλμα.

Το πρώτο πείραμα το παρουσιάζουμε σε printscreen για να φανεί ο τρόπος παρουσίασης των αποτελεσμάτων στο τερματικό. Έντεκα από τα είκοσι κατηγοριοποιήθηκαν σωστά με ποσοστό επιτυχίας 55%.

```

11.txt negative
12.txt negative
13.txt negative
14.txt negative
15.txt negative
16.txt negative
17.txt negative
18.txt negative
19.txt negative
20.txt negative
21.txt negative
22.txt negative
23.txt negative
24.txt negative
25.txt negative
26.txt negative
27.txt negative
38.txt negative
39.txt negative
40.txt negative
Overall accuracy: 0.55 and 11 objects were classified correctly

```

Εικόνα 3.1.2-1: Πρώτο πείραμα βιβλιοθήκης με το αρνητικό αρχείο

Συγκεντρωτικός Πίνακας Αποτελεσμάτων

Λόγω εξοικονόμησης χώρου αντί για την προβολή στιγμιότυπων οθόνης κάθε πειράματος ξεχωριστά επιλέξαμε να καταγράψουμε όλα τα αποτελέσματα της βιβλιοθήκης Sentimental σε έναν συγκεντρωτικό πίνακα που παρουσιάζονται στην συνέχεια.

Αλγόριθμος Sentimental		
Όνομα Αρχείου	Ποσοστό %	Σωστά ταξινομημένα αρχεία
Negative.csv	55 %	22/40
Negative_1.csv	60 %	24/40
Positive.csv	80 %	32/40
thetika_1.csv	75 %	30/40
Neutral.csv	0	0/40

Πίνακας 3.1.2-1: Αποτελέσματα του αλγορίθμου Sentimental

Όπως ήταν αναμενόμενο ο αλγόριθμος αυτός δεν κατηγοριοποιεί κανένα ουδέτερο κείμενο το οποίο οφείλεται στην πολικότητα, αν ανοίγαμε το εύρος της πολικότητας τα αποτελέσματα θα ήταν πιο ικανοποιητικά. Παρατηρούμε ότι τα θετικά κατηγοριοποιούνται καλύτερα από ότι τα αρνητικά.

3.2 Ο αλγόριθμος “Sad Panda” βασισμένος σε λεξικό

Αυτός είναι ο δεύτερος αλγόριθμος για συναισθηματική ανάλυση βασισμένος σε λεξικό που χρησιμοποιήσαμε, γραμμένος στη γλώσσα προγραμματισμού Ruby. Ο αλγόριθμος εκτός από την θετική και αρνητική ταξινόμηση κειμένου κάνει και κατηγοριοποίηση συναισθήματος, αναγνωρίζοντας τα εξής συναισθήματα: "Οργή", "αποστροφή", "χαρά", "έκπληξη", "φόβο" και "θλίψη". Οι λέξεις παίρνουν πολικότητα (polarity), με εύρος από το 1 έως το 10 και εκφράζει την υποκειμενική άποψη, θετική ή αρνητική, που αποτυπώνεται στο σύνολο δεδομένων.

Το εκτελέσιμο πρόγραμμα που χρησιμοποιήσαμε για να τρέξουμε τα πειράματα:

```
1 require 'csv'
2 require 'sad_panda'

3 csv_filename = ARGV[0]

4 def load_class_table(csv_filename)
5   csv_file = File.new(csv_filename, "r")
6   # file name , class
7   table = CSV.read(csv_file, headers: true, header_converters: :symbol ,
8     col_sep: " , ")
9   csv_file.close
10  return table
11 end
```



```

11 def sentiment_classify(table , classifier)
12 cor = 0
13 table.each do |tp|
14 test_object = open(tp[:filename]).read

15 #test_class = classifier.get_sentiment(test_object)
16 test_class = SadPanda.polarity(test_object)
17 real_class = tuple[:class].to_sym

18 acc = 0.0
19 acc = cor / Float(table.count)
20 return accuracy , correct
21 end
22 def setup_sentimental_classifier(dataset_filename = nil)
23 Sentimental.load_defaults
24 Sentimental.load_senti_file(dataset_filename) if dataset_filename
25 classifier = Sentimental.new
26 return classifier
27 end
28 if real_class == test_class
29   a. correct += 1
30 end
31 t_objects = load_class_table(csv_filename)

32 t_objects.each do |tuple|
33 puts "#{tp[:filename]} #{tp[:class]}"
34 end

35 #classifier = setup_classifier()
36 #acc, tp = sentiment_classify(test_objects,classifier)
37 acc, tp = sentiment_classify(t_objects,nil)

38 puts "Overall accuracy: #{acc} and #{tp} objects were classified correctly"

```

3.2.1 Τα ενσωματωμένα συναισθηματικά λεξικά

Στη συνέχεια παρουσιάζεται η συναισθηματική προσέγγιση των λέξεων που χρησιμοποιήθηκαν στον συγκεκριμένο αλγόριθμο και η πολικότητα που τους δόθηκε. Ο αλγόριθμος ακόμα χρησιμοποιεί και μια μέθοδο που βρίσκει την υποκειμενικότητα δίνοντας βαρύτητα στις λέξεις (weaksubjectivity και strongsubjectivity), ανάλογα με το ποσο υποκειμενική είναι η λέξη. Τέλος μπορεί να αναγνωρίσει σε ένα κείμενο τα εξής συναισθήματα: "Οργή", "αποστροφή", "χαρά", "έκπληξη", "φόβο" και "θλίψη".

Κάθε λέξη που χρησιμοποιείται κατατάσσεται σε ένα από τα παραπάνω συναισθήματα.

```

1 module TermPolarities
2   # this method reads a csv file containing 'word,severity,positive/negative'
3   # triplits, and returns a giant hash where the keys are individual words
4   # and the values range between 0 and 2 (0 being more negative, 2 being most positive)
5   def self.get_term_polarities
6     @polarities = {"abandoned"=>2.5, "abandonment"=>2.5, "abandon"=>2.5, "abase"=>0,
7                   "abacement"=>0, "abash"=>0, "abate"=>2.5, "abdicate"=>2.5, "aberration"=>0,
8                   "abhor"=>0, "abhorred"=>0, "abhorrence"=>0, "abhorrent"=>0,
9                   "abhorrently"=>0, "abhors"=>0, "abidance"=>10, "abide"=>10, "abject"=>0,
10                  "abjectly"=>0, "abjure"=>2.5, "abilities"=>7.5, "ability"=>7.5, "able"=>7.5,
11                  "abnormal"=>2.5, "abolish"=>2.5, "abominable"=>0, "abominably"=>0,
12                  "abominate"=>0, "abomination"=>0, "above"=>7.5, "aboveaverage"=>7.5,
13                  "abound"=>7.5, "abrade"=>2.5, "abrasive"=>0, "abrupt"=>2.5, "abscond"=>0,
14                  "absence"=>2.5, "absentee"=>2.5, "absentminded"=>0, "absolve"=>10, "absurd"=>0,
15                  "absurdity"=>0, "absurdly"=>0, "absurdness"=>0, "abundant"=>7.5,
16                  "abundance"=>7.5, "abuse"=>0, "abuses"=>2.5, "abusive"=>0, "abysmal"=>0,
17                  "abysmally"=>0, "abyss"=>0, "accede"=>10, "accept"=>7.5, "acceptance"=>7.5,
18                  "acceptable"=>7.5, "accessible"=>7.5, "accidental"=>2.5, "acclaim"=>10,
19                  "acclaimed"=>10, "acclamation"=>10, "accolade"=>10, "accolades"=>10,
20                  "accommodative"=>7.5, "accomplish"=>7.5, "accomplishment"=>7.5,

```

Εικόνα 3.2.1-1: Απόδοση πολικότητας στις λέξεις απο 1-10

Line	Word	Severity	Polarity
1	abandoned	weaksbj	negative
2	abandonment	weaksbj	negative
3	abandon	weaksbj	negative
4	abase	strongsubj	negative
5	abacement	strongsubj	negative
6	abash	strongsubj	negative
7	abate	weaksbj	negative
8	abdicate	weaksbj	negative
9	aberration	strongsubj	negative
10	abhor	strongsubj	negative
11	abhorred	strongsubj	negative
12	abhorrence	strongsubj	negative
13	abhorrent	strongsubj	negative
14	abhorrently	strongsubj	negative
15	abhors	strongsubj	negative

Εικόνα 3.2.1-2: Λεξικό υποκειμενικότητας του αλγορίθμου

file 1543 lines (1542 sloc) 23.03 kb		
Open Edit Raw Blame History Delete		
<input type="text" value="Search this file..."/>		
1	abhor	anger
2	abhor	anger
3	abhor	disgust
4	abhorrence	anger
5	abhorrent	disgust
6	abomin	anger
7	abomin	disgust
8	abominably	disgust
9	abominate	anger
10	abomination	anger
11	admir	joy
12	admir	surprise
13	admirable	joy
14	admirably	joy
15	admiration	joy
16	admiration	surprise

Εικόνα 3.2.1-3: Λεξικό συναισθημάτων του αλγορίθμου

3.2.2 Αποτελέσματα πειραμάτων

Παρακάτω παρουσιάζονται όλα τα πειράματα σε έναν συγκεντρωτικό πίνακα για όλα τα σύνολα δεδομένων.

Στην εικόνα 3.2.2-4 βλέπουμε ένα πείραμα με το σύνολο κειμένων positive.csv το οποίο αποτελείται

απο είκοσι αρνητικά κείμενα και τα 18 από τα 20 κατηγοριοποιήθηκαν σωστά με ποσοστό επιτυχίας 90%

```
1.txt positive
2.txt positive
3.txt positive
4.txt positive
5.txt positive
6.txt positive
7.txt positive
8.txt positive
9.txt positive
10.txt positive
28.txt positive
29.txt positive
30.txt positive
31.txt positive
32.txt positive
33.txt positive
34.txt positive
35.txt positive
36.txt positive
37.txt positive
Overall accuracy: 0.9 and 18 objects were classified correctly
```

Εικόνα 3.2.2-1: Πείραμα του αλγορίθμου όπως αποτυπώνεται στην οθόνη

Πίνακας Αποτελεσμάτων

Αλγόριθμος Sad Panda		
Όνομα Αρχείου	Ποσοστό %	Σωστά ταξινομημένα αρχεία
Negative.csv	75 %	30/40
Negative_1.csv	60 %	24/40
Positive.csv	90 %	36/40
thetika_1.csv	70 %	28/40
Oudetera.csv	15.8%	6/40

Πίνακας 3.2.2-1: Τα αποτελέσματα του αλγορίθμου

Αυτός ο αλγόριθμος αποδίδει καλύτερα στην ταξινόμηση των θετικών και των αρνητικών και κατηγοριοποίησε ουδέτερα κείμενα αν και με πολύ χαμηλό ποσοστό, μόλις 15.5%. Ο μόνος τρόπος για να έχουμε καλύτερη απόδοση στις παραπάνω βιβλιοθήκες είναι να αλλάξουμε τα λεξικά τους προσθέ-

τοντας άλλες πιο βαρυσήμαντες λέξεις που πιθανόν να έχουν περισσότερο συναισθηματικό φορτίο. Επειδή όμως αυτός ο τρόπος είναι πολύ χρονοβόρος και δεν είναι σίγουρο ότι θα αποδώσει καλύτερα καταφεύγουμε στην επόμενη λύση που είναι ο κατηγοριοποιητής Naïve Bayes.

3.3 Κατηγοριοποιητής “Classifier”

Ο αλγόριθμος Classifier είναι προγραμματισμένος για Bayesian και άλλα είδη κατηγοριοποιήσεων. Ο αλγόριθμος Bayes (Naïve Bayes) είναι ένας κατηγοριοποιητής βασισμένος στις πιθανότητες. Η κατασκευή του εξαρτάται από ένα σύνολο εκπαίδευσης για να εκτιμήσει τις παραμέτρους μιας κατανομής πιθανότητας, δεδομένων των τιμών των χαρακτηριστικών ενός νέου εγγράφου. Αυτές οι πιθανότητες εκτιμώνται με τη βοήθεια του θεωρήματος Bayes.

Το εκτελέσιμο πρόγραμμα που χρησιμοποιήσαμε για να τρέξουμε τα πειράματα.

```
1 require 'csv'
2 require 'fast_stemmer'
3 require 'classifier'
4 csv_filename = ARGV[0]

5 def load_class_table(csv_filename)
6   csv_file = File.new(csv_filename, "r")
7   # file name , class
8   table = CSV.read(csv_file, headers: true, header_converters: :symbol ,
9     col_sep: " , ")
10  csv_file.close
11  return table
12 end

12 def sentiment_classify(table , classifier)
13   cor_count = 0
14   table.each do |tp , assessor = classifier, count = cor_count|
```

```

15 acc = 0.0
16 acc= cor_count / Float(table.count)
17 return acc , cor_count
18 end

19 def setup_classifier(dtset = {})
20 thetika_files = load_class_table(dataset["positive"])
21 arnitika_files = load_class_table(dataset["negative"])

22 clas = Classifier::Bayes.new('negative', 'positive')

23 thetika_files.each do |tp|
24 train_object = open(tp[:filename]).read
    a. clas.train_positive(train_object)
25 end
26 neg_files.each do |tp|
27 train_object = open(tp[:filename]).read
    a. clas.train_negative(training_object)
28 end

29 return classifier
30 end

31 test_objects = load_class_table(csv_filename)

32 test_objects.each do |tuple|
33 puts "#{tp[:filename]} #{tp[:class]}"
34 end

35 clas = setup_classifier({"positive" => "positive.csv", "negative" => "nega-
    tive.csv"})
36 acc, tp = sentiment_classify(test_objects, classifier) do |klasifier,
    test_object , tp, cor_count|
37 puts klasifier
38 puts correct_count
39 test_class = klasifier.classify(test_object)
40 puts "#{real_class} #{test_class}"
41 correct_count
42 end

43 puts "Overall accuracy: #{acc} and #{tp} objects were classified correctly"

```

3.3.1 Data set - Συλλογές δεδομένων

Απο την θεωρία των Naïve Bayes είναι γνωστό οτι απαιτούν ένα ικανοποιητικό σύνολο εκπαίδευσης για να εξάγουν ικανοποιητικά αποτελέσματα. Από αυτή την θεωρία συνεπάγεται ότι η ποιότητα αλλά και η ποσότητα των συλλογών εκπαίδευσης (corpus) που χρησιμοποιούνται παίζουν πολύ σημαντικό ρόλο στην ακρίβεια των αποτελεσμάτων της μεθόδου. Ο εντοπισμός συναισθήματος σε ένα κείμενο απαιτεί την ανθρώπινη επίβλεψη λόγω των αμφίσημων νοημάτων της φυσικής γλώσσας για αυτό απαιτείται κατηγοριοποίηση με επίβλεψη για τον Naïve Bayes κατηγοριοποιητή. Για να πετύχουμε το σύνολο εκπαίδευσης που θα έχει τα πιο καλά αποτελέσματα πρέπει να χωρίζουμε το σύνολο κειμένων σε δύο κατηγορίες, αυτό του συνόλου εκπαίδευσης και του συνόλου δοκιμής. Όποιο σύνολο εκπαίδευσης έχει τα καλύτερα αποτελέσματα θα επιλεγεί και θα χρησιμοποιηθεί για τον κατηγοριοποιητή. Κατά την διάρκεια πειραμάτων παρατηρήθηκε ότι ο κατηγοριοποιητής εξάγει καλύτερα αποτελέσματα όταν τροφοδοτείται με πολλά είδη κειμένων. Έκτός απο τα κείμενα που συλλέχθηκαν χειροκίνητα από το διαδίκτυο χρησιμοποιήθηκε το polarity dataset v2.0 από το Movie Review Data [13] που αποτελείται απο 1000 θετικές και 1000 αρνητικές κριτικές ταινιών.

- thetika_web.csv (20 θετικά κείμενα τεχνολογικού περιεχομένου)
- arnitika_web.csv (20 αρνητικά κείμενα τεχνολογικού περιεχομένου)
- thetika_web_1.csv (60 θετικά κείμενα τεχνολογικού περιεχομένου)
- arnitika_web_1.csv (30 αρνητικά κείμενα τεχνολογικού περιεχομένου)
- thetika_last.csv (30 θετικά κείμενα τεχνολογικού περιεχομένου)
- cornell_thetika.csv (1000 θετικές κριτικές)
- cornell_arnitika.csv (1000 αρνητικές κριτικές)
- thetika_health.csv (30 θετικά κείμενα ιατρικού περιεχομένου)
- arnitika_health.csv (30 αρνητικά κείμενα ιατρικού περιεχομένου)

3.3.2 Πειράματα - Αποτελέσματα

Ακολουθεί όλη η πειραματική διαδικασία με τα σύνολα εκπαίδευσης για τον κατηγοριοποιητή Naïve Bayes, για να καταλήξουμε στα πιο ικανοποιητικά αποτελέσματα.

