



**ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ**  
**ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ**  
**ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ**

**ΤΕΧΝΙΚΕΣ ΑΝΑΚΤΗΣΗΣ ΠΕΡΙΕΧΟΜΕΝΟΥ**  
**ΚΑΙ ΑΝΑΛΥΣΗΣ ΔΙΑΧΥΣΗΣ ΤΗΣ ΕΠΙΡΡΟΗΣ**  
**ΣΤΑ ΚΟΙΝΩΝΙΚΑ ΔΙΚΤΥΑ**

**ΔΙΔΑΚΤΟΡΙΚΗ ΔΙΑΤΡΙΒΗ**

**ΜΑΓΔΑΛΗΝΗΣ Δ. ΚΑΡΔΑΡΑ**

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

**ΕΠΙΒΛΕΠΟΥΣΑ:**

**Θ. ΒΑΡΒΑΡΙΓΟΥ**

Καθηγήτρια Ε.Μ.Π.

**ΑΘΗΝΑ, Ιούλιος 2015**





# ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ

ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ

ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ

## ΤΕΧΝΙΚΕΣ ΑΝΑΚΤΗΣΗΣ ΠΕΡΙΕΧΟΜΕΝΟΥ ΚΑΙ ΑΝΑΛΥΣΗΣ ΔΙΑΧΥΣΗΣ ΤΗΣ ΕΠΙΡΡΟΗΣ ΣΤΑ ΚΟΙΝΩΝΙΚΑ ΔΙΚΤΥΑ

ΔΙΔΑΚΤΟΡΙΚΗ ΔΙΑΤΡΙΒΗ

**ΜΑΓΔΑΛΗΝΗΣ Δ. ΚΑΡΔΑΡΑ**

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

**Συμβουλευτική Επιτροπή:**

1. Θ. ΒΑΡΒΑΡΙΓΟΥ, Καθ. Ε.Μ.Π. (Επιβλέπουσα)
2. Μ. ΘΕΟΛΟΓΟΥ, Καθ. Ε.Μ.Π.
3. Σ. ΠΑΠΑΒΑΣΙΛΕΙΟΥ, Καθ. Ε.Μ.Π.

Εγκρίθηκε από την επταμελή εξεταστική επιτροπή την 6η Ιουλίου, 2015.

.....  
Θ. Βαρβαρίγου,  
Καθ. Ε.Μ.Π.

.....  
Μ. Θεολόγου,  
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, Ιούλιος 2015**

.....  
Μαγδαληνή Δ. Καρδαρά

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

Copyright © Μαγδαληνή Δ. Καρδαρά, 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 γνώμών του συγγραφέα (Ν. 5343/1932, Άρθρο 202).

## *Περίληψη*

Η παρούσα διδακτορική διατριβή ασχολείται με την ανάπτυξη τεχνολογιών ανάλυσης δεδομένων από κοινωνικά δίκτυα. Βασικό αντικείμενο της διατριβής είναι η ανάπτυξη ενός πλαισίου αξιολόγησης για την εκτίμηση της σχετικής επίδοσης θεωριών τοπικής επιρροής. Οι θεωρίες επιρροής είναι τυποποιημένα μοντέλα που εξετάζουν ένα κοινωνικό δίκτυο και εντοπίζουν άτομα που μπορούν να επηρεάσουν και να καθοδηγήσουν τα υπόλοιπα μέλη του δικτύου με τις πράξεις τους. Το ενδιαφέρον γύρω από την ανάπτυξη τέτοιων θεωριών είναι μεγάλο, καθώς μπορούν να ωφελήσουν σημαντικά πλήθος εμπορικών εφαρμογών όπως καμπάνιες προώθησης ή συστήματα προτάσεων. Στις περισσότερες ερευνητικές εργασίες, ωστόσο, υπάρχει έλλειψη μιας τυποποιημένης μεθοδολογίας για την αξιολόγηση των αποτελεσμάτων που προκύπτουν από τις αντίστοιχες θεωρίες σε αρκετά μεγάλη κλίμακα ώστε να επιτρέπει την εξαγωγή ασφαλών συμπερασμάτων.

Στην παρούσα εργασία, εισάγουμε ένα πλαίσιο που επιτρέπει την αυτοματοποιημένη αξιολόγηση της επίδοσης των θεωριών τοπικής επιρροής οι οποίες εφαρμόζονται σε κοινότητες χρηστών που σχηματίζονται γύρω από θεματικές κατηγορίες. Στόχος του πλαισίου είναι να εξετάσει αν οι θεωρίες αυτές εντοπίζουν ως επιδραστικές μικρές ομάδες χρηστών που είναι ικανοί με ελάχιστη προσπάθεια να επηρεάσουν κάποιες αντικειμενικά μετρήσιμες πλευρές μιας κοινότητας. Το πλαίσιο ορίζει πέντε αναγκαίες και ικανές συνθήκες τις οποίες μια αποδοτική θεωρία επιρροής θα πρέπει να ικανοποιεί για ένα σύνολο

θεματικών κατηγοριών. Θέσαμε το πλαίσιο μας σε εφαρμογή χρησιμοποιώντας ένα σύνολο δεδομένων μεγάλης κλίμακας με πραγματικά δεδομένα από 75 κοινότητες του Twitter και εξετάζοντας πάνω σε αυτό πέντε διαδοσόμενες θεωρίες τοπικής επιρροής. Τα αποτελέσματα της ανάλυσης αποκαλύπτουν σημαντικές διαφορές στην επίδοσή τους. Αναπτύξαμε επίσης μεθοδολογία για τη διερεύνηση της εσωτερικής δυναμικής των επιδραστικών ομάδων που ορίζονται από τις θεωρίες επιρροής η οποία χρησιμοποιήθηκε για την ανάλυση των επιλεγμένων θεωριών και των αποτελεσμάτων τους καθώς και για την κατηγοριοποίησή τους με βάση τα ευρήματα της ανάλυσης αυτής.

