



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ
ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ
ΤΟΜΕΑΣ ΕΠΙΚΟΙΝΩΝΙΩΝ, ΗΛΕΚΤΡΟΝΙΚΗΣ ΚΑΙ
ΣΥΣΤΗΜΑΤΩΝ ΠΛΗΡΟΦΟΡΙΚΗΣ

**ΤΕΧΝΙΚΕΣ ΑΝΑΛΥΣΗΣ ΚΟΙΝΩΝΙΚΩΝ
ΔΙΚΤΥΩΝ, ΜΕ ΕΜΦΑΣΗ ΣΕ ΓΡΑΦΟΥΣ
ΕΜΠΙΣΤΟΣΥΝΗΣ**

ΔΙΔΑΚΤΟΡΙΚΗ ΔΙΑΤΡΙΒΗ

ΑΘΑΝΑΣΙΟΥ Δ. ΠΑΠΑΟΙΚΟΝΟΜΟΥ

Διπλωματούχου Ηλεκτρολόγου Μηχανικού &
Μηχανικού Υπολογιστών Ε.Μ.Π.

ΕΠΙΒΛΕΠΟΥΣΑ:

Θ. ΒΑΡΒΑΡΙΓΟΥ

Καθηγήτρια Ε.Μ.Π.

ΑΘΗΝΑ, Απρίλιος 2015



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ
ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ
ΤΟΜΕΑΣ ΕΠΙΚΟΙΝΩΝΙΩΝ, ΗΛΕΚΤΡΟΝΙΚΗΣ ΚΑΙ
ΣΥΣΤΗΜΑΤΩΝ ΠΛΗΡΟΦΟΡΙΚΗΣ

ΤΕΧΝΙΚΕΣ ΑΝΑΛΥΣΗΣ ΚΟΙΝΩΝΙΚΩΝ ΔΙΚΤΥΩΝ, ΜΕ ΕΜΦΑΣΗ ΣΕ ΓΡΑΦΟΥΣ ΕΜΠΙΣΤΟΣΥΝΗΣ

ΔΙΔΑΚΤΟΡΙΚΗ ΔΙΑΤΡΙΒΗ

ΑΘΑΝΑΣΙΟΥ Δ. ΠΑΠΑΟΙΚΟΝΟΜΟΥ

Διπλωματούχου Ηλεκτρολόγου Μηχανικού &
Μηχανικού Υπολογιστών Ε.Μ.Π.

Συμβουλευτική Επιτροπή : 1. Θ. ΒΑΡΒΑΡΙΓΟΥ, Καθ. Ε.Μ.Π. (Επιβλέπουσα)
2. Ν. ΚΟΖΥΡΗΣ, Καθ. Ε.Μ.Π.
3. Σ. ΠΑΠΑΒΑΣΙΛΕΙΟΥ, Καθ. Ε.Μ.Π.

Εγκρίθηκε από την επταμελή εξεταστική επιτροπή την 24η Απριλίου, 2015.

.....
Θ.Α.Βαρβαρίγου,
Καθ. Ε.Μ.Π.

.....
Ν. Κοζύρης,
Καθ. Ε.Μ.Π.

.....
Σ. Παπαβασιλείου,
Καθ. Ε.Μ.Π.

.....
Β. Λούμος,
Καθ. Ε.Μ.Π.

.....
Α. Σταφυλοπάτης,
Καθ. Ε.Μ.Π.

.....
Δ. Ασκούνης,
Αν. Καθ. Ε.Μ.Π.

.....
Α. Δουλάμης,
Λέκτορας Ε.Μ.Π.

ΑΘΗΝΑ, Απρίλιος 2015

.....

Αθανασίου Δ. Παπαοικονόμου

Διδάκτωρ Ηλεκτρολόγος Μηχανικός και Μηχανικός Υπολογιστών Ε.Μ.Π.

Copyright © Αθανάσιος Δ. Παπαοικονόμου, 2015.

Με επιφύλαξη παντός δικαιώματος. All rights reserved.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, εκπαιδευτικής ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα.

Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και δεν πρέπει να ερμηνευθεί ότι αντιπροσωπεύουν τις επίσημες θέσεις του Εθνικού Μετσόβιου Πολυτεχνείου.

Η έγκριση της διδακτορικής διατριβής από την Ανώτατη Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών του Ε.Μ. Πολυτεχνείου δεν υποδηλώνει αποδοχή των γνώμων του συγγραφέα (Ν. 5343/1932, Άρθρο 202).

Πρόλογος

Η διδακτορική διατριβή που παρουσιάζεται στις επόμενες σελίδες εκπονήθηκε από το Ιανουάριο του 2010 μέχρι τον Ιανουάριο του 2015, στο εργαστήριο Τηλεπικοινωνιών του τομέα Επικοινωνιών, Ηλεκτρονικής και Συστημάτων Πληροφορικής, στη Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών του Εθνικού Μετσόβιου Πολυτεχνείου υπό την επίβλεψη της κ. Θεοδώρας Βαρβαρίγου.

Θα ήθελα να ευχαριστήσω την καθηγήτριά μου κ. Θεοδώρα Βαρβαρίγου για την υποστήριξη, και την καθοδήγηση που μου παρείχε από την αρχή ως το τέλος της προσπάθειάς μου, καθώς επίσης τους καθηγητές της τριμελούς συμβουλευτικής επιτροπής κ. κ. Συμεών Παπαβασιλείου και Νεκτάριο Κοζύρη.

Επίσης, θα ήθελα να ευχαριστήσω όλους τους συναδέλφους με τους οποίους συνεργάστηκα κατά την διάρκεια της εκπόνησης της διατριβής μου. Ιδιαίτερες ευχαριστίες ωστόσο θα ήθελα να απευθύνω στους συναδέλφους και φίλους Μαγδαληνή Καρδαρά , Κωνσταντίνο Τσερπέ και Φώτη Αίσωπο.

Τέλος, θα ήθελα να ευχαριστήσω την οικογένειά μου και τους φίλους μου για την στήριξη τους όλα αυτά τα χρόνια.

Αθανάσιος Δ. Παπαοικονόμου

Απρίλιος 2015

Πίνακας περιεχομένων

1	Εισαγωγή	1
1.1	Ανάλυση προσημασμένων κοινωνικών δικτύων	3
1.1.1	<i>Θεωρία δομικής ισορροπίας</i>	3
1.1.2	<i>Θεωρία της κοινωνικής Θέσης (Status Theory)</i>	9
1.2	Συστήματα Προτάσεων και Εμπιστοσύνη	10
1.3	Οργάνωση του εγγράφου	10
2	Η επίδραση των αρνητικών απόψεων	13
2.1	Βασική Ορολογία.....	13
2.2	Παρουσίαση της μεθόδου	17
2.3	Αξιολόγηση της μεθόδου.....	19
2.3.1	<i>Περιγραφή των δεδομένων</i>	19
2.3.2	<i>Προβλεπτική απόδοση</i>	20
2.3.3	<i>Δυνατότητες Γενίκευσης της μεθόδου σε πολλαπλά δίκτυα</i>	23
2.4	Η δύναμη των Αρνητικών Ακμών	24
2.5	Σύνοψη Κεφαλαίου.....	27
3	Πρόβλεψη Συναισθήματος με την χρήση Συχνών Υπογράφων	29
3.1	Βασική Ορολογία.....	29
3.2	Προτεινόμενη Προσέγγιση	33
3.3	Πειραματική αξιολόγηση.....	36
3.3.1	<i>Περιγραφή των δεδομένων</i>	36
3.3.2	<i>Εργαλεία και μετρικές αξιολόγησης</i>	37
3.3.3	<i>Ακρίβεια Μοντέλου και Ανάλυση</i>	38
3.4	Σύνοψη Κεφαλαίου.....	41
4	Εύρεση Δεσμών Εμπιστοσύνης με χρήση Τεχνικών Deep Learning.....	43
4.1	Διατύπωση του προβλήματος	44

4.2	Μέθοδος αναφοράς.....	45
4.3	Αλγόριθμοι και Εργαλεία	46
4.3.1	<i>Restricted Boltzmann Machines (RBMs)</i>	47
4.3.2	<i>Αυτοκωδικοποιητές (Autoencoders)</i>	53
4.4	Προτεινόμενη μέθοδος.....	56
4.5	Αξιολόγηση των αποτελεσμάτων	59
4.5.1	<i>Περιγραφή των δεδομένων</i>	59
4.5.2	<i>Εργαλεία και μετρικές αξιολόγησης</i>	59
4.5.3	<i>Ανάλυση ακρίβειας μοντέλου</i>	60
4.5.4	<i>Αυτοκωδικοποιητές versus Support Vector Machines (SVMs)</i>	62
4.6	Σύνοψη κεφαλαίου.....	63
5	Εύρεση κοινοτήτων σε κοινωνικά δίκτυα	65
5.1	Εισαγωγή	66
5.2	Διατύπωση του Προβλήματος	68
5.3	Μέθοδοι Αναφοράς.....	69
5.4	Προτεινόμενη Προσέγγιση	70
5.4.1	<i>Χώροι σημασιολογικών διανυσμάτων και GloVe</i>	70
5.4.2	<i>Από το Glove στο GloveGraph</i>	72
5.5	Πειραματική Επαλήθευση	77
5.5.1	<i>Περιγραφή των δεδομένων</i>	77
5.5.2	<i>Μέθοδοι Αναφοράς</i>	78
5.5.3	<i>Μετρικές αξιολόγησης</i>	78
5.5.4	<i>Αξιολόγηση του μοντέλου</i>	79
5.6	Σύνοψη του κεφαλαίου	84
6	Σύνοψη	85
6.1	Συνεισφορά – Καινοτομία	85
6.2	Μελλοντική Εργασία	87
	Βιβλιογραφικές Αναφορές.....	90

Σχήματα - Εικόνες

Εικόνα 3.1 α) Αριστερά – Ένα κομμάτι του κοινωνικού γράφου. Η ακμή $s \rightarrow t$ θα λειτουργήσει ως ακμή-γεννήτορας β) Δεξιά – Το εξαγόμενο εγω-γράφημα. Το εγω-γράφημα είναι ‘θετικό’ όπως το πρόσημο της ακμής-γεννήτορα	31
Εικόνα 3.2 Λειτουργία της canonical labeling συνάρτησης. Αριστερά φαίνεται ο γράφος που δίνεται σαν είσοδος στην συνάρτηση. Στο πρώτο βήμα, παράγονται όλοι οι ισοδύναμοι πίνακες γειτνίασης και υπολογίζονται οι αντίστοιχοι κώδικες. Δεξιά, γίνεται η επιλογή του κωδικού	32
Εικόνα 3.3 Τα τέσσερα βήματα του αλγορίθμου. Πρώτα παραγάγουμε το εγω-γράφημα κάθε ακμής και βρίσκουμε τους συχνούς υπογράφους. Έπειτα, δημιουργούμε έναν πίνακα δεδομένων, όπου κάθε γραμμή δείχνει τις συχνότητες εμφάνισης κάθε υπογράφου ανά εγω-γράφημα.	34
Εικόνα 3.4 Ο αριθμός των διαφορετικών υπογράφων για κάθε δίκτυο και διαφορετική τιμή κατωφλίου. Αριστερά, οι γραφικές παραστάσεις για υπογράφους μεγέθους 3 και δεξιά για μέγεθος 4	40
Εικόνα 3.5 Χαρακτηριστικοί υπογράφοι που συνδέονται με θετικές και αρνητικές ακμές-γεννήτορες	42
Εικόνα 4.1 Αναπαράσταση ενός Restricted Boltzmann Machine	48
Εικόνα 4.2 Εκμάθηση RBM μέσω Contrastive Divergence	51
Εικόνα 4.3 Δίκτυο αυτοκωδικοποιητή. Τα δεδομένα \mathbf{x}^i περνούν μέσα από το κρυμμένο επίπεδο και αναδύονται εκ νέου στην έξοδο (\mathbf{x}^i)	54
Εικόνα 4.4 Η μέθοδος μας εκτελείται σε δύο στάδια. α) Εύρεση κωδικών για τους χρήστες: Εκπαιδεύουμε το CF-RBM βάσει των κριτικών των χρηστών για τα προϊόντα. Ένα ξεχωριστό RBM αντιστοιχεί σε κάθε χρήστη. Μετά το πέρας της εκπαίδευσης, υπολογίζουμε την ενεργοποίηση των μονάδων στον κρυμμένο επίπεδο, β) Φάση ταξινόμησης: Λαμβάνουμε τον κωδικό κάθε ακμής, συνδυάζοντας τους κώδικες των χρηστών που βρίσκονται στα άκρα της. Το αποτέλεσμα τροφοδοτείται στο δίκτυο του αυτοκωδικοποιητή για την δημιουργία του ταξινομητή	56
Εικόνα 4.5 Παράδειγμα manifold χαμηλών διαστάσεων, στο οποίο αναπαριστώνται περιοχές θετικού και αρνητικού προσήμου. Κάθε κουκκίδα αντιστοιχεί στο κώδικα ακμής που βρέθηκε με την μέθοδο μας	63
Εικόνα 5.1 Στατιστικά των δεδομένων	77

Πίνακες

Πίνακας 3-1 Στατιστικά των συλλογών δεδομένων.....	37
Πίνακας 3-2 Μετρήσεις AUC για όλα τα πειράματα και για τα τρία δίκτυα, με δύο τιμές για το μέγεθος του υπογράφου και επτά διαφορετικές τιμές για το κατώφλι ελάχιστης συχνότητας.....	39
Πίνακας 4-1 Μετρήσεις AUC για τις διάφορες αρχιτεκτονικές των στοιβαγμένων αυτοκωδικοποιητών, και για διαφορετικό ποσοστό των εξ' αρχής γνωστών προσημασμένων ακμών . Η αρχιτεκτονική αναφέρεται στον κομμάτι του 'κωδικοποιητή'. Η τελευταία σειρά του πίνακα δείχνει του πίνακα δείχνει τις επιδόσεις της μεθόδου-αναφοράς	60
5-1 Στατιστικά στοιχεία των δεδομένων.....	77
5-2 Μετρήσεις Micro-F1 στο BlogCatalog.....	80
5-3 Μετρήσεις Macro-F1 στο BlogCatalog.....	81
5-4 Μετρήσεις Micro-F1 στο Flickr	81
5-5 Μετρήσεις Macro-F1 στο Flickr	82

Περίληψη

Η παρούσα διατριβή προτείνει τεχνικές για την ανάλυση κοινωνικών δικτύων δίνοντας ιδιαίτερη έμφαση σε δίκτυα στα οποία οι χρήστες μπορούν να εκφράζουν εμπιστοσύνη ή δυσπιστία μεταξύ τους. Η ανάλυση τέτοιων γράφων εμπιστοσύνης είναι ένα ενδιαφέρον πρόβλημα με ευρύ φάσμα εφαρμογών όπως η ανάλυση γεωπολιτικών σχέσεων και η εύρεση κοινοτήτων χρηστών.

Στα πρώτα τρία κεφάλαια εξετάζουμε το πρόβλημα της πρόβλεψης της προδιάθεσης ενός χρήστη για έναν άλλο, αντλώντας τεχνικές από τρεις διαφορετικούς τομείς. Αρχικά, χρησιμοποιούμε κλασικές και διαδοσμένες τεχνικές από τον χώρο της Ανάλυσης Κοινωνικών Δικτύων (Social Network Analysis) με σκοπό να ερευνήσουμε τους μηχανισμούς διάδοσης θετικών και αρνητικών απόψεων στο δίκτυο. Έπειτα, ενσωματώνουμε τεχνικές από τον τομέα της Βιοστατιστικής, ώστε να αναλύσουμε μεγάλα κοινωνικά δίκτυα από μικροσκοπική σκοπιά. Στη συνέχεια, με χρήση τεχνικών deep learning δείχνουμε πως είναι δυνατόν να "κατασκευαστεί" ένας γράφος εμπιστοσύνης αξιοποιώντας δεδομένα φαινομενικά άσχετα με αυτόν τον σκοπό, όπως οι κριτικές των χρηστών για διάφορα προϊόντα. Στο τελευταίο κεφάλαιο, παρουσιάζουμε έναν αλγόριθμο εύρεσης κοινοτήτων σε κοινωνικά δίκτυα, βασιζόμενοι σε πρόσφατες προόδους στον τομέα της Ανάλυσης Φυσικής Γλώσσας (Natural Language Processing).

Η σειρά των κεφαλαίων αποτυπώνει την χρονική σειρά των πειραμάτων που εφαρμόσαμε αλλά κάθε κεφάλαιο είναι γραμμένο ώστε να μην έχει σημαντικές συσχετίσεις με τα προηγούμενα και έτσι να μπορεί να διαβαστεί αυτόνομα.

Abstract

This PhD thesis presents novel techniques in the domain of social network analysis, by putting more emphasis on networks in which users can express trust or distrust on each other. Mining such networks is an interesting problem with a variety of application domains like the analysis of international relationships and the detection of user communities.

The first three chapters examine the sign prediction problem by borrowing techniques from three different domains. At first, we employed standard , well-defined measures from the field of Social Network Analysis to study the propagation patterns of the positive and negative opinions in the social graph. Then, we adapted state of the art techniques from the domain of Biostatistics in order to perform microscopic network analysis of large social networks. Next, we investigated the possibility to reconstruct a social trust graph by exploiting external data such as the ratings of the users for certain items, using algorithms from the domain of deep learning. Finally, the last chapter presents a community detection algorithm based on recent advances in Natural Language Processing.

The sequence of the chapters depicts the chronological order of the experiments. I tried to minimize the dependencies among the chapters in order to facilitate the interested reader.

Η σελίδα αυτή είναι σκόπιμα λευκή.

1

Εισαγωγή

Τα τελευταία χρόνια, η χρήση των εφαρμογών κοινωνικής δικτύωσης έχει σημειώσει σημαντική άνοδο με τους χρήστες να αφιερώνουν ένα σημαντικό κομμάτι του ελεύθερου χρόνου τους ώστε να σχολιάσουν θέματα της επικαιρότητας, να μοιραστούν τις ιδέες τους και να δημιουργήσουν δεσμούς φιλίας μεταξύ τους. Για ένα μεγάλο μέρος των χρηστών η ενασχόληση με τους ιστοτόπους κοινωνικής δικτύωσης ξεπερνά το όριο της συνήθειας και σε αρκετές περιπτώσεις, λαοφιλείς εφαρμογές όπως το Facebook¹, αποτελούν την βασική οδό μετάδοσης της πληροφορίας. Η ηλεκτρονική βιομηχανία δεν θα μπορούσε να παραμείνει απαθής σε αυτή την εξέλιξη, γνωρίζοντας ότι οι πελάτες τους βασίζονται στις αποφάσεις τους για μελλοντικές αγορές στα σχόλια και τις κριτικές άλλων χρηστών. Οι επιχειρήσεις ηλεκτρονικού εμπορίου (online retailers) κατάλαβαν γρήγορα τον ρόλο της επιρροής των ομότιμων χρηστών (peer influence) στην διαμόρφωση των αγοραστικών συνηθειών των πελατών τους. Το προσωπικό δίκτυο ενός χρήστη θεωρείται πλέον ισοδύναμο με το ιστορικό αγορών του, και για αυτό το λόγο εταιρίες όπως η Amazon² ψάχνουν τρόπους ώστε να το ενσωματώσουν στα προγνωστικά τους μοντέλα (predictive models). Τέτοιοι

¹ <https://www.facebook.com>

² <http://www.amazon.com/>

μηχανισμοί βοηθούν ώστε να αντιμετωπιστούν προβλήματα όπως το cold-start problem : Στην περίπτωση ενός καινούριου χρήστη , ο οποίος δεν έχει ιστορικό αγορών, το σύστημα προτάσεων (recommender systems) θα μπορούσε να προτείνει προϊόντα βάσει των προτιμήσεων των «φίλων» του σε κάποιο κοινωνικό δίκτυο.

Από τα παραπάνω, γίνεται εμφανές ότι το προσωπικό κοινωνικό δίκτυο ενός χρήστη είναι ένα παραπάνω στοιχείο της προσωπικότητας του, το οποίο μπορεί να βοηθήσει σημαντικά στην μοντελοποίηση της συμπεριφοράς του. Σε αυτό το πλαίσιο, ο ρόλος της εμπιστοσύνης, ανάμεσα στους χρήστες, αναδύεται φυσικά. Τα περισσότερα σύγχρονα κοινωνικά δίκτυα, όπως το Facebook και το Twitter³, απεικονίζουν τις σχέσεις μεταξύ των χρηστών μέσω απλών γράφων γνωριμιών (acquaintance networks), όπου κάθε άτομο "συνδέεται" με κάποιο άλλο με τη δημιουργία του αντίστοιχου δεσμού. Σε αυτά τα δίκτυα, κάθε ακμή του κοινωνικού γράφου θεωρείται ένδειξη φιλίας ή εμπιστοσύνης, αφού δίνει δικαιώματα πρόσβασης σε μεγαλύτερο όγκο προσωπικών δεδομένων (π.χ. φάκελοι φωτογραφιών που είναι ορατοί μόνο σε φίλους κ.α.). Παρ' όλα αυτά υπάρχουν κοινωνικά δίκτυα που επιτρέπουν στους χρήστες να υποδείξουν άλλους χρήστες είτε ως φίλους είτε ως ανεπιθύμητους. Χαρακτηριστικά παραδείγματα σε αυτήν την κατηγορία είναι το τεχνολογικό blog Slashdot , το οποίο επιτρέπει στους χρήστες του να «χρωματίζουν» τις μεταξύ τους σχέσεις με ταμπέλες όπως friend or foe, καθώς και το Epinions , το οποίο είναι ένας ιστοτόπος ανταλλαγής κριτικών για προϊόντα, στο οποίο οι χρήστες μπορούν να κατευθύνουν δεσμούς εμπιστοσύνης/δυσπιστίας δημιουργώντας ένα προσημασμένο κοινωνικό γράφο.

³ <https://twitter.com/>

Για το υπόλοιπο του κεφαλαίου θα παρουσιάσουμε τις δύο δημοφιλέστερες θεωρίες για την επεξήγηση δικτύων φιλίας/εχθρότητας: τη θεωρία δομική ισορροπίας (structural balance theory) και την θεωρία κοινωνικής θέσης (social status theory).

1.1 Ανάλυση προσημασμένων κοινωνικών δικτύων

1.1.1 Θεωρία δομικής ισορροπίας

Η θεωρία της δομικής ισορροπίας είναι η παλιότερη και πιο θεμελιωμένη πειραματικά θεωρία για την επεξήγηση δικτύων φιλίας/εχθρότητας. Καθιερώθηκε από τον Heider [1] το 1940, ενώ αργότερα εκφράστηκε σε γραφο-θεωρητική γλώσσα από τους Cartwright και Harary [2]. Αναπτύχθηκε αρχικά για πλήρη κοινωνικά δίκτυα (δηλαδή δίκτυα όπου υπάρχει πάντα ένας δεσμός μεταξύ δύο οποιονδήποτε χρηστών) των οποίων οι ακμές είναι συμμετρικές. Βάσει της θεωρίας οι χρήστες χωρίζονται σε μικρές ομάδες και εξετάζονται οι διαφορετικές διατάξεις που προκύπτουν βάσει των προτιμήσεων τους για άλλους χρήστες εντός της ίδιας ομάδας. Ο σκοπός της θεωρίας της δομικής ισορροπίας είναι να δείξει ότι κάποιες διατάξεις είναι πιο πιθανές από άλλες. Συγκεκριμένα, υποστηρίζει ότι οι τριάδες χρηστών με περιττό αριθμό θετικών ακμών είναι πιο σταθερές και γι' αυτό το λόγο θα πρέπει να εμφανίζονται με μεγαλύτερη συχνότητα στα δεδομένα. Η Εικόνα 1 παρουσιάζει τις τέσσερις δυνατές διατάξεις μεταξύ τριών χρηστών. Η περίπτωση α) θεωρείται σταθερή διάταξη γιατί αντανακλά την περίπτωση μιας παρέας τριών ατόμων. Κατά τον ίδιο τρόπο και η διάταξη β) θεωρείται φυσική μιας και απεικονίζει δύο ομάδες ατόμων ($\{A,B\}$ και $\{C\}$) που είναι εχθρικές μεταξύ τους. Κατά την θεωρία, η περίπτωση γ) με δύο θετικές και μία αρνητική ακμή είναι η πιο ασταθής γιατί εισάγει ένα «άγχος» ή «αγωνία» στις σχέσεις των χρηστών. Ο χρήστης A είναι φίλος με τον

B και τον C, αλλά οι B και C δεν τα πάνε καλά μεταξύ τους, πράγμα που θα αναγκάσει τον A είτε να προσπαθήσει να τους συμβιβάσει, μετατρέποντας την ακμή B-C σε θετική, είτε να επιλέξει ως φίλο μόνο έναν από τους δύο καταλήγοντας εκ νέου σε μια σταθερότερη διάταξη. Τέλος η περίπτωση δ) χαρακτηρίζεται και αυτή ως ασταθής, στη λογική ότι για δύο άτομα με κοινό έχθρο υπάρχουν οι προϋποθέσεις ώστε να συνεργαστούν εναντίον του. Η θεωρία της δομικής ισορροπίας προβλέπει διασθητικά έγκυρους κοινωνικές κανόνες όπως «ο φίλος του φίλου μου είναι φίλος» και «ο φίλος του εχθρού μου είναι εχθρός».

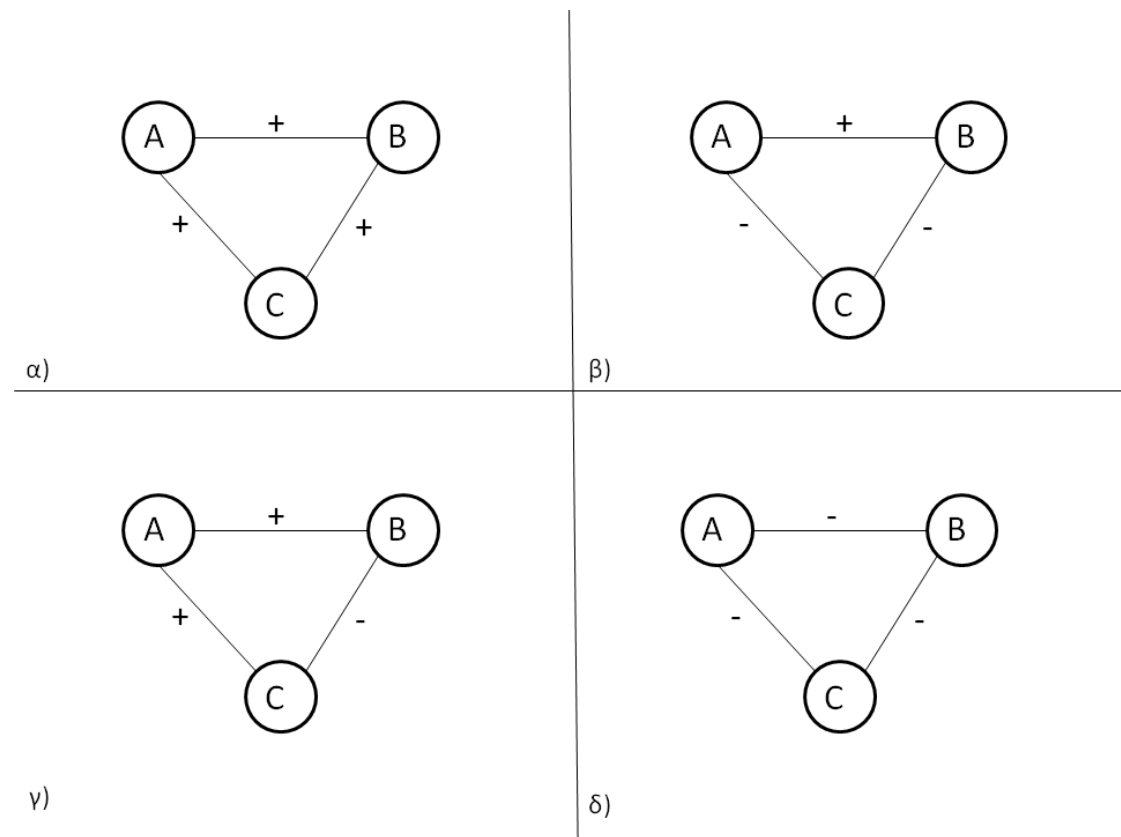
1.1.1.1 Δομή ισορροπημένων κοινωνικών γράφων

Μια ενδιαφέρουσα πτυχή της θεωρίας δομικής ισορροπίας αφορά στον τρόπο πως απλοί τοπικοί κανόνες (τριάδες χρηστών) μπορούν να προβλέψουν την δομή του δικτύου μακροσκοπικά. Ένας προσημασμένος κοινωνικός γράφος θεωρείται ισορροπημένος αν όλες οι δυνατές τριάδες μεταξύ των χρηστών είναι σταθερές, δηλαδή περιέχουν ένα περιττό αριθμό θετικών ακμών. Όπως αναφέρθηκε στην προηγούμενη ενότητα, υπάρχουν δύο σταθερές διατάξεις τριάδων χρηστών. Στην πρώτη (Εικόνα 1.1 α) οι χρήστες είναι όλοι οι φίλοι μεταξύ τους και έτσι για να προκύψει ένα ισορροπημένο δίκτυο συνολικά, θα πρέπει κάθε χρήστης να είναι φίλος με οποιονδήποτε άλλο χρήστη. Η περίπτωση αυτή βέβαια είναι απλή και λιγότερο ρεαλιστική αφού δεν υπάρχει καμμία αρνητική ακμή στο δίκτυο. Πιο ενδιαφέρουσες προεκτάσεις παίρνει η δεύτερη σταθερή διάταξη (Εικόνα 1.1β) όπου σύμφωνα με την θεωρία, το αντίστοιχο κοινωνικό δίκτυο για να είναι ισορροπημένο θα πρέπει να χαρακτηρίζεται από δύο μεγάλες ομάδες χρηστών X και Y , όπου θα υπάρχουν θετικοί δεσμοί μεταξύ όλων των ατόμων εντός ομάδος και αρνητικοί δεσμοί μεταξύ χρηστών διαφορετικής ομάδας. Το Θεώρημα της Ισορροπίας από τον Hagary [2] [3]

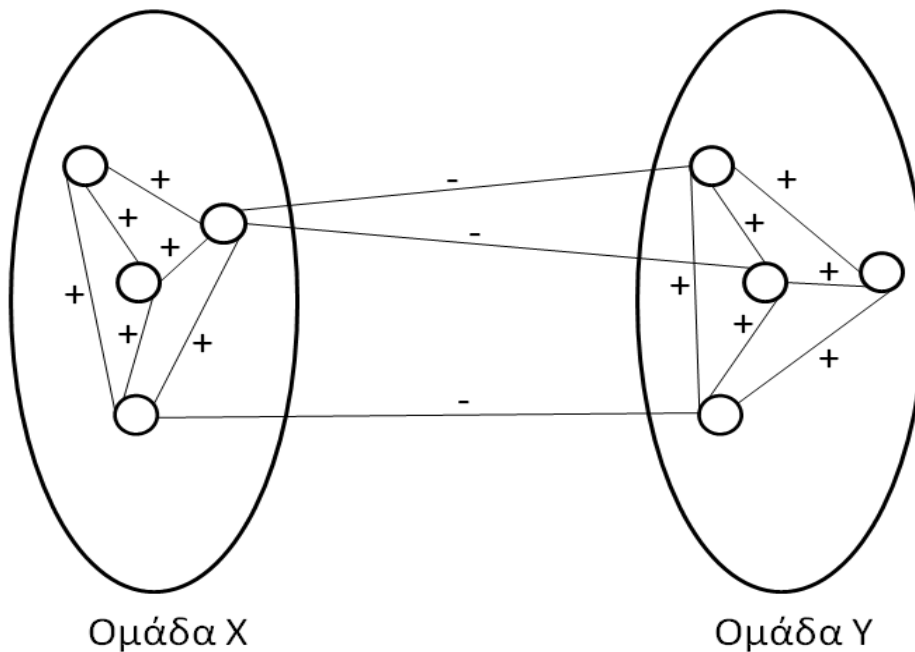
το 1953 ενσωματώνει σε έναν κοινό ορισμό όλες τις δυνατές εκφάνσεις ενός ισορροπημένου, προσημασμένου, κοινωνικού δικτύου και διατυπώνεται ως εξής:

Ορισμός (Ισορροπημένο κοινωνικό δίκτυο)

Ένας προσημασμένος πλήρης γράφος είναι ισορροπημένος αν, είτε όλοι οι χρήστες είναι φίλοι μεταξύ τους είτε μπορούν να χωριστούν σε δύο ομάδες X και Y, τέτοια ώστε όλα τα μέλη του X να είναι φίλοι μεταξύ τους, όλα τα μέλη του Y να είναι φίλοι μεταξύ τους και κάθε μέλος του X να τάσσεται αντίθετα σε κάθε μέλος του Y.