Η λογική της εντολής για να «τρέξει» το πρόγραμμα:

rescue clas.rb TestSet.csv Positive TrainingSet.csv Negative TrainingSet.csv

Σύνολο Δοκιμής (Test set)	Σύνολο Εκπαίδευσης (Training Set)	
	Θετικά	Αρνητικά
thetika_last.csv	thetika_web_1.csv	arnitika_web_1.csv
	18 / 30	
arnitika_1.csv	thetika_web_1.csv	arnitika_web_1.csv
	49/60	
arnitika_1.csv	positive.csv	arnitika_web_1.csv
	0 / 60	
arnitika_1.csv	thetika_1.csv	arnitika_web_1.csv
	3 / 60	
arnitika_1.csv	thetika_web_1.csv	negative.csv
	60/60	
thetika_1.csv	positive.csv	arnitika_web_1.csv
	19 / 20	
thetika_1.csv	thetika_web_1.csv	arnitika_web_1.csv

Πίνακας 3.3.2-1: Πρώτο πείραμα

Σύνολο Δοκιμής (Test Set)	Σύνολο Εκπαίδευσης (Training Set)	
	Θετικά	Αρνητικά
arnitika.csv	thetika_web_1.csv	arnitika_web_1.csv
	58 / 60	
thetika.csv	thetika_web_1.csv	arnitika_web_1.csv
	8 / 60	

Πίνακας 3.3.2-2: Δεύτερο πείραμα με ίδιο training set

Σύνολο Δοκιμής (Test Set)	Σύνολο Εκπαίδευσης (Training Set)	
	Θετικά	Αρνητικά
arnitika_web_1.csv	thetika_web_1.csv	negative.csv
	3 / 30	
thetika_web_1.csv	positive.csv	arnitika_web_1.csv
	1 / 82	
arnitika.csv	positive.csv	arnitika_web_1.csv
	0 / 20	
thetika.csv	thetika_web_1.csv	negative.csv
	2 / 20	

Πίνακας 3.3.2-3: Τρίτο πείραμα

Όπως ήταν αναμενόμενο δέν ήταν όλα τα αποτελέσματα ικανοποιητικά, για αυτό έγινε αλλαγή στο σύνολο εκπαίδευσης για νέο σέτ πειραμάτων.

- all_positive.csv (thetika_1.csv+thetika_web_1.csv)
- all_negative.csv (arnitika_1.csv+ arnitika_web_1.csv)
- all_positive_1.csv (thetika.csv+ thetika_web_1.csv)
- all_negative_1.csv (arnitika.csv+ arnitika_web_1.csv)

Σύνολο Δοκιμής	Σύνολο Εκπαίδευσης	
	Θετικά	Αρνητικά
thetika.csv	all_positive.csv	all_negative.csv
	1 / 20	
thetika.csv	all_positive.csv	arnitika_web_1.csv
	10 / 20	
thetika.csv	all_positive.csv	Negative.csv
	18 / 20	

Πίνακας 3.3.2-4: Πρώτο πείραμα με θετικά σύνολα δοκιμής

Σύνολο Δοκιμής	Σύνολο Εκπαίδευσης	
	Θετικά	Αρνητικά
thetika_1.csv	all_positive_1.csv	all_negative_1.csv
	0 / 20	
thetika_1.csv	all_positive_1.csv	arnitika_web_1.csv
	5 / 20	
thetika_1.csv	positive	all_negative_1.csv
	19 / 20	

Πίνακας 3.3.2-5: Δεύτερο πείραμα με θετικά σύνολα δοκιμής

Σύνολο Δοκιμής	Σύνολο Εκπαίδευσης	
	Θετικά	Αρνητικά
arnitika_1.csv	all_positive_1.csv	all_negative_1.csv
	20 / 20	
arnitika_1.csv	all_positive_1.csv	arnitika_web_1.csv
	17 / 20	
arnitika_1.csv	positive.csv	all_negative_1.csv
	1 / 20	

Πίνακας 3.3.2-6: Πείραμα με αρνητικό σύνολο

Οι παρατηρήσεις που έχουμε να κάνουμε απο το πρώτο σετ πειραμάτων είναι οι εξής:

- Παρατηρούμε ότι τα θετικά εξαρτώνται από το σύνολο εκπαίδευσης τους και αποδίδουν πολυ καλά όταν έχουν κοινή θεματολογία (domain specific), δηλαδή έχουμε καλύτερα αποτελέσματα όταν το σύνολο εκπαίδευσης ανήκει στην ίδια κατηγορία με το σύνολο δοκιμής. Έχουμε χειρότερα αποτελέσματα όταν διαφέρει η κατηγορία του συνόλου εκπαίδευσης είναι απο αυτή του συνόλου δοκιμής.
- Τα αρνητικά σύνολα από την άλλη δεν έχουν το ίδιο πρόβλημα, ο αλγόριθμος δεν είναι domain specific.

Για να δοκιμάσουμε το έτοιμο σύνολο κειμένων του Cornell [13] πρέπει να συνδυάσουμε όλα τα κείμενα που έχουμε συλλέξει. Έτσι δημιουργήθηκαν τα :

- finalpos.csv (περιέχει συνολικά 150 αρχεία)
- finalneg.csv (περιέχει συνολικά 90 αρχεία)

Σύνολο Δοκιμής	Σύνολο Εκπαίδευσης		Απόδοση %	Σύνολο
	Θετικά	Αρνητικά		
finalneg.csv	cornell_thetika.csv	cornell_arnitika.csv	63.3%	57/90
finalpos.csv	cornell_thetika.csv	cornell_arnitika.csv	75,3%	113/150
cornell_arnitika.csv	finalpos.csv	finalneg.csv	97.2%	972/1000
cornell_thetika.csv	finalpos.csv	finalneg.csv	17.3%	173/1000

Πίνακας 3.3.2-7: Πείραμα κειμένων Cornell ως training set

Από τους παραπάνω πίνακες παρατηρούμε ότι το σύνολο δοκιμής του Cornell αποδίδει λίγο καλύτερα στο δετικό σύνολο δοκιμής απο το αρνητικό. Ό μόνος τρόπος να καταλάβουμε γιατί έχει πιο εξασθενημένη απόδοση το σύνολο του Cornell είναι να διασπάσουμε τα σύνολα δοκιμών και να τα δοκιμάσουμε διαδοχικά.

Σύνολο Δοκιμής	Σύνολο Εκπαίδευσης		Απόδοση %	Σύνολο
	Θετικά	Αρνητικά		
thetika_web.csv	cornell_thetika.csv	cornell_arnitika.csv	96.6 %	29/30
positive.csv	cornell_thetika.csv	cornell_arnitika.csv	90 %	18/20
thetika_1.csv	cornell_thetika.csv	cornell_arnitika.csv	83.3 %	25/30

thetika_health.csv	cornell_thetika.csv	cornell_arnitika.csv	70 %	21/30
thetika_web_1.csv	cornell_thetika.csv	cornell_arnitika.csv	81.4%	57/70

Πίνακας 3.3.2-8: Πρώτο πείραμα θετικών συνόλων δοκιμών

Χρησιμοποιώντας το σύνολο Cornell όλα τα υπόλοιπα σύνολα έχουν καλή απόδοση. Τα αρνητικά όπως και στο προηγούμενο σετ πειραμάτων αποδίδουν χειρότερα. Η θεωρία λέει ότι τα αρνητικά πάντα αποδίδουν χειρότερα από τα αρνητικά επειδή είναι πιο δύσκολο να “πιασει” ο αλγοριθμος τις αρνητικές έννοιες. Άλλοι λόγοι μπορεί να είναι η ανομοιότητα μεταξύ συνόλου δοκιμής και εκπαίδευσης στον η στον μικρό αριθμό του αρνητικού συνόλου εκπαίδευσης.

Σύνολο Δοκιμής	Σύνολο Εκπαίδευσης		Απόδοση %	Σύνολο
	Θετικά	Αρνητικά		
arnitika_web.csv	cornell_thetika.csv	cornell_arnitika.csv	25 %	5/20
negative.csv	cornell_thetika.csv	cornell_arnitika.csv	85 %	17/20
arnitika_1.csv	cornell_thetika.csv	cornell_arnitika.csv	100 %	20/20
arnitika_health.csv	cornell_thetika.csv	cornell_arnitika.csv	36.6 %	11 / 30
arnitika_web_1.csv	cornell_thetika.csv	cornell_arnitika.csv	13%	4 /30

Πίνακας 3.3.2-9: Δεύτερο πείραμα αρνητικών συνόλων δοκιμών

Με τις αλλαγές που έγιναν τα ποσοστά επιτυχίας στα αρνητικά αυξήθηκαν ενώ στα θετικά μειώθηκαν, με αυτά τα σύνολα εκπαίδευσης. Τα αποτελέσματα αυτά δεν μπορούν να είναι τα τελικά

που θα χρησιμοποιήσουμε αφού πρέπει να είναι πιο ισορροπημένα τα θετικά με τα αρνητικά.

Σύνολο Δοκιμής	Σύνολο Εκπαίδευσης	
	cornell_thetika	finalneg_test
arnitika_web.csv	20/20	
negative.csv	3/20	
arnitika_1.csv	9/20	
arnitika_web_1.csv	29/30	

Πίνακας 3.3.2-10: Πείραμα θετικού συνόλου με το νέο σύνολο εκπαίδευσης

Testset (Σύνολο Δοκιμής)	Trainingset (Σύνολο Εκπαίδευσης)	
	cornell_thetika	golneg_test
thetika_web.csv		0/20
positive.csv		16/20
thetika_1.csv		19/20
thetika_web_1.csv		0/90

Πίνακας 3.3.2-11: Δοκιμή αρνητικού συνόλου με νέο σύνολο εκπαίδευσης

Μετρώντας τα πειράματα με ποσοστο ακρίβειας μόνο δεν μπορούμε να κατανοήσουμε ακριβώς που δεν λειτουργεί καλά ο αλγόριθμός. Μετά απο καποιές αλλαγές στο πρόγραμμα η απόδοση του αλγορίθμου θα είναι ως προς το precision και το recall, ορισμοί που αναλύσαμε και στο δεύτερο κεφάλαιο.

Ακολουθεί το τελικό σύνολο πειραμάτων με τις εντολές που χρησιμοποιήθηκαν καθώς και τα

αποτελέσματα σε πίνακες που εξήγαγε ο αλγόριθμος Naïve Bayes, τα οποία αποτυπώνονται ως προς την ακρίβεια (recall) και την ανάκληση (precision), καθώς επίσης μας δίνονται και τα λάθος θετικά (false positives) και λάθος αρνητικά (false negatives) κατηγοριοποιημένα κείμενα.

Εντολή: rescue test_oudetera.rb comb_test.csv comb_thetika.csv comb_arnitika.csv 0.04

Overall clas. accuracy: 0.63 - 133/210

ΘΕΤΙΚΑ

ΑΡΝΗΤΙΚΑ

Precision: 0.7

Precision: 0.6

Recall: 0.7

Recall: 0.9

false positives 21.0

false positives 27.0

false negatives 23.0

false negatives 4.0

Πίνακας 3.3.2-12: Πρώτο πείραμα με μετρήσεις ως προς την ακρίβεια και την ανάκληση

Εντολή: rescue test_oudetera.rb comb_test.csv comb_thetika_oudeteral.csv comb_arnitika.csv 0.04

Overall clas. accuracy: 0.6 - 132/210

ΘΕΤΙΚΑ

ΑΡΝΗΤΙΚΑ

Precision: 0.67

Precision: 0.6

Recall: 0.7

Recall: 0.85

false positives 27

false positives 23

false negatives 22

false negatives 7

Πίνακας 3.3.2-13: Δεύτερο πείραμα με μετρήσεις ως προς την ακρίβεια και την ανάκληση

Εντολή: rescue test_oudetera.rb comb_test.csv comb_thetika.csv comb_aritika_oudetera.csv 0.04

Overall clasaccuracy: 0.58- 125/210

POSITIVE

NEGATIVE

Precision: 0.74

Precision: 0.5

Recall: 0.52

Recall: 0.95

false positives 14

false positives 46

false negatives 37

false negatives 2

Πίνακας 3.3.2-14: Τρίτο πείραμα με μετρήσεις ως προς την ακρίβεια και την ανάκληση

Εντολή: rescue test_oudetera.rb comb_test.csv comb_thetika_oudetera.csv comb_arnitika_oudetera.csv 0.04

Overall clas. accuracy: 0.60- 129/210

ΘΕΤΙΚΑ

ΑΡΝΗΤΙΚΑ

Precision: 0.74

Precision: 0.5

Recall: 0.5

Recall: 0.95

false positives 14

false positives 46

false negatives 37

false negatives 2

Πίνακας 3.3.2-15: Τέταρτο πείραμα με μετρήσεις ως προς την ακρίβεια και την ανάκληση

Εντολή: rescue test_oudetera.rb comb_test.csv comb_thetika_oudetera.csv comb_arnitika_oudetera.csv 0.03

Overall clas accuracy: 0.60 - 129/210

ΘΕΤΙΚΑ

ΑΡΝΗΤΙΚΑ

Precision: 0.7

Precision: 0.5

Recall: 0.6

Recall: 0.9

false positives 20

false positives 34

false negatives 29

false negatives 4

Πίνακας 3.3.2-16: Πέμπτο πείραμα με μετρήσεις ως προς την ακρίβεια και την ανάκληση

3.4 Προβλήματα - Συμπεράσματα

Όπως ήταν φυσικό και γνωστό και από τις έρευνες άλλων επιστημόνων για να έχουν ικανοποιητικά αποτελέσματα οι κατηγοριοποιητές με λεξικά πρέπει να υπάρχουν τεράστιες συλλογές δεδομένων. Υπάρχουν όμως και τρόποι να κάνουμε και τα λεξικά να λειτουργήσουν καλύτερα χωρίς να τα εμπλουτίσουμε με νέες λέξεις. Ένας τρόπος είναι να αλλάξουμε την προεπιλεγμένη πολικότητα η οποία απλά μπορεί να βοηθήσει στην κατηγοριοποίηση και των μη συναισθηματικά φορτισμένων κειμένων (ουδέτερων). Παρατηρείται όμως και το φαινόμενο ακόμα και τα λεξικά που έχουν βασιστεί σε κατηγορία λέξεων ενός πεδίου να μην μπορούν να αποδώσουν εξίσου όταν χρησιμοποιούνται σε διαφορετικό πεδίο εφαρμογής. Επειδή όμως υποστηρίζεται ότι οι κατηγοριοποιητές λεξικών βασίζονται σε στατιστικούς τρόπους εξαγωγής αποτελεσμάτων προτιμήσαμε την χρήση του Naive Bayes κατηγοριοποιητή.

Το κυριότερο πρόβλημα που αντιμετωπίσαμε ήταν αυτό της συλλογής μεγάλου όγκου δεδομένων, αφού η χρήση του Naive Bayes κατηγοριοποιητή απαιτεί κείμενα και για εκπαίδευση και για δοκιμή. Τα δεδομένα που συλλέξαμε προορίζονται το μεγαλύτερο μέρος τους για εκπαίδευση και το υπόλοιπο για δοκιμή. Το σύνολο δοκιμής που αποτελεί το 20-30% του συνόλου δεδομένων δεν αποτελεί πρόβλημα διότι ο αλγόριθμος πρέπει να εκπαιδευτεί με τον καταλληλότερο τρόπο για να επιτυγχάνει την καλύτερη κατηγοριοποίηση. Μία ιδιαιτερότητα που παρατηρήθηκε ήταν ότι κατά την εκπαίδευση του αλγορίθμου παρατηρήσαμε ότι το αρνητικό σύνολο εκπαίδευσης πρέπει να είναι μικρότερο από αυτό του θετικού συνόλου για την επίτευξη του καλύτερου δυνατού αποτελέσματος. Μια ακόμη σημαντική παρατήρηση είναι ότι τα σύνολα εκπαίδευσης πρέπει να αποτελούνται από διαφορετικούς τομείς για να μην αποδίδει καλά ο κατηγοριοποιητής σε έναν συγκεκριμένο τομέα (domain specific). Σε γενικές γραμμές η κατηγοριοποίηση με επίβλεψη των κειμένων πρέπει να γίνεται προσεκτικά για να μην οδηγήσουμε τον αλγόριθμο σε εσφαλμένα αποτελέσματα.

Κατά την πειραματική διαδικασία παρατηρήσαμε ότι η ακρίβεια δεν ήταν αρκετό μέτρο σύγκρισης, για αυτό τροποποιήσαμε τον εκτελέσιμο κώδικα του αλγορίθμου έτσι ώστε μετράμε την απόδοση σε ακρίβεια (precision) και ανάκλήση (recall) καθώς και την εμφάνιση των λάθος θετικών και λάθος αρνητικών κατηγοριοποιημένων κειμένων, για να έχουμε καλύτερη γενική εικόνα των αποτελεσμάτων.

Οί λόγοι επιλογής του Naive bayes κατηγοριοποιητή ήταν όχι μόνο λόγω της δυνατότητας παραμετροποίησης του αλλά και της αξιοπιστίας των αποτελεσμάτων του.