*Λέξεις Κλειδιά:* ανάλυση κοινωνικών δικτύων, κοινωνική επιρροή, θεματικές κοινότητες

## *Abstract*

This PhD thesis addresses the issue of influence diffusion in social networks by introducing a novel evaluation framework for evaluating and comparing the performance of topic-specific influence theories. Influence theories constitute formal models that identify those individuals that are able to affect and guide their peers through their activity. There is a large body of work on developing such theories, as they have important applications in viral marketing and systems. Most works, however, lack a formal methodology for evaluating the results produced by their theories in a large enough scale to draw safe, representative conclusions.

In this thesis, we introduce a formalized framework for large-scale, automatic evaluation of topic-specific influence theories that are specialized in Twitter. The objective of the framework is to determine whether these theories identify as influencers small groups of users who are able with minimal effort to affect some objectively measured aspects of a community. The framework consists of five conjunctive conditions that are indicative of real influence exertion for a set of topic categories. We put our framework into practice using a large-scale test-bed with real data from 75 Twitter communities and examining five established topic-specific theories that are applicable to our settings. The outcomes of our analysis reveal significant differences in their performance. To explain them, we introduce a novel methodology for delving into the internal dynamics of the groups of influencers they define. We use it to analyze the implications of the selected theories and their categorisation based on the resulting evidence.

**Keywords:** social network analysis, social influence, topic communities

## Πρόλογος

*Η διδακτορική διατριβή που παρουσιάζεται στις επόμενες σελίδες εκπονήθηκε από τον Ιανουάριο του 2007 μέχρι τον Ιούλιο 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 ιδιαίτερα 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. Μιχαήλ Θεολόγου και Συμεών Παπαβασιλείου.*

*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 Γιώργο Παπαδάκη, Θάνο Παπαϊκονόμου, Κλεοπάτρα Κωνσταντέλη, Βασίλη Καλογήρου, Φώτη Αίσωπο και Κωνσταντίνο Τσερπέ.*

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

*Μαγδαληνή Δ. Καρδαρά*

*Ιούλιος 2015*



## Πίνακας περιεχομένων

1	Εισαγωγή .....	7
1.1	Οργάνωση του εγγράφου .....	14
2	Ανάλυση Κοινωνικών Δικτύων .....	17
2.1	Μελέτη κοινωνικών δικτύων .....	17
2.2	Ερευνητικές περιοχές .....	19
3	Επιρροή στα κοινωνικά δίκτυα .....	25
3.1	Θεωρίες επιρροής .....	25
3.2	Αξιολόγηση θεωριών τοπικής επιρροής .....	32
3.3	Πίεση από τον κοινωνικό περίγυρο .....	33
3.4	Εφαρμογές των θεωριών επιρροής .....	34
3.5	Ομοφιλία ή κοινωνική επιρροή .....	36
4	Ορισμός Πλαισίου Αξιολόγησης .....	39
4.1	Παρουσίαση του Twitter .....	39
4.2	Βασικοί ορισμοί .....	42
4.3	Θεωρίες τοπικής επιρροής .....	46
4.4	Πλαίσιο αξιολόγησης .....	50
5	Αξιολόγηση θεωριών τοπικής επιρροής .....	60
5.1	Twitter Dataset .....	61
5.2	Συνθήκες όγκου .....	64
5.3	Συνθήκη συσχέτισης .....	69
5.4	Συνθήκη χρόνου .....	72
5.5	Συζήτηση .....	74
6	Εσωτερική ανάλυση θεωριών θεματικής επιρροής .....	77
6.1	Ομοφιλία .....	78

6.2	Επικοινωνία .....	81
6.3	Πολυπλευρικότητα.....	86
6.4	Συζήτηση .....	89
7	Ανάσχυση δεδομένων από πολλαπλά κοινωνικά δίκτυα .....	97
7.1	Εργασία με δεδομένα από πολλαπλά κοινωνικά δίκτυα.....	97
7.2	Σχετικές ερευνητικές εργασίες και υπάρχουσες εφαρμογές .....	100
7.3	Μοντέλο δεδομένων .....	101
7.4	Λεπτομέρειες υλοποίησης.....	105
8	Σύνοψη .....	110
	Βιβλιογραφικές Αναφορές.....	113