Εικόνα 1.1 Τέσσερις διαφορετικές διατάξεις τριάδων. Οι α) β) διαθέτουν έναν περιττό αριθμό θετικών ακμών και θεωρούνται "σταθερές" βάσει της θεωρίας δομικής ισορροπίας σε αντίθεση με τις γ) δ) που θεωρούνται "ασταθείς"



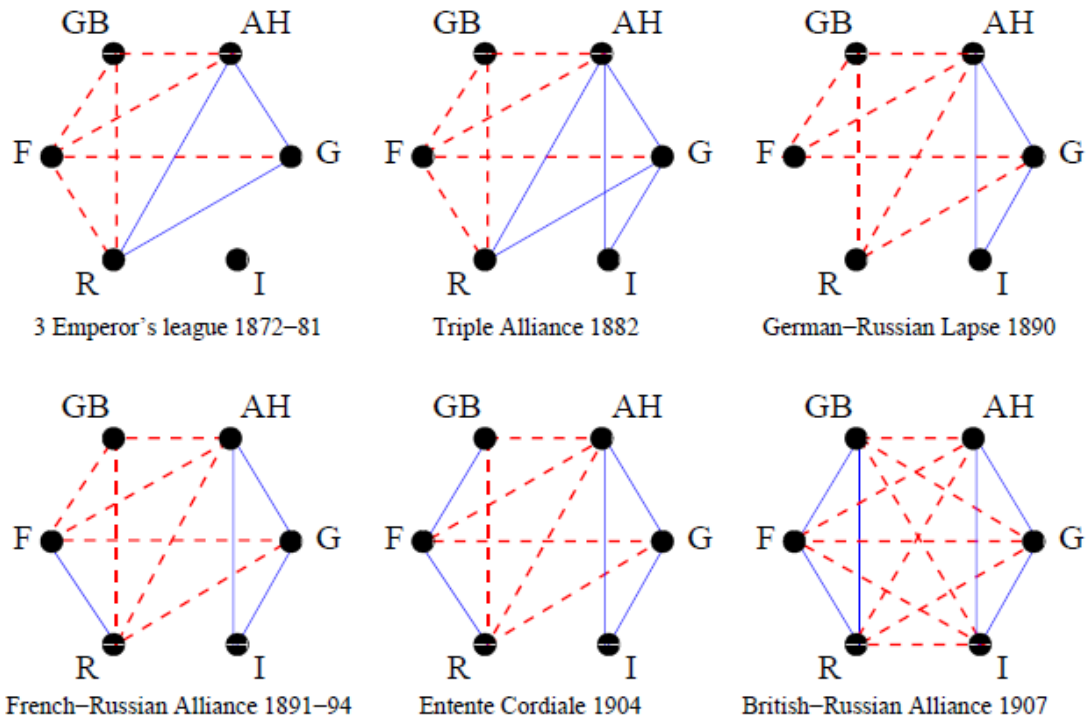
Εικόνα 1.2. Ένας ισορροπημένος κοινωνικός γράφος μπορεί να χωριστεί σε δύο αντιμαχόμενες ομάδες χρηστών. Οι σχέσεις ενός ομάδες είναι (αποκλειστικά) σχέσεις εμπιστοσύνης ενώ αρνητικές ακμές συνδέουν χρήστες διαφορετικών ομάδων

Η θεωρία της δομικής ισορροπίας είναι δύσκολο να εφαρμοστεί με απόλυτη καθαρότητα σε πραγματικά κοινωνικά δίκτυα, ώστε να τα χαρακτηρίσουμε ισορροπημένα ή όχι, λόγω του είναι δύσκολο να πληρωθούν όλες οι προϋποθέσεις της. Πιο συγκεκριμένα:

- Οι γράφοι από πραγματικά κοινωνικά δίκτυα δεν είναι πλήρεις, δηλαδή δεν υπάρχει πάντα μια ακμή να συνδέει δύο οποιουδήποτε χρήστες. Αντίθετα, ο πίνακας γειτνίασης που αντιστοιχεί σε πραγματικά δίκτυα είναι αρκετά *αραιός*.
- Είναι στατιστικά απίθανο να περιμένουμε ότι ένα πραγματικό, προσημασμένο δίκτυο θα περιέχει μόνο σταθερές τριάδες.

Για να ξεπεράσουμε το πρώτο πρόβλημα, είτε παίρνουμε υπ' όψιν μόνο τις τριάδες των χρηστών που εμφανίζονται στα δεδομένα είτε ακολουθούμε μια διαδικασία δημιουργίας δεσμών έτσι ώστε το τελικό αποτέλεσμα να είναι ένας πλήρης και ισορροπημένος κοινωνικός γράφος. Η πρώτη λύση είναι υπολογιστικά "φθηνότερη" και γι' αυτό συνήθως χρησιμοποιείται. Για το δεύτερο ζήτημα, χαλαρώνουμε τον περιορισμό που θέλει κάθε δυνατή τριάδα χρηστών να είναι σταθερή και στη θέση του ορίζουμε ένα κατώφλι $k\%$ και εξετάζουμε εάν η αναλογία των τριάδων είναι πάνω από αυτό το όριο.

Παρόλο που η θεωρία της δομικής ισορροπίας προτάθηκε το 1940, έχει εφαρμοστεί με κάποιες μετατροπές σε σύγχρονα κοινωνικά δίκτυα με μεγάλη επιτυχία. Ίσως η πιο ενδιαφέρουσα εφαρμογή της θεωρίας δομικής ισορροπίας είναι η εξήγηση γεωπολιτικών συσχετισμών όπως αυτή αναλύεται στην εργασία των Antal et al. στο [4] και παρουσιάζεται στην Εικόνα 1.3. Οι συγγραφείς εξετάζουν το "κοινωνικό δίκτυο" των πρωταγωνιστριών-χωρών του Πρώτου Παγκοσμίου Πολέμου και πως οι σχέσεις του μεταβάλλονταν με τον καιρό, ορίζοντας νέες συμμαχίες κάθε φορά.



Εικόνα 1.3 Εξέλιξη των σχέσεων μεταξύ των πρωταγωνιστών του Πρώτου Παγκοσμίου Πολέμου (GB = Μεγάλη Βρετανία, AH = Αυστροουγγαρία, G = Γερμανία, I = Ιταλία, R = Ρωσία, F = Γαλλία). Οι συνεχείας γραμμές αφορούν θετικές ακμές ενώ οι διακεκομμένες αντιστοιχούν σε αρνητικές. Η εικόνα έχει παρθεί από την δημοσίευση των Antal et al. [4]

Η μελέτη ξεκινάει με την υπογραφή της συμφωνίας των τριών αυτοκρατοριών (Three Emperors' League, αρχικά το 1872 και ανανέωση το 1881) η οποία συνδέει τις χώρες της Γερμανίας, της Αυστροουγγαρίας και της Ρωσίας. Έπειτα, η συμφωνία της Τριπλής Συμμαχίας (Triple Alliance), η οποία υπεγράφη το 1882, βάσει της οποίας η Γερμανία, η Αυστροουγγαρία και η Ιταλία συνασπίστηκαν σε ένα κοινό μπλοκ το οποίο συνεχίστηκε μέχρι και τον Πρώτο Παγκόσμιο Πόλεμο. Το 1890, η διμερής συμφωνία μεταξύ Γερμανίας και Ρωσίας έληξε οδηγώντας σε μια Γάλλο-Ρωσική συμμαχία την περίοδο 1891-1894. Στη συνέχεια, η σύναψη της Αντάντ μεταξύ της Γαλλίας και της Μεγάλης Βρετανίας το 1904, και η διμερής συμφωνία Μεγάλης

Βρετανίας-Ρωσίας το 1907 οδήγησε στην διαμόρφωση της Τριπλής Ανταντ που συνέδεε την Γαλλία , την Μεγάλη Βρετανία και την Ρωσία σε ένα κοινό μέτωπο.

1.1.2 Θεωρία της κοινωνικής Θέσης (Status Theory)

Μια εναλλακτική θεωρία για την ανάλυση δικτύων εμπιστοσύνης εισήχθη από τους Leskovec et al στο [5] και στηρίζεται στην θεωρία της κοινωνικής θέσης (status theory). Αναλυτικότερα, οι συγγραφείς θεωρούν ότι οι προσημασμένες ακμές σε έναν γράφο εμπιστοσύνης δεν είναι απλά ενδείξεις φιλίας ή εχθρότητας αλλά ταυτόχρονα περιγράφουν μια ιεραρχία μεταξύ των χρηστών η οποία μπορεί να εξηγηθεί βάσει του status του κάθε χρήστη: Μια θετική (αρνητική) ακμή από έναν χρήστη u σε κάποιον χρήστη v , σημαίνει ότι ο χρήστης u θεωρεί τον v υψηλότερου (χαμηλότερου) status από εκείνον. Έτσι το πρόβλημα της πρόβλεψης του προσήμου μιας ακμής μεταξύ δύο οποιονδήποτε χρηστών ανάγεται στην εύρεση του κατάλληλου status για κάθε χρήστη. Οι ίδιοι συγγραφείς στο [6] επεκτείνουν αυτήν την ιδέα, παρουσιάζοντας ένα μοντέλο λογαριθμιστικής παλινδρόμησης το οποίο αναπαριστά κάθε ακμή του γράφου σε έναν πολυδιάστατο χώρο μεταβλητών που κατατάσσονται σε δύο κλάσεις: Στην πρώτη, εξετάζονται οι συσχετίσεις μεταξύ του εξερχόμενου βαθμού (out-degree) του κόμβου που απευθύνει έναν προσημασμένο δεσμό και του εισερχόμενου βαθμού (in-degree) του κόμβου προορισμού, ενώ στην δεύτερη, εξετάζονται ο αριθμός και ο τύπος των τριάδων που ορίζονται από τους εμπλεκόμενους κόμβους και τους κοινούς τους γείτονες.

1.2 Συστήματα Προτάσεων και Εμπιστοσύνη

Η εισαγωγή της έννοιας της εμπιστοσύνης σε συστήματα προτάσεων (recommender systems) [7] [8] [9] [10] είναι μια από τις τελευταίες εξελίξεις στο συγκεκριμένο χώρο με σκοπό να βελτιώσει την προβλεπτική ικανότητα των συστημάτων αυτών. Στην βιβλιογραφία υπάρχουν αναφορές [11] [12] [13] [14] που δείχνουν προσπάθειες να ενσωματωθεί η έννοια της εμπιστοσύνης.

Στο πλαίσιο αυτής της διδακτορικής διατριβής ασχολήθηκα με το πρόβλημα της εμπιστοσύνης είτε μέσω της απ' ευθείας πρόβλεψης του επιπέδου εμπιστοσύνης μεταξύ δύο χρηστών αξιοποιώντας κατάλληλα τον προσημασμένο κοινωνικό γράφο είτε μέσω της εκτίμησης εμπιστοσύνης βασισμένος σε "εξωτερικά" δεδομένα όπως οι κριτικές των χρηστών.

1.3 Οργάνωση του εγγράφου

Η διατριβή αποτελείται από έξι (6) κεφάλαια. Στις ενότητες των κεφαλαίων αυτών παρουσιάζεται ουσιαστικά και με αναλυτικό τρόπο το αντικείμενο της διδακτορικής διατριβής.

Το Κεφάλαιο 1, είναι εισαγωγικό και παρουσιάζει τις δύο κυρίαρχες θεωρίες στην ανάλυση προσημασμένων κοινωνικών δικτύων: τη θεωρία δομικής ισορροπίας ("structural balance theory") και τη θεωρία της Κοινωνικής θέσης ("status theory"). Επίσης παρουσιάζεται η βασική ορολογία για την ανάλυση κοινωνικών δικτύων.

Στο Κεφάλαιο 2 παρουσιάζεται ένας αλγόριθμος πρόβλεψης συναισθήματος με την χρήση παραδοσιακών μεθόδων από τον χώρο της Ανάλυσης Κοινωνικών Δικτύων (Social Network Analysis). Ως βάση του προτεινόμενου αλγορίθμου διακρίνεται η ομοφιλία [15] μεταξύ των χρηστών, δηλαδή το πως χρήστες επηρεάζουν και

επηρεάζονται από άλλους χρήστες που είναι όμοιοι με αυτούς. Το κεφάλαιο ολοκληρώνεται συγκρίνοντας την σχετική δυναμική μεταξύ των θετικών και αρνητικών απόψεων.

Στην συνέχεια το Κεφάλαιο 3 παρουσιάζεται μια νέα προσέγγιση στο πρόβλημα της πρόβλεψης προσήμου σε δίκτυα εμπιστοσύνης, η οποία βασίζεται στην εύρεση των συχνών υπογράφων. Η τεχνική έχει επιρροές από τον χώρο της Βιοστατιστικής και πιο συγκεκριμένα από την υπόθεση QSAR (Quantitative structure-activity relationship), η οποία υποστηρίζει ότι 'όμοια στοιχεία' έχουν παρόμοιες συμπεριφορές π.χ. η τοξικότητα ενός στοιχείου μπορεί να προβλεφθεί από την κατανομή των μορίων του στον τρισδιάστατο χώρο, συγκρινόμενο με στοιχεία με παρόμοια δομή και γνωστή συμπεριφορά.

Στο Κεφάλαιο 4 εξετάζεται εκ νέου το ζήτημα της πρόβλεψης προσήμου αλλά από την σκοπιά ενός προβλήματος ημι-εποπτευόμενης μάθησης. Η προτεινόμενη προσέγγιση στηρίζεται σε αλγορίθμους από τον χώρο του deep learning (Restricted Boltzmann Machines, Stacked Autoencoders). Στο κεφάλαιο αυτό παρουσιάζουμε έναν πρακτικό τρόπο κατασκευής ενός προσημασμένου κοινωνικού γράφου αξιοποιώντας εξωτερικά δεδομένα όπως οι κριτικές των χρηστών για διάφορα προϊόντα.

Στο Κεφάλαιο 5 παρουσιάζεται ένας γενικός τρόπος εύρεσης κοινοτήτων σε κοινωνικά δίκτυα αξιοποιώντας πρόσφατες προόδους στον χώρο της Επεξεργασίας Φυσικής Γλώσσας (Natural Language Processing).

Τέλος, στο Κεφάλαιο 6 περιλαμβάνεται η σύνοψη της διατριβής και τα συμπεράσματα που εξήχθησαν κατά την εκπόνησή της, η συνεισφορά και η καινοτομία που επιδεικνύει στον αντίστοιχο ερευνητικό χώρο, και συζητούνται θέματα μελλοντικής εργασίας και επέκτασης των ερευνητικών αποτελεσμάτων.

Σημειώνεται ότι στο τέλος του κειμένου της διατριβής παρουσιάζεται ένα γλωσσάριο το οποίο περιέχει τους την αντιστοιχία των βασικών όρων που χρησιμοποιήθηκαν στην διατριβή με την συντομογραφία αυτών για την διευκόλυνση του αναγνώστη. Η διατριβή ολοκληρώνεται με βιβλιογραφικές αναφορές.

2

Η επίδραση των αρνητικών απόψεων

Αυτό το κεφάλαιο αποτελεί την πρώτη από τις προσεγγίσεις που αναπτύχθηκαν στο πλαίσιο της Διδακτορικής διατριβής και αφορά στην πρόβλεψη εμπιστοσύνης/δυσπιστίας σε κοινωνικά δίκτυα. Πιο συγκεκριμένα, παρουσιάζουμε μια ανάλυση σχετικά με την επίδραση των θετικών και αρνητικών ακμών . Ο ρόλος του κεφαλαίου είναι διττός: Πρώτον, παρουσιάζουμε έναν αλγόριθμο πρόβλεψης προσήμου σε κοινωνικά δίκτυα εμπιστοσύνης/δυσπιστίας, χρησιμοποιώντας κλασσικές τεχνικές από τον χώρο της Ανάλυσης Κοινωνικών Δικτύων (Social Network Analysis). Η μέθοδος μας εμφανίζει καλύτερη απόδοση σε σχέση με τις τρέχουσες τεχνικές της βιβλιογραφίας. Δεύτερον, εξετάζουμε την σχετική επιρροή μεταξύ των θετικών και αρνητικών ακμών. Με άλλα λόγια, διερευνούμε κατά πόσο η ύπαρξη ενός τύπου ακμής επηρεάζει την ακρίβεια του μοντέλου.

2.1 Βασική Ορολογία

Για την ανάπτυξη της μεθόδου μας αργότερα θα εισάγουμε χρήσιμη ορολογία από τον τομέα της Ανάλυσης Κοινωνικών Δικτύων (Social Network Analysis) . Ξεκινάμε

με την έννοια του εγω-δικτύου(ego-network) [16] ενός χρήστη κοινωνικού δικτύου, το οποίο είναι το "κομμάτι" του κοινωνικού γράφου που βρίσκεται κοντά στον χρήστη:

Ορισμός (εγω-δίκτυο ενός χρήστη)

Το εγω-δίκτυο ενός χρήστη περιλαμβάνει τον κόμβο του γράφου που αντιστοιχεί σε αυτόν, τους άμεσους γείτονές του, καθώς και το σύνολο των δεσμών που τον συνδέουν με αυτούς.

Για λόγους που θα γίνουν κατανοητοί αργότερα, επεκτείνουμε τον ορισμό αυτό ώστε να περιλαμβάνει γείτονες-κόμβους που βρίσκονται μέχρι και k βήματα "μακριά" από τον υπό εξέταση κόμβο.

Ορισμός (Διευρυμένο εγω-δίκτυο ενός χρήστη)

Το Διευρυμένο εγω-δίκτυο ενός χρήστη περιλαμβάνει τον κόμβο του γράφου που αντιστοιχεί σε αυτόν, τους γείτονές που βρίσκονται k ή λιγότερα βήματα από αυτόν, καθώς και το σύνολο των δεσμών που συνδέουν όλους αυτούς του κόμβους μεταξύ τους. Το διευρυμένο εγω-δίκτυο ενός χρήστη u συμβολίζεται ως $ego_k(u)$.

Το κόστος εξαγωγής ενός διευρυμένου εγω-δικτύου είναι ισοδύναμο με την κόστος εκτέλεσης μιας Breadth First Algorithm (BFS) αναζήτησης, ορίζοντας ως ρίζα τον υπο-εξέταση κόμβο και εκτελώντας μέχρι και βάθος k , και συνεπώς είναι ίσο με $O(b^k)$ όπου b ο μέσος βαθμός των κόμβων του γράφου. Ο επόμενος ορισμός εισάγει τέσσερις γραφικούς τελεστές, μέσω των οποίων μπορούμε να ομαδοποιήσουμε τους γείτονες ενός κόμβου σε έναν κατευθυνόμενο προσημασμένο γράφο.

Ορισμός (Τελεστές)

Έστω u ένας χρήστης κοινωνικού δικτύου. Δηλώνουμε ως $\Gamma_{in}^+(u)$ το σύνολο των κόμβων που απευθύνουν μια θετική ακμή στον κόμβο u , και ως $\Gamma_{out}^+(u)$ το σύνολο των κόμβων που δέχονται μια θετική ακμή από τον κόμβο u . Παρόμοια ορίζουμε και τους τελεστές $\Gamma_{in}^-(u)$ και $\Gamma_{out}^-(u)$ για τις περιπτώσεις των αρνητικών ακμών.

Για το υπόλοιπο αυτής της ενότητας εισαγάγουμε μετρικές ομοιότητας (similarity metrics) μεταξύ χρηστών κοινωνικών δικτύων. Τέτοιες μετρικές χρησιμοποιούνται κατά κόρον στον χώρο της Ανάλυσης Κοινωνικών Δικτύων γιατί εκφράζουν με ποσοτικό τρόπο την έννοια της ομοφιλίας (homophily) στα κοινωνικά δίκτυα, η οποία δηλώνει ότι τα άτομα σε ένα κοινωνικό δίκτυο τείνουν να συγκεντρώνονται κοντά σε άλλα άτομα με παρόμοια συμπεριφορά. Μια ισοδύναμη διατύπωση θα ήταν ότι άτομα με παρόμοιο ιστορικό προτιμήσεων, τείνουν να εμφανίζουν συσχετίσεις και στο μέλλον, το οποίο αποτελεί την λογική των σύγχρονων συστημάτων προτάσεων (recommender systems). Οι μετρικές που διερευνηθούν στα πλαίσια αυτού του κεφαλαίου είναι οι εξής: Κοινοί Γείτονες [17], συντελεστής Jaccard [17] και Adamic/Adar [18]. Ξεκινούμε με τον ορισμό της μετρικής Κοινών Γειτόνων:

Ορισμός (Μετρική - Κοινοί Γείτονες)

Έστω u, v δυο κόμβοι του κοινωνικού δικτύου. Η μετρική των Κοινών Γειτόνων, που συμβολίζεται με $common(u, v)$ ορίζεται ως:

$$common(u, v) = \sum_{\Gamma} |\Gamma(u) \cap \Gamma(v)|,$$

όπου $\Gamma = \{\Gamma_{in}^+, \Gamma_{out}^+, \Gamma_{in}^-, \Gamma_{out}^-\}$ και $|\cdot|$ η πληθικότητα του συνόλου.

Ο συντελεστής Jaccard είναι μια διαδεδομένη μετρική στον χώρο της Ανάκτησης Πληροφορίας, και εκφράζει την πιθανότητα δύο οντότητες x, y να μοιράζονται ένα χαρακτηριστικό f το οποίο διαθέτει ο x ή ο y . Για τα προσημασμένα κοινωνικά δίκτυα που εξετάζουμε, η μετρική αυτή μεταφράζεται ως:

Ορισμός (Μετρική - Συντελεστής Jaccard)

Εστω u, v δυο κόμβοι του κοινωνικού δικτύου. Ο συντελεστής Jaccard, που συμβολίζεται με $Jaccard(u, v)$, ορίζεται ως:

$$Jaccard(u, v) = \frac{\sum_r |\Gamma(u) \cap \Gamma(v)|}{\sum_r |\Gamma(u) \cup \Gamma(v)|}$$

όπου $\Gamma = \{\Gamma_{in}^+, \Gamma_{out}^+, \Gamma_{in}^-, \Gamma_{out}^-\}$ και $|\cdot|$ η πληθικότητα του συνόλου.

Τέλος, παραθέτουμε και την μετρική Adamic/Adar [18] η οποία, όπως και οι προηγούμενες, λαμβάνει υπόψη τα κοινά χαρακτηριστικά δύο οντοτήτων, δίνοντας μεγαλύτερο βάρος στα πιο σπάνια και ορίζεται ως εξής :

Ορισμός (Μετρική - Adamic/Adar)

Εστω u, v δυο κόμβοι του κοινωνικού δικτύου. Ο συντελεστής Adamic/Adar, συμβολίζεται ως $AdamicAdar(u, v)$, και υπολογίζεται από την σχέση:

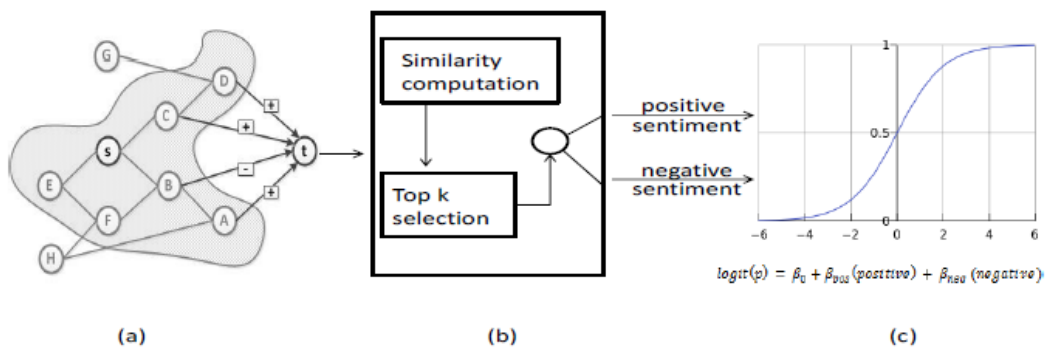
$$AdamicAdar(u, v) = \sum_{t \in \Gamma(u) \cap \Gamma(v)} 1/\log (degree(t))$$

όπου $degree(t) = |\Gamma_{in}^+(t)| + |\Gamma_{out}^+(t)| + |\Gamma_{in}^-(t)| + |\Gamma_{out}^-(t)|$ είναι ο συνολικός βαθμός του κόμβου t .

Αξίζει να σημειωθεί ότι ο υπολογισμός και των τριών μετρικών απαιτεί την εύρεση των κοινών γειτόνων των υπό εξέταση κόμβων. Ισοδύναμα, ένας κόμβος έχει **μη**

μηδενική τιμή "ομοιότητας" μόνο για τους κόμβους που βρίσκονται στο διευρυμένο εγω-δίκτυο βάθους δύο. Επιπλέον, οι μετρικές είναι συμμετρικές, μια ιδιότητα χρήσιμη για την ανάλυση μεγάλης κλίμακας δεδομένων, αφού μετά τον υπολογισμό τους μπορούν να χρησιμοποιηθούν και για τους δύο υπό εξέταση κόμβους. Τέλος, σημειώνεται ότι αυτές οι μετρικές επιλέχθηκαν λόγω της καλής τους ισορροπίας μεταξύ απλότητας και ικανοποιητικής απόδοσης στο πρόβλημα της πρόβλεψης προσήμου που εξετάζουμε.

2.2 Παρουσίαση της μεθόδου



Εικόνα 2.1. Σύνοψη της μεθόδου μας: a) Εξαγωγή του ego^k υπογράφου του κόμβου-πηγή και εύρεση των κόμβων που απευθύνουν μια ακμή(θετική ή αρνητική) προς τον κόμβο-προορισμό. b) Υπολογισμός ομοιοτήτων. c) Εκπαίδευση ενός ταξινομητή λογαριθμιστικής παλινδρόμησης με μεταβλητές τα θετικά και αρνητικά συναισθήματα

Σε αυτή την ενότητα , παρουσιάζουμε την προσέγγιση μας χρησιμοποιώντας την ορολογία που αναπτύχθηκε στην προηγούμενη ενότητα. Η μέθοδος μας στηρίζεται στην λογική ότι για να προβλέψουμε την διάθεση ενός χρήστη A προς κάποιον άλλο

χρήστη B, θα πρέπει να αναζητήσουμε τους χρήστες εκείνους που είναι όσο το δυνατόν πιο "όμοιοι" με τον A και οι οποίοι έχουν απευθύνει μια ακμή προς τον B. Μια σύνοψη της προσέγγισης μας φαίνεται στην Εικόνα 2.1: Πρώτα, εξάγουμε το διευρυμένο εγω-δίκτυο βάθους δύο για τον χρήστη A και έπειτα κρατάμε το σύνολο των κόμβων του υπογράφου αυτού οι οποίοι απευθύνουν μια ακμή προς τον B. Έπειτα, εφαρμόζουμε τις μετρικές ομοιότητας της προηγούμενης ενότητας μεταξύ του χρήστη A και των κόμβων που εξήχθησαν από το πρώτο βήμα και τους κατατάσσουμε σε φθίνουσα σειρά (οι πιο "όμοιοι" θα βρίσκονται ψηλότερα στη λίστα). Στη συνέχεια επιλέγουμε τους k πρώτους χρήστες και έτσι δημιουργούμε μια ομάδα χρηστών τους οποίους καλούμε "υποψήφιους επηρεάζοντες χρήστες" (candidate influencers) και τους συμβολίζουμε με C_i . Το σύνολο C_i θεωρούμε ότι είναι το μοναδικό κριτήριο βάσει του οποίου ο χρήστης A θα κληθεί να αποφασίσει το πρόσημο της ακμής που θα απευθύνει στον χρήστη B. Πειραματιστήκαμε σε ένα μεγάλο εύρος τιμών για το k και διαπιστώσαμε ότι η τιμή του k ίση με 15 πετυχαίνει την καλύτερη ισορροπία μεταξύ ακρίβειας και κάλυψης των δεδομένων (dataset coverage).