Κατά την πειραματική διαδικασία είδαμε ότι η ακρίβεια δεν ήταν αρκετό μέτρο σύγκρισης, για αυτό τροποποιήσαμε τον κώδικα εκτέλεσης του αλγορίθμου έτσι ώστε να μας εξάγει τα αποτελέσματα σε ακρίβεια και ανάκλήση καθώς και την εμφάνιση των λάθος θετικών και λάθος αρνητικών κατηγοριοποιημένων κειμένων, για να έχουμε καλύτερη γενική εικόνα των αποτελεσμάτων.

Για τόν λόγο οτι μπορούμε να επέμβουμε στην πειραματική διαδικασία και τα αποτελέσματα της είναι πιο αξιόπιστα, ο αλγόριθμος ταξινόμησης που επιλέχθηκε για το χτίσιμο της εφαρμογής ήταν ο Naïve Bayes με σκοπό να υπολογιστεί η υπο συνθήκη πιθανότητα ένα μήνυμα d να είναι μέλος μίας κλάσης c , όταν οι πιθανές κλάσεις ήταν οι «αρνητικό» και «θετικό».

4 ΚΕΦΑΛΑΙΟ : Γενική Αρχιτεκτονική

4.1 Έρευνα Τεχνολογιών

4.1.1 Ruby [1]

Ο Yukihiro Matsumoto δημοσίευσε την πρώτη έκδοση της Ruby το 1995, χρησιμοποίησε κομμάτια από τις αγαπημένες του γλώσσες (Perl, Smalltalk, Eiffel, Ada και Lisp). Το 1999 άρχισε να γίνεται γνωστή επειδή δημιουργήθηκε το documentation στα αγγλικά. Είναι γλώσσα ανοιχτού κώδικα (open source) και είναι απλή στην κατανόηση και στη συγγραφή κώδικα. Ο βασικότερος λόγος της σημερινής της επιτυχίας στο διαδίκτυο, είναι το framework Ruby on Rails που είναι



γραμμένο σε Ruby. Συνεχίζει να αναπτύσσεται λόγω της μεγάλης κοινότητας που έχει, και σήμερα βρίσκεται στην έκδοση 2.2.2. Ο δημιουργός της αναφέρει ότι η Ruby φτιάχτηκε για να κάνει ευτυχισμένους τους προγραμματιστές λόγω της απλότητας της. Κατά την συγγραφή κώδικα υπάρχουν λίγοι περιορισμοί και μπορούν να αναπτυχθούν μεγάλα και πολύπλοκα προγράμματα.

Το κυριότερο χαρακτηριστικό της είναι η εξ' ολοκλήρου αντικειμενοστρέφεια της, τα πάντα στην Ruby θεωρούνται αντικείμενα ακόμα και οι αριθμοί και όλες της οι κλάσεις. Επειδή είναι επηρεασμένη από την Lisp έχει στοιχεία αναδρομικών γλωσσών (functional). Είναι πολύ ευέλικτη επειδή είναι διερμηνευτική (interpreted) γλώσσα. Ένα παράδειγμα της ευελιξίας της είναι οι ανοιχτές κλάσεις, δηλαδή η προσθήκη μεθόδων σε υπάρχουσες κλάσεις την ώρα εκτέλεσης ενός προγράμματος. Αυτή η δυνατότητα, θεωρείται από κάποιους αδυναμία γιατί ενδέχεται να γίνει μη ηθελημένη αντικατάσταση κάποιας μεθόδου. Ένα ακόμα από τα πιο γνωστά της στοιχεία, είναι η δυνατότητα «μετα-προγραμματισμού» ή αλλιώς η συγγραφή κώδικα που γράφει κώδικα.

4.1.2 Ruby On Rails [2]

Η Ruby on Rails δημιουργήθηκε το 2004, από την 37signals και τον David Heinemeier Hanson για την ανάπτυξη των εφαρμογών ιστού τους, είναι ένα από τα πρώτα frameworks ιστού τα οποία ακολούθησαν την αρχιτεκτονική Μοντέλο-Προβολή-Ελεγκτής MVC (Model-View-Controller).



Εικόνα 4.1.2-1: Λογότυπο RoR

Η RoR βασίστηκε σε ισχυρές αρχές και θεμέλια τα οποία ακολουθούνται πιστά από την κοινότητα η οποία βοηθάει στην ανάπτυξη της. Οι βασικότερες αρχές είναι:

- Ευέλικτες τεχνικές υλοποίησης (agile development)
- Σύμβαση αντί για παραμετροποίηση (convention over configuration)
- Ανάπτυξη εφαρμογών οδηγούμενες από δοκιμές (tests-driven development)

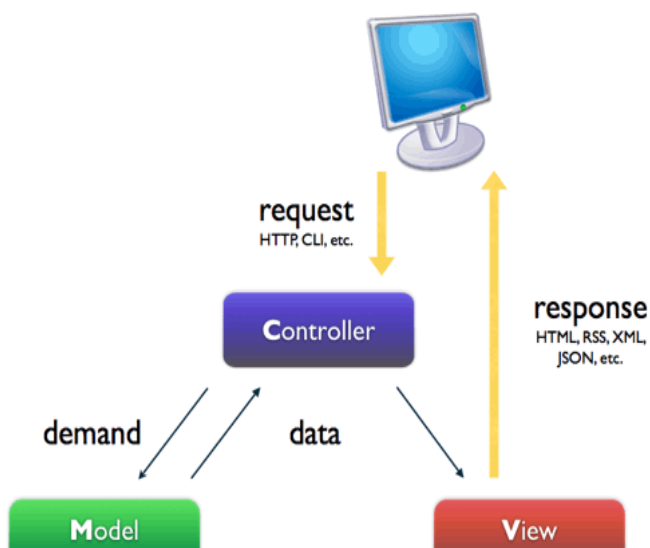
Στην Rails, το HTTP πρωτοκόλλο (GET, POST κλπ) σχετίζονται άμεσα με τα URLs και τις λειτουργίες

HTTP μέθοδοι	Διαδρομή	Λειτουργία βάσης	Τί κάνει
GET	/experiment	Read	Εμφανίζει μια λίστα πειραμάτων
GET	/experiment/new	-	Επιστρέφει μια HTML φόρμα για τη δημιουργία νέου πειραματος
POST	/experiment	Create	Δημιουργεί νέο πειραμα
GET	/experiment/:id	Read	Εμφανίζεισυγκεκριμένο πειραμα
GET	/experiment/:id/edit	-	Επιστρέφει μια HTML φόρμα για τη επεξεργασία συγκεκριμένου πειραματος
PUT	/experiment/:id	Update	Ενημερώνει τη βάση για τις αλλαγές στα δεδομένα συγκεκριμένου πειραματος
DELETE	/experiment/:id	Delete	Διαγράφει συγκεκριμένο πειραμα

Πίνακας 4.1.2-1: Αντιστοίχιση HTTP μεθόδων, Διαδρομών (paths), Λειτουργιών βάσης

4.1.2.1 Μοντέλο – Προβολή – Ελεγκτής

Πως εφαρμόζεται το μοντέλο (MVC) Μοντέλο – Προβολή – Ελεγκτής (Model – View – Controller) στο Ruby on Rails framework.



Εικόνα 4.1.2-2:Αναπαράσταση MVC αρχιτεκτονικής

Με αυτόν τον τρόπο η εφαρμογή διατηρείται καλά οργανωμένη χωρίζοντας τα δεδομένα (model) από τη λογική (controller) και την προβολή (view) του προγράμματος.

Τα μοντέλα περιέχουν τον κώδικα που ελέγχει και χειρίζεται τα δεδομένα ή ορίζει τη λογική. Τα μοντέλα είναι κλάσεις που επικοινωνούν με την βάση δεδομένων.

Οι ελεγκτές (controllers) απαντούν σε όλες τις αιτήσεις των χρηστών. Παίρνουν την είσοδο, την

διαχειρίζονται ανάλογα με το πώς έχουν προγραμματιστεί, καλούν μεθόδους, αλληλεπιδρούν με τα μοντέλα και προωθούν τα δεδομένα εξόδου στις προβολές (view). Οι προβολές εμφανίζουν τα δεδομένα εξόδου συνήθως σε HTML.

4.1.2.2 Δομή φακέλων της εφαρμογής RoR

Ένα από τα πλεονεκτήματα των εφαρμογών σε Ruby on Rails είναι ότι από προεπιλογή με τη δημιουργία της είναι οργανωμένη σε φακέλους. Οι εικόνες, το σχήμα της βάσης δεδομένων, ο κώδικας HTML, CSS και κάθε άλλο διαφορετικό κομμάτι της εφαρμογής, βρίσκονται στην δική τους προκαθορισμένη θέση. Αυτή η λειτουργία συμβάλει στην εύκολη συντήρηση και επέκταση του κώδικα. Στην παρακάτω εικόνα παρουσιάζονται οι σημαντικότεροι φακελοί και αναλύονται οι λειτουργίες τους.

app/ Ο φάκελος που χρησιμοποιείται περισσότερο από τους προγραμματιστές. Περιέχει τον κυρίως κώδικα της εφαρμογής (μοντέλο-προβολή-ελεγκτής).

config/ Οι επιπλέον ρυθμίσεις που απαιτούνται για την κάλυψη των αναγκών της εφαρμογής.

config.ru Αρχείο για την παραμετροποίηση της διεπαφής με τον Rack εξυπηρετητή.

db/ Σχήμα βάσης δεδομένων και πληροφορίες για τα migrations.

Gemfile Καθορίζονται τα gems τα οποία θα χρησιμοποιήσει η εφαρμογή. Τα gems είναι βιβλιοθήκες της Ruby.

lib/ Εδώ υπάρχει κώδικας που είτε δεν ανήκει καπου ή χρησιμοποιείται από περισσότερα από ένα εκ των μοντέλων, προβολών, ελεγκτών (Model-View-Controller).

log/ Όσο λειτουργεί μια Rails εφαρμογή δημιουργούνται καταγραφές (logs). Υπάρχουν τρεις φακελοί για την ανάπτυξη, τις δοκιμές και την κατάσταση που η εφαρμογή έχει δημοσιευτεί.

public/ Ο φάκελος αυτός είναι προσπελάσιμος από το διαδίκτυο και θεωρείται από τον εξυπηρετητή ο βασικός φάκελος της εφαρμογής και περιέχει στατικές σελίδες.

Rakefile Εδώ ορίζονται εργασίες που μπορεί να εκτελούν δοκιμές (tests), να δημιουργούν

τεκμηρίωση, να εκτυπώσουν το σχήμα της βάσης κ.α.

README Οδηγίες εγκατάστασης και χρήσης.

script/ Εδώ βρίσκονται τα scripts της RoR, τα οποία εκτελούνται αν στη κονσόλα γράψουμε την εντολή rails και το όνομα της συνάρτησης που θέλουμε να καλέσουμε π.χ *rails console*

test/ Σ' αυτό το φάκελο βρίσκονται όλες οι δοκιμές, όπως δοκιμές ενσωμάτωσης (integration tests), δοκιμές λειτουργικότητας (functional tests), δοκιμές επανάληψης (iteration tests).

tmp/ Ένας φάκελος για τα προσωρινά αρχεία, όπως τα περιεχόμενα της κρυφής μνήμης.

vendor/ Όλες σχεδόν οι εφαρμογές χρησιμοποιούν έτοιμο κώδικα ο οποίος προσθέτει λειτουργικότητα.



Εικόνα 4.1.2-3: Δομή φακέλων των εφαρμογών Ruby on Rails

4.1.2.3 Υποστήριξη τριών διαφορετικών περιβαλλόντων

Η Ruby on Rails υποστηρίζει τρία ανεξάρτητα περιβάλλοντα: ανάπτυξης, ελέγχων και παραγωγής. Οι ρυθμίσεις γίνονται μια φορά στην αρχική παραμετροποίηση και η μετάβαση από στάδιο σε στάδιο γίνεται με την αλλαγή μεταξύ περιβαλλόντων.

4.1.2.4 Active Record

Το Active Record είναι μια βιβλιοθήκη ORM της Ruby που επιτρέπει τη μεταφορά δεδομένων και συνεργάζεται με την υπόλοιπη εφαρμογή εξυπηρετώντας την καταχώρηση δεδομένων στη βάση, η οποία συνήθως είναι σχεσιακή. Το Active Record περιλαμβάνεται στη Ruby on Rails αλλά επίσης είναι διαθέσιμο και σαν Ruby Gem. Ακολουθώντας την αρχή των συμβάσεων αντί παραμετροποίησης ελαχιστοποιεί τις ρυθμίσεις που χρειάζεται να γίνουν, κάτι που δεν συναντάται συχνά στις υπόλοιπες ORM βιβλιοθήκες.

4.1.3 MySQL

Η βάση δεδομένων MySQL είναι μια απ' τις πιο δημοφιλείς βάσεις δεδομένων. Είναι ανοιχτού κώδικα, σχεσιακή, αξιόπιστη, εύκολη στη χρήση, υψηλών επιδόσεων συμβατή με τις περισσότερες πλατφόρμες όπως τα Linux, Mac OS, Windows κ.α.

4.1.3.1 Σχεσιακή βάση δεδομένων (Relational database)

Με τον όρο σχεσιακή βάση δεδομένων εννοείται μία συλλογή δεδομένων οργανωμένη σε συσχετισμένους πίνακες που παρέχει ταυτόχρονα ένα μηχανισμό για ανάγνωση, εγγραφή, τροποποίηση ή και πιο πολύπλοκες διαδικασίες πάνω στα δεδομένα. Ο σκοπός μιας βάσης δεδομένων είναι η οργανωμένη αποθήκευση πληροφορίας και η δυνατότητα εξαγωγής της πληροφορίας αυτής, ιδίως σε πιο οργανωμένη μορφή, σύμφωνα με ερωτήματα που τίθενται στη σχεσιακή βάση δεδομένων. Τα δεδομένα είναι δυνατόν να αναδιοργανώνονται με πολλούς διαφορετικούς τρόπους, σε νοητούς πίνακες, χωρίς να είναι απαραίτητη η αναδιοργάνωση των φυσικών πινάκων που τα αποθηκεύουν.

Επίσης, σε κάθε πίνακα υπάρχει μια στήλη το περιεχόμενο της οποίας είναι μοναδικό για κάθε γραμμή και ονομάζεται πρωτεύον κλειδί (primary key). Σύμφωνα με τις συμβάσεις της Rails, το πρωτεύον κλειδί είναι ακέραιος αριθμός και ονομάζεται id. Το σύνολο των πινάκων και η δομή τους αποθηκεύονται στο σχήμα της βάσης (database schema).

4.1.3.2 Αντικείμενο-σχεσιακή χαρτογράφηση (ORM – Object Relational Mapping)

Οι βιβλιοθήκες αντικειμενοσχεσιακής χαρτογράφησης είναι εργαλεία που επιτρέπουν την εύκολη και αυτοματοποιημένη αποθήκευση εγγραφών σε μια σχεσιακή βάση δεδομένων. Η λειτουργία τους είναι να κάνουν την αντιστοίχιση πινάκων με κλάσεις, γραμμών με αντικείμενα και στηλών με ορίσματα. Οι μέθοδοι των κλάσεων (class methods) επιδρούν σε πίνακες ενώ οι μέθοδοι των στιγμιοτύπων (instance methods) επιδρούν σε κάποια γραμμή του πίνακα.

Αποτέλεσμα της χρήσης ORM σε μια εφαρμογή είναι η μείωση του χρόνου ανάπτυξης του λογισμικού, του κόστους ανάπτυξης και συντήρησης, η σύνταξη απλούστερου και λιγότερου κώδικα αντί πολύπλοκων SQL ερωτημάτων.

4.1.4 Git

Μια εφαρμογή κατά την ανάπτυξή της χωρίζεται από

τον προγραμματιστή ή την ομάδα ανάπτυξης της σε

στάδια – κατηγορίες. Το **Git** είναι ένα σύστημα ελέγχου διανεμόμενης έκδοσης και διαχείρισης κώδικα (SCM) με έμφαση στην ταχύτητα, στην ακεραιότητα των δεδομένων και στην υποστήριξη για κατανεμόμενες μη γραμμικές ροές εργασίας. Το Git σχεδιάστηκε και αναπτύχθηκε αρχικά από τον Λίνους Τορβαλντς για τη ανάπτυξη του πυρήνα του Λίνουξ το 2005 και έχει γίνει από τότε το πιο πλατιά διαδεδομένο σύστημα ελέγχου εκδόσεων για ανάπτυξη λογισμικού.



Εικόνα 4.1.3-1: Λογότυπο Git

Όπως τα περισσότερα άλλα διανεμόμενα συστήματα ελέγχου εκδόσεων αναθεώρησης και αντίθετα με τα περισσότερα συστήματα πελάτη-διακομιστή κάθε κατάλογος εργασίας του Git είναι ένα ολοκληρωμένο αποθετήριο με πλήρες ιστορικό και δυνατότητες πλήρους παρακολούθησης της έκδοσης, ανεξάρτητα από την πρόσβαση δικτύου ή ενός κεντρικού διακομιστή.