## Σχήματα- Εικόνες

Εικόνα 6-1. Εσωτερική δυναμική των διακεκριμένων ομάδων για Fmnt και Frtw. 90

Εικόνα 6-2. Εσωτερική δυναμική των διακεκριμένων ομάδων για Find και Ftwt.. 90

Εικόνα 6-3. Εσωτερική δυναμική των διακεκριμένων ομάδων για Frnd..... 90

## Πίνακες

Πίνακας 3-1: Δισδιάστατη ταξινόμια των θεωριών επιρροής του Twitter .....	29
Πίνακας 4-1. Σύνοψη των βασικών εννοιών που χρησιμοποιούνται στους ορισμούς .....	43
Πίνακας 4-2: Αναπαράσταση της CAC με ενδεικτικά δεδομένα από τρεις κοινότητες G1, G2 και G3.....	57
Πίνακας 5-1. Τεχνικά χαρακτηριστικά των 75 θεματικών κοινοτήτων που αποτελούν το σύνολο δεδομένων αναφοράς για την πειραματική μελέτη.....	63
Πίνακας 5-2: Ποσοστό του πολωμένου περιεχομένου που παράγεται από τους διακεκριμένους χρήστες για όλες τις θεωρίες επιρροής.....	65
Πίνακας 5-3: Ποσοστό συνολικού περιεχομένου που παράγεται από τις διακεκριμένες ομάδες για όλες τις θεωρίες επιρροής.....	68
Πίνακας 5-4. Συσχέτιση Pearson μεταξύ του συνολικού συναισθήματος των διακεκριμένων ομάδων και των υπόλοιπων μελών της κοινότητας.....	71
Πίνακας 5-5: Τιμές ΔCAC για τις 30 επιλεγμένες κοινότητες.....	73
Πίνακας 6-1. Βαθμός αμοιβαιότητας διακεκριμένων ομάδων για όλες τις θεωρίες επιρροής και τα μεγέθη ομάδων.....	79
Πίνακας 6-2. Πιθανότητα εσωτερικής αναφοράς (αριστερά) και πιθανότητα εσωτερικής αναδημοσίευσης (δεξιά) για όλες τις θεματικές κατηγορίες σε όλες τις θεωρίες επιρροής και τα μεγέθη επιδραστικής ομάδας. Οι ταυτολογίες εμφανίζονται με έντονα γράμματα.....	85
Πίνακας 6-3. Ομοιότητες Jaccard μεταξύ των διακεκριμένων ομάδων που ορίζονται από διαφορετικές θεωρίες επιρροής πάνω στην ίδια θεματική κατηγορία για τα συνήθη μεγέθη διακεκριμένης ομάδας.....	88
Πίνακας 6-4. Εσωτερική δυναμική των διακεκριμένων ομάδων για Fmnt και Frtw.....	90
Πίνακας 6-5. Εσωτερική δυναμική των διακεκριμένων ομάδων για Find και Ftwt.....	90
Πίνακας 6-6. Εσωτερική δυναμική των διακεκριμένων ομάδων για Frnd.....	90
Πίνακας 6-7. Η Gr@20 που ορίζεται από την κάθε θεωρία επιρροής για τη θεματική κοινότητα #iranelection.....	92

Πίνακας 6-8. Η Gr@20 που ορίζεται από την κάθε θεωρία επιρροής για τη θεματική κοινότητα #noh8. .... 94

Πίνακας 7-1: Κύρια αντικείμενα του μοντέλου δεδομένων της εφαρμογής ανάσυρσης δεδομένων από ετερογενή κοινωνικά δίκτυα. ....103

Η σελίδα αυτή είναι σκόπιμα λευκή.

# 1

## *Εισαγωγή*

Ένα κοινωνικό δίκτυο ορίζεται ως ένα δίκτυο σχέσεων ή αλληλεπιδράσεων, στο οποίο οι κόμβοι αποτελούνται από άτομα και οι ακμές είναι οι σχέσεις ή οι αλληλεπιδράσεις μεταξύ αυτών των ατόμων. Αν και η μελέτη των κοινωνικών δικτύων προηγείται χρονικά της εμφάνισης διαδικτυακών κοινωνικών δικτύων, η εμφάνιση και η εξάπλωση των ιστοσελίδων κοινωνικής δικτύωσης τα τελευταία χρόνια επέτρεψε τη μελέτη των δικτύων και της εσωτερικής δυναμικής τους σε πολύ μεγαλύτερη κλίμακα. Αυτό συμβαίνει επειδή στο πλαίσιο τέτοιων δικτύων, οι σχέσεις μεταξύ των ατόμων μοντελοποιούνται με σαφή τρόπο που διευκολύνει την εξαγωγή συμπερασμάτων, ενώ, λόγω της μεγάλης δημοτικότητάς και της πολιτικής ανοιχτών δεδομένων που ακολουθούν οι περισσότερες τέτοιες ιστοσελίδες, ο όγκος δεδομένων που έχουν στη διάθεση τους οι ερευνητές είναι τεράστιος. Η δυνατότητα μελέτης των κοινωνικών δικτύων σε πολύ μεγαλύτερη από την έως τώρα κλίμακα έχει δώσει νέα ώθηση σε παραδοσιακούς ερευνητικούς τομείς που σχετίζονται με την ανακάλυψη γνώσης και την εξόρυξη δεδομένων.