Στο δεύτερο βήμα, αθροίζουμε τις συνεισφορές των χρηστών στο C_i (δηλαδή τις θετικές/αρνητικές τους γνώμες), οι οποίες βαρύνονται σε ποσοστό ανάλογο με του δείκτη ομοιότητας τους με τον χρήστη A. Έτσι, καταλήγουμε με δύο ποσότητες οι οποίες αντιπροσωπεύουν το συνολικό θετικό και αρνητικό συναίσθημα των χρηστών στο C_i προς τον κόμβο-προορισμό βάσει των σχέσεων:

$$pos = \sum_{u \in C_i} score(u, source) * I\{sgn(u, target) = +1\}$$

$$neg = \sum_{u \in C_i} score(u, source) * I\{sgn(u, target) = -1\}$$

όπου I είναι η συνάρτηση-δείκτης, $I\{\chi\} = 1$, αν χ είναι αληθές και 0 σε άλλη περίπτωση.

Το τρίτο και τελευταίο βήμα, περιλαμβάνει την εκπαίδευση ενός ταξινομητή λογαριθμιστικής παλινδρόμησης ο οποίος χρησιμοποιεί ως μεταβλητές τις ποσότητες που αντιστοιχούν στο θετικό και αρνητικό συναίσθημα του προηγούμενου βήματος.

Το τελικό μοντέλο θα έχει την μορφή $logit(p) = \beta_0 + \beta_{pos} * pos + \beta_{neg} * neg$, όπου p είναι η πιθανότητα να είναι θετική η ακμή, β_0 η προδιάθεση (bias) και β_{pos}, β_{neg} οι συντελεστές του θετικού και αρνητικού συναισθήματος αντίστοιχα. Για την αξιολόγηση του μοντέλου χρησιμοποιήσαμε το 60% των δεδομένων για εκπαίδευση και το υπόλοιπο 40% για δοκιμή.

2.3 Αξιολόγηση της μεθόδου

2.3.1 Περιγραφή των δεδομένων

Επαληθεύσαμε την μέθοδο μας σε τρία μεγάλης κλίμακας, πραγματικά, κοινωνικά δίκτυα τα οποία αλιεύσαμε από τον ιστοτόπο SNAP⁴ (Stanford Network Analysis Project). Το πρώτο δίκτυο είναι από το Epinions⁵, έναν ιστοτόπο αξιολόγησης προϊόντων, όπου οι χρήστες έχουν την δυνατότητα να χαρακτηρίζουν άλλους χρήστες ως έμπιστους (ή μη) δημιουργώντας με αυτό τον τρόπο έναν προσημασμένο κοινωνικό γράφο με θετικές και αρνητικές ακμές. Τα δεδομένα καλύπτουν μια περίοδο από το 1999 μέχρι και το 2003, περιλαμβάνοντας περίπου 130 χιλιάδες χρήστες οι οποίοι συνδέονται με 840 χιλιάδες δεσμούς. Το δεύτερο σύνολο

⁴ <http://snap.stanford.edu/>

⁵ <http://www.epinions.com/>

δεδομένων προέρχεται από το Slashdot⁶, έναν τεχνολογικό ιστοτόπο, στο οποίο οι χρήστες μπορούν να ορίσουν ετικέτες για άλλους χρήστες δηλώνοντας την εμπιστοσύνη τους ("friend") ή την δυσπιστία τους ("foe"). Τα δεδομένα αποτυπώνουν την κατάσταση του κοινωνικού γράφου τον Φεβρουάριο του 2009, και περιλαμβάνουν περί τους 82 χιλιάδες χρήστες και 550 χιλιάδες ακμές. Τέλος, το τρίτο δίκτυο αποτελεί την καταγραφή των ψήφων στις εκλογές για την ανάδειξη διαχειριστών στον ιστοτόπο Wikipedia⁷, και περιλαμβάνει δεδομένα από 7 χιλιάδες χρήστες και περίπου 100 χιλιάδες ψήφους.

2.3.2 Προβλεπτική απόδοση

Η ενότητα αυτή παρουσιάζει τα αποτελέσματα από την πειραματική επαλήθευση της μεθόδου μας. Το μεγαλύτερο πρόβλημα που συναντήσαμε κατά αυτή την διαδικασία ήταν η μεγάλη ανισοκατανομή μεταξύ των δύο κλάσεων, αφού περίπου το 80% των ακμών ήταν θετικές. Αντιμετωπίσαμε το πρόβλημα αυτό με δύο τρόπους: Πρώτον, αξιολογήσαμε την μέθοδο μας, με μια μετρική που δεν επηρεάζεται από την ανισοκατανομή και συγκεκριμένα το εμβαδόν της επιφανείας (AUC) κάτω από την συνάρτηση ROC (Receiver Operator Characteristic) [19]. Δεύτερον, εκτιμήσαμε την ακρίβεια της προσέγγισης μας σε πιο *ισορροπημένα* δεδομένα, δηλαδή σε δίκτυα με ίσο αριθμό θετικών και αρνητικών ακμών. Στην βιβλιογραφία, η τεχνική είναι γνωστή και ως υποδειγματοληψία (undersampling) και αποτελεί χρήσιμο οδηγό σε περιπτώσεις μεγάλης ανισορροπίας στα δεδομένα [20] [21]. Για να το επιτύχουμε ακολουθήσαμε την στρατηγική των Guha et. al [22]: Για κάθε δίκτυο κρατήσαμε όλες τις αρνητικές ακμές και επιλέξαμε τυχαία έναν ίδιο αριθμό θετικών ακμών.

⁶ <http://slashdot.org/>

⁷ <https://www.wikipedia.org/>

Επαναλάβαμε την ίδια διαδικασία 10 φορές για να ελαχιστοποιήσουμε την τυχαιότητα των αποτελεσμάτων. Ο Πίνακας παρουσιάζει τα αποτελέσματα για όλες τις παραμέτρους του προβλήματος. Οι πρώτες τρεις στήλες αντιστοιχούν στους συντελεστές του μοντέλου λογαριθμιστικής παλινδρόμησης ενώ οι δύο τελευταίες στήλες στις μετρήσεις ROC-AUC και λάθους 0/1 (ποσοστό των σωστών προβλέψεων).

Μετρική	Κοινός Γείτονας				
	β_0	β_{pos}	β_{neg}	AUC	0/1 loss
Epinions-όλο	2.158	0.002	-0.007	91.43%	93.96%
Epinions-Ισορροπ.	0.525	0.001	-0.009	92.72%	85.00%
Slashdot-όλο	1.412	0.013	-0.052	89.72%	86.88%
Slashdot- Ισορροπ.	0.525	0.010	-0.063	89.82%	82.38%
Wikipedia - όλο	1.049	0.002	-0.005	82.76%	83.64%
Wikipedia – Ισορρ.	-0.197	0.002	-0.006	84.52%	77.70%

Μετρική	Adamic/Adar				
	β_0	β_{pos}	β_{neg}	AUC	0/1 loss
Epinions-όλο	2.155	0.015	-0.043	91.98%	94.08%
Epinions-Ισορροπ.	0.525	0.011	-0.061	92.96%	85.56%
Slashdot-όλο	1.372	0.087	-0.291	90.04%	87.13%
Slashdot- Ισορροπ.	0.500	0.063	-0.350	90.31%	82.85%
Wikipedia - όλο	1.064	0.011	-0.022	82.91%	83.55%
Wikipedia – Ισορρ.	-0.160	0.010	-0.031	83.94%	77.16%

Μετρική	Jaccard				
	β_0	β_{pos}	β_{neg}	AUC	0/1 loss
Erinions-όλο	1.560	3.014	-10.240	94.50%	95.52%
Erinions-Ισορροπ.	0.072	2.216	-13.761	94.89%	88.38%
Slashdot-όλο	0.922	5.653	-14.644	91.61%	87.99%
Slashdot- Ισορροπ.	0.067	4.372	-17.321	91.90%	84.14%
Wikipedia - όλο	0.785	1.472	-2.776	86.69%	85.82%
Wikipedia – Ισορρ.	-0.429	1.393	-3.538	86.52%	78.90%

2.1 Πειραματικές μετρήσεις για τους τρεις αλγορίθμους ομοιότητας.

Κάποια συμπεράσματα είναι εύκολο να εξαχθούν: Πρώτον, φαίνεται ότι ο συντελεστής Jaccard δίνει τα καλύτερα αποτελέσματα και για τις τρεις περιπτώσεις δικτύων, και γι' αυτό τον λόγο χρησιμοποιείται στα επόμενα πειράματα. Δεύτερον, οι μετρήσεις ROC-AUC είναι περίπου ίδιες για τα πλήρη αλλά και τα ισορροπημένα δίκτυα, όπως αναμενόταν, σε αντίθεση με την πιο απλή μετρική του 0/1 loss που μεταβάλλεται σημαντικά. Τρίτον, η απόδοση του ταξινομητή κρίνεται ικανοποιητική. Η βέλτιστη απόδοση παρατηρείται στα δεδομένα του Erinions με τιμή 95,52% και η χειρότερη στην περίπτωση του Wikipedia με ακρίβεια 86,69%. Σε σχέση με την εργασία αναφοράς [6], πετυχαίνουμε μια αύξηση στην ακρίβεια (σε απόλυτες τιμές) 2% και 6% για τις περιπτώσεις των Erinions και Wikipedia αντίστοιχα, ενώ παρατηρείται χειρότερη απόδοση κατά 2% στην περίπτωση του Slashdot. Τέλος, αξίζει να σημειωθεί ότι οι συντελεστές για το αρνητικό και θετικό συναίσθημα διαφέρουν ελάχιστα ανάμεσα στις δύο περιπτώσεις δικτύων (πλήρων και ισορροπημένων). Η μετάβαση από το πλήρες δίκτυο στο ισορροπημένο "απορροφάται" κυρίως από τον συντελεστή προδιάθεσης (bias) β_0 , ο οποίος αρχικά λαμβάνει μεγάλες τιμές, λόγω της ανισοκατανομής των κλάσεων, ενώ στη συνέχεια

μειώνεται σημαντικά προσεγγίζοντας το μηδέν στην περίπτωση των ισορροπημένων δικτύων.

2.3.3 Δυνατότητες Γενίκευσης της μεθόδου σε πολλαπλά δίκτυα

Δίκτυο	Epinions	Slashdot	Wikipedia
Epinions	94.50%	91.68%	87.01%
Slashdot	94.12%	91.61%	86.79%
Wikipedia	93.70%	91.56%	86.69%

Πίνακας 2.2. Δυνατότητες γενίκευσης του ταξινομητή. Εκπαιδεύουμε στα δεδομένα του δικτύου κάθε γραμμής και αξιολογούμε στο δίκτυο της στήλης

Ένα ενδιαφέρον χαρακτηριστικό της προσέγγισης που παρουσιάζουμε είναι οι ιδιότητες γενίκευσης που διαθέτει. Με λίγα λόγια, εκπαιδεύοντας το μοντέλο μας σε ένα σύνολο δεδομένων μπορούμε αργότερα να το εφαρμόσουμε αυτούσιο σε άλλα δεδομένα προερχόμενα από διάφορα κοινωνικά δίκτυα. Στην συγκεκριμένη ενότητα εξετάζουμε την πτυχή αυτή ακολουθώντας την μεθοδολογία των συγγραφέων στο [6]. Ο Πίνακας 2.2 παρουσιάζει τα αποτελέσματα τα οποία πήραμε κάνοντας χρήση του συντελεστή Jaccard. Κάθε γραμμή του Πίνακα 2.2 υποδεικνύει το σύνολο των δεδομένων πάνω στο οποίο εκπαιδεύσαμε το μοντέλο μας ενώ κάθε στήλη παρουσιάζει την ακρίβεια του μοντέλου όταν εφαρμόστηκε στο αντίστοιχο δίκτυο. Αναλύοντας τα περιεχόμενα του Πίνακα 2.2 καταλήγουμε στα ακόλουθα συμπεράσματα. Πρώτον, η χειρότερη ακρίβεια παρουσιάζεται στα δεδομένα του Wikipedia, ανεξαρτήτως του δικτύου που χρησιμοποιήσαμε για εκπαίδευση. Δεύτερον, είναι εντυπωσιακό το γεγονός ότι όταν εκπαιδεύουμε στα δεδομένα του Wikipedia παρατηρούμε ακρίβεια 93,70% και 91,56% στα Epinions και Slashdot

αντίστοιχα, παρόλο που στο ίδιο παίρνουμε μόνο 86,69%. Τέλος, παρατηρούμε ότι τα αποτελέσματα είναι λίγο πολύ συνεπή μεταξύ τους ανεξαρτήτως των δεδομένων εκπαίδευσης. Συγκεκριμένα, στο δίκτυο του Epinions πετυχαίνουμε ακρίβεια γύρω στο 94%, στο Slashdot 91,6% και στο Wikipedia περίπου 87%.

2.4 Η δύναμη των Αρνητικών Ακμών

Στην ενότητα αυτή εξετάζουμε μια διαφορετική πτυχή του προβλήματος ταξινόμησης προσήμου. Πιο συγκεκριμένα, εξετάζουμε την σχετική δυναμική των θετικών και αρνητικών ακμών και πώς αυτή επηρεάζει την προβλεπτική ικανότητα του μοντέλου. Με άλλα λόγια, αναζητούμε εάν ένας συγκεκριμένο τύπος προσήμου φέρει περισσότερη πληροφορία από τον άλλο.

Ξεκινάμε την ανάλυση μας κάνοντας μια πρώτη ανάγνωση των συντελεστών του μοντέλου από τον Πίνακα 2.1. Το μοντέλο λογαριθμιστικής παλινδρόμησης εκφράζει τον *logit* μετασχηματισμό της μεταβλητής εξόδου ως ένα γραμμικό συνδυασμό των μεταβλητών εισόδου. Στην περίπτωση μας, οι ανεξάρτητες μεταβλητές είναι το θετικό και αρνητικό συναίσθημα που παράγεται από το "κοντινούς φίλους" του κόμβου-πηγή οι οποίοι πολλαπλασιάζονται με τους συντελεστές β_{pos} και β_{neg} αντίστοιχα. Γίνεται εύκολα αντιληπτό ότι για όλα τα δίκτυα που εξετάστηκαν αλλά και για όλες τις μετρικές ομοιότητας ο συντελεστής του αρνητικού συναισθήματος β_{neg} είναι πάντα μεγαλύτερος (σε απόλυτες τιμές) του άλλου. Αυτό το γεγονός από μόνο του δεν φαίνεται στατιστικά σημαντικό αλλά αποτελεί μια πρώτη ένδειξη.

Για να επιβεβαιώσουμε τον ισχυρισμό μας, δηλαδή ότι οι δύο συντελεστές διαφέρουν στατιστικά σημαντικά, κάνουμε χρήση κλασικών στατιστικών εργαλείων όπως ο έλεγχος των υποθέσεων (hypothesis testing). Συγκεκριμένα, θα εξετάσουμε δύο

περιπτώσεις: την μηδενική υπόθεση (null hypothesis) που λέει ότι οι δύο συντελεστές δεν διαφέρουν (στατιστικά) σημαντικά και την εναλλακτική υπόθεση η οποία είναι συμβατή με την διαίσθηση μας ότι οι δύο συντελεστές διαφέρουν σημαντικά, με τον αρνητικό συντελεστή να είναι μεγαλύτερος σε απόλυτες τιμές. Στην εναλλακτική υπόθεση αντιστοιχεί το μοντέλο που ήδη έχουμε από τις προηγούμενες ενότητες ενώ για την μηδενική υπόθεση δημιουργούμε ένα νέο μοντέλο:

$$\text{logit}(p) = \gamma_0 + \gamma_{\text{common}} * (\text{pos} + \text{neg})$$

το οποίο προκύπτει από το προηγούμενο με τον επιπλέον περιορισμό ότι $\beta_{\text{pos}} = \beta_{\text{neg}}$. Με χρήση ενός *likelihood ratio test* καταλήξαμε ότι η μηδενική υπόθεση απορρίπτεται σε επίπεδο σημαντικότητας 0.01

Μέχρι τώρα έχουμε ισχυρές ενδείξεις ότι η ύπαρξη ενός αρνητικού δεσμού επηρεάζει περισσότερο από έναν αντίστοιχο θετικό. Σε αυτό το τελευταίο κομμάτι του κεφαλαίου θα προσπαθήσουμε να ποσοτικοποιήσουμε την ανωτερότητα των αρνητικών ακμών και πως αυτό επηρεάζει την συνολική προβλεπτική απόδοση του μοντέλου. Διαμορφώνουμε το ακόλουθο πείραμα: Διαλέγουμε έναν θετικό ακέραιο αριθμό m και για κάθε σύνολο δεδομένων, δημιουργούμε δύο καινούρια

- Στο πρώτο σύνολο, το οποίο θα συμβολίσουμε ως S^- , αφαιρούμε m αρνητικές ακμές, κρατώντας όλες τις θετικές, ενώ
- στο δεύτερο, το οποίο θα συμβολίσουμε ως S^+ αφαιρούμε m θετικές ακμές, και κρατάμε όλες τις αρνητικές.

Για κάθε καινούριο σύνολο επανεκπαιδεύουμε τον ταξινομητή μας και αξιολογούμε την ακρίβεια του. Είναι προφανές ότι τα S^- και S^+ έχουν το ίδιο μέγεθος αλλά διαφορετικές αναλογίες προσήμων.

Επαναλάβαμε το πείραμα μας για $k = 10$ επαναλήψεις, όπου κάθε φορά αφαιρούσαμε $k \cdot 10\% \cdot |E^-|$ ακμές του ίδιου τύπου, όπου $|E^-|$ είναι ο συνολικός

αριθμός των αρνητικών ακμών στον γράφο. Τα αποτελέσματα φαίνονται στο Σχήμα

2.2



Σχήμα 2.2. Προβλεπτική απόδοση των τροποποιημένων δικτύων. Η διακεκομμένη γραμμή αντιστοιχεί στο δίκτυο, από το οποίο έχουν αφαιρεθεί θετικές ακμές (S^+) ενώ η συνεχής γραμμή στο δίκτυο από το έχουν διαγραφεί οι αρνητικές (S^-).

Η διακεκομμένη γραμμή αντιστοιχεί στο S^+ σύνολο ενώ η συνεχής στο S^- . Στην περίπτωση του Epinions οι δύο γραφικές παραστάσεις εκκινούν από το ίδιο σημείο (95,52% ακρίβεια) αλλά ακολουθούν διαφορετική πορεία. Η παράσταση του S^+ συνόλου παραμένει σταθερή καθ' όλη την διάρκεια του πειράματος, ενώ η άλλη μειώνεται φτάνοντας τελική ακρίβεια 93,54%. Παρόμοια χαρακτηριστικά παρατηρούνται και στην περίπτωση του Wikipedia, όπου η αρχική απόδοση 85,83% γίνεται 85,98% για το S^+ σύνολο και 82,22% για το S^- σύνολο. Μεγαλύτερη διακύμανση παρατηρείται στην περίπτωση του Slashdot, όπου η αρχική ακρίβεια 88,01% καταλήγει σε 88,27% (για το S^+) και 82,54% (για το S^-). Τα αποτελέσματα είναι συμβατά με την διαίσθηση μας. Η αφαίρεση έστω και ενός μικρού αριθμού αρνητικών ακμών οδηγεί σε σημαντική μείωση της ακρίβειας, το οποίο μπορεί να είναι κρίσιμο για διάφορες εφαρμογές όπως τα συστήματα προτάσεων (recommender systems). Λαμβάνοντας υπόψη πέρα από τις προτιμήσεις και τις αρνητικές γνώμες των χρηστών μπορούμε να μοντελοποιήσουμε αποδοτικότερα την κοινωνική τους

συμπεριφορά. Τα ευρήματά μας συμφωνούν με την εργασία των Garcia et al. [23], όπου οι συγγραφείς μελέτησαν την χρήση συναισθηματικά φορτισμένων εκφράσεων στα Αγγλικά, Γερμανικά και Ισπανικά και διαπίστωσαν ότι οι εκφράσεις με θετικά ορισμένο περιεχόμενο είναι πολύ πιο συχνές από τις αντίστοιχες αρνητικές, και άρα ότι υπάρχει μια θετική προδιάθεση στην ανθρώπινη επικοινωνία (Pollyanna hypothesis [24]). Τέλος, καταλήγουν πως οι αρνητικές λέξεις φέρουν στατιστικά περισσότερη πληροφορία από τις θετικές.

2.5 Σύνοψη Κεφαλαίου

Ο σκοπός αυτού του κεφαλαίου ήταν διπλός: Πρώτον, κατασκευάσαμε έναν μηχανισμό πρόβλεψης προσήμου για δίκτυα εμπιστοσύνης χρησιμοποιώντας ως βασικά συστατικά

- την αρχή της *ομοφιλίας*, η οποία θέτει ότι τα άτομα ενός κοινωνικού δικτύου τείνουν να επηρεάζονται περισσότερο από άλλα όμοια άτομα και
- τρεις παραδοσιακές μετρικές ομοιότητας (Κοινοί Γείτονες, Jaccard, Adamic/Adar)

Η προσέγγισή μας εμφάνισε βελτιώσεις σε σχέση με την υπάρχουσα βιβλιογραφία. Επίσης, εξετάσαμε την σχετική δυναμική των θετικών και αρνητικών απόψεων και διαπιστώσαμε πως οι αρνητικές γνώμες τείνουν να έχουν μεγαλύτερη ισχύ.

3

Πρόβλεψη Συναισθήματος με την χρήση Συχνών Υπογράφων

Σε αυτό το κεφάλαιο παρουσιάζουμε μια διαφορετική προσέγγιση για την πρόβλεψη προσήμου σε κοινωνικά δίκτυα. Βάση της μεθόδου που παρουσιάζουμε αποτελεί η εύρεση των συχνών υπογράφων που παρατηρούνται σε γράφους εμπιστοσύνης, τους οποίους χρησιμοποιούμε ως μεταβλητές για περαιτέρω στατιστική ανάλυση. Η διάρθρωση του κεφαλαίου είναι ως εξής: Πρώτα παρουσιάζουμε βασικούς ορισμούς και έννοιες από τον χώρο της γραφο-θεωρίας. Στη συνέχεια παραθέτουμε τις σχετικές εργασίες στην βιβλιογραφία, ενώ στο τέλος παρουσιάζουμε την προσέγγιση μας αξιολογώντας την σε τρεις, μεγάλης κλίμακας, γράφους εμπιστοσύνης από online κοινωνικά δίκτυα.

3.1 Βασική Ορολογία

Όλη η προσέγγιση που παρουσιάζουμε σε αυτό το κεφάλαιο στηρίζεται στην έννοια του υπογράφου ενός κοινωνικού δικτύου. Ως υπογράφο $g(V', E')$ ενός μεγαλύτερου γράφου $G(V, E)$ θεωρούμε οποιονδήποτε υπογράφο του οποίου οι ακμές και κόμβοι

περιέχονται στον G , και συμβολίζεται ως $g \subseteq G$. Εισάγουμε επίσης την έννοια του *συχνού υπογράφου* ο οποίος ορίζεται ως εξής :

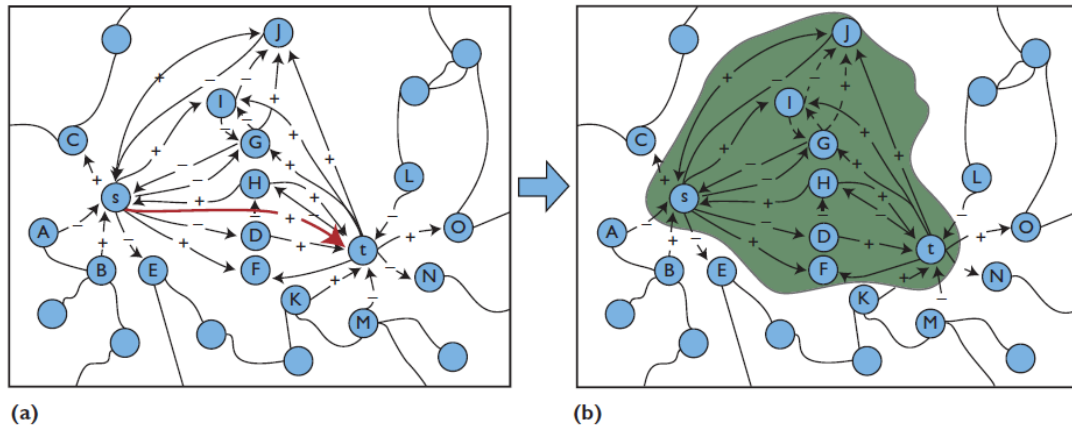
Ορισμός (Συχνός ε-υπογράφος)

*Δεδομένης μιας συλλογής γράφων D , και ενός υπογράφου g , ορίζουμε ως υποστηρίζουν σύνολο D_g το σύνολο $D_g := \{ d \mid g \subseteq d, d \in D \}$. Η υποστήριξη του g δίνεται από την σχέση $support(g) = |D_g|/|D|$. Ένας υπογράφος καλείται **συχνός ε-υπογράφος** αν η υποστήριξη του είναι πάνω από ένα προκαθορισμένο κατώφλι ϵ .*

Το κατώφλι ϵ εξαρτάται σε μεγάλο βαθμό από το πεδίο εφαρμογής και συνήθως για την εύρεση της βέλτιστης τιμής πραγματοποιείται μια γραμμική αναζήτηση σε ένα ικανό εύρος τιμών. Ο δεύτερος ορισμός που παραθέτουμε αφορά τον *επαγόμενο υπογράφο* (induced subgraph). Ένας υπογράφος H καλείται *επαγόμενος* ενός γράφου G αν για όλους τους κόμβους $u_i \in V(H) \subseteq V(G)$, ο H περιέχει όλες τις ακμές ανάμεσα στα u_i από τον αρχικό γράφο δηλαδή $E(H) = \{(u_i, u_j) \mid (u_i, u_j) \in E(G)\}$. Θέλοντας να βρούμε έναν τρόπο ώστε να εξάγουμε το κομμάτι του γράφου που βρίσκεται ‘κοντά’ σε μια ακμή , εισάγουμε την έννοια του *εγω-γραφήματος* (ego-graph) μιας ακμής.

Ορισμός (εγω-γράφημα)

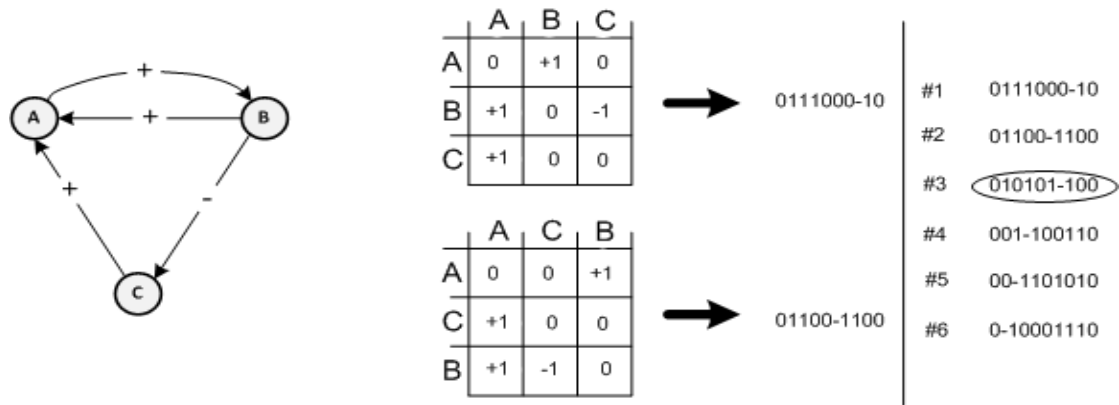
*Δεδομένου ενός γράφου G και μιας ακμής $e \in E(G)$, ορίζουμε ως εγω-γράφημα της ακμής e , τον επαγόμενο υπογράφο του G , του οποίου το σύνολο των κόμβων περιέχει τους κόμβους της ακμής e καθώς και τους κοινούς τους γείτονες. Η ακμή e καλείται **ακμή-γεννήτορας** του ego-graph.*



Εικόνα 3.1 α) Αριστερά - Ένα κομμάτι του κοινωνικού γράφου. Η ακμή $s \rightarrow t$ θα λειτουργήσει ως ακμή-γεννήτορας β) Δεξιά - Το εξαγόμενο εγω-γράφημα. Το εγω-γράφημα είναι 'θετικό' όπως το πρόσημο της ακμής-γεννήτορα

Η Εικόνα 3.1α παρουσιάζει ένα παράδειγμα προσημασμένου γράφου όπου η ακμή ενδιαφέροντος είναι σχεδιασμένη σε πιο έντονο φόντο. Είναι μια θετική ακμή που ξεκινάει από τον κόμβο s και καταλήγει στον κόμβο t . Χρησιμοποιώντας τη ως ακμή-γεννήτορα κατασκευάζουμε το αντίστοιχο εγω-γράφημα, το οποίο φαίνεται στην πράσινη περιοχή της εικόνας 3.1β. Εξ ορισμού, όλοι οι κοινοί γείτονες των s και t περιλαμβάνονται στο εγω-γράφημα, φέροντας και τις μεταξύ τους ακμές από τον αρχικό γράφο. Η μόνη εξαίρεση αφορά στην ακμή-γεννήτορα, η οποία δεν συμπεριλαμβάνεται στο σύνολο των ακμών του εγω-γραφήματος, αλλά χρησιμοποιούμε το πρόσημο της για να χαρακτηρίσουμε συνολικά το εγω-γράφημα π.χ. το εξαγόμενο εγω-γράφημα εδώ είναι θετικό γιατί η ακμή $s \rightarrow t$ είναι θετική.

Ολοκληρώνουμε αυτή την ενότητα με την περιγραφή μιας συνάρτησης **κανονικής επισήμανσης** (*canonical labeling*). Εν συντομία, οι συναρτήσεις αυτές δέχονται ως είσοδο γραφικά δεδομένα και παράγουν αλφαριθμητικούς κώδικες, οι οποίοι αντανakλούν τις δομικές τους ιδιότητες. Χρησιμοποιούνται κυρίως για να αντιμετωπίσουν το πρόβλημα του **ισομορφισμού**, αφού δύο ισομορφικοί γράφοι θα καταλήξουν να έχουν τον ίδιο κωδικό. Η μέθοδος μας εκτελείται ως εξής:



...