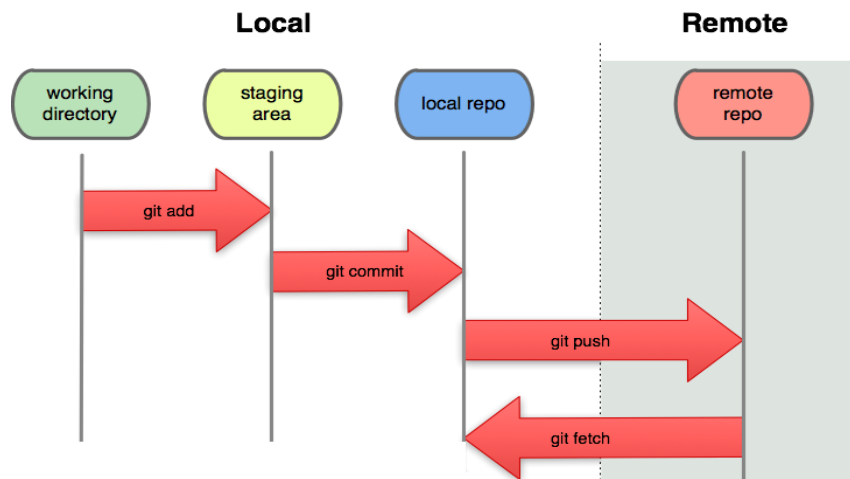
Οι προγραμματιστές να μπορούν να ανακαλέσει τις συγκεκριμένες εκδόσεις αργότερα. Από τη γέννησή του το 2005, το Git έχει εξελιχθεί, για να είναι εύκολο στη χρήση αλλά και να διατηρεί τις αρχικές του ιδιότητες. Είναι γρήγορο, πολύ αποτελεσματικό με τα μεγάλα έργα, και έχει ένα σύστημα διακλάδωσης για μη-γραμμική ανάπτυξη και χρησιμοποιείται από μεγάλες εταιρείες και μεγάλα έργα όπως οι Google, Microsoft, gnome, linux, twitter, facebook, LinkedIn, PostgreSQL, android, eclipse.

Η εγκατάσταση είναι εύκολη σε κάθε σύστημα εκτελώντας:

- ✦ `yum install git-core` (CentOS, fedora)
- ✦ `apt-get install git-core` (debian, ubuntu)
- ✦ `brew install git` (OS X)

Στα Windows, όπως ένα κοινό πρόγραμμα, μεταφορτώνεται και εγκαθίσταται το msysGit package.

Τις περισσότερες φορές υπάρχουν ένα ή περισσότερα απομακρυσμένα αποθετήρια (remote repositories), εκτός από το τοπικό (local repository). Για να ενημερωθεί ένα απομακρυσμένο αποθετήριο για τις τελευταίες εκδόσεις εκτελείται η εντολή `git push`. Και αντίστροφα για την ενημέρωση του τοπικού αποθετηρίου από το απομακρυσμένο, εκτελείται η εντολή `git fetch`. Όμως με την εντολή `git fetch` ενημερώνεται το τοπικό αποθετήριο χωρίς να ενσωματώνει τις αλλαγές. Για την ενσωμάτωση η κατάλληλη εντολή είναι `git merge`.



Εικόνα 4.1.4-2: Ενημέρωση τοπικού και απομακρυσμένου repository

Όταν θέλουμε να πάρουμε ένα αντίγραφο τότε μπορούμε να κλωνοποιήσουμε ένα απομακρυσμένο αποθετήριο με την εξής εντολή: `git clone https://github.com`

Το git, επίσης υποστηρίζει διακλαδώσεις (branches). Συνηθίζεται να υπάρχει μια κύρια διακλάδωση (master) η οποία περιλαμβάνει μόνο το κώδικα που προορίζεται για την τελική έκδοση της εφαρμογής κι άλλες διακλαδώσεις στις οποίες μπορεί να δοκιμάζεται κώδικας κι αν πληρεί τις προϋποθέσεις στη να προστεθεί στη κύρια διακλάδωση.



Εικόνα 4.1.4-3: Διακλαδώσεις στο Git

4.2 Επιλεγμένες Τεχνολογίες

4.2.1 Βάση δεδομένων – PostgreSQL

Η PostgreSQL αποτελεί μια ανοιχτού κώδικα σχεσιακή βάση δεδομένων. Η ανάπτυξη της ήδη διαρκεί πάνω από 20 χρόνια και βασίζεται σε μια αποδεδειγμένα καλή αρχιτεκτονική η οποία έχει δημιουργήσει μια ισχυρή αντίληψη των χρηστών της γύρω από την αξιοπιστία, την ακεραιότητα δεδομένων και την ορθή λειτουργία.

Η PostgreSQL τρέχει σε όλα τα βασικά λειτουργικά συστήματα, περιλαμβάνοντας **Linux**, **UNIX** (AIX, BSD, Mac OS X, Solaris), και **Windows**.



Εικόνα 4.2.1-1: Λογότυπο της PostgreSQL

Μερικές Γενικές οριακές Τιμές συμπεριλαμβάνονται στον παρακάτω πίνακα:

Limit	Value
Maximum Database Size	Unlimited
Maximum Table Size	32 TB
Maximum Row Size	1.6 TB
Maximum Field Size	1 GB
Maximum Rows per Table	Unlimited
Maximum Columns per Table	250 - 1600 depending on column types
Maximum Indexes per Table	Unlimited

Πίνακας 4.2.1-1: Οριακές τιμές ενεργών εγκαταστάσεων PostgreSQL

4.2.2 Twitter Bootstrap

Το Bootstrap είναι μια συλλογή εργαλείων ανοιχτού κώδικα (Ελεύθερο λογισμικό) για τη δημιουργία ιστοσελίδων και διαδικτυακών εφαρμογών. Περιέχει HTML και CSS για τις μορφές τυπογραφίας, κουμπιά πλοήγησης και άλλων στοιχείων του περιβάλλοντος, καθώς και προαιρετικές επεκτάσεις JavaScript. Έχει το πιο δημοφιλές πρόγραμμα στο GitHub και έχει χρησιμοποιηθεί από τη NASA και το MSNBC, μεταξύ άλλων.



Εικόνα 4.2.2-1: Λογότυπο Bootstrap

Το Bootstrap έχει σχετικά ελλιπή υποστήριξη για HTML5 και CSS, αλλά είναι συμβατό με όλους τους φυλλομετρητές (browsers). Βασικές πληροφορίες συμβατότητας των ιστοσελίδων ή εφαρμογές είναι διαθέσιμες για όλες τις συσκευές και τα προγράμματα περιήγησης.

Το Bootstrap είναι σπονδυλωτό και αποτελείται ουσιαστικά από μια σειρά στυλ(stylsheets) που εφαρμόζουν τα διάφορα συστατικά του πακέτου εργαλείων. Ένα στυλ που ονομάζεται bootstrap.less περιλαμβάνει τα συστατικά stylesheets. Οι προγραμματιστές μπορούν να προσαρμόσουν το αρχείο Bootstrap, επιλέγοντας τα στοιχεία που θέλουν να χρησιμοποιήσουν στο έργο τους.

Προσαρμογές είναι δυνατές σε περιορισμένη έκταση μέσω ενός κεντρικού στυλ διαμόρφωσης. Η χρήση γλώσσας στυλ επιτρέπει τη χρήση για μεταβλητές, λειτουργίες και φορείς (operators), ένθετους επιλογείς, γνωστά και ως μείγματα mixin.

Από την έκδοση 2.0, η διαμόρφωση του Bootstrap έχει επίσης μία ειδική επιλογή "Προσαρμογή" στην τεκμηρίωση (documentation). Επιπλέον, ο σχεδιαστής του έργου επιλέγει σε μια φόρμα τα επιθυμητά συστατικά και τα προσαρμόζει, εάν είναι αναγκαίο, σε τιμές διαφόρων εναλλακτικών λύσεων για τις ανάγκες του. Στη συνέχεια δημιουργείται ένα πακέτο που περιλαμβάνει ήδη το προχτισμένο CSS στυλ.

Εκτός από τα βασικά HTML στοιχεία, το Bootstrap περιέχει και άλλα στοιχεία περιβάλλοντος που χρησιμοποιούνται συχνά. Αυτά περιλαμβάνουν κουμπιά με προηγμένα χαρακτηριστικά (π.χ. ομαδοποίηση κουμπιών ή drop-down επιλογή, οριζόντιες και κάθετες καρτέλες, πλοήγηση, σελιδοποίηση, κ.λπ.), ετικέτες, προηγμένες τυπογραφικές δυνατότητες, εικονίδια, προειδοποιητικά μηνύματα και μια γραμμή προόδου.

4.2.3 Highcharts

Το Highcharts είναι μια βιβλιοθήκη διαγραμμάτων γραμμένη σε καθαρή JavaScript, προσφέροντας στον προγραμματιστή έναν εύκολο τρόπο να προσθέσει διδραστικά διαγράμματα στην ιστοσελίδα ή στην διαδικτυακή



Εικόνα 4.2.3-1: Λογότυπο

εφαρμογή του. Το Highcharts υποστηρίζει τους εξής τύπους διαγραμμάτων : line, spline, area, areaspline, column, bar, pie, Angular gauges polar chart types.

Οι χρήστες έχουν την δυνατότητα να εξαγάγουν το γράφημα απευθείας από την ιστοσελίδα. Είναι πλήρως συμβατή με όλους τους περιηγητές συμπεριλαμβάνοντας και τα iPhone/iPad και τον Internet Explorer από την έκδοση 6.

4.2.4 RVM [18]

Το RVM (Ruby Version Manager) δημιουργήθηκε από τον Seguin Wayne με σκοπό την ταυτόχρονη εγκατάσταση πολλών εκδόσεων της γλώσσας Ruby. Επίσης υποστηρίζει την καλύτερη οργάνωση των βιβλιοθηκών σε gemsets. Κάθε gemset αποτελείται από ένα σύνολο βιβλιοθηκών (gems) της Ruby. Ανάλογα με την εφαρμογή χρησιμοποιείται η

κατάλληλη έκδοση βιβλιοθηκών και οργανώνονται καλύτερα οι αλληλεξαρτήσεις των βιβλιοθηκών.

4.2.5 Bundler

Ο Bundler αναπτύχθηκε από έναν από τους πρωτοπόρους της βασικής ομάδας της Rails και της jQuery, τον Yehuda Katz. Ελέγχει τις αλληλεξαρτήσεις μεταξύ των βιβλιοθηκών της Ruby και φροντίζει ώστε να γίνεται με αυτοματοποιημένο τρόπο η μεταφόρτωση όπως και η εγκατάσταση όλων των βιβλιοθηκών μιας εφαρμογής σε διαφορετικά συστήματα.



Εικόνα 4.2.4-1: : Λογότυπο του RVM

4.2.6 Λόγοι χρήσης επιμέρους τεχνολογιών

Ο *κάθε προγραμματιστής επιλέγει τις τεχνολογίες και τα εργαλεία που κάνουν πιο εύκολη την δουλειά του. Σε αυτό το σημείο αναφέρουμε τους λόγους που κάνουν την Ruby on Rails ένα ισχυρό web framework γραμμένο σε Ruby.*

Το γεγονός ότι το Ruby on Rails framework είναι καινούργιο όμως με τεράστια κοινότητα που το βοηθάει στην ανάπτυξη αποδείχτηκε εξαιρετική επιλογή καθώς έδωσε τη δυνατότητα στους προγραμματιστές να μάθουν τα βασικά χαρακτηριστικά σε πολύ μικρό χρονικό διάστημα και να μπορέσουν να δημιουργίσουν τα βασικά στοιχεία της εφαρμογής.

Κατά τη διάρκεια δημιουργίας της εφαρμογής αποδείχτηκε μεγάλη βοήθεια η τεράστια ποικιλία σε βιβλιοθήκες Gems που βοηθούν τον

προγραμματιστή στην ανάπτυξη σημαντικών χαρακτηριστικών του συστήματος.

Όσον αφορά τη βάση δεδομένων στα αρχικά στάδια, χρησιμοποιήθηκε η πιο διαδεδομένη σχεσιακή βάση δεδομένων, η MySQL. Όμως όταν έφτασε η στιγμή η εφαρμογή να περάσει στο στάδιο παραγωγής (development), η επιλογή της χρήσης του “υπολογιστικού νέφους” προέκυψε η ανάγκη για τη μετάβαση στη βάση δεδομένων PostgreSQL.

Αυτή η αλλαγή ήταν η καλύτερη επιλογή όσον αφορά το κομμάτι της αποθήκευσης των δεδομένων. Η εύκολη και ταχύτατη διαχείριση από τους προγραμματιστές, η σταθερότητα και αξιοπιστία μέσω του Heroku αλλά και η δυνατότητα επεκτασιμότητας είναι μερικοί από τους λόγους που κάνουν την PostgreSQL άριστη επιλογή πέρα από αναγκαία.

Λόγω του όγκου εργασίας που είχε η εφαρμογή απο μόνη της έπρεπε να βρεθεί μια ευκολη λύση

για τη δημιουργία ενός εύχρηστου και όμορφου περιβάλλοντος. Ύστερα από αρκετή μελέτη βρέθηκε η βιβλιοθήκη Twitter Bootstrap, η οποία μπορεί να βοηθήσει τον αρχάριο σχεδιστή εφαρμογών διαδικτύου να αναπτύξει με λίγη προσπάθεια ένα καλαίσθητο αποτέλεσμα. Το Twitter Bootstrap προσφέρει μια συλλογή από εργαλεία που χρησιμοποιούν τεχνολογίες όπως HTML, CSS και JavaScript.

Ενα ακόμα χαρακτηριστικό της εφαρμογής είναι τα διαγράμματα. Μέσα από ένα διάγραμμα μπορεί ο χρήστης να κατανοήσει πώς απέδωσε ο αλγόριθμος σε σύγκριση με τις δικές του βαθμολογίες. Για την απεικόνιση των διαγραμμάτων χρησιμοποιήθηκε η βιβλιοθήκη Highcharts προσφέροντας γρήγορη απόδοση και όμορφο αποτέλεσμα. Η βιβλιοθήκη Highcharts παίρνει τις πληροφορίες που χρειάζεται από τη βάση δεδομένων.

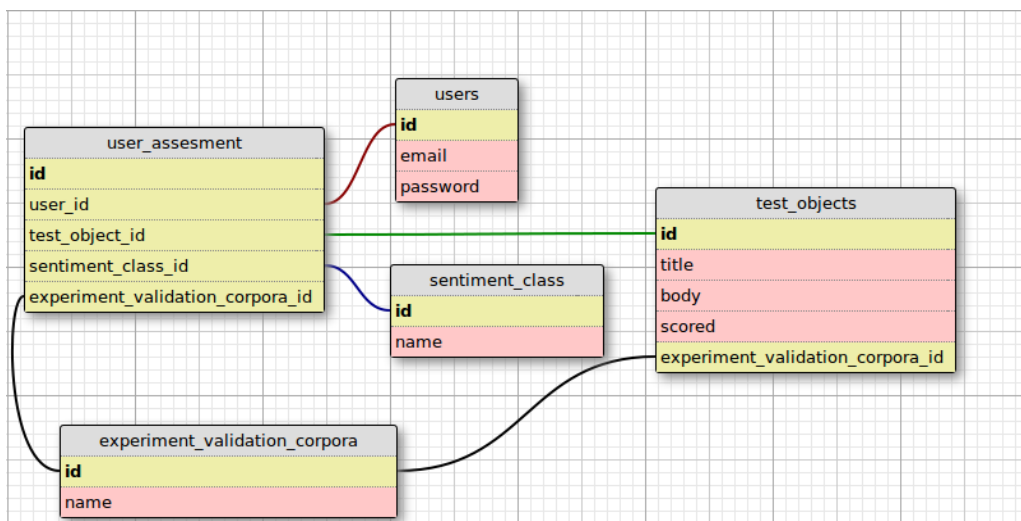
5 ΚΕΦΑΛΑΙΟ : Η Εφαρμογή αξιολόγησης Analyment

5.1 Η δημιουργία του σκελετού της εφαρμογής

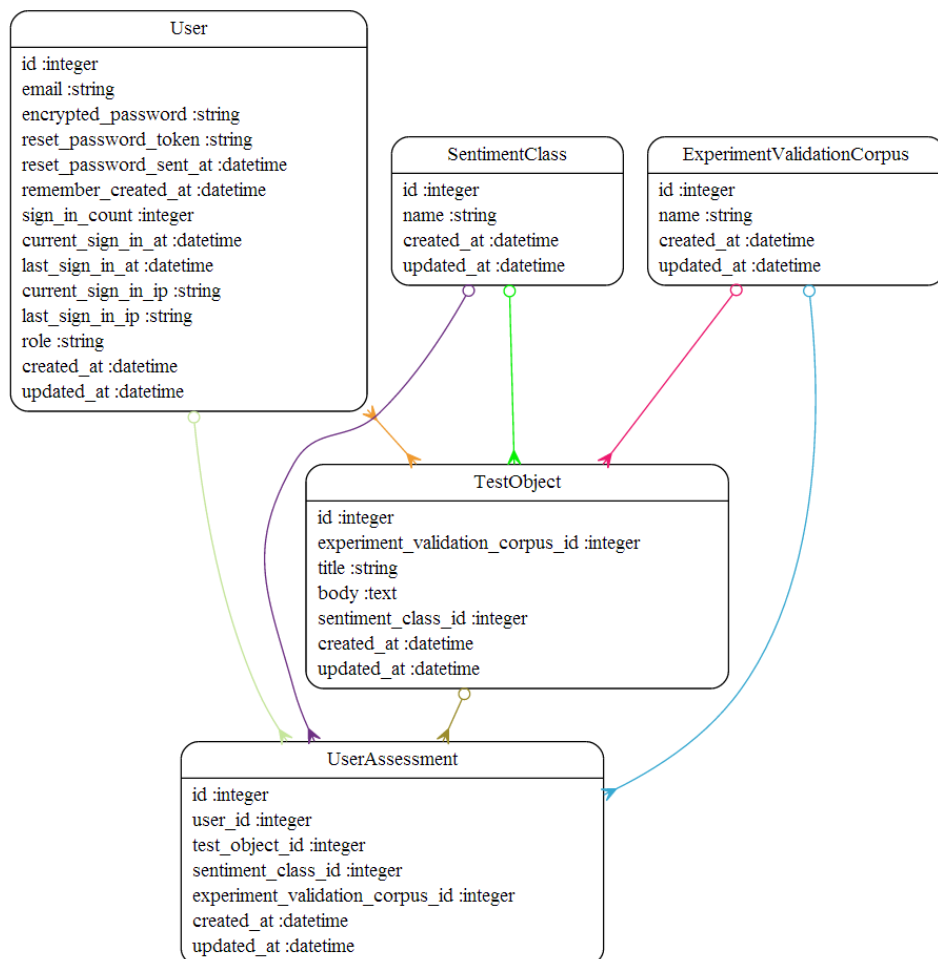
Η ιδέα για τη δημιουργία της εφαρμογής η οποία στοχεύει στην αξιολόγηση των αλγορίθμων που κάνουν συναισθηματική ανάλυση προήλθε από το γεγονός ότι ο ανθρώπινη νοϋμοσύνη είναι ακόμα σε πολύ υψηλότερα επίπεδα από αυτά των μηχανών. Έτσι ο άνθρωπος καλείται να βαθμολογήσει

κείμενα τα οποία έχουν βαθμολογηθεί απο έναν αλγόριθμο πρώτα.