Μια από τις ερευνητικές περιοχές που ωφελήθηκε από την άνθιση των διαδικτυακών κοινοτήτων είναι η μελέτη της επιρροής με στόχο τον εντοπισμό επιδραστικών χρηστών σε ένα κοινωνικό δίκτυο. Στο πλαίσιο ενός κοινωνικού δικτύου, οι

επιδραστικοί χρήστες είναι ξεχωριστά άτομα με κάποια ειδικά χαρακτηριστικά που τους επιτρέπουν να επηρεάζουν ένα δυσανάλογα μεγάλο αριθμό άλλων χρηστών με τις πράξεις τους. Αυτά τα ειδικά χαρακτηριστικά σχετίζονται με την ατομική τους δραστηριότητα και το κοινωνικό τους υπόβαθρο καθώς και με τη θέση τους στο δίκτυο, δηλαδή τις διασυνδέσεις τους με τα άλλα μέλη. Οι επιδραστικοί χρήστες συνήθως διαδραματίζουν καθοδηγητικό ρόλο σε μια κοινότητα και η δυνατότητα επηρεασμού της συμπεριφοράς τους μπορεί να είναι χρήσιμη για ένα πλήθος επιστημονικών και επιχειρησιακών εφαρμογών [1]. Για παράδειγμα, οι καμπάνιες προώθησης θα μπορούσαν να αποκομίσουν σημαντικό όφελος από αυτή τη διαδικασία, καθώς οι πελάτες τείνουν να ακολουθούν τους χρήστες με υψηλή επιρροή στην προτίμηση συγκεκριμένων προϊόντων [2]. Αντί να κατακλύζουν την καταναλωτική βάση με μαζικές αλλά τυφλές διαφημίσεις, οι προωθητικές καμπάνιες θα μπορούσαν να στοχεύουν σε ένα μικρό αριθμό επιδραστικών ανθρώπων. Αυτή η οικονομική εναλλακτική ονομάζεται *viral marketing* και δύναται να επιτύχει παρόμοια επίπεδα διάχυσης προϊόντος με τις παραδοσιακές τακτικές [3].

Η βάση για τη δημιουργία τέτοιων εφαρμογών είναι ο εντοπισμός των επιδραστικών χρηστών, μια ερευνητική περιοχή που συγκεντρώνει έντονο ενδιαφέρον από ερευνητές. Το ενδιαφέρον αυτό είχε ως αποτέλεσμα την ανάπτυξη *θεωριών επιρροής (influence theories)*, δηλαδή τυποποιημένων μοντέλων τα οποία αξιολογούν την επιρροή που ασκεί κάθε μέλος ενός κοινωνικού δικτύου στους υπόλοιπους χρήστες με βάση ένα ή περισσότερα κριτήρια. Κάποιες τέτοιες θεωρίες μελετούν την επιρροή σε καθολική κλίμακα, λαμβάνοντας δηλαδή υπόψη τη δραστηριότητα και τη βάση χρηστών ολόκληρου του κοινωνικού δικτύου. Είναι ωστόσο πιο σύνθητες η επιρροή ενός χρήστη να είναι τοπική: ο χρήστης μπορεί να θεωρείται αυθεντία σε ένα συγκεκριμένο τομέα αλλά η γνώμη του να μην έχει ιδιαίτερη σημασία εκτός αυτού



του τομέα. Με βάση αυτή την αρχή, οι θεωρίες τοπικής επιρροής έχουν ως στόχο να εντοπίσουν επιδραστικούς χρήστες ανάμεσα στα μέλη συγκεκριμένων κοινοτήτων που σχηματίζονται γύρω από ένα συγκεκριμένο θέμα.

Μια βασική πρόκληση σχετική με τη μελέτη των θεωριών τοπικής επιρροής είναι η αξιολόγησή τους. Τα τελευταία χρόνια τα δικτυακά κοινωνικά δίκτυα πρόσφεραν στους ερευνητές πολύτιμα εργαλεία για τη μελέτη της δυναμικής της διάχυσης επιρροής. Περιέχουν μεγάλο όγκο δεδομένων που δημιουργείται από τους χρήστες καθώς και σαφείς διασυνδέσεις μεταξύ των μελών, επιτρέποντας την ανάλυση της κοινωνικής επιρροής σε πολύ μεγαλύτερη κλίμακα από ό,τι στο παρελθόν. Παρόλα αυτά, η επιρροή αποτελεί μια υποκειμενική έννοια και ως τέτοια είναι δύσκολο να μετρηθεί και να σκιαγραφηθεί. Στις περισσότερες ερευνητικές εργασίες υπάρχει έλλειψη μιας τυποποιημένης μεθοδολογίας για την αξιολόγηση των αποτελεσμάτων που προκύπτουν από τις αντίστοιχες θεωρίες. Αντί γι' αυτό, συχνά καταφεύγουν στην επιλογή ενός μικρού δείγματος από τους κορυφαίους σε κατάταξη χρήστες και επιχειρούν να εκτιμήσουν την αξία της γνώμης τους στον πραγματικό κόσμο (π.χ. τη φήμη τους ή την ποιότητα του περιεχομένου που παράγουν) [4] [5]. Αυτή η πρακτική όμως, δε μπορεί να αναχθεί σε κλίμακα μεγάλου όγκου δεδομένων και είναι συνεπώς αδύνατο να χρησιμοποιηθεί για την εξαγωγή αντιπροσωπευτικών, αναπαράξιμων και γενικεύσιμων αποτελεσμάτων.