Εικόνα 3.2 Λειτουργία της canonical labeling συνάρτησης. Αριστερά φαίνεται ο γράφος που δίνεται σαν είσοδος στην συνάρτηση. Στο πρώτο βήμα, παράγονται όλοι οι ισοδύναμοι πίνακες γειτνίασης και υπολογίζονται οι αντίστοιχοι κώδικες. Δεξιά, γίνεται η επιλογή του κωδικού

1. Δεδομένου ενός γράφου, κατασκευάζουμε τον πίνακα γειτνίασης του καθώς και όλους τους ισοδύναμους με αυτόν πίνακες που προκύπτουν με μετάθεση γραμμών/στηλών.
2. Για κάθε πίνακα, παράγουμε τον αντίστοιχο κωδικό, συγχωνεύοντας τις στήλες του πίνακα από αριστερά προς τα δεξιά όπως φαίνεται στην Εικόνα 3.2 Κάθε κωδικός χαρακτηρίζεται ως υποψήφιος για τον γράφο.
3. Επιλέγουμε τον βέλτιστο κωδικό βάσει μιας σειράς κανόνων. Κάθε κωδικός είναι μια ακολουθία από '0', '1' και '-1' και ο σκοπός μας είναι να διαλέξουμε εκείνον, του οποίου οι μη-μηδενικές συνιστώσες εμφανίζονται πρώτες, ενώ οι μηδενικές τιμές τελευταίες. Σε περίπτωση που κάποιος κανόνας δεν μπορεί να επιλέξει έναν μοναδικό κωδικό, προωθεί τους ισοδύναμους κωδικούς στο επόμενο κατά σειρά κανόνα. Οι κανόνες είναι οι εξής:
 - a. Επιλογή του κωδικού με την μεγαλύτερη ακολουθία μηδενικών στο τέλος π.χ. ο κωδικός '0101010100' είναι προτιμότερος του

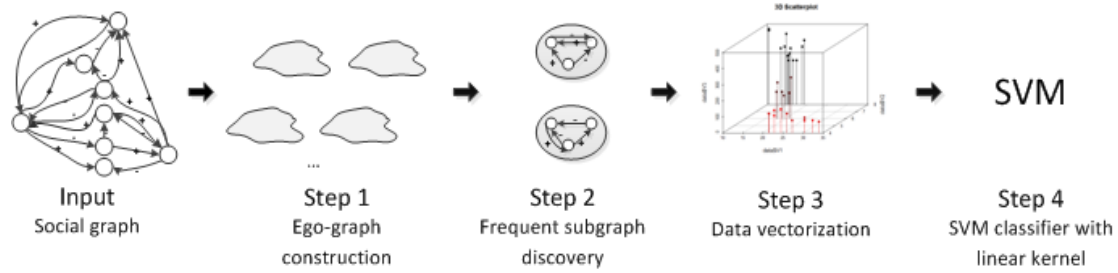
‘0111000-10’ γιατί έχει δύο μηδενικά στο τέλος έναντι ενός του δεύτερου.

- b. Επιλέγεται ο κωδικός με το μικρότερο δείκτη για την τελευταία εγγραφή που περιέχει ‘1’. Αν ακόμα δεν μπορεί να επιλεγεί μοναδικός κωδικός, συνεχίζουμε με το ‘1’ που είναι δεύτερο από το τέλος κ.ο.κ.
- c. Επαναλαμβάνουμε το βήμα b για τις εγγραφές που περιέχουν ‘-1’.

Ένα πρόβλημα που έχει η συγκεκριμένη μέθοδος είναι ότι έχει εκθετική πολυπλοκότητα σε σχέση με το μέγεθος του γράφου, δεδομένων των ισοδύναμων πινάκων γεινίασης που πρέπει να ληφθούν υπόψη. Αυτό όμως δεν επηρεάζει αισθητά την δική μας προσέγγιση, καθώς κατά τα πειράματά μας η συνάρτηση εφαρμόζεται σε μικρούς υπογράφους με 3 ή 4 κόμβους, κάτι που καθιστά το κόστος σταθερό.

3.2 Προτεινόμενη Προσέγγιση

Σε αυτή την ενότητα παρουσιάζουμε την προσέγγιση μας για την πρόβλεψη προσήμου των ακμών κοινωνικών γράφων. Ο βασικός σκοπός μας είναι η μετατροπή του προβλήματος σε ένα ισοδύναμο πρόβλημα ταξινόμησης γράφων (graph classification) και η επίλυση του με την χρήση μοντέρνων τεχνικών ανάλυσης δικτύων. Η μετατροπή επιτυγχάνεται μέσω της συνάρτησης εγω-γραφήματος που παρουσιάστηκε στην προηγούμενη ενότητα, η οποία δεχόμενη μια ακμή, δίνει ως έξοδο το κομμάτι του γράφου που ‘βρίσκεται κοντά’ σε αυτή. Το πρόσημο της ακμής-γεννήτορα χρησιμοποιείται για να χαρακτηρίσει το κάθε εγω-γράφημα ως θετικό ή αρνητικό. Η προσέγγιση μας βασίζεται στην πεποίθηση ότι η δομή των εγω-



Εικόνα 3.3 Τα τέσσερα βήματα του αλγορίθμου. Πρώτα παραγάγουμε το εγω-γράφημα κάθε ακμής και βρίσκουμε τους συχνούς υπογράφους. Έπειτα, δημιουργούμε έναν πίνακα δεδομένων, όπου κάθε γραμμή δείχνει τις συχνότητες εμφάνισης κάθε υπογράφου ανά εγω-γράφημα.

γραφημάτων επιδεικνύει διαφορετικά μοτίβα (graph patterns) ανάλογα με το είδος του προσήμου της ακμής-γεννήτορα. Στον τομέα της βιοστατιστικής αυτή η προσέγγιση είναι γνωστή ως υπόθεση QSAR (Quantitative structure-activity relationship), η οποία υποστηρίζει ότι ‘όμοια στοιχεία’ έχουν παρόμοιες συμπεριφορές π.χ. η τοξικότητα ενός στοιχείου μπορεί να προβλεφθεί από την κατανομή των μορίων του στον τρισδιάστατο χώρο, συγκρινόμενο με στοιχεία με παρόμοια δομή και γνωστή συμπεριφορά. Η Εικόνα 3.4 απεικονίζει τα τέσσερα στάδια της μεθόδου μας και κάθε βήμα εξηγείται αναλυτικά.

Βήμα 1ο - Δημιουργία εγω-γραφημάτων

Το πρώτο βήμα πραγματοποιεί την μεταφορά από τον χώρο των ακμών του κοινωνικού γράφου σε έναν ισοδύναμο χώρο υπογράφων, εφαρμόζοντας σε κάθε ακμή την συνάρτηση του εγω-γραφήματος. Για να εξασφαλίσουμε ότι υπάρχει αρκετή πληροφορία στην ‘γειτονιά’ της κάθε ακμής, απαιτούμε ο αριθμός των κοινών γειτόνων των κόμβων κάθε ακμής να είναι πάνω από ένα προκαθορισμένο όριο. Επιλέξαμε την τιμή 25, αφενός γιατί πετυχαίνει καλή απόδοση καλύπτοντας την πλειονότητα των περιπτώσεων, αφετέρου γιατί η ίδια τιμή χρησιμοποιείται και στην μέθοδο αναφοράς [6] με την οποία συγκρινόμαστε.

Βήμα 2ο - Εύρεση συχνών υπογράφων

Στο δεύτερο βήμα πραγματοποιείται η εύρεση των συχνών υπογράφων [25] των εγω-γραφημάτων που παρήχθησαν από το πρώτο βήμα. Οι υπονήφιοι υπογράφοι έχουν ένα καθορισμένο μέγεθος k , το οποίο δίνεται από τον χρήστη. Περιορίσαμε την μελέτη μας σε υπογράφους οι οποίοι περιέχουν και τους δύο κόμβους της ακμής-γεννήτορα. Με άλλα λόγια, από τους k κόμβους του υπογράφου δύο από αυτούς είναι οι κόμβοι της ακμής-γεννήτορα ενώ οι υπόλοιπες $k - 2$ θέσεις καλύπτονται από το σύνολο των κοινών τους γειτόνων. Για να αντιμετωπίσουμε το πρόβλημα των ισομορφισμών, εφαρμόσαμε την συνάρτηση της κανονικής επισήμανσης της προηγούμενης ενότητας. Τέλος, για να κρατήσουμε το κόστος υπολογισμού χαμηλά, πειραματιστήκαμε με δύο μικρές τιμές για το k , και πιο συγκεκριμένα τις τιμές 3 και 4 καθώς και επτά διαφορετικές τιμές για το κατώφλι ελάχιστης συχνότητας στο διάστημα από 0.5% μέχρι 10%. Οι συχνόι υπογράφοι περνάνε στο επόμενο βήμα για περαιτέρω επεξεργασία.

Βήμα 3ο - Διανυσματοποίηση των δεδομένων

Σε αυτό το βήμα δημιουργούμε έναν πίνακα δεδομένων από το σύνολο των εγω-γραφημάτων του πρώτου βήματος, χρησιμοποιώντας ως μεταβλητές τους συχνούς υπογράφους του δεύτερου βήματος. Πιο συγκεκριμένα, έστω G και F τα σύνολα των εγω-γραφημάτων και συχνών υπογράφων αντίστοιχα. Για κάθε εγω-γράφημα g που ανήκει στο G , δημιουργούμε μια γραμμή στον πίνακα δεδομένων, η οποία απεικονίζει την συχνότητα του αντίστοιχου συχνού υπογράφου της αντίστοιχης στήλης, δηλ. $x(g) = \text{frequency}(f_i), i = 1..|F|$. Ο πίνακας δεδομένων δίνεται ως είσοδος στο επόμενο βήμα.

Βήμα 4ο - Εκπαίδευση ταξινομητή

Σε αυτό το τελευταίο βήμα, εφαρμόζουμε έναν Support Vector Machine (SVM) ταξινομητή με γραμμικό πυρήνα, ο οποίος προβλέπει το πρόσημο του εγωγραφήματος (άρα και της αντίστοιχης ακμής-γεννήτορα) ψάχνοντας για ένα υπερ-επίπεδο μέγιστου περιθωρίου (maximum margin hyper-plane). Το γεγονός ότι χρησιμοποιήσαμε γραμμικό πυρήνα κάνει την μέθοδο μας δυνητικά ικανή να αντιμετωπίσει περιπτώσεις ακόμα και μεγαλύτερων δικτύων.

3.3 Πειραματική αξιολόγηση

Σε αυτό το κεφάλαιο, παρουσιάζονται τα πειραματικά αποτελέσματα της μεθόδου μας και γίνεται η σύγκριση με την τρέχουσα καλύτερη προσέγγιση στην βιβλιογραφία [6]. Αρχικά, δίνουμε μια περιγραφή των δεδομένων που χρησιμοποιήσαμε και στη συνέχεια εξάγουμε συμπεράσματα για την αποτελεσματικότητα της δουλειάς μας.

3.3.1 Περιγραφή των δεδομένων

Πραγματοποιήσαμε τα πειράματα μας σε τρία μεγάλης κλίμακας πραγματικά, κοινωνικά δίκτυα. Όλες οι πληροφορίες αντλήθηκαν από το Stanford Network Analysis Project (SNAP) . Η πρώτη συλλογή προέρχεται από το Epinions, που είναι ένας ιστότοπος αξιολόγησης προϊόντων. Τα μέλη του Epinions έχουν την δυνατότητα να εκφράζουν δεσμούς εμπιστοσύνης ή/και δυσπιστίας προς άλλα μέλη, δημιουργώντας έτσι έναν προσημασμένο κοινωνικό γράφο με θετικές και αρνητικές ακμές. Τα δεδομένα καλύπτουν μια περίοδο από το 1999 έως το 2003 και περιλαμβάνουν περίπου 130 χιλιάδες χρήστες με 840 χιλιάδες μεταξύ τους δεσμούς. Η δεύτερη συλλογή δεδομένων προέρχεται από το τεχνολογικό blog Slashdot , στο

οποίο οι χρήστες έχουν την δυνατότητα να χρησιμοποιούν τις ετικέτες "friend" και "foe" για να εκφράσουν το επίπεδο της εμπιστοσύνης τους προς τους υπόλοιπους χρήστες. Τα δεδομένα αποτελούν ένα στιγμιότυπο του κοινωνικού γράφου, το οποίο δημιουργήθηκε τον Φεβρουάριο του 2009 και περιέχει περίπου 82 χιλιάδες χρήστες με 549.000 ακμές. Τέλος, αναλύσαμε δεδομένα προερχόμενα από το Wikipedia, και συγκεκριμένα από τις εκλογές για την ανάδειξη διαχειριστών. Για να γίνει κάποιος διαχειριστής πρέπει να κάνει το αντίστοιχο αίτημα (Request for Adminship – RfA) και έπειτα το αίτημα του να εξεταστεί από την αντίστοιχη επιτροπή. Τα δεδομένα περιέχουν πληροφορίες για 7 χιλιάδες χρήστες και περίπου 100 χιλιάδες ψήφους. Στατιστικά για τις τρεις συλλογές φαίνονται στον Πίνακα 1.

<i>Δίκτυο</i>	<i>#χρηστών</i>	<i>#δεσμών</i>	<i>% θετικών ακμών</i>
<i>Epinions</i>	131,828	841,200	85.0
<i>Slashdot</i>	81,867	549,202	77.4
<i>Wikipedia</i>	7,194	103,747	78.7

Πίνακας 3-1 Στατιστικά των συλλογών δεδομένων

3.3.2 Εργαλεία και μετρικές αξιολόγησης

Ο κώδικας για την εξαγωγή των γράφων υλοποιήθηκε σε Java και για την στατιστική ανάλυση χρησιμοποιήσαμε το περιβάλλον της R καθώς και τα πακέτα ROCR και LibLineR. Το μεγαλύτερο πρόβλημα που αντιμετωπίσαμε κατά την αξιολόγηση έχει να κάνει με την μεγάλη ασυμμετρία των δεδομένων, αφού περίπου 80% των ακμών κατά μέσο όρο είναι θετικές, επιτρέποντας σε ένα οποιονδήποτε θετικά προκατειλημμένο (biased) ταξινομητή να επιτύχει μεγάλη ακρίβεια. Αντιμετωπίσαμε αυτό το πρόβλημα με δύο τρόπους: Πρώτον, επιλέξαμε ως μετρική αξιολόγησης την επιφάνεια (Area Under Curve – AUC) κάτω από την γραφική παράσταση της Receiver Operator Characteristic (ROC), η οποία δεν επηρεάζεται από την ασυμμετρία στην κατανομή των κλάσεων. Δεύτερον, για να βελτιώσουμε την

απόδοση του ταξινομητή μας, διερευνήσαμε την συνεισφορά των βαρών κόστους για τις θετικές και αρνητικές ακμές του Support Vector Machine (SVM). Ακολουθήσαμε μια συνήθη πρακτική στην βιβλιογραφία, η οποία απαιτεί οι παράμετροι για το βάρος να είναι αντιστρόφως ανάλογες της αντιπροσωπευτικότητας της αντίστοιχης κλάσης στα δεδομένα, έτσι ώστε να προβλέπεται μεγαλύτερη ποινή στα λάθη που αφορούν την λιγότερη συχνή κλάση. Εφαρμόσαμε μια αναζήτηση πλέγματος (grid-search) στον δι-διάστατο χώρο των παραμέτρων κόστους και καταλήξαμε ότι ένα σχήμα με τιμές 1.0 και 0.1 για τις αρνητικές και θετικές ακμές, αντίστοιχα, είναι το βέλτιστο. Τέλος, αξιολογήσαμε τα αποτελέσματα μέσω 5-απλής διασταυρωμένης επικύρωσως (5-fold cross validation).

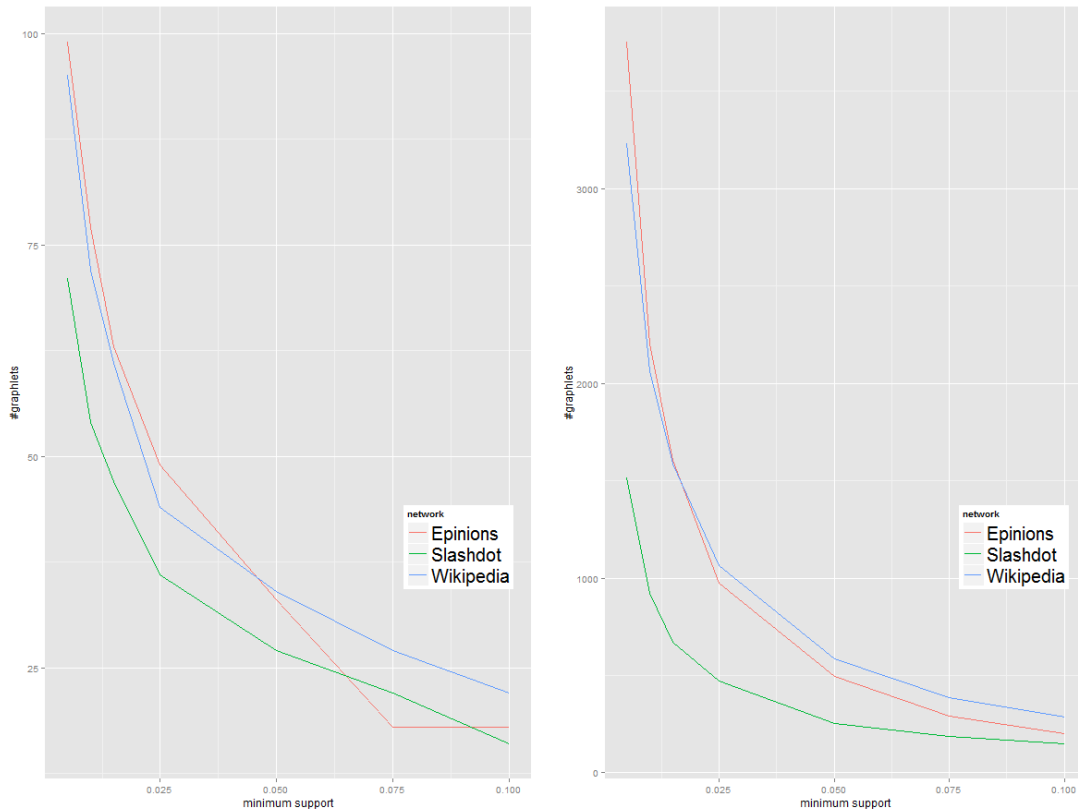
3.3.3 Ακρίβεια Μοντέλου και Ανάλυση

Στον Πίνακα 3.2 παρατίθενται οι μετρήσεις του AUC για τις διαφορετικές διατάξεις των παραμέτρων (μέγεθος υπογράφου, κατώφλι ελάχιστης συχνότητας). Κάποια συμπεράσματα μπορούν να εξαχθούν άμεσα: Πρώτον, η απόδοση της προσέγγισης μας κρίνεται ικανοποιητική, δεδομένου ότι εμφανίζει βελτίωση σε σχέση με την μέθοδο αναφοράς κατά 4-15% σε απόλυτες τιμές. Αυτά τα αποτελέσματα επιβεβαιώνουν την αρχική μας υπόθεση ότι οι συχνοί υπογράφοι διαθέτουν αρκετή ισχύ ώστε να ταξινομήσουν με ακρίβεια τα εγω-γραφήματα σε θετικά και αρνητικά. Στις περιπτώσεις των Slashdot και Wikipedia, πετυχαίνουμε την βέλτιστη απόδοση (97.83% και 96.84% αντίστοιχα) για μέγεθος υπογράφου ίσο με τρία και ελάχιστο κατώφλι ίσο με 0.5%, ενώ στην περίπτωση του Epinions, η βέλτιστη κατανομή αφορά μέγεθος υπογράφου ίσο με τέσσερα και τιμή κατωφλίου στο 1.5%. Γενικά, μικρότερο κατώφλι για την ελάχιστη συχνότητα οδηγεί σε καλύτερα αποτελέσματα αλλά συνεπάγεται και μεγαλύτερο αριθμό διαστάσεων. Δεύτερον, παρατηρήσαμε ότι

ο αλγόριθμος μας λειτουργεί ως *αραιός κωδικοποιητής (sparse coder)* , αφού ο πίνακας δεδομένων που προκύπτει από το τρίτο στάδιο της μεθόδου έχει το 70 με 95 τοις εκατό των εγγραφών του μηδενικά. Αυτή είναι μια ενδιαφέρουσα ιδιότητα που μας επιτρέπει να αναλύσουμε μεγαλύτερα δίκτυα χωρίς μεγάλο επιπλέον κόστος. Η Εικόνα 3.4, δείχνει τον αριθμό των διαφορετικών υπογράφων ανά δίκτυο, για μέγεθος υπογράφου τρία και τέσσερα.

Δίκτυο	Μέγεθος υπογράφου	Κατώφλι ελάχιστης συχνότητας						
		10.0%	7.5%	5.0%	2.5%	1.5%	1.0%	0.5%
Slashdot	3	96.25%	97.15%	97.45%	97.50%	97.69%	97.78%	97.83%
Wikipedia	3	96.12%	96.24%	96.18%	96.29%	96.78%	96.84%	96.84%
Epinions	3	97.20%	97.20%	98.54%	99.32%	99.47%	99.68%	98.79%
Slashdot	4	95.32%	95.88%	95.85%	96.10%	95.66%	94.84%	92.67%
Wikipedia	4	95.73%	95.87%	95.70%	95.18%	94.27%	94.33%	94.04%
Epinions	4	98.38%	99.41%	99.67%	99.78%	99.80%	99.71%	99.70%

Πίνακας 3-2 Μετρήσεις AUC για όλα τα πειράματα και για τα τρία δίκτυα, με δύο τιμές για το μέγεθος του υπογράφου και επτά διαφορετικές τιμές για το κατώφλι ελάχιστης συχνότητας



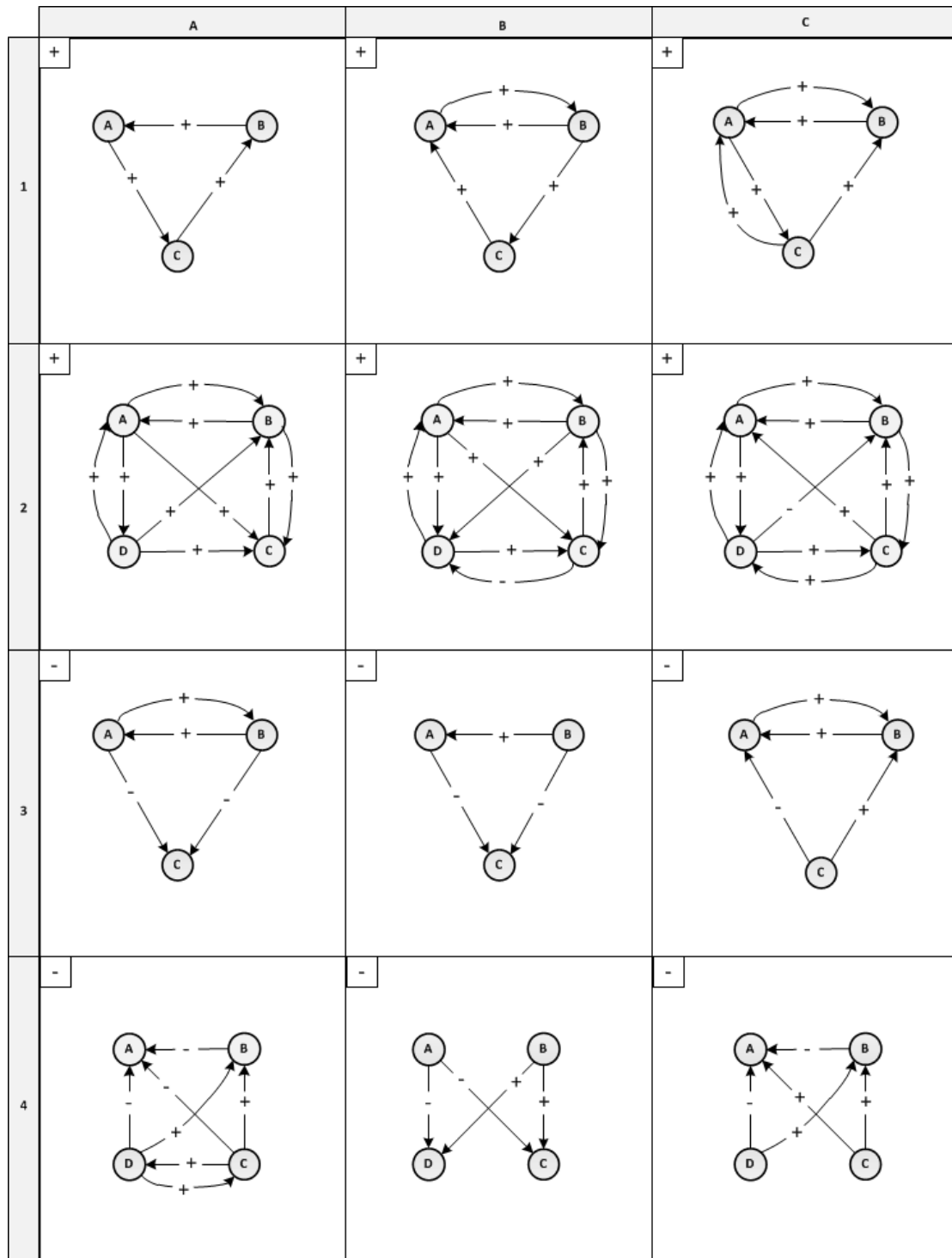
Εικόνα 3.4 Ο αριθμός των διαφορετικών υπογράφων για κάθε δίκτυο και διαφορετική τιμή κατωφλίου. Αριστερά, οι γραφικές παραστάσεις για υπογράφους μεγέθους 3 και δεξιά για μέγεθος 4

Για το υπόλοιπο αυτής της ενότητας θα εστιάσουμε σε κάποιους χαρακτηριστικούς τύπους συχνών υπογράφων οι οποίοι συνδέονται άμεσα με θετικές ή αρνητικές ακμές. Το Σχήμα 3.5 παραθέτει τις περιπτώσεις με το μεγαλύτερο ενδιαφέρον. Οι πρώτες δύο γραμμές απεικονίζουν συχνούς υπογράφους που ‘εντοπίζονται’ κοντά σε θετικές ακμές, ενώ οι δύο τελευταίες γραμμές συνδέονται με αρνητικές. Ως γενική παρατήρηση, συμπεραίνουμε ότι οι θετικές ακμές εμφανίζονται σε κομμάτια του κοινωνικού γράφου που κυριαρχούνται από θετικούς (σχεδόν) πλήρεις υπογράφους. Ακόμα και στις περιπτώσεις 2B και 2C όπου υπάρχουν αρνητικές ακμές, η επίδρασή τους ‘απορροφάται’ από την πυκνή σύνδεση μεταξύ των κόμβων. Από την άλλη πλευρά, οι αρνητικές ακμές εμφανίζουν πιο πολύπλοκη δυναμική και γι αυτό είναι δυσκολότερο να προβλεφθούν. Στις περισσότερες περιπτώσεις των αρνητικών ακμών,

παρατηρούμε δύο διακριτές και αντιτιθέμενες ομάδες χρηστών, όπου υπάρχουν θετικές ακμές ανάμεσα στα μέλη της ίδιας ομάδας και αρνητικές μεταξύ χρηστών διαφορετικών ομάδων. Οι περιπτώσεις των 3A, 3B, 4A και 4B ανήκουν σε αυτήν την κατηγορία. Επιπλέον, κοιτώντας πιο προσεκτικά τις περιπτώσεις 4A και 4B παρατηρούμε την ύπαρξη δύο ομάδων, $\{B, C, D\}$ και $\{A\}$, οι οποίες σχηματίζονται κάτω από διαφορετικές συνθήκες: Στην περίπτωση του 4A, οι χρήστες B, C και D απορρίπτουν τον χρήστη A, ενώ στην περίπτωση 4B ο χρήστης A 'απομονώνεται' στρέφοντας αρνητικές ακμές προς την ομάδα $\{B, C, D\}$. Τέλος, υπάρχουν συχνοί υπογράφοι που δεν μπορούν να εξηγηθούν διαισθητικά όπως οι περιπτώσεις 3C και 4C.

3.4 Σύνοψη Κεφαλαίου

Στο κεφάλαιο αυτό παρουσιάσαμε μια διαφορετική προσέγγιση στην πρόβλεψη προσήμου σε κοινωνικούς γράφους, η οποία περιλαμβάνει την αναζήτηση μοτίβων-υπογράφων που εμφανίζονται συχνά στο κοινωνικό δίκτυο. Συμπεράναμε ότι οι συχνοί υπογράφοι διαθέτουν αρκετή ισχύ ώστε να προβλέψουν με ακρίβεια τις σχέσεις μεταξύ των χρηστών. Αξιολογήσαμε την μέθοδο μας μέσα από μια σειρά πειραμάτων, τα οποία περιελάμβαναν τρία μεγάλης κλίμακας κοινωνικά δίκτυα και επιβεβαιώσαμε ότι εμφανίζει καλύτερη επίδοση σε σχέση με τις τρέχουσες τεχνικές.