Ο κορμός της εφαρμογής είναι η βάση δεδομένων της, δηλαδή οι πίνακες που περιέχει. Αρχικά δημιουργήθηκαν οι πίνακες: Πίνακας των κατηγοριών, της συναισθηματικής κλάσης, κειμένου δοκιμής , και ανάθεσης χρήστη. Ο χρήστης έχει τη δυνατότητα να επιλέγει την κατηγορία του πειράματος και να αξιολογήσει τα κείμενα της κατηγορίας αυτής.



Εικόνα 4.2.6-1: Το αρχικό Schema βάσης δεδομένων



Εικόνα 4.2.6-3: Schema βάσης δεδομένων προγραμματιστικά

Στο δεύτερο βήμα δημιουργήθηκαν οι αλγόριθμοι για να υπολογίζουν τις πληροφορίες που εισάγει ο χρήστης. Οι πληροφορίες αυτές είναι οι εξής:

- Το σύνολο κειμένων που αξιολόγησε ο χρήστης
- Σύγκριση των αποτελεσμάτων
- Τελικά αποτελέσματα με διαγράμματα

Τέλος δημιουργήθηκε το περιβάλλον διεπαφής των χρηστών. Όλα τα παραπάνω προεκτάθηκαν για πολλούς χρήστες μέσω του πίνακα χρηστών και της διαδικασίας διαπίστευσης χρηστών. Επίσης δεν δώθηκαν δικαιώματα χρήσης των παραπάνω στους μη συνδεδεμένους χρήστες.

5.2 Τα αρχεία που απαρτίζεται η εφαρμογή

Τα gems που χρησιμοποιήθηκαν είναι τα παρακάτω:

```
source 'https://rubygems.org'  
# Bundle edge Rails instead: gem 'rails', github: 'rails/rails'  
gem 'rails', '4.0.2'  
# Use PostgreSQL as the database  
gem 'mysql2'  
# Application server  
gem 'puma'  
# Use SCSS for stylesheets  
gem 'sass-rails', '~> 4.0.0'  
# Use Uglifier as compressor for JavaScript assets  
gem 'uglifier', '>= 1.3.0'  
gem "therubyracer"  
gem "less-rails" #Sprockets (what Rails 3.1 uses for its asset pipeline) supports LESS  
gem "twitter-bootstrap-rails"  
# Imports csv files  
gem 'roo'  
# Use CoffeeScript for .js.coffee assets and views  
gem 'coffee-rails', '~> 4.0.0'  
# See https://github.com/sstephenson/execjs#readme for more supported runtimes  
# gem 'therubyracer', platforms: :ruby  
# Use jquery as the JavaScript library  
gem 'jquery-rails'  
# Turbolinks makes following links in your web application faster. Read  
more: https://github.com/rails/turbolinks  
gem 'turbolinks'  
# Build JSON APIs with ease. Read more: https://github.com/rails/jbuilder  
gem 'jbuilder', '~> 1.2'
```

```
# User authentication
gem 'devise'
gem 'bcrypt-ruby', '~> 3.1.2'
#User authorization
gem 'cancan'
# Pagination
gem 'kaminari'
group :doc do
  # bundle exec rake doc:rails generates the API under doc/api.
  gem 'sdoc', require: false
end
group :development do
  # Debugging tool
  gem 'pry'
end
# Use ActiveRecord has_secure_password
# gem 'bcrypt-ruby', '~> 3.1.2'

# Use unicorn as the app server
# gem 'unicorn'

# Use Capistrano for deployment
# gem 'capistrano', group: :development

# Use debugger
# gem 'debugger', group: [:development, :test]
ruby "2.0.0"
```

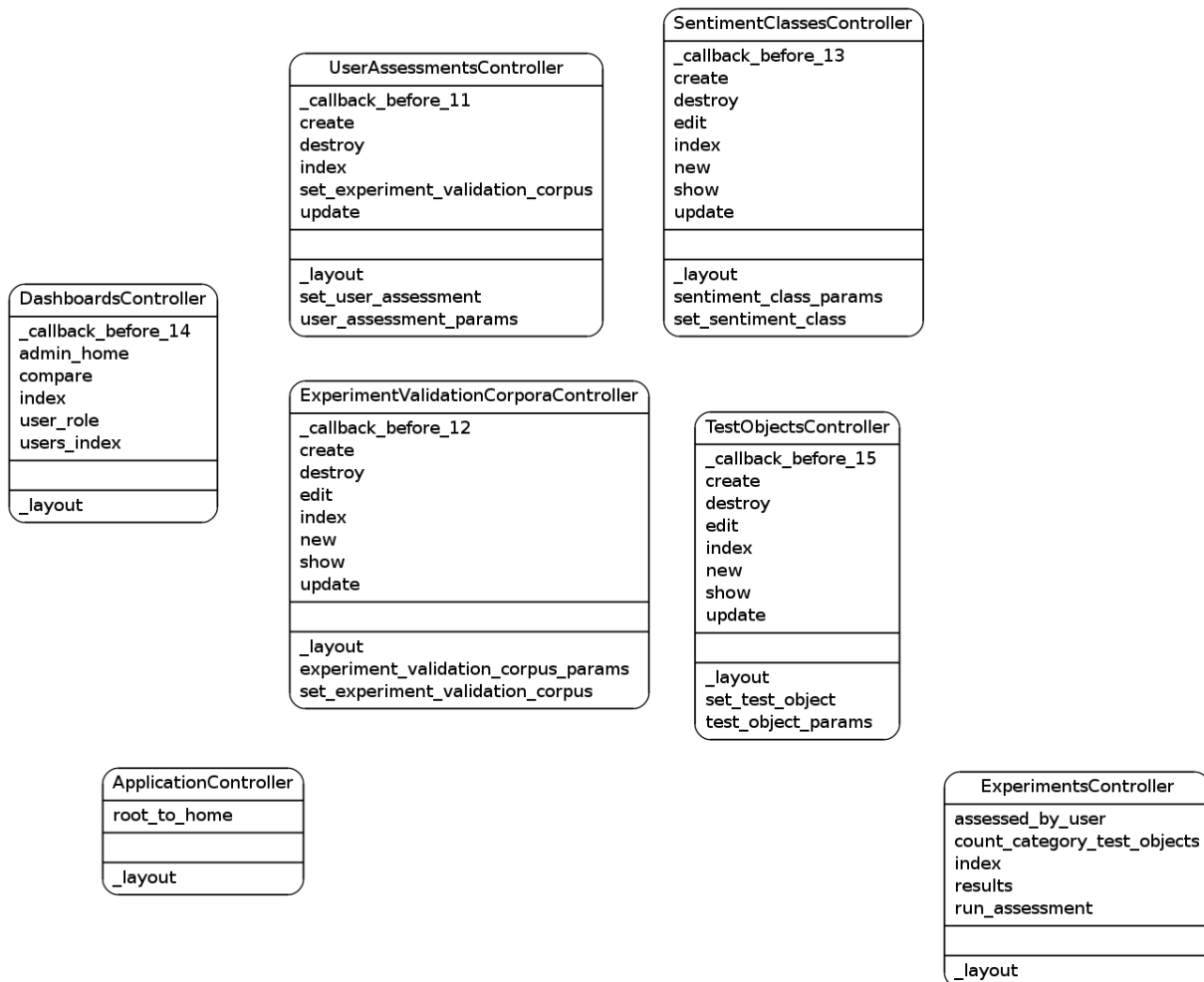
Στην παρακάτω εικόνα φαίνονται τα αρχεία της εφαρμογής:



Εικόνα 4.2.6-1: Δομή φακέλων της εφαρμογής *Analyment*

5.3 Πώς σχεδιάστηκε η εφαρμογή

Αφού οι controller αποτελούν την το κέντρο της εφαρμογής μίας και συνδέουν όλες τις λειτουργίες όπως τη βάση με τους προβολείς παραθέτουμε την εικόνα με τον σχεδιασμό τους. Το διάγραμμα κλάσεων για τους controller δημιουργήθηκε χρησιμοποιώντας μόνο μία βιβλιοθήκη της Ruby η οποία μόλις την εκτελέσεις παράγει την εικόνα.



Εικόνα 4.2.6-1: Δομή και σύνδεση των controllers

5.4 Λειτουργίες χρηστών

Η αρχική σελίδα που εμφανίζεται σε εγγεγραμένους και μη εγγεγραμμένους, και τους δίνει την δυνατότητα εγγραφής (sign up) καθώς και η δυνατότητα σύνδεσης (sign in) στο σύστημα από την μπάρα που βρίσκεται στην κορυφή της εφαρμογής.

Analymment Sign In

Analymment

Sentiment Evaluation.

Signed out successfully.

Analymment by N.T - © TeiHer 2015

5.4.1 Δημιουργία – σύνδεση χρήστη

Ολόκληρη η εφαρμογή είναι σχεδιασμένη και φτιαγμένη να είναι όσο πιο εύκολη στην χρήση της. Ο απλός χρήστης μπορεί να ανοίξει πολύ εύκολα λογαριασμό στην εφαρμογή Analyment. Στην κορυφή της αρχικής σελίδας επιλέγοντας το κουμπί sign up εμφανίζεται η φόρμα συμπλήρωσης με τρία υποχρεωτικά πεδία:

- E-mail. Υποχρεωτικά ένα έγκυρο e-mail το οποίο θα χρησιμοποιείται για τη σύνδεση του χρήστη.
- Password (κωδικός). Ο μυστικός κωδικός τουλάχιστον 8 χαρακτήρων που θα χρησιμοποιηθεί σε συνδυασμό με το e-mail για τη σύνδεση του χρήστη στο σύστημα.
- Password Confirmation (επιβεβαίωση κωδικού). Για την επανάληψη του κωδικό του για να αποφευχθούν τυχόν λάθη.



Analyment

Sign in

Email

Password

Remember me

Sign in

[Sign up](#)

[Forgot your password?](#)

Analyment by N.T - © TeiHer 2015



Sign up

Email

Password

Password confirmation

Sign up

[Sign in](#)

Analymment by N.T - © TeiHer 2015

Μόλις κάποιος χρήστης δημιουργήσει τον λογαριασμό του μπορεί να συνδεθεί στην εφαρμογή Analymnet μόνο με το e-mail και τον κωδικό του.

Επιλέγοντας το Remember_me αποθηκεύεται στην σελίδα και δεν εισάγει κάθε φορά τα στοιχεία του. Υπάρχει και η επιλογή του forgot your password όπου ο χρήστης έχει τη δυνατότητα να ανακτήσει τον κώδικό του εισάγοντας μόνο το email του.



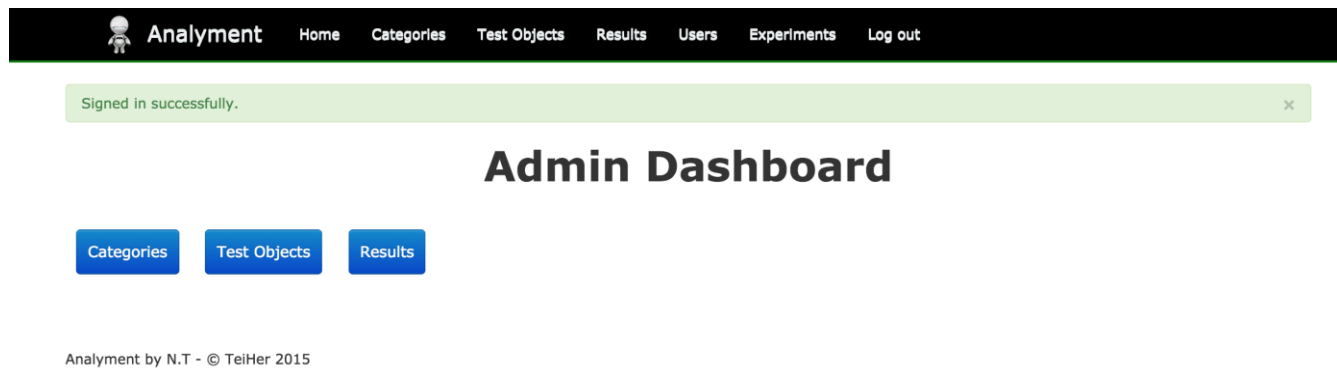
Εικόνα 5.4.1-3: Λειτουργία ανάκτησης κωδικού

5.4.2 Λειτουργίες διαχειριστών - απλών χρηστών

Στην εφαρμογή Analymnet όπως και σε όλες τις διαδικτυακές εφαρμογές υπάρχουν δύο είδη χρηστών, απλός χρήστης και ο διαχειριστής. Εδώ ο απλός χρήστης τα μόνα δικαιώματα που έχει είναι αυτό της αξιολόγησης και της προβολής των αποτελεσμάτων αφού έχει τελειώσει ολόκληρη την κατηγορία με τα κείμενα. Ο διαχειριστής από την άλλη έχει όλα τα δικαιώματα ακόμα και να κάνει τον απλό χρήστη διαχειριστή και το αντίθετο.

Αρχική σελίδα διαχειριστή

Η πρώτη σελίδα περιέχει μόνο ένα γρήγορο μενού που οδηγεί σε όλες τις λειτουργίες, και υπάρχει σε όλες τις σελίδες για να κάνει την πλοήγηση πιο εύκολη.



Εικόνα 5.4.2-1: Αρχική σελίδα admin

Categories-Κατηγορίες

Στο κουμπί Categories φαίνονται όλες οι κατηγορίες που υπάρχουν για “βαθμολόγηση”. Ακόμα δίνεται η δυνατότητα να επεξεργαστούν να δημιουργηθούν νέες αλλά και να διαγραφούν μαζί με όλα τα κείμενα τους.

The screenshot shows the 'Analymnt' application interface. At the top, there is a navigation bar with the following items: Home, Categories, Test Objects, Results, Users, Experiments, and Log out. The main heading is 'Listing Categories'. Below this, there is a table with the following structure:

Name			
Movies	Show	Edit	Destroy
Quotes	Show	Edit	Destroy

Below the table, there is a green button labeled 'New Category'. At the bottom left, there is a small text: 'Analymnt by N.T - © TeiHer 2015'.

Εικόνα 5.4.2-2: Όλες οι κατηγορίες των πειραμάτων

Αντικείμενα προς βαθμολόγηση

Σε αυτό το σημείο είναι συγκεντρωμένα όλα τα κείμενα που είναι προς κατηγοριοποίηση. Εδώ μπορεί ο διαχειριστής να εισάγει νέα κείμενα είτε μαζικά από ένα csv αρχείο είτε ένα-ένα χειροκίνητα. Ακόμα μπορεί το κάθε κείμενο να επεξεργαστεί και να διαγραφεί. Όπως είναι φυσικό τα κείμενα μίας κατηγορίας μπορεί να είναι πολλά και να μην αρκεί μια σελίδα να εμφανιστούν όλα, για αυτό προστεθηκε σελιδοποίηση στην εφαρμογή για να τα ταξινομήι ανά ομάδες κειμένων.