Στην παρούσα διατριβή, επιχειρούμε να ξεπεράσουμε αυτή την έλλειψη, ορίζοντας ένα καθορισμένο πλαίσιο που μπορεί να χρησιμοποιηθεί για την αξιολόγηση θεωριών τοπικής επιρροής σε κοινωνικά δίκτυα σε μεγάλη κλίμακα. Το πλαίσιο παίρνει ως δεδομένα εισόδου τις ομάδες επιδραστικών χρηστών που ορίζει η κάθε θεωρία – τις οποίες στην παρούσα εργασία ονομάζουμε διακεκριμένες ομάδες – μαζί με την υπόλοιπη κοινότητα και την αντίστοιχη δραστηριότητα. Ο στόχος του πλαισίου είναι

να αξιολογήσει την σχετική ακρίβεια των θεωριών επιρροής στην πρόβλεψη προτύπων συμπεριφοράς που υποδηλώνουν μίμηση από την υπόλοιπη κοινότητα. Στο εσωτερικό του, το πλαίσιο μας συμπεριλαμβάνει πέντε προϋποθέσεις που θα πρέπει να ικανοποιούνται από μια διακεκριμένη ομάδα ώστε να έχει πραγματική επιρροή στα μέλη της υπόλοιπης κοινότητας. Οι προϋποθέσεις αυτές συνοψίζονται ως εξής:

1. οι πραγματικά επιδραστικοί χρήστες αποτελούν ένα μικρό υποσύνολο της κοινότητας,
2. δύνανται να επηρεάσουν τα υπόλοιπα μέλη με ελάχιστο κόστος, παράγοντας δηλαδή ένα μικρό ποσοστό της συνολικής δραστηριότητας της κοινότητας,
3. η δραστηριότητα τους παρουσιάζει υψηλή συσχέτιση με αυτή της υπόλοιπης κοινότητας σε σχέση με μια αντικειμενικά υπολογίσιμη μετρική,
4. η σχετική με τη μετρική αυτή δραστηριότητά τους προηγείται χρονικά αυτήν των υπόλοιπων μελών της κοινότητας, και
5. ο όγκος της σχετικής με τη μετρική δραστηριότητας τους αντιστοιχεί σε ένα μικρό υποσύνολο της συνολικής δραστηριότητας που ορίζει αυτή η μετρική.

Οι προϋποθέσεις 1, 2 και 5 αποτελούν τις απαιτήσεις προ-επεξεργασίας του πλαισίου μας. Στόχος τους είναι να διασφαλίζουν ότι μια διακεκριμένη ομάδα είναι συμβατή με το πλαίσιο, από την άποψη ότι σχετίζεται με ένα περιορισμένο υποσύνολο δραστηριότητας και χρηστών της υποκείμενης κοινότητας. Είναι θεμελιώδεις προϋποθέσεις, χωρίς τις οποίες δεν είναι δυνατή η εξαγωγή ασφαλών συμπερασμάτων από το πλαίσιο μας. Οι δύο προϋποθέσεις που απομένουν εμπεριέχουν την πραγματική λειτουργία του πλαισίου. Βασίζονται σε μία

αντικειμενικά υπολογίσιμη μετρική που συσχετίζει την δραστηριότητα μιας διακεκριμένης ομάδας με αυτή της υπόλοιπης κοινότητας. Για την κατανόηση της λειτουργίας της, ας θεωρήσουμε μια μετρική που αξιολογεί το συνολικό συναίσθημα μιας ομάδας χρηστών: η υψηλή συσχέτιση μεταξύ της διακεκριμένης ομάδας και των υπόλοιπων μελών της κοινότητας είναι ένδειξη ότι η στάση της πρώτης συμπίπτει με τη συνολική διάθεση της κοινότητας. Μια θεωρία επιρροής που αδυνατεί να ικανοποιήσει κάποια από τις δύο προϋποθέσεις που ορίσαμε είναι ανεπαρκής για να εντοπίσει πραγματικά επιδραστικούς χρήστες. Αντιθέτως, μια θεωρία επιρροής είναι επιτυχής αν οι χρήστες τους οποίους αξιολογεί ως επιδραστικούς ικανοποιούν και τις δύο προϋποθέσεις. Όσο περισσότερο ικανοποιούνται οι συνθήκες αυτές, τόσο πιο επιτυχής είναι μια θεωρία, δηλαδή ανάμεσα σε δύο θεωρίες με παρόμοια επίδοση, προτιμάται αυτή που επιτυγχάνει τη μεγαλύτερη συσχέτιση μεταξύ διακεκριμένων χρηστών και της υπόλοιπης κοινότητας.