Εικόνα 3.5 Χαρακτηριστικοί υπογράφοι που συνδέονται με θετικές και αρνητικές ακμές-γεννήτορες

4

Εύρεση Δεσμών Εμπιστοσύνης με χρήση Τεχνικών Deep Learning

Σε αυτή την ενότητα, αναλύουμε τεχνικές και αλγορίθμους προερχόμενα από τον χώρο του deep learning [26], μια νέα γενιά αλγορίθμων μηχανικής μάθησης οι οποίοι βασίζονται σε νευρωνικά δίκτυα πολλών επιπέδων . Οι αλγόριθμοι αυτοί έχουν επιτύχει αξιοσημείωτες επιδόσεις σε ένα μεγάλο φάσμα προβλημάτων κυρίως βάσει της ικανότητας τους να παράγουν καλύτερες και αντιπροσωπευτικότερες μεταβλητές περιγραφής (feature engineering). Στο πλαίσιο του διδακτορικού επικεντρωθήκαμε σε δύο βασικούς deep learning αλγορίθμους : Restricted Boltzmann Machine [27], και Autoencoders [28]. Ο σκοπός του κεφαλαίου αυτού είναι να παρουσιάσει μια καινούρια οπτική στο πρόβλημα της πρόβλεψης εμπιστοσύνης σε ένα κοινωνικό δίκτυο, ιδωμένο στα πλαίσια ενός προβλήματος ημί-εποπτευόμενης μάθησης (semi-supervised learning). Η δομή του κεφαλαίου είναι η εξής: Πρώτα παρουσιάζεται η διατύπωση του προβλήματος όπου διαφαίνονται οι διαφορές με τα προηγούμενα κεφάλαια και έπειτα εξετάζεται η σχετική βιβλιογραφία. Το κύριο μέρος αφορά στην

περιγραφή των αλγορίθμων deep learning που χρησιμοποιήσαμε καθώς και την δική μας προσέγγιση, η οποία αξιολογείται μέσω ενός μεγάλου όγκου συνόλου δεδομένων.

4.1 Διατύπωση του προβλήματος

Σκοπός του κεφαλαίου είναι να παρουσιάσει έναν καινοτόμο αλγόριθμο για την ανίχνευση σχέσεων εμπιστοσύνης σε online κοινωνικά δίκτυα. Η διαφορά με τα προηγούμενα κεφάλαια έγκειται στο ότι εδώ θεωρούμε ότι γνωρίζουμε τις σχέσεις μεταξύ των χρηστών μόνο για ένα σχετικά μικρό κομμάτι του κοινωνικού γράφου και επιχειρούμε να καλύψουμε το έλλειμμα της πληροφορίας χρησιμοποιώντας "εξωτερικά" δεδομένα, όπως οι κριτικές των χρηστών για διάφορα αντικείμενα. Πιο συγκεκριμένα, τα δεδομένα που εξετάζουμε είναι της ακόλουθης μορφής:

- Ένας προσημασμένος κοινωνικός γράφος $G(V, E, L)$, όπου V είναι το σύνολο των κόμβων που αντιστοιχούν στους χρήστες, E είναι το σύνολο των (κατευθυνόμενων) ακμών και L το σύνολο των ετικετών που αποδίδονται στις ακμές
- ένα σύνολο αντικειμένων M ,
- ένα σύνολο από βαθμολογίες $R: V \times M \rightarrow \{1..K\}$, που δίνονται από έναν χρήστη $u \in V$ για ένα αντικείμενο $m \in M$

Όπως και στα προηγούμενα κεφάλαια, δύο χρήστες αντιστοιχούν σε γειτονικούς κόμβους εάν τουλάχιστον ένας από τους δύο έχει εκφράσει την άποψη του για τον άλλο με την ετικέτα $l \in L$ να υποδεικνύει το είδος της σχέσης. Στην περίπτωση μας αρκεί να ορίσουμε το σύνολο $L = \{-1, +1\}$, κατηγοριοποιώντας τις ακμές του γράφου σε θετικές (φιλικές) και αρνητικές (ανταγωνιστικές).

Ο σκοπός του παρόντος κεφαλαίου είναι να δει το αντικείμενο της πρόβλεψης προσήμου ως ένα πρόβλημα ημί-εποπτευόμενης μάθησης (semi-supervised learning).

Πιο συγκεκριμένα:

Ημι-εποπτευόμενη μάθηση προσήμων: Δεδομένου ενός μερικώς προσημασμένου κοινωνικού γράφου G_{known} , στον οποίο μόνο ένα ποσοστό $\epsilon\%$ των ακμών είναι προσημασμένες, ο στόχος μας είναι να "προσημάνουμε" και τον υπόλοιπο γράφο εκμεταλλευόμενοι τις κριτικές των χρηστών R για τα αντικείμενα $m \in M$.

Η προσέγγιση αυτή είναι σχετική με την έννοια του *transfer learning* [29] [30], όπου η γνώση που αποκτάται σε έναν τομέα, χρησιμοποιείται για να λύσει ένα πρόβλημα σε έναν άλλο. Για το υπόλοιπο του κεφαλαίου θεωρούμε ότι το δίκτυο γνωριμιών (acquaintance network) μεταξύ όλων των χρηστών είναι γνωστό, γνωρίζοντας όμως τις ετικέτες για ένα μικρό ποσοστό των ακμών αυτών. Με άλλα λόγια, ξέρουμε ότι δύο χρήστες έχουν αλληλεπιδράσει μεταξύ τους, αλλά δεν ξέρουμε πάντα το είδος της αλληλεπίδρασης (συμπάθεια ή αντιπάθεια).

4.2 Μέθοδος αναφοράς

Σε αυτή την ενότητα παρουσιάζουμε την μέθοδο αναφοράς, η οποία πετυχαίνει την μεγαλύτερη ακρίβεια στην βιβλιογραφία. Πρόκειται για την δουλειά των Yang et al. [31] οι οποίοι για πρώτη φορά δείχνουν έναν πρακτικό τρόπο προσήμανσης ενός κοινωνικού γράφου που κυριαρχείται από ακμές αγνώστου ετικέτας. Με άλλα λόγια, δείχνουν πως ένα απλό δίκτυο γνωριμιών (π.χ. Facebook) μπορεί να μετατραπεί σε ένα προσημασμένο γράφο εμπιστοσύνης (π.χ. Epinions). Για να το επιτύχουν συνδυάζουν τεχνικές μη- και ημι-εποπτευόμενης μάθησης, μοντελοποιώντας την συμπεριφορά των χρηστών στην λήψη αποφάσεων. Οι συγγραφείς επαληθεύουν την ορθότητα του μοντέλου τους σε δύο μεγάλα σύνολα δεδομένων από το Epinions και το Yahoo! Pulse.

Κάθε χρήστης u αναπαριστάται μέσω ενός διανύσματος λανθανόντων παραγόντων φ_u και κάθε αντικείμενο i με ένα διάνυσμα ψ_i . Η βαθμολογία ενός χρήστη u προς το αντικείμενο i είναι ανάλογη του γινομένου $\varphi_u^T \psi_i$ ενώ οι σχέσεις εμπιστοσύνης μεταξύ δύο χρηστών u, v χαρακτηρίζονται από το γινόμενο $\varphi_u^T \varphi_v$, όπου μεγαλύτερες τιμές συνεπάγονται μεγαλύτερη πιθανότητα θετικής ακμής. Τέλος, ένας χρήστης u λαμβάνει απόφαση για ένα καινούριο προϊόν βάσει της ακόλουθης διαδικασίας: Με πιθανότητα p επηρεάζεται από τις δική του "πεποίθηση" $\varphi_u^T \psi_i$ για το προϊόν, ενώ με πιθανότητα $1 - p$ λαμβάνει υπόψη τις προτιμήσεις ενός φίλου του v , με πιθανότητα $\varphi_u^T \varphi_v$. Οι συγγραφείς δείχνουν έναν αποτελεσματικό τρόπο εκπαίδευσης του μοντέλου μέσω *βαθμιαίας καθόδου* (gradient descent) [32] [33] [34] σε μια κατάλληλη συνάρτηση κόστους, η οποία μετράει τις αποκλίσεις των προβλέψεων από τα πειραματικά δεδομένα.

Οι διαφορές της προτεινόμενης από εμάς μεθόδου σε σχέση με αυτή της αναφοράς είναι οι ακόλουθες:

- Στην δική μας προσέγγιση μοντελοποιούμε απ' ευθείας τους χρήστες (μέσω διανυσμάτων λανθανόντων παραγόντων) αποφεύγοντας έτσι το επιπλέον κόστος μοντελοποίησης των αντικειμένων, μιας και δεν σχετίζονται με το πρόβλημα προσήμανσης του κοινωνικού γράφου.
- Η μέθοδος αναφοράς επιλύει το πρόβλημα κάνοντας αποκλειστικά γραμμικές μετατροπές στα δεδομένα. Αντίθετα, η δική μας προσέγγιση περιλαμβάνει και μη-γραμμικούς μετασχηματισμούς μέσω των σιγμοειδών (sigmoids) συναρτήσεων που υλοποιούνται στις μονάδες των κρυμμένων επιπέδων.

4.3 Αλγόριθμοι και Εργαλεία

Στην ενότητα αυτή παρουσιάζουμε αλγόριθμους και τεχνικές από τον τομέα του deep learning, τα οποία χρησιμοποιούμε αργότερα ως βασικά συστατικά της δικής μας προσέγγισης. Συγκεκριμένα, ασχολούμαστε με τους αλγόριθμους Restricted Boltzmann Machine και Autoencoder.

4.3.1 Restricted Boltzmann Machines (RBMs)

Τα RBMs ανήκουν στην οικογένεια των ενεργειακών (energy-based) [28] [35] μοντέλων όπου για κάθε διάταξη των μεταβλητών του υπό εξέταση προβλήματος αντιστοιχεί μια πεπερασμένη τιμή ενέργειας. Για την εκπαίδευση τέτοιων μοντέλων απαιτείται ο ορισμός μιας συνάρτησης ενέργειας (energy function), η οποία να κατέχει κάποιες επιθυμητές ιδιότητες. Για παράδειγμα, ένα ενεργειακό μοντέλο θα μπορούσε να εκπαιδευτεί ώστε να δίνει ως έξοδο μια χαμηλή τιμή ενέργειας όταν η διάταξη των μεταβλητών εμφανίζεται συχνά στα πειραματικά δεδομένα, δίνοντας έτσι έναν μοντέλο γεννήτορα (generative model). Επιπλέον, η συνάρτηση ενέργειας δίνει την δυνατότητα να οριστούν κατανομές πιθανότητας της μορφής :

$$P(x) = \frac{e^{-Energy(x)}}{Z} \quad (4.1)$$

όπου x είναι το σύνολο των μεταβλητών και Z ένας παράγοντας κανονικοποίησης με

$$\text{τιμή } Z = \sum_x e^{-Energy(x)} \quad (4.2)$$

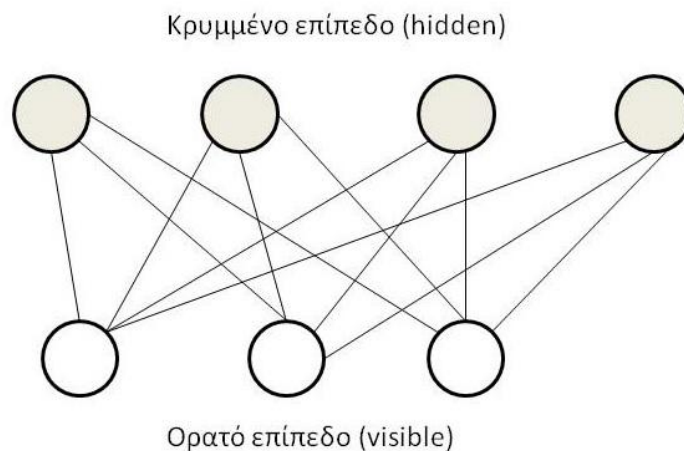
(Το Z στην βιβλιογραφία καλείται και *partition function*).

Σε αρκετά προβλήματα δεν υπάρχει η δυνατότητα να παρατηρηθούν ταυτόχρονα όλες οι τιμές των μεταβλητών x_i . Μια δημοφιλής επιλογή σε αυτές τις περιπτώσεις είναι να διαχωρίζονται οι μεταβλητές σε δύο ομάδες: την ορατή (visible) και την κρυφή (hidden). Η κοινή κατανομή πιθανότητας ορίζεται τότε ως :

$$P(x, h) = \frac{e^{-Energy(x,h)}}{Z} \quad (4.3)$$

όπου x, h είναι η ορατή και κρυφή ομάδα αντίστοιχα.

Ο αλγόριθμος Restricted Boltzmann Machine ακολουθεί μια αντίστοιχη φιλοσοφία. Πιο συγκεκριμένα, πρόκειται για ένα συμμετρικό γραφικό μοντέλο, με ένα ορατό και ένα κρυμμένο επίπεδο, και την επιπλέον ιδιαιτερότητα ότι οι μονάδες (units) του ενός επιπέδου συνδέονται (και άρα εξαρτώνται) *μόνο* με μονάδες του άλλου επιπέδου. Ο περιορισμός αυτός καθιστά πιο εύκολη την εκπαίδευση τέτοιων διατάξεων. Η Εικόνα δίνει μια γραφική απεικόνιση ενός RBM : Το επάνω, γκρι σκιασμένο επίπεδο είναι αντιστοιχεί στις κρυμμένες μονάδες ενώ το κάτω στις φανερές.



Εικόνα 4.1 Αναπαράσταση ενός Restricted Boltzmann Machine

Η συνάρτηση ενέργειας ενός RBM με V ορατές μονάδες και H κρυφές ορίζεται ως :

$$E(v, h) = - \sum_{i=1}^V \sum_{j=1}^H v_i h_j w_{ij} - \sum_{i=1}^V v_i b_i^v - \sum_{j=1}^H h_j b_j^h \quad (4.4)$$

όπου,

- v είναι ένα δυαδικό (binary) διάνυσμα κατάστασης των ορατών μονάδων
- h είναι ένα δυαδικό (binary) διάνυσμα κατάστασης των κρυμμένων μονάδων
- v_i είναι η κατάσταση της φανεράς μονάδας i
- h_j είναι η κατάσταση της κρυμμένης μονάδας j

- w_{ij} είναι το βάρος της σύνδεσης μεταξύ της φανεράς μονάδας i και της κρυμμένης μονάδας j
- b_i^v και b_j^h είναι η προδιάθεση (bias) της φανεράς μονάδας i και της κρυμμένης μονάδας j αντίστοιχα.

Η δεσμευμένη πιθανότητα $p(v|h)$ δίνεται από την σχέση:

$$p(v|h) = \frac{e^{-E(v,h)}}{\sum_g e^{-E(g,h)}} \quad (4.5)$$

ενώ στην ειδική περίπτωση που εξετάζουμε μια συγκεκριμένη μονάδα i του ορατού επιπέδου, η δεσμευμένη κατανομή πιθανότητας που της αποδίδεται, γνωρίζοντας την κατάσταση του το κρυφού επιπέδου h είναι :

$$p(v_k = 1|h) = \frac{1}{1+e^{-\sum_{j=1}^H h_j w_{kj} + b_k^v}} \quad (4.6)$$

Αντίστοιχα, οι δεσμευμένες πιθανότητες $p(h|v)$ και $p(h_k = 1|v)$ ορίζονται από τις σχέσεις :

$$p(h|v) = \frac{e^{-E(v,h)}}{\sum_g e^{-E(v,g)}} \quad \text{και} \quad p(h_k = 1|v) = \frac{1}{1+e^{-\sum_{i=1}^V v_i w_{kj} + b_k^h}} \quad (4.7)$$

Οι σχέσεις αυτές εκφράζουν την ανεξαρτησία των μονάδων των δύο επιπέδων από τις μονάδες του ίδιου επιπέδου.

4.3.1.1 Εκπαίδευση ενός RBM

Όπως επισημάνθηκε και στην εισαγωγή του κεφαλαίου, τα ενεργειακά μοντέλα μπορούν να εκπαιδευτούν εάν δοθούν οι κατάλληλοι περιορισμοί που αφορούν την ενέργεια κάθε διάταξης. Στην περίπτωση των RBMs και δεδομένου ενός συνόλου παρατηρήσεων, αυτό που επιζητούμε είναι να μεγιστοποιήσουμε την πιθανότητα της εμφάνισης αυτών των δεδομένων από το μοντέλο μας, δίνοντας τους χαμηλότερη ενέργεια. Πιο συγκεκριμένα, έστω ότι δίνεται ένα σύνολο $|C|$ περιπτώσεων εκπαίδευσης (training cases) $\{v^c | c \in C\}$, τότε η εκπαίδευση του RBM συνίσταται

στην εύρεση τιμών για τις παραμέτρους του, οι οποίες να μεγιστοποιούν την μέση λογαριθμική πιθανότητα εμφάνισης του συνόλου C :

$$\sum_{c=1}^C \log p(v^c) = \sum_{c=1}^C \log \frac{\sum_g e^{-E(v^c, g)}}{\sum_u \sum_g e^{-E(u, g)}} \quad (4.8)$$

Διάφορα προβλήματα μηχανικής μάθησης συνήθως βελτιστοποιούνται μέσω *gradient descent*, όπου σε κάθε βήμα εκτέλεσης οδηγούμαστε την τελική λύση αξιοποιώντας τις πληροφορίες που δίνει η *gradient* συνάρτηση σε εκείνο το σημείο. Έστω ότι θέλουμε να βρούμε την εξίσωση για την ανανέωση των τιμών των βαρών w_{ij} . Παίρνοντας την μερική παράγωγο της συνάρτησης κόστους ως προς w_{ij} και με χρήση της (4.5) παίρνουμε :

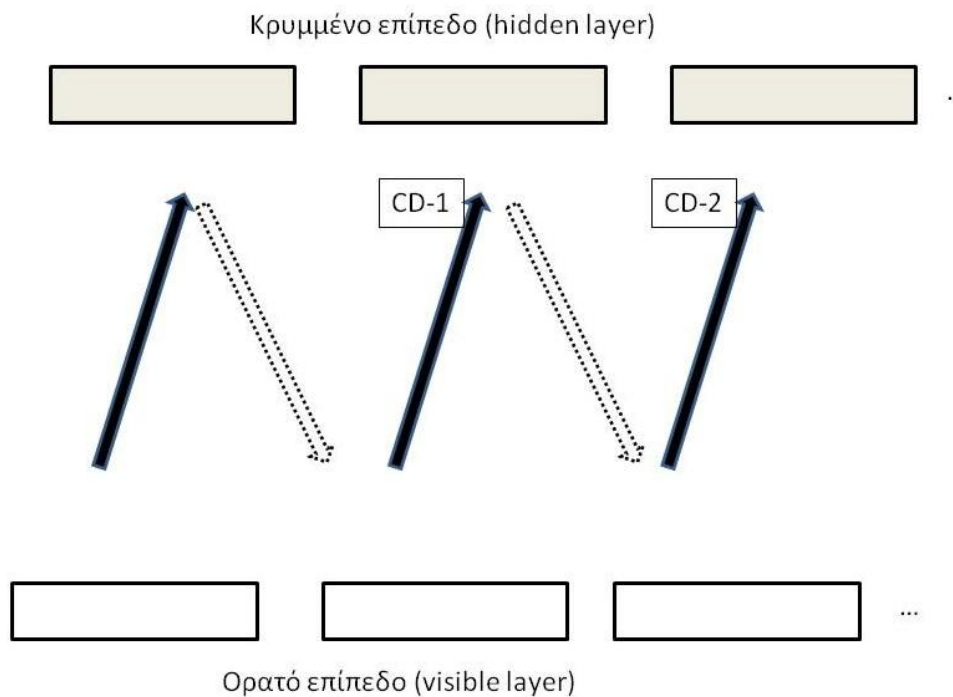
$$\frac{\partial}{\partial w_{ij}} \sum_{c=1}^C \log p(v^c) = \frac{\partial}{\partial w_{ij}} (\sum_{c=1}^C \log \sum_g e^{-E(v^c, g)} - \log \sum_u \sum_g e^{-E(u, g)}) \quad (4.9)$$

Ο πρώτος όρος υπολογίζει με την μέση τιμή της παράστασης $v_i^c g_j$ όταν το ορατό επίπεδο του RBM "οδηγείται" από τα δεδομένα v^c , ενώ ο δεύτερος όρος αντιστοιχεί στην τιμή της παράστασης $v_i g_j$ όταν τα δεδομένα "παράγονται" από το μοντέλο. Ένας ισοδύναμος τρόπος διατύπωσης θα έλεγε ότι κάθε βάρος w_{ij} θα πρέπει να υποστεί μεταβολή ίση με :

$$\Delta w_{ij} = \varepsilon_w (E_{data}[v_i h_j] - E_{model}[v_i h_j]) \quad (4.10)$$

ώστε το RBM μοντέλο να αρχίζει να προσεγγίζει την πραγματική κατανομή των δεδομένων. Ο πρώτος όρος $E_{data}[v_i h_j]$ στην (4.10) είναι εύκολο να υπολογιστεί αφού ξέροντας τις τιμές των μονάδων στο ορατό επίπεδο, μπορούμε μέσω των εξισώσεων (4.7) να υπολογίσουμε την δεσμευμένη πιθανότητα για κάθε μονάδα στο κρυμμένο επίπεδο. Πιο δύσκολος είναι ο υπολογισμός του δεύτερου όρου, ο οποίος προϋποθέτει την ύπαρξη δειγμάτων από το μοντέλο αυτό καθ' αυτό. Ο πιο σωστός

τρόπος είναι να αρχικοποιήσουμε τις μονάδες στο ορατό επίπεδο σε τυχαίες τιμές και μετά να πάρουμε δείγματα μέσω Gibbs sampling μέσω των εξισώσεων (4.6) και (4.7), επαναλαμβάνοντας για αρκετές φορές, αναγκάζοντας το μοντέλο να "ξεχάσει" την αρχική (τυχαία) κατάσταση. Δυστυχώς, αυτή η μέθοδος δεν είναι αποδοτική γιατί απαιτεί μεγάλο χρόνο εκτέλεσης. Η λύση που προκρίνεται σε αυτές τις περιπτώσεις περιλαμβάνει βελτιστοποίηση μέσω Contrastive Divergence (CD) [36], η οποία απεικονίζεται στην Εικόνα 4.2



Εικόνα 4.2 Εκμάθηση RBM μέσω Contrastive Divergence

Η μέθοδος CD επιχειρεί να προσεγγίσει την ποσότητα $E_{model}[v_i h_j]$, εκτελώντας την δειγματοληψία για έναν μικρό αριθμό επαναλήψεων. Οι μονάδες στο ορατό επίπεδο αρχικοποιούνται σε ένα δείγμα από τα υπάρχοντα πραγματικά δεδομένα και μέσω των εξισώσεων (4.6) και (4.7) εκτελούμε N επαναλήψεις, παίρνοντας κάποια "ανακατασκευασμένα" δεδομένα λόγω της συνεισφοράς του μοντέλου. Η μέθοδος CD θα αναλάβει να δώσει χαμηλότερη ενέργεια στα πραγματικά δεδομένα και αρκετά υψηλότερη ενέργεια στις "ανακατασκευές" που προκύπτουν από αυτά,

βοηθώντας έτσι το μοντέλο να προσεγγίσει την πραγματική κατανομή των δεδομένων.

4.3.1.2 RBMs σε περιπτώσεις συνεργατικού φιλτραρίσματος

Για την ανάπτυξη της μεθόδου μας, η οποία θα παρουσιαστεί στην επόμενη ενότητα, χρησιμοποιήσαμε μια επέκταση του RBM μοντέλου όπως αυτή παρουσιάστηκε από τους Salakhutdinov et al. στο [37]. Πρόκειται για μια παραλλαγή προσαρμοσμένη σε προβλήματα συνεργατικού φιλτραρίσματος (collaborative filtering) [7], με σκοπό να προβλέψει τις βαθμολογίες των χρηστών για διάφορα αντικείμενα όπως ταινίες. Οι συγγραφείς ασχολήθηκαν με δεδομένα από το Netflix⁸, μια ηλεκτρονική πλατφόρμα ενοικίασης ταινιών πετυχαίνοντας 6% καλύτερη ακρίβεια σε σχέση με τους αλγόριθμους που χρησιμοποιούνταν από την εταιρία. Στο [37] οι συγγραφείς εξηγούν ότι είναι αναγκαίες κάποιες τροποποιήσεις στο παραδοσιακό μοντέλο του RBM ώστε να αντιμετωπιστούν χαρακτηριστικά ζητήματα των μηχανισμών συνεργατικού φιλτραρίσματος όπως αυτό της αραιότητας (*sparseness*) των δεδομένων, όπου ο χρήστης αξιολογεί μόνο έναν πολύ μικρό αριθμό των διαθέσιμων αντικειμένων. Συγκεκριμένα, ορίζουν ένα διαφορετικό RBM για κάθε χρήστη, επιβάλλοντας παράλληλα τους αντίστοιχους περιορισμούς:

- Κάθε RBM έχει τον ίδιο αριθμό μονάδων στο κρυμμένο επίπεδο, αλλά το RBM κάθε χρήστη έχει ενεργές ορατές μονάδες **μόνο** για τις ταινίες που αξιολογήθηκαν από τον συγκεκριμένο χρήστη.
- Σε κάθε RBM αντιστοιχεί μόνο ένα παράδειγμα εκπαίδευσης, αλλά τα βάρη των συνδέσεων w_{ij} και τα biases b_i^v, b_j^h είναι κοινά για όλα τα RBMs.

⁸ www.netflix.com

Το δεύτερο σημείο υποδηλώνει ότι εάν δύο χρήστες αξιολόγησαν την ίδια ταινία, τότε τα RBMs που τους αντιστοιχούν θα πρέπει να έχουν τα ίδια βάρη μεταξύ της ορατής μονάδας (που αντιστοιχεί στην ταινία) και των μονάδων του κρυμμένου επιπέδου. Από την άλλη, η κατάσταση του κρυμμένου επιπέδου μπορεί να είναι διαφορετική για κάθε RBM. Κάθε ορατή μονάδα μοντελοποιείται ως μια δεσμευμένη πολυωνυμική (multinomial) κατανομή ενώ κάθε μονάδα του κρυμμένου επιπέδου θεωρείται, όπως και πριν, ως μεταβλητή Bernoulli. Πιο συγκεκριμένα, για κάθε χρήστη που έχει αξιολογήσει m προϊόντα με βαθμολογίες από το σύνολο $\{1..K\}$, κατασκευάζουμε έναν δυαδικό ενδεικτικό πίνακα V μεγέθους $K \times m$ όπου κάθε κελί παίρνει τιμή $v_i^k = 1$, αν ο χρήστης βαθμολόγησε την ταινία i με k και 0 σε διαφορετική περίπτωση. Οι εκφράσεις για τις δεσμευμένες πιθανότητες των δύο επιπέδων γίνονται:

$$p(h_j = 1|V) = 1/(1 + \exp(a_j + \sum_{i=1}^m \sum_{j=1}^d v_i^k W_{ij}^k)), j = 1..d \quad (1)$$

και

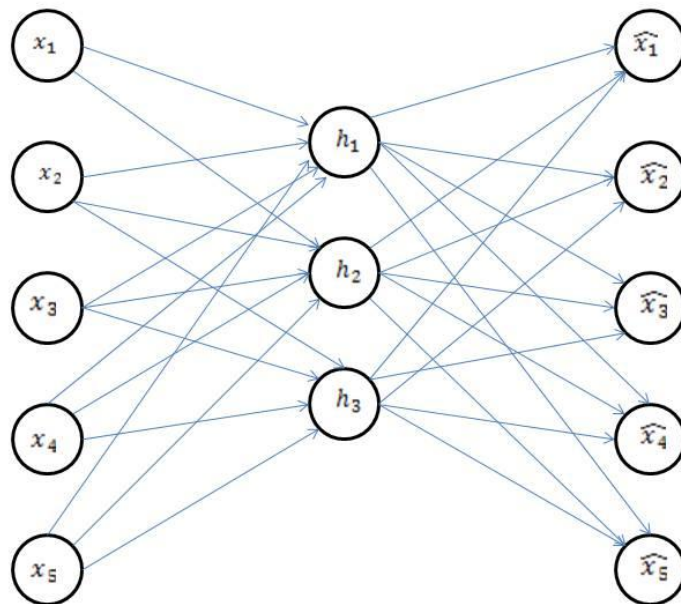
$$p(v_i^k = 1|h) = \frac{\exp(b_i^k + \sum_{j=1}^d h_j W_{ij}^k)}{\sum_{l=1}^K \exp(b_i^l + \sum_{j=1}^d h_j W_{ij}^l)}, i = 1..m \quad (2),$$

Η εκπαίδευση το μοντέλου γίνεται και πάλι μέσω Contrastive Divergence. Για λόγους απλότητας από εδώ και στο εξής, θα αναφερόμαστε στο γραφικό μοντέλο που παρουσιάζεται στο [37] με τον όρο *CF-RBM (Collaborative Filtering RBM)*.

4.3.2 Αυτοκωδικοποιητές (Autoencoders)

Ο δεύτερος αλγόριθμος με τον οποίο ασχοληθήκαμε είναι ο αυτοκωδικοποιητής [28]. Πρόκειται για ένα νευρωνικό δίκτυο με τρία επίπεδα (είσοδο, έξοδο και κρυμμένο) το

οποίο προορίζεται για περιπτώσεις μη επιτευόμενης μάθησης (unsupervised learning). Ένας αυτοκωδικοποιητής μπορεί να εκπαιδευτεί με χρήση back-propagation θέτοντας το επίπεδο εξόδου ίσο με το επίπεδο εισόδου. Με άλλα λόγια προσπαθεί να προσομοιώσει την ταυτοτική συνάρτηση (identity function), ασκώντας μη γραμμικές μεταμορφώσεις στα δεδομένα. Η Εικόνα 4.3 δίνει ένα παράδειγμα αυτοκωδικοποιητή με 5 μονάδες στα επίπεδα εισόδου και εξόδου καθώς και 3 μονάδες στο κρυμμένο επίπεδο.



Εικόνα 4.3 Δίκτυο αυτοκωδικοποιητή. Τα δεδομένα x_i περνούν μέσα από το κρυμμένο επίπεδο και αναδύονται εκ νέου στην έξοδο (\hat{x}_i)

Η επίδραση του αυτοκωδικοποιητή πάνω στα δεδομένα καθορίζεται από την αρχιτεκτονική του. Για παράδειγμα, εάν το μέγεθος του κρυμμένου επιπέδου είναι μικρότερο από το επίπεδο εισόδου τότε ο αυτοκωδικοποιητής πραγματοποιεί μια μη γραμμική μείωση διαστάσεων (dimensionality reduction) [38] . Από την άλλη

πλευρά, εάν το μέγεθος του κρυμμένου επιπέδου είναι μεγαλύτερο τότε συνήθως οριοθετούνται επιπλέον περιορισμοί που αποτρέπουν το δίκτυο από το να φτάσει σε μια τετριμμένη λύση. Ένα τέτοιο παράδειγμα είναι του αραιού αυτοκωδικοποιητή (sparse autoencoder) [39], όπου η επιπλέον συνθήκη που επιβάλλεται αναγκάζει τις μονάδες στο κρυμμένο επίπεδο να παραμένουν αδρανείς τις περισσότερες φορές. Άλλη περίπτωση τέτοιου δικτύου είναι ο στοιβαγμένος αυτοκωδικοποιητής, όπου πολλαπλά δίκτυα αυτοκωδικοποιητών στοιβάζονται σε επίπεδα, κατά τρόπο ώστε η έξοδος του ενός να αποτελεί είσοδο για το επόμενο.