Listing Test Objects

Category	Title	Body	Classifier Result			
Quotes	Be strong	Work hard for what you want because it won't come to you without a fight. You...	positive	Show	Edit	Destroy
Quotes	Sun	Keep your face to the sunshine and you cannot see a shadow.	positive	Show	Edit	Destroy
Movies	Bad movie	This is perhaps the most appalling piece of art (lol) ever committed to...	negative	Show	Edit	Destroy
Movies	Worst film ever	Terrible over-the-top acting, poor filming/photography, lame storyline and...	negative	Show	Edit	Destroy
Quotes	Joy	Find a place inside where there's joy, and the joy will burn out the pain.	positive	Show	Edit	Destroy

1 2 3 ... Next > Last »

New Test object

Επιλογή αρχείου

Δεν επιλέχθηκε κανένα αρχείο.

Import

Analymet by N.T - © TeiHer 2015

Εικόνα 5.4.2-3: Σελίδα διαχείρισης των Κειμένων

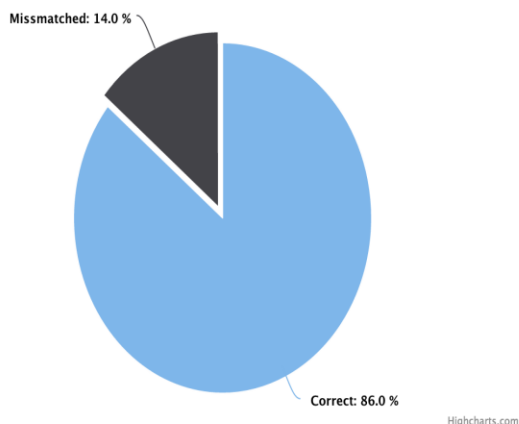
Compare Results

Η σελίδα με τα συγκεντρωτικά αποτελέσματα των χρηστών. Μπορεί να την δει μόνο ο διαχειριστής καθώς οι χρήστες μπορούν να δούν μόνο τα δικά τους.



Results - Category: Quotes

Comparison between algorithm and user results ☰



Test Object	Algorithm Result	Your Result
Be strong	Positive	Positive
Sun	Positive	Positive
Joy	Positive	Positive
Positive side of things	Positive	Positive
Positive anything	Positive	Positive
Instincts	Positive	Positive
Worth living	Positive	Positive
The life to come	Positive	Positive
Focus on	Positive	Positive
Positive events	Positive	Positive

Εικόνα 5.4.2-4: Σύγκριση αποτελεσμάτων ανά κατηγορία

Αν επιλέξει μία κατηγορία του εμφανίζονται σε κάθε σελίδα τα κείμενά της καθώς επίσης την βαθμολογία του αλγορίθμου, τα email των χρηστών και πώς το βαθμολόγησε ο κάθε χρήστης το συγκεκριμένο κείμενο.

Anlyment Home Categories Test Objects Results Users Experiments Log out

Compare Results - Movies

Test Object: Bad movie
Algorithm Assessment: Negative

User email	User assessment
trullis@gmail.com	Negative
nt@hotmail.com	Positive

1 2 3 ... Next > Last >

Anlyment by N.T - © TeiHer 2015

Εικόνα 5.4.2-5: Σελίδα σύγκρισης αποτελεσμάτων ανά Test Object

Χρήστες

Σε περίπτωση που υπάρξουν πολλοί χρήστες προστέθηκε η λειτουργία στην εφαρμογή όλοι οι χρήστες να μπορούν να γίνουν διαχειριστές καθώς και το αντίστροφο, δηλαδή ένας διαχειριστής να μπορεί να κάνει τον εαυτό του απλό χρήστη.

Anlyment Home Categories Test Objects Results Users Experiments Log out

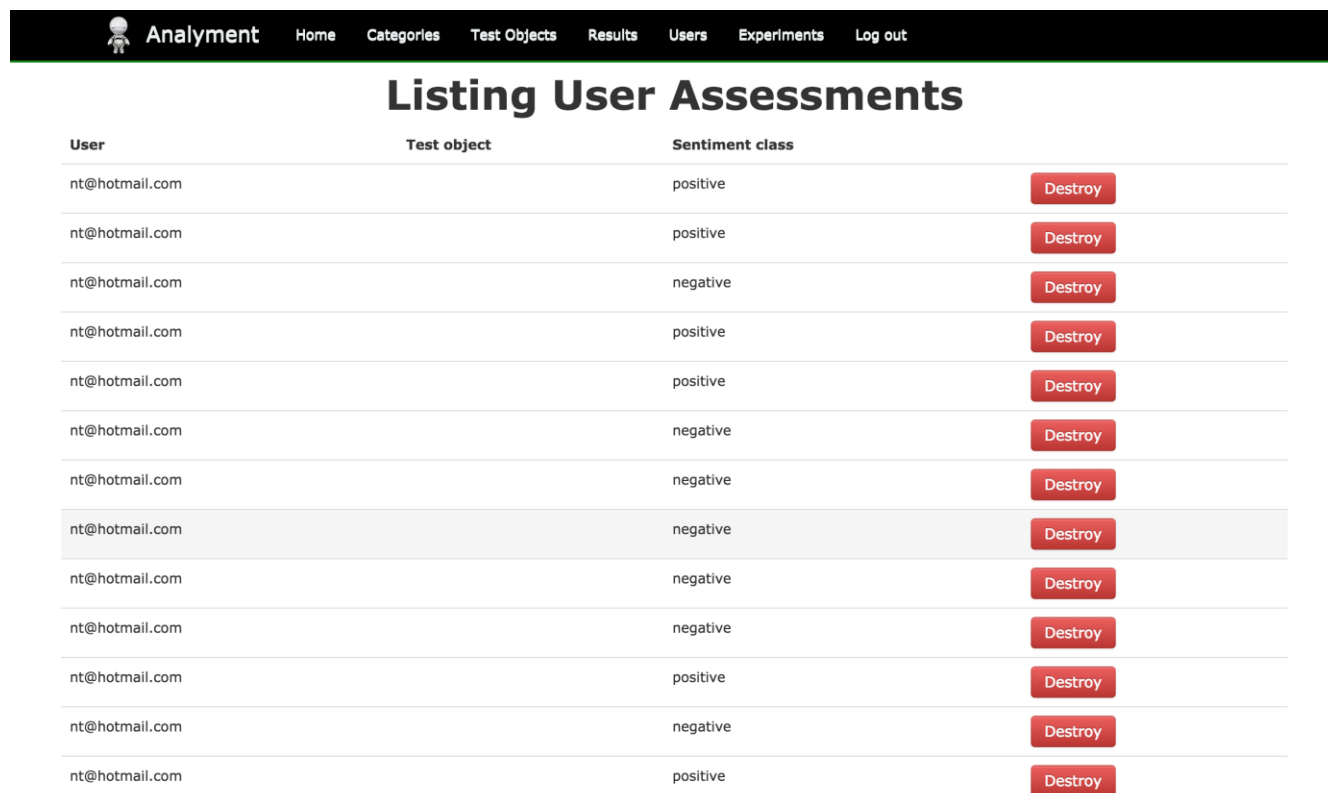
Listing users

Email	
nt@hotmail.com	Make end-user
aineias@hotmail.com	Make end-user

Anlyment by N.T - © TeiHer 2015

Εικόνα 5.4.2-6: Σελίδα χρηστών

Με το πράσινο κουμπί ο διαχειριστής μπορεί να δει συγκεντρωτικά τα αποτελέσματα ακόμα και να τα διαγραφει όταν έχει ολοκληρωθεί η πειραματική διαδικασία.



User	Test object	Sentiment class	
nt@hotmail.com		positive	Destroy
nt@hotmail.com		positive	Destroy
nt@hotmail.com		negative	Destroy
nt@hotmail.com		positive	Destroy
nt@hotmail.com		positive	Destroy
nt@hotmail.com		negative	Destroy
nt@hotmail.com		negative	Destroy
nt@hotmail.com		negative	Destroy
nt@hotmail.com		negative	Destroy
nt@hotmail.com		negative	Destroy
nt@hotmail.com		positive	Destroy
nt@hotmail.com		negative	Destroy
nt@hotmail.com		positive	Destroy

Εικόνα 5.4.2-7: Εμφάνιση όλων των κειμένων

Πειράματα

Η τελευταία και βασική λειτουργία της εφαρμογής είναι το Experiment, η οποία είναι και η μοναδική κοινή λειτουργία των διαχειριστών με τους απλούς χρήστες.

Εδώ παρουσιάζονται όλα τα πειράματα και ο συνολικός αριθμός των κειμένων που περιέχουν και πόσα απο αυτά ο χρήστης έχει βαθμολογήσει. Όταν κατηγοριοποιήσει όλα τα κείμενα εμφανίζεται ένα κουμπί που οδηγεί στην σελίδα των αποτελεσμάτων μαζί με ένα γράφημα δείχνοντας τι ποσοστό απόκλισης είχε ο αλγόριθμος από τον χρήστη.

Category	Assessed	Total
Movies	1	15
Quotes	15	15

View Results

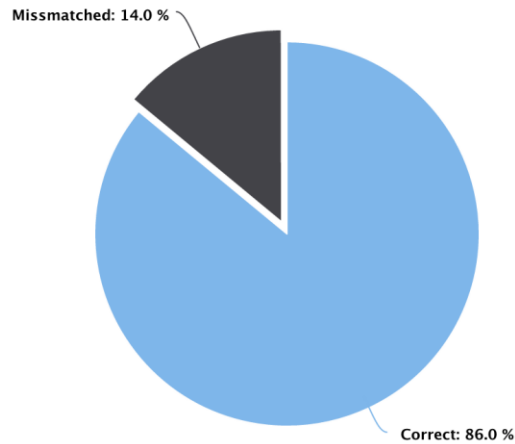
Analymet by N.T - © TeiHer 2015

Εικόνα 5.4.2-8: Πειραματική διαδικασία για τον διαχειριστή



Results - Category: Quotes

Comparison between algorithm and user results



Highcharts.com

Εικόνα 5.4.2-9: Γράφημα με την ολοκλήρωση της πειραματικής διαδικασίας



Assess Movies

Bad movie

This is perhaps the most appalling piece of art (lol) ever committed to celluloid, the acting, if you can call it that, consists of inane lines punctuated by silence while the mannequin imitating cast catch up and spout their increasingly dire reply. The editing is a case history in how modern technology can be abused to turn a screenplay into a horror story (not for the content of the movie but for anyone viewing it who has in their life time graduated beyond watching a spinning toy above a child's cot (crib)). From the total lack of acting skills this "thing" seems to be a product of someone who found a camera, asked a couple of friends to join in and then put together a script while eating at Chucky Cheese (or equivalent). If those responsible ever read these reviews in future if your intending to make another movie, assuming you are not now working in a mall or a drive-thru, try to ensure the lighting is balanced, it looked like Jill (?) was having a torch shone in her eyes in the boat, and for the two girls at the start there is no need to have them splashing water at each other for 20 minutes and generally if your going to film in a stairwell again at least balance the light. If you have nothing better to do for an hour or so (I am only 25 minutes in but stopped to write this)it is highly recommended as it truly has to be seen to be believed.

Sentiment class

Submit assessment

Exit experiment

1 2 3 ... Next > Last >

Εικόνα 5.4.2-10: Ανάθεση βαθμολογίας στα κείμενα

5.5 Προβλήματα - Συμπεράσματα και μελλοντική εξέλιξη της εφαρμογής

Αυτό το κεφάλαιο αφιερώνεται στην καταγραφή των προβλημάτων που αντιμετώπισε η ομάδα μας καθώς και τον τρόπο που έφτασε στην επίλυσή τους. Αξίζει να αναφερθεί ότι και κατά την σχεδίαση και κατά την υλοποίηση της εφαρμογής Analymnet η αναζήτηση στα κατάλληλα φόρουμ στο διαδίκτυο μας οδήγησε στην εύρεση της λύσης, καθώς όλοι οι προγραμματιστές έχουν αντιμετωπίσει κάποιο πρόβλημα και το έχουν μοιραστεί στα ειδικά φόρουμ.

Το πρώτο πρόβλημα που αντιμετωπίστηκε ήταν η εγκατάσταση του περιβάλλοντος εργασίας της γλώσσας προγραμματισμού Ruby αλλά και του framework της Ruby on Rails, καθώς ήταν απαραίτητη η εργασία απο διαφορετικά λειτουργικά συστήματα. Χρησιμοποιήθηκαν οι εκδόσεις λειτουργικών Ubuntu 12.04 και Mac OSX 10.10. Μετά από αναζήτηση στο διαδίκτυο (κυρίως στο stackoverflow.com) βρέθηκε τρόπος εγκατάστασης της πλατφόρμας στα διαφορετικά λειτουργικά συστήματα. Ο λόγος ύπαρξής τους ήταν για να δοκιμάσουμε πώς λειτουργεί η εφαρμογή αλλά και επειδή η ομάδα είχε διαφορετικά λειτουργικά συστήματα.

Αναφέρθηκε σε προηγούμενο κεφάλαιο η ευκολία χρήσης του twitter bootstrap. Αξίζει να αναφέρουμε οτι αντιμετωπίστηκαν αρκετά προβλήματα κατά την εγκατάστασή της βιβλιοθήκης μέχρι να καταλήξουμε στην πιο λειτουργική, η οποία προτεινόταν και στα εγχειρίδια της βιβλιοθήκης

και της κοινότητας της Ruby on Rails.

Μια πολύ χρήσιμη βιβλιοθήκη, η `kaminari` βρέθηκε για να μας βοηθήσει στην αντιμετώπιση του προβλήματος της σελιδοποίησης. Λόγω του όγκου δεδομένων που έχει η εφαρμογή έπρεπε να τα ταξινομήσουμε για να είναι πιο εύκολα στην αναζήτηση αλλά και την πλοήγηση για τον χρήστη.

Μία από τις πιο χρήσιμες βιβλιοθήκες είναι η `device` που διαχωρίζει δικαιώματα μεταξύ απλού χρήστη και διαχειριστή, ήταν δύσκολη η απόφαση με ποιόν τρόπο θα διαχωριστούν τα δικαιώματα μέχρι την εγκατάστασή της. Έπρεπε να την χρησιμοποιήσουμε αφού η εφαρμογή βασίζεται κυρίως στις λειτουργίες του διαχειριστή.

Το πιο δύσκολο πρόβλημα που αντιμετωπίστηκε είναι αυτό του `deployment` στο υπολογιστικό νέφος το οποίο προκαλούσε δυσλειτουργία στην εφαρμογή. Η εφαρμογή για να δώσει τα επιθυμητά αποτελέσματα, εκτός από τα αρχεία που περιέχουν κώδικα σε Ruby έχει και αρχεία με κώδικα άλλων γλωσσών όπως `javascript`, `CSS` και εικόνες τα `assets`. Τα αρχεία αυτά περνάνε από σύνταξη (`compilation`) για να τρέξουν παράλληλα με τον κύριο κορμό της εφαρμογής και δημιουργήθηκαν προβλήματα στην σωστή λειτουργία της εφαρμογής.

Οι χρήστες μπορούν να κάνουν χειροκίνητη αξιολόγηση του συναισθήματος των κειμένων βοηθώντας την τελική εξαγωγή συμπεράσματος για την καλή λειτουργία των `classifiers`. Εάν και καταλήξαμε στην υλοποίηση ενός εύχρηστου και ευκολού "εργαλείου" ο σκοπός είναι το εξελίξουμε στο μέλλον σε μία πιο αυτοματοποιημένη εφαρμογή. Το πρώτο βήμα είναι να είναι συνδεδεμένος ο `classifier` με την εφαρμογή και να κάνει `real time` εκπαίδευση. Αφού εκπαιδευτεί να εισάγει ο χρήστης URL της επιθυμίας του και ο αλγόριθμος να εμφανίζει το αποτέλεσμα. Η πιο χρήσιμη λειτουργία για τον αλγόριθμο είναι τα κείμενα τα οποία κατηγοριοποιεί σωστά να τα παίρνει ως ανατροφοδότηση "feedback" και να αποτελούν παραπάνω σύνολα εκπαίδευσης. Οι λειτουργίες αυτές δεν υλοποιήθηκαν λόγω έλλειψης χρόνου και αυξημένων υποχρεώσεων, αλλά το βασικότερο λόγω αυξημένων προγραμματιστικών αναγκών.