Δεδομένου ότι και οι πέντε προϋποθέσεις βασίζονται σε αντικειμενικά υπολογίσιμες μετρικές, το πλαίσιο μας επιτρέπει τη σύγκριση θεωριών τοπικής επιρροής σε μεγάλη κλίμακα χωρίς περαιτέρω παρέμβαση. Για να το θέσουμε σε εφαρμογή, συγκροτούμε ένα σύνολο πειραματικών δεδομένων μεγάλης κλίμακας το οποίο αποτελείται από πραγματικά δεδομένα. Χρησιμοποιούμε δεδομένα από το Twitter , το οποίο επιλέξαμε για αρκετούς λόγους [1] [4] [5]: είναι ένα από τα πιο δημοφιλή κοινωνικά δίκτυα, διέπεται από δυναμικούς αλλά εύληπτους και καθορισμένους κανόνες για την κοινωνική αλληλεπίδραση μεταξύ των μελών του, περιέχει πληθώρα δυναμικών θεματικών κοινοτήτων και τέλος παρέχει εύκολη πρόσβαση σε μεγάλο όγκο περιεχομένου χρηστών. Συνολικά, τα πειραματικά δεδομένα μας αποτελούνται από 75 θεματικές κοινότητες από το Twitter με πάνω από 600,000 δραστήριους χρήστες που έχουν δημοσιεύσει πάνω από 6 εκατομμύρια μηνύματα σε μια χρονική περίοδο 7

μηνών. Είναι συνεπώς κατάλληλα για την πραγματοποίηση ποιοτικών και ποσοτικών αναλύσεων μεγάλης κλίμακας με το πλαίσιο μας.

Η πειραματική μελέτη μας συμπεριλαμβάνει επίσης αρκετές θεωρίες τοπικής επιρροής από τη βιβλιογραφία. Για να διευκολύνουμε την κατανόηση της λειτουργίας τους, εισάγουμε μια δισδιάστατη ταξινόμια που τις κατατάσσει σε κατηγορίες με βάση το εύρος και τη μετρική που χρησιμοποιούν για την εκτίμηση της επιρροής. Το πρώτο κριτήριο διαιρεί τις θεωρίες σε καθολικές, τοπικές (δηλαδή βάσει θέματος) και υβριδικές, ενώ το δεύτερο εξετάζει τις παραμέτρους που λαμβάνουν υπόψη και τις κατηγοριοποιεί σε θεωρίες βασισμένες στο περιεχόμενο, βασισμένες στο γράφο και υβριδικές. Στη συνέχεια τοποθετούμε τις βασικές θεωρίες επιρροής για το Twitter στην ταξινόμια μας και εξηγούμε ποιες από αυτές είναι συμβατές ώστε να αξιολογηθούν με το πλαίσιο που εισάγουμε. Η ανάλυση μας καταλήγει στην επιλογή πέντε καθιερωμένων και αντιπροσωπευτικών θεωριών επιρροής οι οποίες συναντώνται ευρέως στη βιβλιογραφία.

Τα αποτελέσματα της αξιολόγησης κατέδειξαν σημαντικές διαφορές στην επίδοση των υπό εξέταση θεωριών επιρροής. Για να εξηγήσουμε τα αποτελέσματα που προκύπτουν, εισάγουμε μια καινοτόμα μεθοδολογία για την επισκόπηση της εσωτερικής λειτουργίας κάθε θεωρίας ώστε να εξετάσουμε τη δυναμική της ομάδας επιδραστικών χρηστών που ορίζει. Η μεθοδολογία αυτή, αποτελείται ουσιαστικά από ένα σύνολο στατιστικών αναλύσεων που καταδεικνύουν τρεις πλευρές κάθε διακεκριμένης ομάδας:

1. τα επίπεδα ομοφιλίας ανάμεσα στα μέλη,
2. η πολλαπλότητα της δραστηριότητάς τους, και
3. η επικοινωνία μεταξύ τους, με βάση τη συχνότητα των αλληλεπιδράσεών τους σε ζεύγη.

Τα αποτελέσματα της μεθοδολογίας αυτής συνηγορούν προς μια τριμερή κατηγοριοποίηση των θεωριών επιρροής: (i) αυτές που σχηματίζουν ομάδες επιδραστικών χρηστών με ισχυρούς δεσμούς μεταξύ τους, (ii) αυτές που επιλέγουν ασύνδετους αλλά ατομικά ισχυρούς επιδραστικούς χρήστες και (iii) αυτές που ορίζουν ως επιδραστικούς συνηθισμένους χρήστες χωρίς καμία αίσθηση ομαδικότητας και με χαμηλά επίπεδα σύμπραξης. Στην πράξη, η τελευταία κατηγορία επιτυγχάνει χαμηλή επίδοση στο πλαίσιο αξιολόγησης μας ενώ η πρώτη εντοπίζει πραγματικά επιδραστικούς χρήστες που συντονίζονται μεταξύ τους ώστε να διαχύσουν την επιρροή τους σε όλη την κοινότητα. Παρόμοια επίδοση επιτυγχάνεται από τους επιδραστικούς χρήστες της δεύτερης κατηγορίας, παρά την περιορισμένη συνεργασία μεταξύ τους, καθώς επωφελούνται της ατομικής τους αξίας.

Εκτός από την παραπάνω ερευνητική μελέτη, η οποία αποτελεί και το βασικό αντικείμενο της εργασίας, στην παρούσα διατριβή εξετάζουμε και το πρόβλημα της εξόρυξης δεδομένων από ετερογενή κοινωνικά δίκτυα. Αν και το περιεχόμενο των κοινωνικών δικτύων είναι άπλετο σε όγκο και εύκολα προσβάσιμο, η αξιοποίησή του συνεχίζει να παρουσιάζει σημαντικές προκλήσεις καθώς, παρά τις ομοιότητες σε έννοιες και βασικές λειτουργίες, η αναπαράσταση των δεδομένων μεταξύ διαφορετικών κοινωνικών δικτύων είναι σημαντικά ετερογενής. Προκειμένου να αντιμετωπίσουμε το πρόβλημα αυτό, αναπτύξαμε μια εφαρμογή που επιτρέπει το συνδυασμό δεδομένων και λειτουργικότητας από πολλαπλές πλατφόρμες κοινωνικής δικτύωσης με ενιαίο τρόπο.