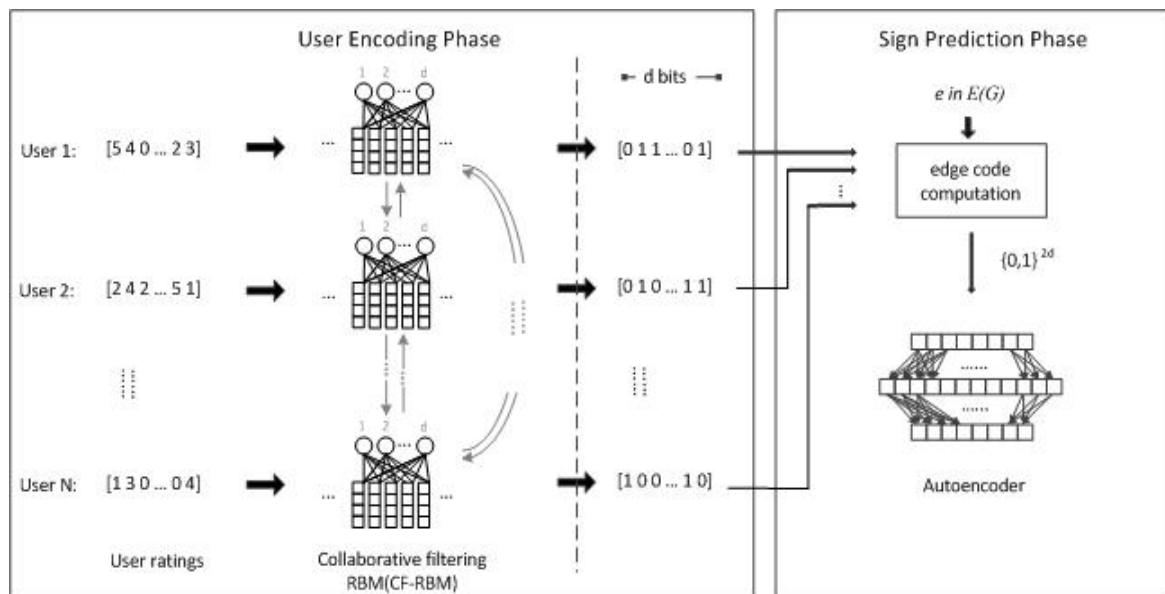
Οι αυτοκωδικοποιητές, όπως και κάθε άλλος deep learning αλγόριθμος, είναι ιδιαίτερα χρήσιμοι σε προβλήματα ημι-επιτευόμενης μάθησης, όπου υπάρχει ένας μεγάλος αριθμός ακατηγοριοποίητων δεδομένων (unlabeled data) και ένα σημαντικά μικρότερος αριθμός εποπτευμένων δεδομένων. Η συνήθης πρακτική είναι να αφήνουμε τον αυτοκωδικοποιητή να «μαθαίνει» την κατανομή των δεδομένων μέσω μη εποπτευόμενης μάθησης και έπειτα να χρησιμοποιούμε τις παραμέτρους του για να αρχικοποιήσουμε ένα κλασικό νευρωνικό δίκτυο ίδιων διαστάσεων το οποίο θα παράξει το τελικό μοντέλο για την κατηγοριοποίηση των δεδομένων. Ο λόγος που πιστεύεται ότι μια τέτοια τεχνική θα είναι αποδοτική οφείλεται στην υποβόσκουσα υπόθεση ότι για μια συνάρτηση με πεδίο ορισμού X και πεδίο τιμών Y , «μαθαίνοντας» την πιθανοτική κατανομή του X θα καταστεί χρήσιμο για την εκτίμηση της δεσμευμένης κατανομής του Y . Με άλλα λόγια, θεωρούμε ότι

$$P(X) \rightarrow P(Y | X)$$

Για την ανάπτυξη της δικής μας μεθόδου χρησιμοποιήσαμε στοιβαγμένους αραιούς αποκωδικοποιητές με στόχο 0.05 (δηλ, οι μονάδες να είναι ενεργές μόνο το 5% του χρόνου) και πειραματιστήκαμε με διάφορες αρχιτεκτονικές μετρώντας τις επιπτώσεις στην απόδοση.

4.4 Προτεινόμενη μέθοδος

Αυτή η ενότητα παρουσιάζει την προτεινόμενη μέθοδο για την επίλυση του ημι-εποπτευόμενου προβλήματος μάθησης όπως αυτό ορίστηκε στην αρχή του κεφαλαίου. Η προσέγγιση μας χωρίζεται σε δύο στάδια: Το πρώτο αφορά σε έναν μηχανισμό, ο οποίος αποδίδει πεπερασμένους δυαδικούς κώδικες σε κάθε χρήστη, ενώ το δεύτερο στάδιο χρησιμοποιεί τους κώδικες αυτούς για να υλοποιήσει τον τελικό ταξινομητή. Η Εικόνα 4.4 δίνει ένα περίγραμμα του αλγορίθμου, ενώ οι επόμενες υπο-ενότητες δίνουν τις απαραίτητες λεπτομέρειες.



Εικόνα 4.4 Η μέθοδος μας εκτελείται σε δύο στάδια. α) Εύρεση κωδικών για τους χρήστες: Εκπαιδεύουμε το CF-RBM βάσει των κριτικών των χρηστών για τα προϊόντα. Ένα ξεχωριστό RBM αντιστοιχεί σε κάθε χρήστη. Μετά το πέρας της εκπαίδευσης, υπολογίζουμε την ενεργοποίηση των μονάδων στον κρυμμένο επίπεδο, β) Φάση ταξινόμησης: Λαμβάνουμε τον κωδικό κάθε ακμής, συνδυάζοντας τους κώδικες των χρηστών που βρίσκονται στα άκρα της. Το αποτέλεσμα τροφοδοτείται στο δίκτυο του αυτοκωδικοποιητή για την δημιουργία του ταξινομητή

1ο βήμα - Κωδικοποίηση χρηστών

Σε αυτό το στάδιο, κάθε χρήστης λαμβάνει έναν δυαδικό κωδικό βάσει των κριτικών του για τα προϊόντα που έχει αξιολογήσει. Για να το πετύχουμε αυτό, βασιζόμαστε στο CF-RBM μοντέλο, που παρουσιάστηκε στην προηγούμενη ενότητα. Εκκινούμε εκπαιδεύοντας το μοντέλο για 30 εποχές (epochs), τοποθετώντας 100 μονάδες στο κρυμμένο επίπεδο. Έπειτα, μετά το πέρας της εκπαίδευσης, "περνάμε" τα δεδομένα εκ νέου από το RBM, υπολογίζοντας αυτή τη φορά την πιθανότητα κάθε μονάδας στο κρυμμένο επίπεδο να είναι ενεργή $p(h_j = 1|V)$ βάσει της (4.7). Εν συνεχεία, κάθε RBM θα πάρει μια στοχαστική απόφαση για το αν θα 'ενεργοποιήσει' ($h_j = 1$) ή όχι ($h_j = 0$), $j = 1..d$, κάθε μονάδα. Εφόσον κάθε χρήστης έχει το δικό του RBM και η κατάσταση του κρυμμένου επιπέδου του κάθε RBM είναι ανεξάρτητη των άλλων, οι προτιμήσεις του κάθε χρήστη θα καθορίσουν το σύνολο των ενεργοποιημένων μονάδων. Τελικά, ο δυαδικό κωδικός για κάθε χρήστη θα προκύψει αντανακλώντας τις καταστάσεις των μονάδων αυτών. Για παράδειγμα, αν οι καταστάσεις των τριών πρώτων μονάδων είναι {"off", "off", "on",...}, ο αντίστοιχος κωδικός για τον χρήστη θα ξεκινά με [0 0 1 ...]. Προφανώς, το μήκος του κώδικα ισούται με το μέγεθος του κρυμμένου επιπέδου, το οποίο εδώ λαμβάνει την τιμή 100. Αυτό το βήμα λειτουργεί ως μετάβαση από ένα χώρο μεγάλο αριθμό διαστάσεων (αν κάθε αντικείμενο θεωρηθεί μια διάσταση) σε ένα πολύ χαμηλότερο χώρο δυαδικών διανυσμάτων.

2ο βήμα - Μοντελοποίηση και ταξινόμηση με αυτοκωδικοποιητές

Το βήμα αυτό είναι υπεύθυνο για την δημιουργία του ταξινομητή χρησιμοποιώντας δίκτυα αυτοκωδικοποιητών και εκμεταλλεύόμενο τους κώδικες των χρηστών από το προηγούμενο βήμα. Θεωρούμε ότι το δίκτυο γνωριμιών (ακμές χωρίς πρόσημο) είναι γνωστό για όλους τους χρήστες του δικτύου και ότι διαθέτουμε ακόμα ένα μικρό

ποσοστό προσημασμένων ακμών. Ξεκινάμε υπολογίζοντας έναν κωδικό για κάθε ακμή του γράφου, τοποθετώντας τον ένα μετά τον άλλο τους κώδικες των κόμβων που την ορίζουν. Εφόσον εξετάζουμε δίκτυα με κατευθυνόμενες ακμές, επιλέγουμε ως σύμβαση να βάζουμε πρώτο τον κωδικό του κόμβου-πηγής και έπειτα του κόμβου προορισμού π.χ. για την ακμή $u \rightarrow v$ ο κωδικός της είναι

$$code(u \rightarrow v) := [[\leftarrow code(u) \rightarrow], [\leftarrow code(v) \rightarrow]].$$

Έχοντας υπολογίσει τον κώδικα κάθε ακμής προχωρούμε στην δημιουργία του ταξινομητή βασιζόμενοι στον αριθμό των προσημασμένων ακμών που διαθέτουμε. Αρχικά, τροφοδοτούμε τους κώδικες των ακμών στο δίκτυο του αυτοκωδικοποιητή και τον εκπαιδεύουμε με μη-εποπτευόμενο τρόπο (unsupervised learning), με την ελπίδα ότι θα εκτιμήσει σωστά την κατανομή των ακμών, το οποίο θεωρούμε κρίσιμο για την επιτυχία της ταξινόμησης. Έπειτα χρησιμοποιούμε τις παραμέτρους του αυτοκωδικοποιητή για να αρχικοποιήσουμε ένα παραδοσιακό νευρωνικό δίκτυο, το οποίο και εκπαιδεύουμε χρησιμοποιώντας μόνο τις προσημασμένες ακμές. Είναι σαφές ότι για να λειτουργήσει το νευρωνικό δίκτυο ως ταξινομητής, χρειάζεται ένα επιπλέον επίπεδο εξόδου, το οποίο θα αναλάβει να κωδικοποιήσει το πρόσημο των ακμών. Για τον σκοπό αυτό χρησιμοποιούμε ‘one-hot encoding’: κάθε αρνητικός δεσμός αντιστοιχεί σε έξοδο [1 0] ενώ ένας θετικός σε [0 1]. Τέλος, ο αλγόριθμος back-propagation μέσω *gradient descent* βρίσκει τις τελικές κατάλληλες τιμές των παραμέτρων του ταξινομητή.

4.5 Αξιολόγηση των αποτελεσμάτων

4.5.1 Περιγραφή των δεδομένων

Εκτελέσαμε τα πειράματά μας σε ένα μεγάλης κλίμακας σύνολο δεδομένων από 130 χιλιάδες χρήστες από τον ιστότοπο αξιολόγησης προϊόντων Epinions, το οποίο συλλέχθηκε από τους Massa et al. [40]. Η περίπτωση του Epinions είναι ενδιαφέρουσα μιας και εκτός της αξιολόγησης προϊόντων, οι χρήστες μπορούν να αλληλεπιδρούν μεταξύ τους: ορισμένοι χρήστες γράφουν κριτικές για κάποια προϊόντα, οι οποίες αργότερα βαθμολογούνται από άλλους χρήστες. Επίσης, οι χρήστες μπορούν να δηλώσουν με σαφήνεια άλλους χρήστες ως *φίλους* ή *εχθρούς*. Τα δεδομένα περιέχουν πληροφορίες και για τις δύο παραμέτρους: Υπάρχει ένας προσημασμένος κοινωνικός γράφος βάσει 840 χιλιάδων δηλώσεων εμπιστοσύνης (και μη), ένα σύνολο από 1.5 εκατομμύριο κριτικές προϊόντων και 13,6 εκατομμύρια βαθμολογίες.

4.5.2 Εργαλεία και μετρικές αξιολόγησης

Ο κώδικας για το CF-RBM υλοποιήθηκε σε Java, ενώ για την ανάλυση με τους αυτοκωδικοποιητές χρησιμοποιήθηκε η “DeepLearnToolbox”, μια βιβλιοθήκη για deep learning γραμμένη σε MATLAB [41]. Πραγματοποιήσαμε περαιτέρω στατιστική ανάλυση με την χρήση του στατιστικού περιβάλλοντος R καθώς και τα πακέτα ROCR [42], LibLinear [43] and e1071 [44]. Το μεγαλύτερο πρόβλημα που αντιμετωπίσαμε κατά την αξιολόγηση είχε να κάνει με την ανισομερή κατανομή των προσήμων, αφού περίπου 80% των ακμών ήταν θετικές επιτρέποντας ακόμα και σε έναν τυχαίο, θετικά προκατειλημμένο ταξινομητή να πετύχει σημαντική ακρίβεια. Για τον λόγο αυτό, επιλέξαμε ως μετρική το εμβαδόν της επιφανείας κάτω από την γραφική παράσταση της ROC (ROC-AUC), η οποία δεν επηρεάζεται από την

ανισοκατανομή των κλάσεων . Δεδομένου ότι το σύνολο των εξ' αρχής γνωστών προσημασμένων ακμών είναι τυχαίο, επαναλάβουμε τα πειράματα μας για 10 φορές και καταγράψαμε την μέση απόδοση, ώστε να περιορίσουμε την ισχύ της τυχειότητας στην πειραματική διαδικασία.

4.5.3 Ανάλυση ακρίβειας μοντέλου

Αρχιτεκτονική	Ποσοστό των προσημασμένων ακμών που είναι εξ' αρχής γνωστές										
	0.1%	0.3%	0.5%	1.0%	5.0%	10%	20%	30%	50%	70%	90%
200-750	0.7544	0.7687	0.7780	0.8016	0.8476	0.8687	0.8846	0.8947	0.9070	0.9153	0.9232
200-1000	0.7690	0.7707	0.7873	0.8060	0.8535	0.8801	0.8988	0.9105	0.9200	0.9284	0.9338
200-2000	0.7345	0.7813	0.7984	0.8156	0.8551	0.8780	0.8955	0.9082	0.9195	0.9295	0.9353
200-500-200	0.7405	0.7731	0.7776	0.8084	0.8491	0.8623	0.8878	0.8968	0.9115	0.9209	0.927
200-1000-500-200	0.7267	0.7802	0.7844	0.8004	0.8401	0.8608	0.8809	0.8936	0.9074	0.9185	0.9280
200-800-800-2000	0.7134	0.7955	0.8029	0.8066	0.8443	0.8563	0.8778	0.8938	0.9053	0.9195	0.9269
200-50	0.7450	0.7564	0.7762	0.7916	0.8360	0.8580	0.8648	0.8683	0.8753	0.8807	0.8853
200-50-50	0.7389	0.7627	0.7847	0.7739	0.8272	0.8345	0.8442	0.8499	0.8624	0.8686	0.8766
<i>Baseline</i> [31]	-	-	-	0.731	0.747	0.776	0.818	0.836	0.869	0.894	0.912

Πίνακας 4-1 Μετρήσεις AUC για τις διάφορες αρχιτεκτονικές των στοιβαγμένων αυτοκωδικοποιητών, και για διαφορετικό ποσοστό των εξ' αρχής γνωστών προσημασμένων ακμών . Η αρχιτεκτονική αναφέρεται στον κομμάτι του 'κωδικοποιητή'. Η τελευταία σειρά του πίνακα δείχνει του πίνακα δείχνει τις επιδόσεις της μεθόδου-αναφοράς

Ο Πίνακας 4-1 παρουσιάζει τις μετρήσεις του ROC-AUC για τις διαφορετικές διατάξεις των παραμέτρων του προβλήματος (αρχιτεκτονική αυτοκωδικοποιητή και ποσοστό των εξ' αρχής γνωστών προσημασμένων ακμών). Η στήλη ' Αρχιτεκτονική' αναφέρεται μόνο στο κομμάτι του κωδικοποιητή, αφού ο αποκωδικοποιητής είναι η συμμετρική του εικόνα π.χ. για μια αρχιτεκτονική κωδικοποιητή '200-750', υπάρχει ένας αντίστοιχος αποκωδικοποιητής '750-200', ο οποίος θα αναδημιουργήσει την

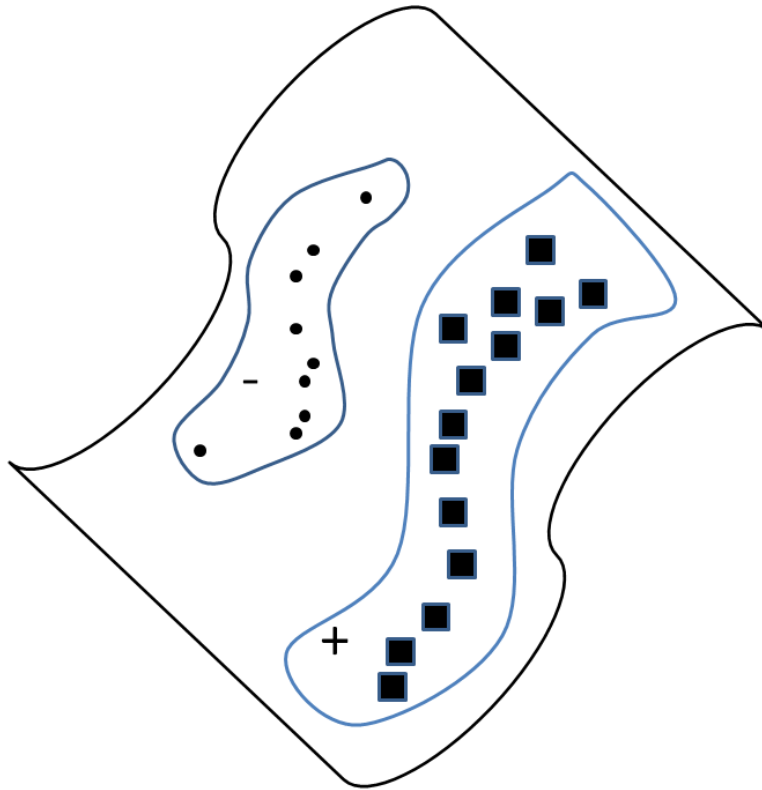
είσοδο. Κάθε αρχιτεκτονική είναι της μορφής $[i - h_1 - h_2 - \dots - o]$, όπου τα i, o αναφέρονται στις διαστάσεις των επιπέδων εισόδου και εξόδου αντίστοιχα, ενώ τα h_i αποτυπώνουν το μέγεθος των κρυμμένων επιπέδων. Εδώ όλες οι αρχιτεκτονικές ξεκινούν με ένα επίπεδο εισόδου με μέγεθος 200. Εξετάσαμε οκτώ διαφορετικές αρχιτεκτονικές, οι οποίες είχαν από 2 έως 4 επίπεδα, με διαφορετικές διαστάσεις επιπέδων. Για τις πρώτες έξι, ορίσαμε ένα *σχετικά μεγάλο* πρώτο κρυμμένο επίπεδο (≥ 500 μονάδες) ενώ για τα δύο τελευταία χρησιμοποιήσαμε ένα αρκετά μικρότερο ('bottleneck') των 50 μονάδων. Εκπαιδεύσαμε το μοντέλο μας με mini-batch gradient descent για 50 epochs και μέγεθος δεσμίδας ίσο με 250. Η τελευταία γραμμή του Πίνακα 4.1 παραθέτει τις μετρήσεις της μεθόδου αναφοράς.

Κάποια συμπεράσματα είναι εύκολο να εξαχθούν: Πρώτον, η απόδοση του αλγορίθμου είναι ικανοποιητική επιβεβαιώνοντας την επιλογή μας να αναλύσουμε βαθμολογίες χρηστών (user ratings) με χρήση αλγορίθμων deep learning, καταλήγοντας σε αποδοτικούς μηχανισμούς εύρεσης κωδικών για τους χρήστες. Πιο συγκεκριμένα, οι έξι πρώτες αρχιτεκτονικές εμφανίζουν 6-10% απόλυτη αύξηση στην ακρίβεια (σε σχέση με την μέθοδο αναφοράς) για μικρότερα ποσοστά γνωστών προσημασμένων ακμών (1% έως 30%) και περίπου 2-4% αύξηση για μεγαλύτερα ποσοστά. Παρόμοια χαρακτηριστικά εμφανίζουν και οι δύο (μικρότερες) τελευταίες αρχιτεκτονικές, οι οποίες είναι πιο αποδοτικές για ποσοστό γνωστών έως και 50%, ενώ για μεγαλύτερες τιμές η απόδοση τους πέφτει. Τα αποτελέσματα υποδεικνύουν ότι η επιπλέον πολυπλοκότητα που εισάγεται λόγω ενός μεγάλου κρυμμένου επιπέδου, γίνεται σημαντική όσο περισσότερα ταξινομημένα δεδομένα γίνονται γνωστά. Δεύτερον, παρατηρούμε ότι μια μοντελοποίηση χαμηλών διαστάσεων για έναν προσημασμένο κοινωνικό γράφο είναι δυνατή, αφού κάθε χρήστης μπορεί να αναπαρασταθεί με μεγάλη ακρίβεια με ένα χαμηλών διαστάσεων δυαδικό διάνυσμα. Σε αυτή την λογική, τα δεδομένα μας είναι συμβατά με την δουλειά στο [45], όπου η

πρόβλεψη προσήμου επιλύεται ως ένα πρόβλημα 'συμπλήρωσης' ενός μικρού βαθμού πίνακα.

4.5.4 Αυτοκωδικοποιητές versus Support Vector Machines (SVMs)

Διεξήγαμε επιπλέον πειράματα, χρησιμοποιώντας τον αλγόριθμο Support Vector Machines (SVMs) ως ταξινομητή τροφοδοτώντας τον με τους κώδικες των ακμών, όπως αυτοί παρήχθησαν από το CF-RBM και για ποσοστά εξ' αρχής γνωστών προσήμων ίδια με του Πίνακα 4.1. Εξετάσαμε γραμμικούς και μη-γραμμικούς (κυρίως Radial Basis Function) πυρήνες και επαναλάβαμε για αρκετές φορές, διαλέγοντας κάθε φορά ένα διαφορετικό σύνολο γνωστών ακμών. Διαπιστώσαμε ότι τα SVMs δεν ήταν εύρωστα όταν το ποσοστό των προσημασμένων ακμών ήταν μικρό: Από επανάληψη σε επανάληψη, η ακρίβεια κυμαινόταν από 25% έως 75%, πράγμα το οποίο αποδεικνύει την υψηλή εξάρτηση των SVMs από την αρχική επιλογή των γνωστών προσημασμένων ακμών. Αυτό είναι το σημείο όπου αποδίδει η δυνατότητα των deep learning αρχιτεκτονικών να εισάγουν ένα στάδιο μη-εποπτευόμενης μάθησης στα δεδομένα. Τα αποτελέσματα φαίνεται να επιβεβαιώνουν την υπόθεση της ύπαρξης *manifolds* [46], βάσει της οποίας η πιθανοτική κατανομή δεδομένων ακόμα υψηλών διαστάσεων τείνει να 'συγκεντρώνεται' γύρω από έναν χώρο λίγων διαστάσεων. Το στάδιο της μη-εποπτευόμενης μάθησης που αναφέρθηκε φαίνεται ότι βρίσκει ένα τέτοιο χώρο για τα δεδομένα μας, και το οποίο χρησιμοποιείται από το επόμενο στάδια της μεθόδου για να καθορίσει 'περιοχές' θετικού και αρνητικού συναισθήματος σε αυτό.



Εικόνα 4.5 Παράδειγμα manifold χαμηλών διαστάσεων, στο οποίο αναπαριστώνται περιοχές θετικού και αρνητικού προσήμου. Κάθε κουκκίδα αντίστοιχα στο κώδικα ακμής που βρέθηκε με την μέθοδο μας

4.6 Σύνοψη κεφαλαίου

Σε αυτό το κεφάλαιο παρουσιάσαμε μια καινοτόμα τεχνική για την πρόβλεψη των δεσμών εμπιστοσύνης ενός κοινωνικού γράφου βασιζόμενοι στις κριτικές των χρηστών για προϊόντα. Σε πρώτη φάση καταφέραμε να αντιστοιχίσουμε έναν δυαδικό κώδικα σε κάθε χρήστη με την χρήση ενός Restricted Boltzmann Machine, λαμβάνοντας υπόψη τις ιδιαιτερότητες των συστημάτων συνεργατικού φιλτραρίσματος όπως η αραιότητα των δεδομένων. Έπειτα, παράξαμε κώδικες για κάθε ακμή του γράφου βάσει των κωδικών των χρηστών που της αντιστοιχούν και στο τέλος τους τροφοδοτήσαμε σε ένα δίκτυο αυτοκωδικοποιητών για την επίλυση του προβλήματος ταξινόμησης. Δείξαμε την υπεροχή των αλγορίθμων deep learning

επί πιο συμβατικών μεθόδων στην ανάλυση δεδομένων όπως τα Support Vector Machines.

5

Εύρεση

κοινοτήτων σε κοινωνικά

δίκτυα

Σε αυτό το κεφάλαιο εξετάζουμε το ζήτημα της εύρεσης κοινοτήτων (community detection) σε κοινωνικά δίκτυα υπό την μορφή της ταξινόμησης ετικετών για τους χρήστες (κάθε ετικέτα είναι ένδειξη συμμετοχής σε κοινότητα). Το ζήτημα ορίζεται στα πλαίσια ενός προβλήματος ημι-εποπτευόμενης μάθησης., όπου γνωρίζοντας τις κοινότητες που ανήκει ένα μικρό υποσύνολο των χρηστών, προσπαθούμε να προβλέψουμε τις ετικέτες των υπολοίπων αξιοποιώντας την δομή του κοινωνικού δικτύου. Η προσέγγιση μας βασίζεται στην έννοια των *σημασιολογικών διανυσματικών χώρων* (semantic vector spaces), μια τεχνική που χρησιμοποιείται ευρέως στην Επεξεργασία Φυσικής Γλώσσας (Natural Language Processing) τα τελευταία χρόνια. Επαληθεύουμε την μέθοδο μας σε δύο μεγάλης κλίμακας, πραγματικά σύνολα δεδομένων και δείχνουμε ότι πετυχαίνει απόδοση συγκρίσιμη ή και καλύτερη σε σχέση με τις βέλτιστες τεχνικές στην τρέχουσα βιβλιογραφία.

5.1 Εισαγωγή

Η θεωρία της κοινωνικής συσχέτισης (social correlation) υποθέτει ότι οι πράξεις ενός ατόμου εντός κοινωνικού δικτύου *δεν* είναι αυτόνομες αλλά αντίθετα εξαρτώνται σε μεγάλο βαθμό από την συμπεριφορά των "οικείων" του. Από την μία πλευρά, η θεωρία της κοινωνικής επιρροής (social influence) υποστηρίζει ότι τα άτομα ενδέχεται να υιοθετούν αντιλήψεις ώστε οι προτιμήσεις τους να συμπίπτουν με αυτές των "κοντινών" τους προσώπων. Από την άλλη, η έννοια της "ομοφιλίας" προϋποθέτει ότι άτομα που μοιράζονται πολλά κοινά χαρακτηριστικά, θα τείνουν να βρίσκονται πιο κοντά στον κοινωνικό γράφο. Και οι δύο θεωρίες καταλήγουν στο γεγονός ότι η ίδια η δομή ενός κοινωνικού δικτύου δίνει τη δυνατότητα να προβλέψουμε τα ενδιαφέροντα ενός χρήστη αξιοποιώντας κατάλληλα το προσωπικό του δίκτυο.

Η δυνατότητα ταυτοποίησης των χρηστών και η κατάταξη τους σε κατηγορίες (*topics*) είναι μεγάλης σημασίας για εμπορικές εφαρμογές όπως η στοχευόμενη διαφήμιση (*targeted advertising*) και τα συστήματα προτάσεων (*recommender systems*) [47]. Όπως είναι φυσικό αυτός ο τομέας παρουσιάζει μεγάλο ερευνητικό ενδιαφέρον, με βασικό ζητούμενο την εύρεση των ενδιαφερόντων του χρήστη βάσει της θέσης του στο κοινωνικό δίκτυο. Μια μεγάλη πρόκληση στο ζήτημα αυτό έχει να κάνει με την μεγάλη ετερογένεια που παρατηρείται στις κοινωνικές σχέσεις. Για παράδειγμα, κάποιος χρήστης Α ενός κοινωνικού δικτύου μπορεί να συνδέεται με κάποιον χρήστη Β λόγω συγγενικής σχέσης, ενώ επιπλέον να συνδέεται με τον κάποιον άλλο χρήστη Γ λόγω του ότι μοιράζονται τον ίδιο εργασιακό χώρο. Έτσι κάθε δεσμός συνήθως είναι διαφορετικός και αντανακλά μια διαφορετική πτυχή της προσωπικότητάς του. Συνήθως όμως, στα περισσότερα κοινωνικά δίκτυα (π.χ. Facebook) οι δεσμοί φιλίας δεν φέρουν κάποια ετικέτα η οποία να υποδεικνύει τον

λόγο της δημιουργίας του δεσμού, κάνοντας έτσι πιο δύσκολη την εύρεση των χαρακτηριστικών που ενώνουν τους δύο χρήστες.