Παράρτημα

Κατάλογος κώδικα

Κώδικας 1 - H:\Analymet\production.rb	127
Κώδικας 2 - H:\Analymet\application.rb	129
Κώδικας 3 - H:\Analymet\results.js.coffee	130
Κώδικας 4 - H:\Analymet\custom.css.scss.....	131
Κώδικας 5 - H:\Analymet\user_assessments_controller.rb.....	133
Κώδικας 6 - H:\Analymet\test_objects_controller.rb.....	135
Κώδικας 7 - H:\Analymet\sentiment_classes_controller.rb.....	137
Κώδικας 8 - H:\Analymet\experiments_controller.rb.....	139
Κώδικας 9 - H:\Analymet\experiments_validation_corpora_controller.rb	140
Κώδικας 10 - H:\Analymet\dashboards_controller.rb	142
Κώδικας 11 - H:\Analymet\application.html.erb.....	143
Κώδικας 12 - H:\Analymet\schema.rb.....	145

Production.rb

```
seeds.rb × index.html.erb × production.rb ●
1 Analymet::Application.configure do
2   # Settings specified here will take precedence over those in config/application.rb.
3
4   # Code is not reloaded between requests.
5   config.cache_classes = true
6
7   # Eager load code on boot. This eager loads most of Rails and
8   # your application in memory, allowing both thread web servers
9   # and those relying on copy on write to perform better.
10  # Rake tasks automatically ignore this option for performance.
11  config.eager_load = true
12
13  # Full error reports are disabled and caching is turned on.
14  config.consider_all_requests_local = false
15  config.action_controller.perform_caching = true
16
17  # Enable Rack::Cache to put a simple HTTP cache in front of your application
18  # Add `rack-cache` to your Gemfile before enabling this.
19  # For large-scale production use, consider using a caching reverse proxy like nginx, varnish or squid.
20  # config.action_dispatch.rack_cache = true
21
22  # Disable Rails's static asset server (Apache or nginx will already do this).
23  config.serve_static_assets = false
24
25  # Compress JavaScripts and CSS.
26  config.assets.js_compressor = :uglifier
27  # config.assets.css_compressor = :sass
28
29  # Do not fallback to assets pipeline if a precompiled asset is missed.
30  config.assets.compile = false
31
32  # Generate digests for assets URLs.
33  config.assets.digest = true
34
35  # Version of your assets, change this if you want to expire all your assets.
36  config.assets.version = '1.0'
37
38  # Specifies the header that your server uses for sending files.
39  # config.action_dispatch.x_sendfile_header = "X-Sendfile" # for apache
40  # config.action_dispatch.x_sendfile_header = 'X-Accel-Redirect' # for nginx
41
42  # Force all access to the app over SSL, use Strict-Transport-Security, and use secure cookies.
43  # config.force_ssl = true
44
45  # Set to :debug to see everything in the log.
46  config.log_level = :info
47
48  # Prepend all log lines with the following tags.
49  # config.log_tags = [ :subdomain, :uuid ]
```



```
48 # Prepend all log lines with the following tags.
49 # config.log_tags = [ :subdomain, :uuid ]
50
51 # Use a different logger for distributed setups.
52 # config.logger = ActiveSupport::TaggedLogging.new(SyslogLogger.new)
53
54 # Use a different cache store in production.
55 # config.cache_store = :mem_cache_store
56
57 # Enable serving of images, stylesheets, and JavaScripts from an asset server.
58 # config.action_controller.asset_host = "http://assets.example.com"
59
60 # Precompile additional assets.
61 # application.js, application.css, and all non-JS/CSS in app/assets folder are already added.
62 # config.assets.precompile += %w( search.js )
63
64 # Ignore bad email addresses and do not raise email delivery errors.
65 # Set this to true and configure the email server for immediate delivery to raise delivery errors.
66 # config.action_mailer.raise_delivery_errors = false
67
68 # Enable locale fallbacks for I18n (makes lookups for any locale fall back to
69 # the I18n.default_locale when a translation can not be found).
70 config.i18n.fallbacks = true
71
72 # Send deprecation notices to registered listeners.
73 config.active_support.deprecation = :notify
74
75 # Disable automatic flushing of the log to improve performance.
76 # config.autoflush_log = false
77
78 # Use default logging formatter so that PID and timestamp are not suppressed.
79 config.log_formatter = ::Logger::Formatter.new
80 end
81
```

Application.rb

```
experiments_controller.rb x sentiment_classes_contro x experiment_validation_co x dashboards_controller.rb x application.html.erb x schema.rb x
1  require File.expand_path('../boot', __FILE__)
2
3  require 'rails/all'
4  require 'csv'
5
6  # Require the gems listed in Gemfile, including any gems
7  # you've limited to :test, :development, or :production.
8  Bundler.require(:default, Rails.env)
9
10 module Anallyment
11   class Application < Rails::Application
12     config.encoding = "utf-8"
13     # Settings in config/environments/* take precedence over those specified here.
14     # Application configuration should go into files in config/initializers
15     # — all .rb files in that directory are automatically loaded.
16
17     # Set Time.zone default to the specified zone and make Active Record auto-convert to this zone.
18     # Run "rake -D time" for a list of tasks for finding time zone names. Default is UTC.
19     # config.time_zone = 'Central Time (US & Canada)'
20
21     # The default locale is :en and all translations from config/locales/*.rb,yml are auto loaded.
22     # config.i18n.load_path += Dir[Rails.root.join('my', 'locales', '*.rb,yml')].to_s
23     # config.i18n.default_locale = :de
24   end
25 end
26
```

Results.js.coffee

```
results.js.coffee *
1  $ ->
2  $("#comparison_chart").highcharts
3  chart:
4  |   plotBackgroundColor: null
5  |   plotBorderWidth: null
6  |   plotShadow: false
7
8  |   title:
9  |     text: "Comparison between algorithm and user results"
10
11 |   tooltip:
12 |     pointFormat: "{series.name}: <b>{point.percentage:.1f}%</b>"
13
14 |   plotOptions:
15 |     pie:
16 |       allowPointSelect: true
17 |       cursor: "pointer"
18 |       dataLabels:
19 |         enabled: true
20 |         format: "<b>{point.name}</b>: {point.percentage:.1f} %"
21 |         style:
22 |           color: (Highcharts.theme and Highcharts.theme.contrastTextColor) or "black"
23
24 |   series: [
25 |     type: "pie"
26 |     name: "User assessments"
27 |     data: [
28 |       {
29 |         name: "Correct"
30 |         y: parseInt($('#user_correct').text())
31 |         sliced: true
32 |         selected: true
33 |       }
34 |       [
35 |         "Missmatched"
36 |         100 - parseInt($('#user_correct').text())
37 |       ]
38 |     ]
39 |   ]
40 ]
41
42 return
43
```

Custom.css.scss

```
custom.css.scss *
1  body {
2    background-image: white;
3  }
4
5  h1 {
6    text-align: center
7  }
8
9  table {
10   margin: 0 auto;
11 }
12
13 footer{
14   margin-top: 40px;
15 }
16 .container-fluid {
17   background-color: Black;
18 }
19 .navbar {
20   overflow: hidden;
21   margin-bottom: 0px;
22 }
23 .main-wrapper {
24   width: 1170px;
25   margin: 0 auto;
26   background-color: white;
27 }
28 .navbar .nav {
29   margin-top: 5px;
30 }
31 .navbar .brand:hover {
32   background-color: #3E4449;
33   text-decoration: none;
34   color: #777777;
35   text-shadow: 0 1px 0 #777777;
36 }
37
38 .navbar .brand {
39   float: left;
40   display: block;
41   padding: 8px 20px 4px;
42   margin-left: 50px;
43   font-size: 20px;
44   font-weight: 200;
45   color: #ecf0f1;
46   text-shadow: 0 1px 0 #ecf0f1;
47 }
48 .navbar-inner {
49
50   padding-left: 0px;
51   padding-right: 0px;
52   background-color: #3E4449;
53   background-image: -webkit-linear-gradient(top, #333333, #222222);
54   border-bottom: 2px solid green;
55 }
```

```

56
57 .navbar .nav>li>a {
58 float: none;
59 padding: 10px 15px 10px;
60 color: #ecf0f1;
61 text-decoration: none;
62 text-shadow: 0 1px 0 #ecf0f1;
63 display: block;
64 }
65
66 .navbar .nav>li>a: hover {
67 color: #777777;
68 text-shadow: 0 1px 0 #777777;
69 }
70 }
71
72
73
74
75
76 .alert {
77 margin-top: 20px;
78 }
79
80 .dashboard {
81 }
82
83 .dashboard a {
84 margin: 30px 10px;
85 padding: 10px;
86 }
87
88 a.btn-primary:visited, a.btn-success:visited, a.btn-danger:visited {color:#fff}
89
90 .btn-default {
91 color: #523352;
92 }
93
94 h4.result_categories {
95 float: left;
96 }
97 .result_categories li {
98 list-style: none;
99 float: left;
100 margin: 5px 10px;
101 }
102
103
104 img.style_image{
105 margin: 20px auto;
106 display: block;
107 width: 75%;
108 }

```

assessments_controller.rb

```
user_assessments_controller.rb x
1  class UserAssessmentsController < ApplicationController
2    authorize_resource class: false
3    before_action :set_user_assessment, only: [:show, :edit, :update, :destroy]
4
5    # GET /user_assessments
6    # GET /user_assessments.json
7    def index
8      @user_assessments = UserAssessment.all.page(params[:page]).per(20)
9    end
10
11   # POST /user_assessments
12   # POST /user_assessments.json
13   def create
14     @user_assessment = UserAssessment.new(user_assessment_params)
15     @user_assessment.user_id = current_user.id
16     set_experiment_validation_corpus
17     respond_to do |format|
18       if @user_assessment.save
19         format.html { redirect_to :back, notice: 'User assessment was successfully created.' }
20         format.json { render action: 'show', status: :created, location: @user_assessment }
21       else
22         if @user_assessment.errors[:user_id] == ["has already been taken"]
23           update and return
24         end
25         format.html { render action: 'new' }
26         format.json { render json: @user_assessment.errors, status: :unprocessable_entity }
27       end
28     end
29   end
30
31   def update
32     @user_assessment = UserAssessment.where(user_id: current_user.id, test_object_id: user_assessment_params[:test_object_id] ).f
33     respond_to do |format|
34       if @user_assessment.update_attributes(sentiment_class_id: user_assessment_params[:sentiment_class_id])
35         format.html { redirect_to :back, notice: 'User assessment was successfully updated.' }
36         format.json { head :no_content }
37       else
38         format.html { render action: 'edit' }
39         format.json { render json: @user_assessment.errors, status: :unprocessable_entity }
40       end
41     end
42   end
43
44   # DELETE /user_assessments/1
45   # DELETE /user_assessments/1.json
46   def destroy
47     @user_assessment.destroy
48     respond_to do |format|
49       format.html { redirect_to user_assessments_url }
50       format.json { head :no_content }
51     end
52   end
53
54   def set_experiment_validation_corpus
55     @user_assessment.experiment_validation_corpus_id = TestObject.find(user_assessment_params[:test_object_id]).experiment_valida
56   end
57 end
```

```
57
58 private
59 # Use callbacks to share common setup or constraints between actions.
60 def set_user_assessment
61   @user_assessment = UserAssessment.find(params[:id])
62 end
63
64 # Never trust parameters from the scary internet, only allow the white list through.
65 def user_assessment_params
66   params.require(:user_assessment).permit(:user_id, :test_object_id, :sentiment_class_id)
67 end
68 end
69
```

objects_controller.rb

```
user_assessments_controller.rb * test_objects_controller.rb *
1 class TestObjectsController < ApplicationController
2   authorize_resource class: false
3   before_action :set_test_object, only: [:show, :edit, :update, :destroy]
4
5   # GET /test_objects
6   # GET /test_objects.json
7   def index
8     @test_objects = TestObject.all.page(params[:page]).per(5)
9   end
10
11  # GET /test_objects/1
12  # GET /test_objects/1.json
13  def show
14  end
15
16  # GET /test_objects/new
17  def new
18    @test_object = TestObject.new
19  end
20
21  # GET /test_objects/1/edit
22  def edit
23  end
24
25  # POST /test_objects
26  # POST /test_objects.json
27  def create
28    @test_object = TestObject.new(test_object_params)
29
30    respond_to do |format|
31      if @test_object.save
32        format.html { redirect_to @test_object, notice: 'Test object was successfully created.' }
33        format.json { render action: 'show', status: :created, location: @test_object }
34      else
35        format.html { render action: 'new' }
36        format.json { render json: @test_object.errors, status: :unprocessable_entity }
37      end
38    end
39  end
40
41  # PATCH/PUT /test_objects/1
42  # PATCH/PUT /test_objects/1.json
43  def update
44    respond_to do |format|
45      if @test_object.update(test_object_params)
46        format.html { redirect_to @test_object, notice: 'Test object was successfully updated.' }
47        format.json { head :no_content }
48      else
49        format.html { render action: 'edit' }
50        format.json { render json: @test_object.errors, status: :unprocessable_entity }
51      end
52    end
53  end
54
55  # DELETE /test_objects/1
56  # DELETE /test_objects/1.json
```



```

56 # DELETE /test_objects/1.json
57 def destroy
58   @test_object.destroy
59   respond_to do |format|
60     format.html { redirect_to test_objects_url }
61     format.json { head :no_content }
62   end
63 end
64
65 def import_csv
66   TestObject.import(params[:file])
67   redirect_to test_objects_path, notice: "Test objects imported."
68 end
69
70 def import_body
71   test_object = TestObject.find(params[:test_object_id])
72   file_name = test_object.body
73   file = File.open("#{Rails.root}/tmp/#{file_name}")
74   test_object.body = file.read()
75   test_object.save
76   redirect_to(:back)
77 end
78
79 private
80 # Use callbacks to share common setup or constraints between actions.
81 def set_test_object
82   @test_object = TestObject.find(params[:id])
83 end
84
85 # Never trust parameters from the scary internet, only allow the white list through.
86 def test_object_params
87   params.require(:test_object).permit(:experiment_validation_corpus_id, :title, :body, :sentiment_class_id)
88 end
89 end
90

```

1

1 classes_controller.rb

```
user_assessments_controller.rb × test_objects_controller.rb × sentiment_classes_controller.rb ×
1 class SentimentClassesController < ApplicationController
2   authorize_resource class: false
3   before_action :set_sentiment_class, only: [:show, :edit, :update, :destroy]
4
5   # GET /sentiment_classes
6   # GET /sentiment_classes.json
7   def index
8     @sentiment_classes = SentimentClass.all
9   end
10
11  # GET /sentiment_classes/1
12  # GET /sentiment_classes/1.json
13  def show
14  end
15
16  # GET /sentiment_classes/new
17  def new
18    @sentiment_class = SentimentClass.new
19  end
20
21  # GET /sentiment_classes/1/edit
22  def edit
23  end
24
25  # POST /sentiment_classes
26  # POST /sentiment_classes.json
27  def create
28    @sentiment_class = SentimentClass.new(sentiment_class_params)
29
30    respond_to do |format|
31      if @sentiment_class.save
32        format.html { redirect_to @sentiment_class, notice: 'Sentiment class was successfully created.' }
33        format.json { render action: 'show', status: :created, location: @sentiment_class }
34      else
35        format.html { render action: 'new' }
36        format.json { render json: @sentiment_class.errors, status: :unprocessable_entity }
37      end
38    end
39  end
40
41  # PATCH/PUT /sentiment_classes/1
42  # PATCH/PUT /sentiment_classes/1.json
43  def update
44    respond_to do |format|
45      if @sentiment_class.update(sentiment_class_params)
46        format.html { redirect_to @sentiment_class, notice: 'Sentiment class was successfully updated.' }
47        format.json { head :no_content }
48      else
49        format.html { render action: 'edit' }
50        format.json { render json: @sentiment_class.errors, status: :unprocessable_entity }
51      end
52    end
53  end
end
```

```

54
55 # DELETE /sentiment_classes/1
56 # DELETE /sentiment_classes/1.json
57 def destroy
58   @sentiment_class.destroy
59   respond_to do |format|
60     format.html { redirect_to sentiment_classes_url }
61     format.json { head :no_content }
62   end
63 end
64
65 private
66 # Use callbacks to share common setup or constraints between actions.
67 def set_sentiment_class
68   @sentiment_class = SentimentClass.find(params[:id])
69 end
70
71 # Never trust parameters from the scary internet, only allow the white list through.
72 def sentiment_class_params
73   params.require(:sentiment_class).permit(:name)
74 end
75 end

```

1 Experiments_controller.rb

```
experiments_controller.rb ✕
1 class ExperimentsController < ApplicationController
2   def index
3     @experiments = []
4     ExperimentValidationCorpus.all.each do |experiment|
5       exp = {}
6       exp[:name] = experiment.name
7       exp[:assessed] = assessed_by_user(experiment).length
8       exp[:total] = count_category_test_objects(experiment)
9       @experiments << exp
10    end
11  end
12
13  def assessed_by_user experiment
14    UserAssessment.where(user_id: current_user.id, experiment_validation_corpus_id: experiment.id)
15  end
16
17  def count_category_test_objects experiment
18    TestObject.where(experiment_validation_corpus_id: experiment.id).length
19  end
20
21  def results
22    @experiment_validation_corpus = ExperimentValidationCorpus.find_by_name(params[:format])
23    # in case someone tries to see results (if he knows the URL) without having assessed all test objects
24    if assessed_by_user(@experiment_validation_corpus).length == count_category_test_objects(@experiment_validation_corpus)
25      @results = UserAssessment.where(experiment_validation_corpus_id: @experiment_validation_corpus.id, user_id: current_user.id)
26
27      @score = UserAssessment.algorithm_comparison(@results, current_user)
28    end
29  end
30
31  def run_assessment
32    @experiment_validation_corpus = ExperimentValidationCorpus.find_by_name(params[:format])
33    @test_objects = TestObject.where(experiment_validation_corpus_id: @experiment_validation_corpus.id).page(params[:page]).per(1)
34    @user_assessment = UserAssessment.new
35  end
36 end
37
```