Συνολικά τα βασικά επιτεύγματα της παρούσας εργασίας είναι τα εξής:

- Μοντελοποιούμε το πρόβλημα της αξιολόγησης της επίδοσης θεωριών τοπικής επιρροής σε μεγάλη κλίμακα. Στην πράξη το περιορίζουμε στον έλεγχο πέντε αντικειμενικά μετρήσιμων συνθηκών που παρέχουν

σημαντικές ενδείξεις πραγματικής επιρροής στο πλαίσιο οποιουδήποτε κοινωνικού δικτύου.

- Θέτουμε το πλαίσιο αξιολόγησης σε πράξη, εξετάζοντας πέντε καθιερωμένες θεωρίες τοπικής επιρροής πάνω από ένα μεγάλο σύνολο δεδομένων το οποίο αποτελείται από 75 κοινότητες του Twitter με πάνω από 600,000 χρήστες και 6 εκατομμύρια tweets.
- Αναλύουμε την επίδοση των επιλεγμένων θεωριών επιρροής μέσω μιας καινοτόμας μεθοδολογίας που επιτρέπει την κατανόηση της λειτουργίας και της δυναμικής των διακεκριμένων ομάδων που ορίζουν.
- Αναλύουμε περαιτέρω την επίδοση των επιλεγμένων θεωριών επιρροής, εισάγοντας μια διδιάστατη ταξινόμια που κατηγοριοποιεί τις θεωρίες επιρροής με βάση το εύρος και τις μετρικές που χρησιμοποιούν. Η ταξινόμια εφαρμόζεται πάνω στις βασικές θεωρίες επιρροής για το Twitter αλλά είναι αρκετά γενικεύσιμη ώστε να μπορεί να συμπεριλάβει θεωρίες για οποιοδήποτε κοινωνικό δίκτυο.
- Παρουσιάζουμε ένα εργαλείο εξόρυξης δεδομένων από πολλαπλές πλατφόρμες κοινωνικής δικτύωσης το οποίο λειτουργεί ως ενιαία διεπαφή στη λειτουργικότητα και τα δεδομένα των υποκείμενων κοινωνικών δικτύων.

## **1.1 Οργάνωση του εγγράφου**

Η παρούσα διατριβή αποτελείται από οκτώ (8) κεφάλαια. Στις ενότητες των κεφαλαίων αυτών παρουσιάζεται με αναλυτικό τρόπο το αντικείμενο της διδακτορικής διατριβής.

Το υπόλοιπο έγγραφο δομείται ως εξής: στην Ενότητα 2 εξετάζουμε τις κυριότερες ερευνητικές περιοχές στις οποίες έχει εφαρμογή η ανάλυση των δεδομένων από

κοινωνικά δίκτυα. Στην Ενότητα 3 εξετάζουμε τα βασικά ερευνητικά θέματα σχετικά με τη μελέτη της διάχυσης της επιρροής και παρουσιάζουμε τις βασικές θεωρίες επιρροής για το Twitter, οργανώνοντας τις σε μία διδιάστατη ταξινόμια. Η Ενότητα 4 μοντελοποιεί τις έννοιες που διέπουν το Twitter και βασιζόμενη σε αυτές εισάγει το πλαίσιο αξιολόγησης μας. Στην Ενότητα 5 εξετάζουμε την επίδοση κάθε θεωρίας επιρροής με βάση τις πέντε συνθήκες του πλαισίου μας, ενώ στην Ενότητα 6 εισάγουμε μια νέα μεθοδολογία για την ανάλυση της εσωτερικής δυναμικής των διακεκριμένων ομάδων. Στην Ενότητα 7 παρουσιάζουμε ένα εργαλείο που επιτρέπει την εξόρυξη δεδομένων από ετερογενή κοινωνικά δίκτυα μέσω μιας ενιαίας διεπαφής. Τέλος η Ενότητα 8 συνοψίζει την εργασία μας και προτείνει κατευθύνσεις για μελλοντική έρευνα. Η διατριβή ολοκληρώνεται με βιβλιογραφικές αναφορές.

Η σελίδα αυτή είναι σκόπιμα λευκή.



# 2

## *Ανάλυση Κοινωνικών Δικτύων*

Όπως αναφέρθηκε προηγουμένως, η ραγδαία αύξηση της δημοφιλίας των κοινωνικών δικτύων είχε ως αποτέλεσμα την αναζωπύρωση του ερευνητικού ενδιαφέροντος γύρω από περιοχές που συνδέονται με την ανάλυση τέτοιων δικτύων. Στο παρόν κεφάλαιο παρουσιάζουμε πώς η εμφάνιση των διαδικτυακών κοινοτήτων γύρω από σελίδες κοινωνικής δικτύωσης επηρέασε τη μελέτη των κοινωνικών δικτύων και εξετάζουμε τις κυριότερες ερευνητικές περιοχές στις οποίες έχει εφαρμογή η ανάλυση των δεδομένων από κοινωνικά δίκτυα.