Ένας μεγάλος αριθμός ερευνητικών προσπαθειών προσπαθεί να επιλύσει το πρόβλημα της ετερογένειας αναπαριστώντας τα γραφικά δεδομένα ως διανύσματα, τα οποία συνέχεια τροφοδοτούνται ως μεταβλητές σε παραδοσιακούς ταξινομητές. Μια από τις πρώτες τεχνικές που παρουσιάζεται στην βιβλιογραφία είναι αυτή των Tang et al. [48] οι οποίοι αντιστοιχούν ένα διάνυσμα ανά κόμβο μέσω της φασματικής ανάλυσης του πίνακα γειτνίασης (και παραλλαγών του) [49]. Μια ενδιαφέρουσα εναλλακτική παρουσιάζεται από τους Perozzi et al. [50], οι οποίοι ανάγουν το πρόβλημα της ταξινόμησης ετικετών των χρηστών ενός κοινωνικού δικτύου σε ένα πρόβλημα Επεξεργασίας Φυσικής Γλώσσας. Οι συγγραφείς παρατηρούν την στις κατανομές των βαθμών γειτνίασης (degree distribution) των κόμβων και αυτή των συχνοτήτων εμφάνισης λέξεων σε κείμενα.

Σε αυτό το κεφάλαιο, προτείνουμε ένα τρόπο επίλυσης του προβλήματος της ταξινόμησης ετικετών ορίζοντας κατάλληλες *λανθάνουσες κοινωνικές διαστάσεις* (latent social dimensions). Πιο συγκεκριμένα, χρησιμοποιούμε τον αλγόριθμο GloVe, ο οποίος προτάθηκε από τους Pennington et al. στο [51]. Πρόκειται για ένα εργαλείο μοντελοποίησης που αναπτύχθηκε στο Stanford NLP Group, το οποίο παρουσιάζει έναν πρακτικό τρόπο αντιστοίχισης διανυσμάτων σε λέξεις βάσει των ποσοστών κοινής τους εμφάνισης τους σε ένα κείμενο εκπαίδευσης. Βασιζόμενοι στην αναλογία μεταξύ των συχνοτήτων εμφάνισης των κόμβων σε 'στοχαστικούς περιπάτους' (random walks) και των λέξεων στην φυσική γλώσσα, όπως σημειώνεται από τον Perozzi [50], προτείνουμε μια νέα προσέγγιση στην διανυσματική αναπαράσταση των κόμβων του κοινωνικού γράφου βάσει των κοινών τους "εμφανίσεων". Επαληθεύουμε την ακρίβεια των αναπαραστάσεων πάνω σε δύο

πραγματικά, μεγάλης κλίμακας σύνολα δεδομένων και κάνουμε μια συγκριτική μελέτη με τις τρέχουσες μεθόδους στην βιβλιογραφία, όπου επαληθεύεται η ακρίβεια της προσέγγισης μας. Το υπόλοιπο του κεφαλαίου είναι οργανωμένο ως εξής: Πρώτα, εισάγεται η διατύπωση του προβλήματος και η βασική ορολογία. Έπειτα παρουσιάζεται η προτεινόμενη μέθοδος και η αξιολόγησή της. Το κεφάλαιο τελειώνει με την σύνοψη και τα συμπεράσματα που αποκομίσαμε.

5.2 Διατύπωση του Προβλήματος

Εξετάζουμε το πρόβλημα της κατηγοριοποίησης χρηστών (*node classification*) ως ένα υπο-πρόβλημα της εύρεσης κοινοτήτων σε κοινωνικά δίκτυα. Τα δεδομένα μοντελοποιούνται ένας γράφος $G(V, E, L)$ όπου V , E είναι το σύνολο των κόμβων και των ακμών αντίστοιχα και L το σύνολο των ετικετών που μπορούν να αποδοθούν στους χρήστες. Με άλλα λόγια θεωρούμε ότι κάθε ετικέτα είναι ένδειξη για την συμμετοχή του χρήστη στην αντίστοιχη κοινότητα και κάθε χρήστης μπορεί να ανήκει σε μία ή και περισσότερες κατηγορίες. Αντιμετωπίζουμε το ζήτημα ως ένα πρόβλημα ημι-εποπτευόμενης μάθησης: Πρώτα, επιχειρούμε εκμάθηση της δικτυακής δομής και έπειτα αξιοποιούμε ένα μικρό υποσύνολο $V_L \subset V$ των κόμβων του δικτύου, για τους οποίους θεωρούμε ότι γνωρίζουμε τις κατηγορίες στις οποίες ανήκουν, ώστε να προβλέψουμε την συμμετοχή των υπόλοιπων κόμβων στις κοινότητες. Ο στόχος αυτού του κεφαλαίου είναι να προτείνει έναν πρακτικό τρόπο αναπαράστασης κάθε κόμβου με ένα χαμηλών διαστάσεων διάνυσμα, το οποίο θα λειτουργήσει εν τέλει ως δείκτης για την συμμετοχή του χρήστη σε κάποια κοινότητα.

5.3 Μέθοδοι Αναφοράς

Το πρόβλημα της ταξινόμησης κόμβων σε γράφους έχει εξεταστεί σε μεγάλο βαθμό από την σύγχρονη βιβλιογραφία [52]. Στο [53] οι Macskassy and Provost εκτιμούν τις πιθανότητες ένταξης κάθε κόμβου σε μια κλάση παίρνοντας τον μέσο όρο των άμεσων γειτόνων τους. Μια παρόμοια προσέγγιση ακολουθείται και από του Lu and Getoor στο [54] οι οποίοι προτείνουν ένα πλαίσιο για την μοντελοποίηση των κατανομών των δεσμών ενός κοινωνικού γράφου ώστε να χτίσουν ταξινομητές οι οποίοι θα ενεργούν πάνω στις ιδιότητες των συνδεδεμένων οντοτήτων. Στο [55] οι Zhu et al. παρουσιάζουν μια μέθοδο βασισμένη σε στοχαστικούς περιπάτους (random walks) ορίζοντας το πρόβλημα εκμάθησης βάσει ενός Gaussian random field, ενώ στο [56] οι Szummer and Jaakola εισάγουν εκ νέου μια προσέγγιση σε Markov στοχαστικούς περιπάτους ώστε να τα ταξινομήσουν μερικώς εποπτευμένο κείμενο. Στα [48] [57] οι Tang and Liu ακολουθούν ένα διαφορετικό μονοπάτι: Πρώτα εξάγουν λανθάνουσες κοινωνικές διαστάσεις για κάθε χρήστη βασιζόμενοι στην δομή του κοινωνικού δικτύου και έπειτα χρησιμοποιούν τα διανύσματα αυτά για να εκπαιδεύσουν κατάλληλους ταξινομητές, οι οποίοι θα μπορούν να προβλέπουν την συμμετοχή του κάθε χρήστη στις κοινότητες. Η παραγωγή των διανυσμάτων γίνεται μέσω της φασματικής ανάλυσης του πίνακα Modularity [48] καθώς και του κανονικοποιημένου πίνακα Laplacian [57]. Στο [49] οι ίδιοι συγγραφείς αναπτύσσουν μια δεσμό-κεντρική μέθοδο ομαδοποίησης χρησιμοποιώντας μια παραλλαγή του k-means αλγορίθμου στον πίνακα γειννίας. Τέλος, στο [50] οι Perozzi et al. παρουσιάζουν ένα πλαίσιο ανάθεσης λανθανόντων αναπαραστάσεων στους κόμβους του δικτύου αξιοποιώντας πρόσφατες εξελίξεις στον χώρο της Επεξεργασίας Φυσικής Γλώσσας. Πιο συγκεκριμένα, εφαρμόζουν τον αλγόριθμο SkipGram [58] σε στοχαστικούς περιπάτους πάνω στον κοινωνικό γράφο

λογίζοντας τους ως ισοδύναμους των προτάσεων.

5.4 Προτεινόμενη Προσέγγιση

Σε αυτή την ενότητα εισάγουμε έννοιες και αλγορίθμους τα οποία θα χρησιμοποιηθούν ως βασικές μονάδες υλοποίησης της μεθόδου και προέρχονται από τον χώρο της Επεξεργασίας Φυσικής Γλώσσας. Μια από τις τεχνικές που θα μας απασχολήσουν είναι τα σημασιολογικά διανυσματικά μοντέλα, τα οποία τα τελευταία χρόνια αποτελούν μια δημοφιλή επιλογή για την στατιστική μοντελοποίηση προβλημάτων Επεξεργασίας Φυσικής Γλώσσας.

5.4.1 Χώροι σημασιολογικών διανυσμάτων και GloVe

Τα σημασιολογικά διανυσματικά μοντέλα λειτουργούν αντιστοιχώντας σε κάθε λέξη του λεξιλογίου ένα διάνυσμα πραγματικών αριθμών το οποίο θα χρησιμοποιηθεί ως είσοδος για περαιτέρω ανάλυση σε παραδοσιακούς αλγορίθμους μηχανικής μάθησης. Η ακρίβεια των μοντέλων αυτών εξαρτάται από την δυνατότητα τους να βρίσκουν λανθάνουσες διαστάσεις με νοηματικό και σημασιολογικό περιεχόμενο. Μια ενδιαφέρουσα εφαρμογή παρουσιάζεται από τους Mikolov et al. [58] όπου οι συγγραφείς ανακαλύπτουν αναλογίες ανάμεσα στις λέξεις εφαρμόζοντας απλές αριθμητικές πράξεις. Για παράδειγμα η φράση "η Αθήνα είναι για την Ελλάδα ότι και το Βερολίνο για την Γερμανία" αντανακλάται στον σημασιολογικό χώρο μέσω της ταυτότητας:

$$v(\text{"Athens"}) - v(\text{"Greece"}) = v(\text{"Berlin"}) - v(\text{"Germany"})$$

όπου $v(\tau)$ είναι το διάνυσμα που αντιστοιχεί στην λέξη τ .

Για την υλοποίηση της μεθόδου μας, χρησιμοποιούμε έναν παρόμοιο αλγόριθμο, τον αλγόριθμο GloVe, ο οποίος προτάθηκε πρόσφατα από τους Pennington et al. [51]. Το υπόλοιπο της ενότητας εστιάζει στις βασικές ιδιότητες του αλγορίθμου και εξηγείται πως μπορεί να χρησιμοποιηθεί στα πλαίσια του προβλήματος κατηγοριοποίησης κόμβων.

Μέσω του GloVe μπορούμε να πάρουμε διανυσματικές αναπαραστάσεις λέξεων, τα οποία αντικατοπτρίζουν τον αριθμό των κοινών τους εμφανίσεων σε ένα κείμενο εκπαίδευσης. Χρησιμοποιώντας την ορολογία των συγγραφέων, ορίζουμε ως X τον πίνακα των κοινών εμφανίσεων των λέξεων, του οποίου οι εγγραφές X_{ij} αντικατοπτρίζουν τον αριθμό που η λέξη i εμφανίζεται στο γενικό πλαίσιο (*context*) της λέξης j : δεδομένης μιας πρότασης από ένα κείμενο εκπαίδευσης και ένα προκαθορισμένο παράθυρο w , οι λέξεις i και j βρίσκονται λιγότερο από w λέξεις μακριά. Τα ζεύγη των λέξεων λαμβάνουν βάρος αντίστοιχο της απόστασης τους και έτσι λέξεις που απέχουν κατά d αποκτούν βάρος $\propto 1/d$. Επίσης, ορίζουμε ως $X_i = \sum_k X_{ik}$ τον αριθμό των φορών που μια οποιαδήποτε λέξη εμφανίζεται στο πλαίσιο της λέξης i και τέλος ως $P_{ji} = P(j|i) = X_{ij}/X_i$ την δεσμευμένη πιθανότητα εμφάνισης της λέξης j στο πλαίσιο της λέξης i . Στο GloVe κάθε λέξη έχει έναν *κύριο* (*main*) και έναν *δευτερεύοντα* (*context*) ρόλο και γι' αυτό το λόγο ο αλγόριθμος δίνει δύο αντίστοιχες αναπαραστάσεις (w_i και \tilde{w}_i) για κάθε κόμβο. Η εκπαίδευση του μοντέλου γίνεται βάσει της ελαχιστοποίησης μιας συνάρτησης κόστους της ακόλουθης μορφής

$$J = \sum_{i,j=1}^V f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log(X_{ij}))^2$$

,όπου V είναι το μέγεθος του λεξιλογίου, b_i, \tilde{b}_j οι σταθερές προδιάθεσης (biases) για τους κύριους και δευτερεύοντες ρόλους αντίστοιχα. Η $f(x)$ είναι μια αύξουσα

συνάρτηση η οποία δρα ως κανονικοποιητής (*regularizer*) απαγορεύοντας σε ζεύγη λέξεων που εμφανίζονται (υπερβολικά) συχνά να λαμβάνουν μεγάλο βάρος. Στην αρχική δημοσίευση του GloVe η συνάρτηση $f(x)$ ορίζεται ως

$$f(x) = \begin{cases} (x/x_{max})^a, & \text{εάν } x < x_{max} \\ 1, & \text{σε άλλη περίπτωση} \end{cases},$$

με $x_{max} = 100$ και $a = 0.75$. Το μοντέλο εκπαιδεύεται με την χρήση του αλγορίθμου Adagrad [59], όπου τα μη-μηδενικά στοιχεία του X δειγματοληπτούνται στοχαστικά με ρυθμό εκμάθησης 0.05.

5.4.2 Από το Glove στο GloveGraph

5.4.2.1 Υπολογισμός ομοιότητας μεταξύ των χρηστών

Ο υπολογισμός της ομοιότητας μεταξύ των χρηστών ενός κοινωνικού δικτύου είναι μια κεντρική έννοια στην μοντέρνα θεωρία της Ανάλυσης Κοινωνικών Δικτύων. Κοινωνικές αρχές όπως η κοινωνική επιρροή και η ομοφιλία καθορίζουν σε μεγάλο βαθμό την συμπεριφορά των ατόμων τα οποία τείνουν να ευθυγραμμίζουν τις προτιμήσεις τους με αυτές των "κοντινών" τους στον κοινωνικό γράφο. Η ποικιλία των μετρικών ομοιότητας στην βιβλιογραφία [52] [17] πιστοποιεί ότι η επιλογή της κατάλληλης μετρικής εξαρτάται από το πεδίο εφαρμογής. Μια πρώτη οικογένεια υπολογίζουν την ομοιότητα δύο χρηστών βάσει του βαθμού επικάλυψης των άμεσων γειτόνων του. Σε αυτή την κατηγορία ανήκουν η μετρική των Κοινών Γειτόνων [60], η οποία μετράει τον απόλυτο αριθμό των κοινών γειτόνων μεταξύ δύο χρηστών, η μετρική Jaccard [61], η οποία μετράει την πιθανότητα δύο χρήστες u και v να μοιράζονται έναν γείτονα για κάποιον τυχαία επιλεγμένο χρήστη που έχει είτε ο u είτε ο v , και τέλος η μετρική Adamic/Adar [18] η οποία δίνει μεγαλύτερο βάρος σε κόμβους μικρότερου βαθμού. Ένα πρόβλημα με αυτές τις μετρικές είναι ότι έχουν

περιορισμένη οπτική του κοινωνικού δικτύου αφού δύο χρήστες δύνανται να έχουν μη μηδενική τιμή ομοιότητας αν βρίσκονται το πολύ δύο βήματα μακριά ο ένας από τον άλλο. Στην δική μου προσέγγιση αποφάσισα να χρησιμοποιήσω μια μετρική, η οποία επιτρέπει στην ομοιότητα να "ταξιδεύει" μέσα στον κοινωνικό γράφο, δίνοντας μεγαλύτερη ποινή σε μακρύτερα μονοπάτια. Πιο συγκεκριμένα, χρησιμοποίησα την μετρική Katz [62], η οποία ορίζεται μεταξύ δύο χρηστών u και v ως ακολούθως:

$$Katz(u, v) = \sum_{l=1}^{\infty} (\beta_{katz})^l \cdot paths^{<l>}(u \rightarrow v)$$

Όπου $\beta_{katz} < 1$ καλείται ο συντελεστής Katz και $paths^{<l>}(u \rightarrow v)$ είναι ο συνολικός αριθμός των μονοπατιών μήκους l τα οποία ξεκινάνε από τον u και καταλήγουν στον v . Μια μικρή τιμή του β_{katz} καθιστά την μετρική Katz ισοδύναμη ε τις προηγούμενες μετρικές.

Ο υπολογισμός της μετρικής Katz για μεγάλα κοινωνικά δίκτυα μπορεί να αποτελέσει πρόκληση. Στην πράξη, υπάρχουν δύο μέθοδοι υπολογισμού: Η πρώτη είναι ακριβής και απαιτεί την εκθετοποίηση του πίνακα γειτνίασης: Ο πίνακας γειτνίασης $A(G)$ ενός γράφου G είναι ο τετράγωνος πίνακας με αριθμό γραμμών (και στηλών) ίσο με τον αριθμό των κόμβων. Η εγγραφή A_{ij} του πίνακα είναι ίση με 1 όταν ο κόμβος i συνδέεται με τον κόμβο j και ίση με το 0 σε κάθε άλλη περίπτωση. Στα δίκτυα που μελετάμε οι πίνακες γειτνίασης έχουν δύο βασικά χαρακτηριστικά: Πρώτον, είναι συμμετρικά και επομένως $A_{ij} = A_{ji}$ και δεύτερον είναι αραιά, δηλαδή η πλειονότητα των εγγραφών A_{ij} είναι ίσες με το μηδέν. Η αραιότητα των δεδομένων είναι σημαντική γιατί υπάρχουν αποδοτικές τεχνικές αναπαράστασης και αποθήκευσης αραιών πινάκων, γεγονός που καθιστά δυνατή την ανάλυση ακόμα και εξαιρετικά μεγάλων κοινωνικών δικτύων. Η l -οστή δύναμη του πίνακα γειτνίασης μας δίνει απ' ευθείας τον συνολικό αριθμό των μονοπατιών μήκους l που συνδέουν

μεταξύ τους δύο οποιουσδήποτε κόμβους. Έτσι, ο υπολογισμός της μετρικής Katz σε αυτή την περίπτωση συνίσταται στον υπολογισμό του βοηθητικού πίνακα K όπου $K = \sum_{l=1}^n (\beta_{katz})^l \cdot A(G)^l$ για κάποια κατάλληλη απόσταση n . Το βασικό πρόβλημα με αυτή την προσέγγιση είναι τα πραγματικά κοινωνικά δίκτυα είναι small-world δίκτυα, δηλαδή δύο οποιοδήποτε κόμβοι συνδέονται με ένα μονοπάτι μικρού μήκους. Έτσι, μετά από κάποιες εκθετοποιήσεις του πίνακα γειτνίασης θα προκύπτουν πυκνοί πίνακες, οι οποίοι είναι δύσκολα διαχειρίσιμοι αφού καθιστούν ανεφάρμοστες τις τεχνικές για την αποθήκευση και ανάλυση αραιών πινάκων που συνήθως εφαρμόζονται στις περιπτώσεις των κοινωνικών δικτύων.

Στις περιπτώσεις που η εκθετοποίηση του πίνακα γειτνίασης δεν είναι δυνατή, η χρήση των *στοχαστικών περιπάτων* (random walks) αποτελεί μια ικανοποιητική προσέγγιση. Ένας στοχαστικός περίπατος είναι μια στοχαστική διαδικασία η οποία παράγει μια ακολουθία κόμβων, όπου ο κάθε τρέχων κόμβος επιλέγει τυχαία ένας από τους άμεσους γείτονες του ως τον επόμενο της αλυσίδας, ώσπου το μήκος της ακολουθίας φτάσει μια επιθυμητή τιμή. Στην περίπτωση μας, η μετρική Katz μεταξύ δύο κόμβων u και v μπορεί να προσεγγιστεί ως εξής: Από τον κόμβο u εκκινούμε έναν προκαθορισμένο αριθμό στοχαστικών περιπάτων και αθροίζουμε τις συνεισφορές του κόμβου v . Με άλλα λόγια, κάθε φορά που ο κόμβος v εμφανίζεται σε κάποιον στοχαστικό περίπατο σε απόσταση m από τον u , τότε η τιμή της μεταξύ τους ομοιότητας αυξάνεται κατά $(\beta_{katz})^m$. Γενικά, η εκθετοποίηση του πίνακα γειτνίασης δίνει πιο ακριβή αποτελέσματα με βασικό μειονέκτημα το υψηλό υπολογιστικό κόστος. Από την άλλη, η λύση που βασίζεται στους στοχαστικούς περιπάτους με 800 περιπάτους ανά κόμβο δίνει ισοδύναμα αποτελέσματα.

5.4.2.2 Προτεινόμενος αλγόριθμος

Επιλέξαμε τον αλγόριθμο GloVe [51] λόγω της απλότητας του και της ικανότητας να προσαρμόζεται με τρόπο διαισθητικό εύκολο στο πρόβλημα της κατηγοριοποίησης κόμβων που εξετάζουμε εδώ.

Η Εικόνα 5.1 δίνει τα βασικά στάδια του αλγορίθμου. Ως είσοδο δεχόμαστε τον κοινωνικό γράφο G , το μέγεθος των παραγόμενων διανυσμάτων d , τον αριθμό των στοχαστικών περιπάτων που θα παράξουμε ανά κόμβο w_{pv} (walks per vertex), το μέγεθος κάθε στοχαστικού περιπάτου wL , καθώς και τον συντελεστή β_{katz} για τον υπολογισμό της συνάρτησης ομοιότητας Katz. Ο αλγόριθμος λειτουργεί ως εξής : Πρώτα, αρχικοποιούνται τυχαία οι παράμετροι του GloVe δηλαδή οι "κύριες" αναπαραστάσεις των κόμβων W , οι "δευτερεύουσες" αναπαραστάσεις των κόμβων \tilde{W} , και οι αντίστοιχοι συντελεστές προδιάθεσης b και \tilde{b} . Στη συνέχεια, αλλάζουμε την σειρά "επίσκεψης" του αλγορίθμου στους κόμβους και για κάθε τρέχοντα κόμβο παράγουμε τον προκαθορισμένο στοχαστικών περιπάτων, μέσω των οποίων θα υπολογίσουμε τις Katz ομοιότητες μεταξύ του τρέχοντα κόμβου και των γειτονικών του. Έπειτα, υπολογίζουμε τις απαραίτητες ανανεώσεις που προκύπτουν από την βελτιστοποίηση της συνάρτησης κόστους μέσω του αλγορίθμου Adagrad [59] (βλέπε Εικόνα 5.2). Αφού διατρέξουμε όλους τους κόμβους, η μέθοδος μας επαναλαμβάνεται για έναν προκαθορισμένο αριθμό επαναλήψεων.

Algorithm GloveGraph

Input :

$G(V,E)$: the social graph
 d : vector size
 wpv : walks per vertex
 wL : walk length
 β_{katz} : damping factor of the Katz measure
 $numIters$: number of iterations

Output : vector representation matrix $W_{|V| \times d}^{out}$

1. Initialize *GloVe* parameters $(W, \tilde{W}, b, \tilde{b}) \sim Uniform(-1/d, 1/d)$
 2. for *iteration* = 1 to *numIters* do :
 3. $all_vertex_ids = G.V.randomizeIds()$;
 4. foreach *vertex_id* in *all_vertex_ids* do :
 5. Generate wpv random walks of length wL starting from *vertex_id*;
 6. compute *List<Katz> scores*;
 7. $updateGloVe (main_id = vertex_id, context_id=katz.id, value = katz.value)$
 8. end foreach;
 9. end for;
 10. output $(W_{|V| \times d}^{out} = W + \tilde{W})$
-

5.1 Αλγόριθμος GloveGraph

Algorithm Glove

Input :

$x_{max}=30.0$
 $a = 0.80$;
 X : *similarity pairs*

1. compute $J = \sum_{i,j=1}^V f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log(X_{ij}))^2$
 2. find new search directions through *Adagrad*
 3. update $W, \tilde{W}, b, \tilde{b}$
-

5.2 Υπολογισμός ανανεώσεων μέσω του αλγορίθμου Adagrad

5.5 Πειραματική Επαλήθευση

Σύνολο Δεδομένων	BlogCatalog	Flickr
$ V $	10,312	80,513
$ E $	333,983	5,899,882
Τύπος δεσμών	39 Κατηγορίες (topics)	195 Ομάδες Ενδιαφερόντων (interest groups)

5-1 Στατιστικά στοιχεία των δεδομένων

5.5.1 Περιγραφή των δεδομένων

Αξιολογήσαμε την μέθοδο μας σε δύο μεγάλης κλίμακας σύνολα δεδομένων από πραγματικά κοινωνικά δίκτυα. Το πρώτο σύνολο προέρχεται από τον ιστοτόπο BlogCatalog [48], ο οποίος αποτελεί έναν online καταχωρητή blogs, στο οποίο ένας χρήστης μπορεί να εγγράψει το blog του χρησιμοποιώντας παράλληλα και μετα-δεδομένα ώστε να χαρακτηρίσει το περιεχόμενο του. Τα μετα-δεδομένα περιλαμβάνουν τις κατηγορίες των θεμάτων τις οποίες καλύπτει το blog καθώς και τις συνδέσεις του blog με άλλα. Το σύνολο δεδομένων αποτελείται από περίπου 10 χιλιάδες χρήστες, 330 χιλιάδες ακμές μεταξύ των χρηστών και 39 κατηγορίες blog. Το δεύτερο σύνολο δεδομένων προέρχεται από τον δημοφιλή ιστοτόπο αποθήκευσης εικόνων και βίντεο Flickr [48]. Οι χρήστες του Flickr μπορούν να συνδέονται μεταξύ τους με δεσμούς φιλίας και να εγγράφονται σε ομάδες ενδιαφερόντων (interest groups). Τα δεδομένα αποτελούνται από 80 χιλιάδες χρήστες, το σύνολο των επαφών τους το οποίο αριθμεί περίπου 6 εκατομμύρια ακμές, καθώς και 195 ομάδες ενδιαφερόντων. Κάποια στατιστικά μπορούν να βρεθούν στον Πίνακα 5.1.

5.5.2 Μέθοδοι Αναφοράς

Για να αξιολογήσουμε την απόδοση της προσέγγισης μας την συγκρίνουμε με κάποιες βασικές μεθόδους αναφοράς. Αυτές είναι:

- [SpectralClustering] : Αυτή η μέθοδος προτάθηκε από τους Tang et al. [57], όπου η αναπαράσταση για κάθε κόμβο προκύπτει μέσω της φασματικής ανάλυσης του κανονικοποιημένου Laplacian πίνακα.
- [DeepWalk] : Αυτή η μέθοδος προτάθηκε από τους Perozzi et al. [50]. Οι συγγραφείς χρησιμοποιούν τον αλγόριθμο SkipGram [58] [63] [64] σε 'στοχαστικούς περιπάτους' ώστε να μάθουν λανθάνουσες αναπαραστάσεις για τους κόμβους.
- [Modularity] : Αυτή η μέθοδος προτείνεται από τους Tang et al. στο [48]. Οι αναπαραστάσεις για τους κόμβους προκύπτουν από την φασματική ανάλυση του Modularity πίνακα.
- [EdgeCluster] : Η μέθοδος αυτή [49] εξετάζει την εφαρμογή του αλγορίθμου k-means στον πίνακα γειτνίασης.
- [wwRN] : Η μέθοδος παρουσιάζεται εδώ [53]. Πρόκειται για έναν *σχεσιακό* ταξινομητή, ο οποίος εκτιμά την πιθανότητα συμμετοχής ενός χρήστη σε μια κοινότητα παίρνοντας τον μέσο όρο των γειτόνων του.
- [Majority] : Η μέθοδος αυτή προτείνει πάντα την ετικέτα της πλειοψηφούσας κοινότητας . Δεν μαθαίνει τίποτα για την δομή του δικτύου.

5.5.3 Μετρικές αξιολόγησης

Για να διευκολύνουμε την σύγκριση με τις υπάρχουσες μεθόδους στην βιβλιογραφία, ακολουθούμε την ίδια πειραματική διάταξη. Πρώτα, διαλέγουμε τυχαία ένα υποσύνολο των κόμβων του γράφου και τους χρησιμοποιούμε για εκμάθηση. Θεωρούμε ότι για κάθε χρήστη γνωρίζουμε όλες τις κοινότητες στις οποίες συμμετέχει. Έπειτα εκπαιδεύουμε έναν ταξινομητή για κάθε κοινότητα και ,τέλος, αξιολογούμε την ακρίβεια τους εφαρμόζοντας τους στους υπόλοιπους χρήστες και υπολογίζοντας τις μέσες αποδόσεις για τις μετρικές Micro-F1 και Macro-F1. Επαναλαμβάνουμε για δέκα φορές ώστε να περιορίσουμε τον ρόλο της τυχαιότητας στην διαδικασία.

Στην περίπτωση του BlogCatalog, το ποσοστό των "γνωστών" κόμβων ξεκινά από 10% και αυξάνεται σταδιακά ανά 10% μέχρι να φτάσει το 90%, ενώ για το Flickr εξετάζουμε το διάστημα [1%-10%] με βήμα 1%. Όλα τα μοντέλα είναι ορισμένα στις βέλτιστες τιμές τους π.χ. για τις μεθόδους 'SpectralClustering','Modularity' και 'EdgeCluster' χρησιμοποιούμε διανύσματα μήκους 500, ενώ για στην περίπτωση του 'DeepWalk' επιλέγουμε διανύσματα μήκους 128. Για τον δικό μας αλγόριθμο επιλέξαμε διανύσματα μήκους 140 στην περίπτωση του BlogCatalog και μήκους 40 στην περίπτωση του Flickr.