1 Validation_Corpora_controller.rb

```
1 class ExperimentValidationCorporaController < ApplicationController
2   authorize_resource class: false
3   before_action :set_experiment_validation_corpus, only: [:show, :edit, :update, :destroy]
4
5   # GET /experiment_validation_corpora
6   # GET /experiment_validation_corpora.json
7   def index
8     @experiment_validation_corpora = ExperimentValidationCorpus.all
9   end
10
11  # GET /experiment_validation_corpora/1
12  # GET /experiment_validation_corpora/1.json
13  def show
14  end
15
16  # GET /experiment_validation_corpora/new
17  def new
18    @experiment_validation_corpus = ExperimentValidationCorpus.new
19  end
20
21  # GET /experiment_validation_corpora/1/edit
22  def edit
23  end
24
25  # POST /experiment_validation_corpora
26  # POST /experiment_validation_corpora.json
27  def create
28    @experiment_validation_corpus = ExperimentValidationCorpus.new(experiment_validation_corpus_params)
29
30    respond_to do |format|
31      if @experiment_validation_corpus.save
32        format.html { redirect_to @experiment_validation_corpus, notice: 'Category was successfully created.' }
33        format.json { render action: 'show', status: :created, location: @experiment_validation_corpus }
34      else
35        format.html { render action: 'new' }
36        format.json { render json: @experiment_validation_corpus.errors, status: :unprocessable_entity }
37      end
38    end
39  end
40
41  # PATCH/PUT /experiment_validation_corpora/1
42  # PATCH/PUT /experiment_validation_corpora/1.json
43  def update
44    respond_to do |format|
45      if @experiment_validation_corpus.update(experiment_validation_corpus_params)
46        format.html { redirect_to @experiment_validation_corpus, notice: 'Category was successfully updated.' }
47        format.json { head :no_content }
48      else
49        format.html { render action: 'edit' }
50        format.json { render json: @experiment_validation_corpus.errors, status: :unprocessable_entity }
51      end
52    end
53  end
54 end
```

```

54
55 # DELETE /experiment_validation_corpora/1
56 # DELETE /experiment_validation_corpora/1.json
57 def destroy
58   @experiment_validation_corpus.destroy
59   respond_to do |format|
60     format.html { redirect_to experiment_validation_corpora_url }
61     format.json { head :no_content }
62   end
63 end
64
65 private
66 # Use callbacks to share common setup or constraints between actions.
67 def set_experiment_validation_corpus
68   @experiment_validation_corpus = ExperimentValidationCorpus.find(params[:id])
69 end
70
71 # Never trust parameters from the scary internet, only allow the white list through.
72 def experiment_validation_corpus_params
73   params.require(:experiment_validation_corpus).permit(:name)
74 end
75 end
76

```

1 File – H:\Analymet\ Dashboards_controller.rb

```
1 class DashboardsController < ApplicationController
2   authorize_resource class: false
3   def admin_home
4     end
5
6   def index
7     @experiment_validation_corpora = ExperimentValidationCorpus.all
8   end
9
10  def compare
11    @experiment_validation_corpus = ExperimentValidationCorpus.find(params[:format])
12    @test_objects_by_corpus = @experiment_validation_corpus.test_objects.page(params[:page]).per(1)
13  end
14
15  def users_index
16    @users = User.all
17  end
18
19  def user_role
20    @user = User.find(params[:format])
21    if @user.role.present?
22      @user.update_attributes(role: "")
23    else
24      @user.update_attributes(role: "admin")
25    end
26    redirect_to :back, notice: 'User role was successfully updated.'
27  end
28 end
29
```

2

3

4

5

6

7

8

9

10

11

12

13

14 **File – H:\Analyment\ application.html.erb**

```

1 <!DOCTYPE html>
2 <html lang="en">
3 <head>
4 <meta charset="utf-8">
5 <meta http-equiv="X-UA-Compatible" content="IE=Edge,chrome=1">
6 <meta name="viewport" content="width=device-width, initial-scale=1.0">
7 <title>{{ content_for?(:title) ? yield(:title) : "Analymnet" }}</title>
8 {{ csrf_meta_tags }}
9
10 <!-- Le HTML5 shim, for IE6-8 support of HTML elements -->
11 <!--[if lt IE 9]>
12 | <script src="//cdnjs.cloudflare.com/ajax/libs/html5shiv/3.6.1/html5shiv.js" type="text/javascript"></script>
13 <![endif]-->
14
15 {{ stylesheet_link_tag "application", :media => "all" }}
16
17 <!-- For third-generation iPad with high-resolution Retina display: -->
18 <!-- Size should be 144 x 144 pixels -->
19 {{ favicon_link_tag 'apple-touch-icon-144x144-precomposed.png', :rel => 'apple-touch-icon-precomposed', :type => 'image/png'
20 , :sizes => '144x144' }}
21
22 <!-- For iPhone with high-resolution Retina display: -->
23 <!-- Size should be 114 x 114 pixels -->
24 {{ favicon_link_tag 'apple-touch-icon-114x114-precomposed.png', :rel => 'apple-touch-icon-precomposed', :type => 'image/png'
25 , :sizes => '114x114' }}
26
27 <!-- For first- and second-generation iPad: -->
28 <!-- Size should be 72 x 72 pixels -->
29 {{ favicon_link_tag 'apple-touch-icon-72x72-precomposed.png', :rel => 'apple-touch-icon-precomposed', :type => 'image/png',
30 :sizes => '72x72' }}
31
32 <!-- For non-Retina iPhone, iPod Touch, and Android 2.1+ devices: -->
33 <!-- Size should be 57 x 57 pixels -->
34 {{ favicon_link_tag 'apple-touch-icon-precomposed.png', :rel => 'apple-touch-icon-precomposed', :type => 'image/png' }}
35
36 <!-- For all other devices -->
37 <!-- Size should be 32 x 32 pixels -->
38 {{ favicon_link_tag 'favicon.ico', :rel => 'shortcut icon' }}
39
40 {{ javascript_include_tag "application" }}
41 <script src="http://code.highcharts.com/highcharts.js"></script>
42 <script src="http://code.highcharts.com/modules/exporting.js"></script>
43 </head>
44 <body>
45
46 <div class="navbar navbar-fluid-top">
47 <div class="navbar-inner">
48 <div class="container-fluid">
49 <a class="btn btn-navbar" data-target=".nav-collapse" data-toggle="collapse">
50 <span class="icon-bar"></span>
51 <span class="icon-bar"></span>
52 <span class="icon-bar"></span>
53 </a>
54 <a class="brand" href="/">{{ image_tag("logo1.png", size: "38") }} Analymnet</a>
55 <div class="container-fluid nav-collapse">
56 <ul class="nav">
57 {{ render '/layouts/menu' if current_user }}

```

```
55     </ul>
56   </div><!--/.nav-collapse -->
57 </div>
58 </div>
59 </div>
60
61 <div class="container-fluid main-wrapper">
62   <div class="row-fluid">
63     <div class="span12">
64       <%= bootstrap_flash %>
65       <%= yield %>
66     </div>
67   </div><!--/row-->
68
69   <footer>
70     <p>Analymnt by N.T - &copy; TeiHer 2015</p>
71   </footer>
72
73 </div> <!-- /container -->
74
75 </body>
76 </html>
77
```

16

17

18

19

20

21

22

23

24

25

26

27

28

29 File – H:\Analymnt\ schema.rb

```

1 # encoding: UTF-8
2 # This file is auto-generated from the current state of the database. Instead
3 # of editing this file, please use the migrations feature of Active Record to
4 # incrementally modify your database, and then regenerate this schema definition.
5 #
6 # Note that this schema.rb definition is the authoritative source for your
7 # database schema. If you need to create the application database on another
8 # system, you should be using db:schema:load, not running all the migrations
9 # from scratch. The latter is a flawed and unsustainable approach (the more migrations
10 # you'll amass, the slower it'll run and the greater likelihood for issues).
11 #
12 # It's strongly recommended that you check this file into your version control system.
13
14 ActiveRecord::Schema.define(version: 20140111133824) do
15
16   create_table "experiment_validation_corpora", force: true do |t|
17     t.string "name"
18     t.datetime "created_at"
19     t.datetime "updated_at"
20   end
21
22   create_table "sentiment_classes", force: true do |t|
23     t.string "name"
24     t.datetime "created_at"
25     t.datetime "updated_at"
26   end
27
28   create_table "test_objects", force: true do |t|
29     t.integer "experiment_validation_corpus_id"
30     t.string "title"
31     t.text "body"
32     t.integer "sentiment_class_id"
33     t.datetime "created_at"
34     t.datetime "updated_at"
35   end
36
37   add_index "test_objects", ["experiment_validation_corpus_id"], name: "index_test_objects_on_experiment_validation_corpus_id", u
38   add_index "test_objects", ["sentiment_class_id"], name: "index_test_objects_on_sentiment_class_id", using: :btree
39
40   create_table "user_assessments", force: true do |t|
41     t.integer "user_id"
42     t.integer "test_object_id"
43     t.integer "sentiment_class_id"
44     t.integer "experiment_validation_corpus_id"
45     t.datetime "created_at"
46     t.datetime "updated_at"
47   end
48
49   add_index "user_assessments", ["experiment_validation_corpus_id"], name: "index_user_assessments_on_experiment_validation_corpu
50   add_index "user_assessments", ["sentiment_class_id"], name: "index_user_assessments_on_sentiment_class_id", using: :btree
51   add_index "user_assessments", ["test_object_id"], name: "index_user_assessments_on_test_object_id", using: :btree
52   add_index "user_assessments", ["user_id", "test_object_id"], name: "index_user_assessments_on_user_id_and_test_object_id", uniq
53   add_index "user_assessments", ["user_id"], name: "index_user_assessments_on_user_id", using: :btree
54
55   create_table "users", force: true do |t|
56     t.string "email", default: "", null: false
57     t.string "encrypted_password", default: "", null: false

```

30
31
32

```
57 t.string "encrypted_password", default: "", null: false
58 t.string "reset_password_token"
59 t.datetime "reset_password_sent_at"
60 t.datetime "remember_created_at"
61 t.integer "sign_in_count", default: 0, null: false
62 t.datetime "current_sign_in_at"
63 t.datetime "last_sign_in_at"
64 t.string "current_sign_in_ip"
65 t.string "last_sign_in_ip"
66 t.string "role"
67 t.datetime "created_at"
68 t.datetime "updated_at"
69 end
70
71 add_index "users", ["email"], name: "index_users_on_email", unique: true, using: :btree
72 add_index "users", ["reset_password_token"], name: "index_users_on_reset_password_token", unique: true, using: :btree
73
74 end
75
```

33

Αναφορές

- [1]: David Thomas, Andrew Hunt, (2001), *The Pragmatic Programmer's Guide*. Διαθέσιμο: www.ruby-doc.org/docs/ProgrammingRuby/html/index.html. Προσπελάστηκε: 20η Απρ. 2012
- [2]: Ruby on Rails community, (2014), *Rails Guides*. Διαθέσιμο: guides.rubyonrails.org/generators.html. Προσπελάστηκε: 1^η Απρ. 2014
- [3]: Oracle Corporation, (2014), *MySQL*. Διαθέσιμο: www.mysql.com/why-mysql. Προσπελάστηκε: 26^η Μαρτίου 2014
- [4]: Christian Neukirchen, (2012), *Rack: a Ruby Webserver Interface*. Διαθέσιμο: rack.github.com. Προσπελάστηκε: 26 Απρ. 2012
- [5]: Michael Hartl, (2012), *Ruby on Rails Tutorial - Learn Rails by Example*. Διαθέσιμο: ruby.railstutorial.org/chapters. Προσπελάστηκε: 21η Απρ. 2012
- [6]: Ruby on Rails community, (2012), *Rails Guides*. Διαθέσιμο: guides.rubyonrails.org/generators.html. Προσπελάστηκε: 21η Απρ. 2012
- [7]: heroku, (2014), *heroku website*. Διαθέσιμο: www.heroku.com. Προσπελάστηκε: 25^η Μαΐου 2014
- [8]: C.Baun et al. (2011), *Cloud Computing*. 2nd ed. London: Springer. Κεφάλαιο 6 pp. 49-62
- [9]: Roy T. Fielding, (2008), *REST APIs must be hypertext-driven*. Διαθέσιμο: roy.gbiv.com/untangled/2008/rest-apis-must-be-hypertext-driven. Προσπελάστηκε: 21^η Απρ. 2014
- [10]: Ruby on Rails Community, (2014), *Active Record Migrations*. Διαθέσιμο: <http://api.rubyonrails.org/classes/ActiveRecord/Migration.html> Προσπελάστηκε: 28^η Μαΐου 2014
- [11]: <http://www.cs.cornell.edu/people/pabo/movie-review-data>
- [12]: What is HTML?, (2014), Διαθέσιμο: www.yourhtmlsource.com. Προσπελάστηκε: 25^η Μαΐου 2014
- [13]: RVM, (2014), Διαθέσιμο: rvm.io. Προσπελάστηκε: 25^η Μαΐου 2014

Βιβλιογραφία

- [1] GAMON, M. 2004. Sentiment classification on customer feedback data: noisy data, large feature vectors, and the role of linguistic analysis. In Proceedings the 20th International Conference on Computational Linguistics, pp. 841–847.
- [2] Pang, B., & Lee, L. (2004, July). A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In Proceedings of the 42nd annual meeting on Association for Computational Linguistics (p. 271). Association for Computational Linguistics.
- [3] Xiaowen Ding, Bing Liu and Philip S. Yu. "A Holistic Lexicon-Based Approach to Opinion Mining." Proceedings of First ACM International Conference on Web Search and Data Mining (WSDM-2008), Feb 11-12, 2008, Stanford University, Stanford, California, USA
- [4] SES: Sentiment Elicitation System for Social Media Data. Kunpeng Zhang, Yu Cheng, Yusheng Xie, Daniel Honbo, Ankit Agrawal, Diana Palsetia, Kathy Lee, Wei-keng Liao, Alok Choudhary. s.l. : 11th IEEE International Conference on Data Mining Workshops, 2011.
- [5] Albert Bifet, Eibe Frank. Sentiment Knowledge Discovery in Twitter Streaming Data. Lecture Notes in Computer Science, Volume 6332, Discovery Science, Pages 1-15. 2010.
- [6] Richard D. Waters, Jia Y. Jamal. Tweet, tweet, tweet: A content analysis of nonprofit organizations' Twitter updates. Public Relations Review. Elsevier Inc, September 2011, Τόμ. Volume 37, Issue 3, Pages 321-324.
- [7] Farah Benamara, Baptiste Chardon, Yvette Yannick Mathieu, Vladimir Popescu: Towards Context-Based Subjectivity Analysis. IJCNLP 2011: 1180-1188
- [8] Thorsten Joachims, Transductive Inference for Text Classification using Support Vector Machines. International Conference on Machine Learning (ICML), 1999.
- [9] RAVI PARIKH - MATIN MOVASSATE, SENTIMENT ANALYSIS OF USERGENERATED TWITTER UPDATES USING VARIOUS CLASSICATION TECHNIQUES, JUNE 2009, STANDFORD UNIVERSITY
- [10] An Introduction to Support Vector Machines and Other Kernel-based Learning Methods, Nello Cristianini and John Shawe-Taylor, Cambridge University Press, 2000
- [11] V. PANDEY, K. IYER, SENTIMENT ANALYSIS OF MICROBLOGS, STANDFORD UNIVERSITY
- [12] PREM MELVILLE - WOJCIECH GRYC - RICHARD D. LAWRENCE , SENTIMENT ANALYSIS OF BLOGS BY COMBINING LEXICAL KNOWLEDGE WITH TEXT CLASSIFICATION, KDD JUNE 2009
- [13]Matsumoto, S., Takamura, H., & Okumura, M. (2005). Sentiment classification using word subsequences and dependency sub-trees. In Advances in Knowledge Discovery and Data Mining (pp.301-311). Springer Berlin Heidelberg.
- [14]Kennedy, A., & Inkpen, D. (2006). Sentiment classification of movie reviews using contextual valence shifters. Computational Intelligence, 22(2), 110-125.
- [15] Qiu, L., Zhang, W., Hu, C., & Zhao, K. (2009, November). Selc: a self-supervised model for sentiment classification. In Proceedings of the 18th ACM conference on Information and knowledge management (pp. 929-936). ACM.
- [16] Hatzivassiloglou, V., Klavans, J. L., Holcombe, M. L., Barzilay, R., Kan, M. Y., & McKeown, K. (2001). Simfinder: A flexible clustering tool for summarization. Proceedings of the NAACL Workshop on Au-

omatic Summarizatio.