5.5.4 Αξιολόγηση του μοντέλου

Οι πίνακες 5.2 έως 5.5 δίνουν τα αποτελέσματα για τα δύο σύνολα δεδομένων και για τις δύο μετρικές αξιολόγησης που επιλέξαμε:

BlogCatalog				Micro-F1 (%)			
Μέθοδος Ποσοστό γνωστών	Προτεινόμενη (GloveGraph)	DeepWalk	SpectralClustering	EdgeCluster	Modularity	wwRN	Majority
10%	25.17	36.00	31.06	27.94	27.35	19.51	16.51
20%	33.73	38.20	34.95	30.76	30.74	24.34	16.66
30%	37.86	39.60	37.27	31.85	31.77	25.62	16.61
40%	40.40	40.30	38.93	32.99	32.97	28.82	16.70
50%	41.23	41.00	39.97	34.12	34.09	30.37	16.91
60%	42.11	41.30	40.99	35.00	36.13	31.81	16.99
70%	42.87	41.50	41.66	34.63	36.08	32.19	16.92
80%	43.47	41.50	42.42	35.99	37.23	33.33	16.49
90%	43.96	42.00	42.62	36.29	38.18	34.28	17.26

5-2 Μετρήσεις Micro-F1 στο BlogCatalog

BlogCatalog				Macro-F1 (%)			
Μέθοδος Ποσοστό γνωστών	Προτεινόμενη (GloveGraph)	DeepWalk	SpectralClustering	EdgeCluster	Modularity	wwRN	Majority
10%	17.00	21.30	19.14	16.16	17.36	6.25	2.52
20%	22.04	23.80	23.57	19.16	20.00	10.13	2.55
30%	24.60	25.30	25.97	20.48	20.80	11.64	2.52
40%	26.07	26.30	27.46	22.00	21.85	14.24	2.58

50%	27.08	27.30	28.31	23.00	22.65	15.86	2.58
60%	27.88	27.60	29.46	23.64	23.41	17.18	2.63
70%	28.98	27.90	30.13	23.82	23.89	17.98	2.61
80%	29.47	28.20	31.38	24.61	24.20	18.86	2.48
90%	29.71	28.90	31.78	24.92	24.97	19.57	2.62

5-3 Μετρήσεις Macro-F1 στο BlogCatalog

Flickr				Micro-F1 (%)			
Μέθοδος	Προτεινόμενη (GloveGraph)	DeepWalk	SpectralClustering	EdgeCluster	Modularity	wwRN	Majority
Ποσοστό γνωστών							
1%	21.14	32.4	27.43	25.75	22.75	17.7	16.34
2%	25.74	34.6	30.11	28.53	25.29	14.43	16.31
3%	29.48	35.9	31.63	29.14	27.3	15.72	16.34
4%	32.19	36.7	32.69	30.31	27.6	20.97	16.46
5%	33.94	37.2	33.31	30.85	28.05	19.83	16.65
6%	35.13	37.7	33.95	31.53	29.33	19.42	16.44
7%	36.10	38.1	34.46	31.75	29.43	19.22	16.38
8%	36.87	38.3	34.81	31.76	28.89	21.25	16.62
9%	37.31	38.5	35.14	32.19	29.17	22.51	16.67
10%	37.63	38.7	35.41	32.84	29.2	22.73	16.71

5-4 Μετρήσεις Micro-F1 στο Flickr

Flickr				Macro-F1 (%)			
Μέθοδος	Προτεινόμενη	DeepWalk	SpectralClustering	EdgeCluster	Modularity	wwRN	Majority
1%	10.39	14.0	13.84	10.52	10.21	1.53	0.45
2%	13.40	17.3	17.49	14.10	13.37	2.46	0.44
3%	15.67	19.6	19.44	15.91	15.24	2.91	0.45
4%	17.18	21.1	20.75	16.72	15.11	3.47	0.46
5%	18.38	22.1	21.60	18.01	16.14	4.95	0.47
6%	19.19	22.9	22.36	18.54	16.64	5.56	0.44
7%	19.78	23.6	23.01	19.54	17.02	5.82	0.45
8%	20.85	24.1	23.36	20.18	17.1	6.59	0.47
9%	21.48	24.6	23.82	20.78	17.14	8.00	0.47
10%	21.73	25.0	24.05	20.85	17.12	7.26	0.47

5-5 Μετρήσεις Macro-F1 στο Flickr

Στην περίπτωση του BlogCatalog η προτεινόμενη μέθοδος δείχνει καλή ακρίβεια ξεπερνώντας σε κάποιες περιπτώσεις τις τρέχουσες τεχνικές στην βιβλιογραφία. Ένα σημαντικό πλεονέκτημα της είναι ότι μπορεί να ενσωματώσει ευκολότερα κάθε καινούρια πληροφορία την στιγμή που οι άλλες μέθοδοι δείχνουν σημαντικά μικρότερη αύξηση στην προβλεπτική τους ικανότητα στην παρουσία νέων δεδομένων. Αυτό είναι σημαντικό για δύο λόγους: Πρώτον, το μοντέλο δείχνει ότι διαθέτει αρκετή πολυπλοκότητα ώστε να εκμεταλλευτεί κάθε διαθέσιμο δεδομένο ώστε να βελτιστοποιήσει την συνάρτηση κόστους. Δεύτερον, με εξαίρεση τον

αλγόριθμο DeepWalk [50] , οι άλλες μέθοδοι οφείλουν να επανεκκινούν όταν παρουσιάζονται καινούρια δεδομένα. Με άλλα λόγια, δεν υπάρχει δηλαδή ένας ασφαλής τρόπος online εκπαίδευσης των αλγορίθμων αυτών, και άρα η εκτέλεση τους επιφέρει μεγαλύτερο υπολογιστικό κόστος. Παρόμοια συμπεράσματα προκύπτουν και από την εξέταση στο σύνολο δεδομένων από το Flickr, όπου η προτεινόμενη μέθοδος εναλλάσσεται στην δεύτερη και τρίτη θέση με κυριότερο "ανταγωνιστή" την μέθοδο 'SpectralClustering' [57].

Άξιο σχολιασμού είναι και τα προτεινόμενα μεγέθη των διανυσματικών αναπαραστάσεων. Οι μέθοδοι που βασίζονται στην φασματική ανάλυση του πίνακα γειτνίασης [57] [48] [49] απαιτούν διανύσματα μήκους πεντακοσίων(500) θέσεων ενώ ο αλγόριθμος DeepWalk έχει σημαντικά μικρότερες απαιτήσεις με μέγεθος διανύσματος ίσο με 128. Στην δική μου προσέγγιση, οι βέλτιστες μετρήσεις παρατηρήθηκαν για μήκος 140 στην περίπτωση του BlogCatalog και μόλις 40 για την περίπτωση του Flickr. Ειδικά για την περίπτωση του Flickr αυτό το γεγονός αποτελεί ένα μεγάλο υπολογιστικό κέρδος δεδομένης της έκτασης του γράφου. Επίσης, είναι κάπως παράδοξο πως το μεγαλύτερο δίκτυο προερχόμενο από το Flickr απαιτεί μικρότερου μήκους διανυσματικές αναπαραστάσεις σε σχέση με το κατά πολύ μικρότερο δίκτυο του BlogCatalog.

Τέλος, αξίζει να σημειωθεί ότι η μέθοδος που παρουσιάζεται εδώ θα μπορούσε να ειδωθεί ως ένα γενικότερο πλαίσιο ανάλυσης όπου η συνάρτηση ανάλυσης θα μπορεί να μεταβάλλεται ανάλογα με τον τομέα εφαρμογής. Για παράδειγμα, στο κεφάλαιο αυτό βασιστήκαμε στην θεωρία της κοινωνικής επιρροής, βάσει της οποίας άτομα που βρίσκονται **κοντά** στον κοινωνικό γράφο αλληλοεπηρεάζονται σε μεγάλο βαθμό στις προτιμήσεις τους, μια δύναμη που όμως εξασθενεί όσο "απομακρύνονται" μεταξύ τους. Έτσι, η επιλογή της μετρικής Katz [62] [65] ως συνάρτησης ομοιότητας

αντανακλά με διαισθητικά προφανή τρόπο την κατανομή της επιρροής σε έναν κοινωνικό γράφο. Στον τομέα της Επεξεργασίας Φυσικής Γλώσσας, απ' όπου προήλθε και η αρχική δημοσίευση για το GloVe [51], χρησιμοποιήθηκε η αρμονική συνάρτηση $f(r) = 1/r$, αποδίδοντας βάρη στα ζεύγη των λέξεων αντιστρόφως ανάλογα της απόστασης τους στο κείμενο εκπαίδευσης. Με άλλα λόγια, η μέθοδος που παρουσιάσαμε μπορεί να λειτουργήσει ως ένα γενικότερο πλαίσιο ανάλυσης, το οποίο προσαρμόζεται σε κάθε πεδίο εφαρμογής με την επιλογή της κατάλληλης συνάρτησης ομοιότητας κάθε φορά.

5.6 Σύνοψη του κεφαλαίου

Σε αυτό το κεφάλαιο παρουσιάσαμε έναν αλγόριθμο εύρεσης κοινοτήτων σε κοινωνικά δίκτυα βασισμένο σε πρόσφατες εξελίξεις στον επιστημονικό τομέα της Επεξεργασίας Φυσικής Γλώσσας. Συγκεκριμένα, προτείναμε έναν αποδοτικό τρόπο αντιστοίχισης διανυσμάτων στους κόμβους του κοινωνικού δικτύου, τα οποία ήταν ενδεικτικά της δομής του κοινωνικού γράφου με αποτέλεσμα να μπορούμε να προβλέψουμε την κατανομή των ετικετών ανά χρήστη με μεγάλη ακρίβεια.

6

Σύνοψη

Στο συγκεκριμένο κεφάλαιο περιλαμβάνεται η σύνοψη της διατριβής και τα συμπεράσματα που εξήχθησαν κατά την εκπόνησή της, αναλύεται η συνεισφορά και η καινοτομία που επιδεικνύει στον αντίστοιχο ερευνητικό χώρο, και συζητούνται θέματα μελλοντικής εργασίας και επέκτασης των ερευνητικών αποτελεσμάτων.

6.1 Συνεισφορά – Καινοτομία

Στο πλαίσιο της παρούσας διδακτορικής διατριβής παρουσιάστηκαν τεχνικές ανάλυσης κοινωνικών δικτύων δίνοντας μεγαλύτερη έμφαση σε γράφους εμπιστοσύνης. Ένας βασικός στόχος που επετεύχθη ήταν να δω το πρόβλημα από πολλές διαφορετικές σκοπιές, λαμβάνοντας υπόψη βέλτιστες τεχνικές από φαινομενικά μη σχετικούς χώρους με την Ανάλυση Κοινωνικών Δικτύων, όπως η Βιοστατιστική ή η Επεξεργασία Φυσικής Γλώσσας. Προσπάθησα να αναλύσω το ζήτημα της μετάδοσης συναισθήματος ως ένα πεδίο εφαρμογής μιας *διεπιστημονικής έρευνας (interdisciplinary research)*, και να προτείνω τρόπους ώστε να είναι δυνατή η μετάβαση μεταξύ των διάφορων επιστημονικών περιοχών.

Αρχικά, στο 1ο κεφάλαιο παρουσιάστηκαν οι επικρατούσες θεωρίες στον χώρο της ανάλυσης προσημασμένων γράφων. Από την μία, η θεωρία της δομικής ισορροπίας, η οποία έλκει την καταγωγή της από τον Heider [1] και χρονολογείται στην δεκαετία του 1940. Έχει εφαρμοστεί με επιτυχία σε ένα μεγάλο εύρος κοινωνικών δικτύων και αποδεικνύεται εξαιρετικά ακριβής ακόμα και στις περιπτώσεις των σύγχρονων online κοινωνικών δικτύων. Το 1ο κεφάλαιο κλείνει με την παρουσίαση της θεωρίας της κοινωνικής θέσης (social status) [22] [5]. Στο 2ο κεφάλαιο εξετάσαμε τους τρόπους διάδοσης του θετικού και αρνητικού συναισθήματος καταλήγοντας στο γεγονός ότι οι αρνητικές απόψεις φέρουν περισσότερη πληροφορία από τις θετικές [66]. Το αποτέλεσμα αυτής της μελέτης υποστηρίζει την ανάγκη εισαγωγής μηχανισμών στα σύγχρονα κοινωνικά δίκτυα ώστε να μπορούν να ενσωματώνονται καλύτερα και οι αρνητικές απόψεις των χρηστών. Κάποια κοινωνικά δίκτυα⁹ κινούνται ήδη προς αυτήν την κατεύθυνση.

Το 3ο κεφάλαιο παρουσιάσαμε μια καινοτόμα τεχνική ανάλυσης προσημασμένων δικτύων μέσω της εύρεσης των συχνών υπογράφων [67], μια ευρέως διαδεδομένη τεχνική στον χώρο της Βιοστατιστικής. Οι συχνοί υπογράφοι λειτουργούν ως γράφοι αναφοράς για την μετατροπή ενός κοινωνικού δικτύου σε ένα ισοδύναμο διάγραμμα, κάθε εγγραφή του οποίου ισούται με την συχνότητα εμφάνισης του συγκεκριμένου συχνού υπογράφου. Η τεχνική παρουσίασε μεγάλη ακρίβεια στην πρόβλεψη συναισθήματος ξεπερνώντας κατά πολύ τις τρέχουσες τεχνικές. Στο 4ο κεφάλαιο προσπάθησα να εισάγω τεχνικές από τον ταχέως εξελισσόμενο χώρο του deep learning, ο οποίος κατά μεγάλο βαθμό στηρίζεται στην κατασκευή μεταβλητών (feature engineering) μέσω αρχιτεκτονικών νευρωνικών δικτύων πολλών επιπέδων.

⁹ <http://www.dailymail.co.uk/sciencetech/article-2871378/Is-Facebook-finally-DISLIKE-button-Mark-Zuckerberg-admits-considering-it.html>

Σε σχέση με την βιβλιογραφία, η μέθοδος που παρουσίασα ήταν η πρώτη προσπάθεια πρόβλεψης συναισθήματος ως ένα πρόβλημα ημι-εποπτευόμενης μάθησης με χρήση τεχνικών deep learning. Συγκεκριμένα προτείναμε μια αρχιτεκτονική 4 έως 6 επιπέδων, στο πρώτο επίπεδο της οποίας τοποθετήσαμε ένα Restricted Boltzmann Machine (RBM) ενώ το υπόλοιπο δίκτυο λειτούργησε ως δίκτυο αυτοκωδικοποιητή. Η μέθοδος εμφάνισε αύξηση 4-10% σε απόλυτες τιμές ακρίβειας σε σχέση με τις τρέχουσες τεχνικές.

Τέλος, στο 5ο κεφάλαιο ασχολήθηκα με το πρόβλημα της ανεύρεσης κοινοτήτων σε μεγάλα κοινωνικά δίκτυα, αξιοποιώντας τις πρόσφατες εξελίξεις στον χώρο της Επεξεργασίας Φυσικής Γλώσσας και πιο συγκεκριμένα των σημασιολογικών διανυσματικών αναπαραστάσεων. Ακολουθώντας τις βέλτιστες πρακτικές [50] στον συγκεκριμένο χώρο, παρουσίασα μια μέθοδο εύκολα υλοποιήσιμη, με μικρό υπολογιστικό κόστος, η οποία παρουσιάζει μεγάλη ακρίβεια στον εντοπισμό *ετερογενών* κοινοτήτων.

6.2 Μελλοντική Εργασία

Στα πλαίσια της διδακτορικής διατριβής θεώρησα ως εξαιρετικό ενδιαφέρον το γεγονός ότι τεχνικές και αλγόριθμοι από διάφορους επιστημονικούς τομείς μπόρεσαν (με μικρές παραλλαγές) να επιλύσουν προβλήματα στον χώρο της Ανάλυσης Κοινωνικών Δικτύων. Η μελλοντική ερευνητική μου εργασία θα έχει σκοπό να συνεχίσει σε αυτό το μονοπάτι, συνδυάζοντας με την σωστή αναλογία τις έννοιες της κοινωνικής θεωρίας με τις προηγμένες τεχνικές από άλλους τομείς. Μεγαλύτερη εξέλιξη προμηνύεται στον χώρο του deep learning, έναν ιδιαίτερο τομέα της μηχανικής μάθησης, ο οποίος δείχνει να έχει ωριμάσει αρκετά, γεγονός που αποδεικνύεται και από το ευρύ πεδίο εφαρμογών του. Βασικό χαρακτηριστικό των

αλγορίθμων αυτών είναι το μεγάλο βάρος που δίνεται στην κατασκευή αξιόπιστων μεταβλητών περιγραφής απ' ευθείας από τα δεδομένα. Σκοπός μου είναι να μεταφέρω τις ιδέες του deep learning στον χώρο της Ανάλυσης Δικτύων, δίνοντας στον τελικό χρήστη/ερευνητή τα κατάλληλα εργαλεία για την μελέτη ενός κοινωνικού δικτύου.

Γλωσσάριο

Στον παρακάτω πίνακα παρατίθενται οι όροι που χρησιμοποιήθηκαν στη διατριβή:

RBM	Restricted Boltzmann Machine
logit	Το αντίστροφο της σιγμοειδούς συνάρτησης
GloVe	Αλγόριθμος για την εύρεση semantic vector representations [51]
SVM	Support Vector Machine

Βιβλιογραφικές Αναφορές

- [1] F. Heider, «Attitudes and Cognitive Organization,» *Journal of Psychology*, τόμ. 21, pp. 107-112, 1946.
- [2] D. Cartwright και F. Harary, «Structural Balance : a generalization of Heider's theory,» *Psychological Review*, τόμ. 63, αρ. 5, pp. 277-293, 1956.
- [3] F. Harary, «On the notion of balance of a signed graph,» *Michigan Mathematical Journal*, pp. 143-146, 1953.
- [4] T. Antal, P. L. Krapivsky και S. Redner, *Social Balance on Networks: The Dynamics of Friendship and Enmity*, 2008.
- [5] J. Leskovec, D. P. Huttenlocher και J. M. Kleinberg, «Signed networks in social media,» σε *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '10)*, Atlanta, Georgia, USA, 2010.
- [6] J. Leskovec, D. Huttenlocher και J. Kleinberg, «Predicting positive and negative links in online social networks,» σε *In Proceedings of the 19th international conference on World wide web (WWW '10)*, 2010.
- [7] X. Su και T. M. Khoshgoftaar, «A Survey of Collaborative Filtering Techniques,» *Advances in Artificial Intelligence*, 2009.
- [8] Y. Hu, Y. Koren και C. Volinsky, «Collaborative filtering for implicit feedback datasets,» σε *In IEEE International Conference on Data Mining (ICDM 2008)*, 2008.
- [9] K. Yu, A. Schwaighofer, V. Tresp, W.-Y. Ma και H. J. Zhang, «Collaborative Ensemble Learning: Combining Collaborative and Content-Based Information Filtering,» σε *In Proceedings of the 19th Conference on Uncertainty in Artificial Intelligence*, 2003.
- [10] A. Paterek, «Improving regularized singular value decomposition for collaborative filtering,» σε *In Proc. KDD Cup Workshop at SIGKDD'07*, 2007.
- [11] P. Massa και B. Bhattacharjee, «Using Trust in Recommender Systems: An

- Experimental Analysis,» σε *In Proceedings of iTrust2004 International Conference*, 2004.
- [12] Y. Karabulut, J. Mitchell, P. Herrmann και C. Jensen, «Trust-Based Collaborative Filtering,» *IFIP – The International Federation for Information Processing*, 2008.
- [13] P. Massa και P. Avesani, «Trust-aware bootstrapping of recommender systems,» σε *ECAI Workshop on Recommender Systems*, 2006.
- [14] J. O'Donovan και B. Smyth, «Trust in recommender systems,» σε *In Proceedings of the 10th international conference on Intelligent user interfaces (IUI '05)*, 2005.
- [15] M. McPherson, L. Smith-Lovin και J. M. Cook, «Birds of a Feather: Homophily in Social Networks,» *Annual Review of Sociology*, August 2001.
- [16] M. Everett και S. P. Borgatti, «Ego network betweenness,» *Social Networks*, τόμ. 27, January 2005.
- [17] D. Liben-Nowell και J. Kleinberg, «The link prediction problem for social networks,» σε *Proceedings of the twelfth international conference on Information and knowledge management (CIKM 2003)*, 2003.
- [18] L. A. Adamic και E. Adar, «Friends and Neighbors on the Web,» *Social Networks*, τόμ. 25, pp. 211-230, 2003.
- [19] T. Fawcett, «An introduction to ROC analysis,» *Pattern Recogn. Lett.*, pp. 861-874, June 2006.
- [20] F. Provost, *Machine Learning from Imbalanced Data Sets 101 (Extended Abstract)*.
- [21] J. Zhang και I. Mani, «KNN Approach to Unbalanced Data Distributions: A Case Study Involving Information Extraction,» σε *Proceedings of the ICML'2003 Workshop on Learning from Imbalanced Datasets*, 2006.
- [22] R. Guha, R. Kumar, P. Raghavan και A. Tomkins, «Propagation of trust and distrust,» σε *Proceedings of the 13th international conference on World Wide Web (WWW '04)*, New York, NY, USA, 2004.
- [23] D. Garcia, A. Garas και F. Schweitzer, «Positive words carry less information than negative words,» *CoRR*, 2011.
- [24] J. Boucher και C. E. Osgood, «The pollyanna hypothesis,» *Journal of Verbal Learning and Verbal Behavior*.
- [25] C. C. Aggarwal και J. Han, Επιμ., *Frequent Pattern Mining*, Springer, 2014.
- [26] J. Bengio, A. Courville και P. Vincent, «Representation Learning: A Review and New

- Perspectives,» *IEEE Trans. Pattern Anal. Mach. Intell.*, 2013.
- [27] G. Hinton, «A Practical Guide to Training Restricted Boltzmann Machines,» 2010.
- [28] Y. Bengio, «Learning Deep Architectures for AI,» *Found. Trends Mach. Learn.*, January 2009.
- [29] S. J. Pan και Q. Yang, «A Survey on Transfer Learning,» *IEEE Trans. on Knowl. and Data Eng.*, October 2010.
- [30] R. Raina, A. Battle, H. Lee, B. Packer και Y. A. Ng, «Self-taught learning: Transfer learning from unlabeled data,» σε *Proceedings of the Twenty-fourth International Conference on Machine Learning*, Corvallis, Oregon, 2007.
- [31] S. H. Yang, A. Smola, B. Long, H. Zha και Y. Chang, «Friend or frenemy?: Predicting signed ties in social networks,» σε *In Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval (SIGIR '12)*, 2012.
- [32] E. K. P. Chong και S. H. Zak, *An Introduction to Optimization*, 4th edition, Wiley, 2013.
- [33] E. Alpaydin, *Introduction to Machine Learning*, 2nd edition, 2009.
- [34] M. P. Kevin, *Machine Learning: a Probabilistic Perspective*, The MIT Press, 2012.
- [35] N. Le Roux και Y. Bengio, «Representational Power of Restricted Boltzmann Machines and Deep Belief Networks,» *Neural Comput.*, June 2008.
- [36] G. E. Hinton, «Training products of experts by minimizing contrastive divergence,» *Neural Computation*, τόμ. 14, αρ. 8, pp. 1771 - 1800, August 2002.
- [37] R. Salakhutdinov, A. Mnih και G. Hinton, «Restricted Boltzmann machines for collaborative filtering,» σε *In Proceedings of the 24th international conference on Machine learning (ICML '07)*, Corvallis, Oregon, 2007.
- [38] G. E. Hinton και R. R. Salakhutdinov, «Reducing the dimensionality of data with neural networks,» *Science*, τόμ. 313, pp. 504 - 507, 28 July 2006.
- [39] P. Baldi, «Autoencoders, Unsupervised Learning, and Deep Architectures,» σε *JMLR: Workshop on Unsupervised and Transfer Learning*, 2012.
- [40] P. Massa και P. Avesani, «Trust-aware bootstrapping of recommender systems,» σε *ECAI Workshop on Recommender Systems*, 2006.
- [41] B. R. Palm, *Prediction as a candidate for learning deep hierarchical models of data - M.S. thesis*, 2012.
- [42] T. Sing, O. Sander, N. Beerenwinkel και T. Lengauer, «ROCR: visualizing classifier

- performance in R,» *Bioinformatics*, τόμ. 21, p. 7881, 2005.
- [43] T. Helleputte, «LiblineaR: Linear Predictive Models Based on the LIBLINEAR C/C++ Library,» 2015.
- [44] D. Meyer και T. U. Wien, *Support Vector Machines. The Interface to libsvm in package e1071. Online-Documentation of the package e1071 for R*, 2001.
- [45] C.-J. Hsieh, K.-Y. Chiang και I. S. Dhillon, «Low rank modeling of signed networks,» σε *In Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '12)*, 2012.
- [46] H. Narayanan και S. Mitter, «Sample Complexity of Testing the Manifold Hypothesis,» σε *Advances in Neural Information Processing Systems*, 2010.
- [47] F. Ricci, L. Rokach, B. Shapira και P. B. Kantor, Επιμ., *Recommender Systems Handbook*, Springer US, 2011.
- [48] L. Tang και H. Liu, «Relational learning via latent social dimensions,» σε *In Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '09*, 2009.
- [49] L. Tang και H. Liu, «Scalable learning of collective behavior based on sparse social dimensions,» σε *In Proceedings of the 18th ACM conference on Information and knowledge management*, 2009.
- [50] B. Perozzi, R. Al-Rfou και S. Skiena, «DeepWalk: online learning of social representations,» σε *In Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '14)*, 2014.
- [51] J. Pennington, R. Socher και C. D. Manning, «Glove: Global vectors for word representation,» σε *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 2014.
- [52] C. C. Aggarwal, *Social Network Data Analytics*, Springer Publishing Company, 2011.
- [53] S. A. Macskassy και F. Provost, «A simple relational classifier,» σε *In Proceedings of the Second Workshop on Multi-Relational Data Mining (MRDM-2003) at KDD-2003*, 2003.
- [54] Q. Lu και L. Getoor, «Link-based Classification,» σε *In Proceedings of the 20th International Conference on Machine Learning (ICML '03)*, 2003.
- [55] X. Zhu, Z. Ghahramani και J. Lafferty, «Semi-Supervised Learning Using Gaussian Fields and Harmonic Functions,» σε *In Proceeding of the 20th International Conference*

- on *Machine Learning*, ICML 2003.
- [56] M. Szummer και T. Jaakkola, «Partially labeled classification with Markov random walks,» σε *In Proceedings of Advances in Neural Information Processing Systems 14 (NIPS 2001)*, 2001.
- [57] L. Tang και H. Liu, «Leveraging social media networks for classification,» *Data Mining and Knowledge Discovery*, pp. 447-478, 2011.
- [58] T. Mikolov, I. Sutskever, K. Chen, G. Corrado και J. Dean, «Distributed Representations of Words and Phrases and their Compositionality,» σε *In Proceedings of NIPS*, 2013.
- [59] J. Duchi, E. Hazan και Y. Singer, «Adaptive Subgradient Methods for Online Learning and Stochastic Optimization,» *J. Mach. Learn. Res.*, 2011.
- [60] E. M. Newman, «Clustering and preferential attachment in growing networks,» *Physical Review Letters E*, 2001.
- [61] G. Salton και M. J. McGill, *Introduction to Modern Information Retrieval*, McGraw Hill, 1983.
- [62] L. Katz, «A new status index derived from sociometric analysis,» *Psychometrika*, 1953.
- [63] T. Mikolov, K. Chen, G. Corrado και J. Dean, «Efficient Estimation of Word Representations in Vector Space,» σε *In Proceedings of Workshop at ICLR*, 2013.
- [64] T. Mikolov, W.-t. Yih και G. Zweig, «Linguistic Regularities in Continuous Space Word Representations,» σε *In Proceedings of NAACL HLT*, 2013.
- [65] M. Newman, *Networks: An Introduction*, 2010.
- [66] A. Papaoikonomou, M. Kardara, K. Tserpes και T. Varvarigou, «The Strength of Negative Opinions,» σε *Engineering Applications of Neural Networks*, 2013.
- [67] A. Papaoikonomou, M. Kardara, K. Tserpes και T. Varvarigou, «Predicting Edge Signs in Social Networks Using Frequent Subgraph Discovery,» *IEEE Internet Computing*, pp. 36-43, Sept-Oct 2014.
- [68] D. Easley και J. Kleinberg, *Networks, Crowds, and Markets: Reasoning about a Highly Connected World*, New York: Cambridge University Press, 2010.
- [69] D. J. Watts και S. H. Strogatz, «Collective dynamics of 'small-world' networks,» *Nature*, pp. 440-442, 1998.
- [70] P. Massa και P. Avesani, «Trust-aware recommender systems,» σε *In Proceedings of the 2007 ACM conference on Recommender systems (RecSys '07)*, 2007.

- [71] G. Pitsilis και L. Marshall, «A trust-enabled P2P recommender system,» σε *Enabling Technologies: Infrastructure for Collaborative Enterprises, 2006. WETICE'06. 15th IEEE International Workshops on*, 2006.
- [72] G. Pitsilis και L. Marshall, «Trust as a key to improving recommendation systems,» σε *In Proceedings of the Third international conference on Trust Management (iTrust'05)*, 2005.
- [73] N. Lathia, S. Hailes και L. Capra, «Trust-Based Collaborative Filtering,» σε *IFIP – The International Federation for Information Processing*, 2008.
- [74] Y. Bengio, R. Ducharme, V. Pascal και C. Janvin, «A neural probabilistic language model,» *The Journal of Machine Learning Research*, 2003.
- [75] R. A. Hanneman και M. Riddle, «Introduction to social network methods,» Riverside, CA: University of California, Riverside.
- [76] H. Goh, N. Thome και M. Cord, «Biasing restricted Boltzmann machines to manipulate latent selectivity and sparsity,» σε *NIPS workshop on deep learning and unsupervised feature learning*, 2010.



Ο Δρ. Αθανάσιος Α. Παπαϊκονόμου γεννήθηκε στην Αθήνα το 1984. Αποφοίτησε από τη Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών του Εθνικού Μετσόβιου Πολυτεχνείου το 2008. Το 2015, ολοκλήρωσε τη διδακτορική του διατριβή στο Εθνικό Μετσόβιο Πολυτεχνείο με θέμα την ανάπτυξη τεχνικών για την ανάλυση κοινωνικών δικτύων με έμφαση σε γράφους εμπιστοσύνης. Από το 2010 εργάζεται ως ερευνητικός συνεργάτης στο Εργαστήριο Τηλεπικοινωνιών της Σχολής Ηλεκτρολόγων Μηχανικών και Μηχανικών Ηλεκτρονικών Υπολογιστών του Εθνικού Μετσόβιου Πολυτεχνείου και συμμετέχει σε ερευνητικά προγράμματα χρηματοδοτούμενα από την Ευρωπαϊκή Ένωση. Τα ερευνητικά του ενδιαφέροντα εστιάζονται στον τομέα της ανάλυσης κοινωνικών γράφων, μηχανικής μάθησης και τεχνητής νοημοσύνης.