### **2.1 Μελέτη κοινωνικών δικτύων**

Ένα κοινωνικό δίκτυο ορίζεται ως ένα δίκτυο σχέσεων ή αλληλεπιδράσεων, στο οποίο οι κόμβοι αντιστοιχούν σε άτομα και οι ακμές στις σχέσεις ή αλληλεπιδράσεις μεταξύ αυτών των ατόμων. Η έννοια των κοινωνικών δικτύων δεν περιορίζεται στην κατηγορία των διαδικτυακών κοινωνικών δικτύων όπως τα Facebook και Twitter αλλά προηγείται χρονικά της εμφάνισης τέτοιων κοινοτήτων. Οι παραδοσιακές μελέτες στην ανάλυση κοινωνικών δικτύων εξετάζουν τις αλληλεπιδράσεις μεταξύ των μελών μιας ομάδας χωρίς να επικεντρώνονται στις διαδικτυακές κοινότητες.

Ένα κλασικό κοινωνικό δίκτυο, όπως εξετάζεται παραδοσιακά από την επιστήμη της κοινωνιολογίας, βασίζεται αποκλειστικά στις ανθρώπινες αλληλεπιδράσεις. Οι μελέτες τέτοιων δικτύων συχνά περιελάμβαναν επίπονες και χρονοβόρες μεθόδους για τη μέτρηση των αλληλεπιδράσεων μεταξύ των ατόμων, καθώς δεν ήταν δυνατή η αυτόματη συλλογή των δεδομένων σχετικών με τις αλληλεπιδράσεις. Ένα παράδειγμα τέτοιας μελέτης είναι το γνωστό πείραμα των έξι βαθμών διαχωρισμού (six degrees of separation) του Milgram [6]. Στο πείραμα αυτό, οι ερευνητές χρησιμοποίησαν ταχυδρομικές επιστολές μεταξύ των συμμετεχόντων για να διαπιστώσουν αν δύο τυχαία άτομα μπορούν να συνδεθούν με μια αλυσίδα 6 ακμών. Η κάθε ακμή αντιπροσώπευε την προώθηση της επιστολής από τον ένα χρήστη στον άλλο. Τα αποτελέσματα του πειράματος αυτού αμφισβητήθηκαν έντονα από την επιστημονική κοινότητα, εξαιτίας του χαμηλού ποσοστού προώθησης των επιστολών [7]. Παρόμοια θέματα εμπόδιζαν συχνά την επιστημονικά τεκμηριωμένη διεξαγωγή τέτοιων πειραμάτων, καθώς ο ρυθμός απόκρισης των συμμετεχόντων δε μπορεί να μοντελοποιηθεί άρτια ως παράγοντας κοινωνικής αλληλεπίδρασης. Επιπλέον, τέτοιου είδους κοινωνικά πειράματα ευνοούν την επιλεκτική επιλογή συμμετεχόντων ώστε να διασφαλιστεί η πιθανότητα μέγιστης συμμετοχής. Παρά τις αδυναμίες αυτές, πάντως, τα αποτελέσματα των πειραμάτων αυτών έγιναν σε μεγάλο μέρος δεκτά από την επιστημονική κοινότητα ως προς την ορθότητα των συμπερασμάτων τους από άποψη ποιότητας.

Τα τελευταία χρόνια, έκανε την εμφάνιση του ένα σημαντικό πλήθος ιστοσελίδων κοινωνικής δικτύωσης, στο πλαίσιο των οποίων οι σχέσεις μεταξύ των ατόμων μοντελοποιούνται με σαφή τρόπο. Παραδείγματα τέτοιων δικτύων είναι τα Facebook, Twitter και LinkedIn. Εκτός από αυτά, διαδικτυακές εφαρμογές όπως το YouTube και το Flickr που ως στόχο έχουν την ανταλλαγή αρχείων πολυμέσων όπως βίντεο

και φωτογραφιών μεταξύ των χρηστών, εμπίπτουν επίσης στη ευρεία κατηγορία των κοινωνικών δικτύων, καθώς εμπεριέχουν μεγάλο βαθμό αλληλεπίδρασης μεταξύ των χρηστών. Στο πλαίσιο των εφαρμογών αυτών, όπου η αλληλεπίδραση προέρχεται κυρίως από τη βασική λειτουργία της εφαρμογής που είναι η ανταλλαγή περιεχομένου, πολλές από τις βασικές αρχές της κοινωνικής δικτύωσης συνεχίζουν να έχουν ισχύ.

## **2.2 Ερευνητικές περιοχές**

Η εμφάνιση και η υψηλή δημοτικότητα των διαδικτυακών κοινοτήτων, επηρέασε σημαντικά τη μελέτη των κοινωνικών δικτύων, καθώς κατέστησε εφικτή τη δοκιμή, έστω και σε ένα διαδικτυακό περιβάλλον, θεωριών και υποθέσεων όπως αυτές που περιγράφηκαν πιο πάνω και οι οποίες για αρκετά χρόνια αποτελούσαν εικασίες για την επιστήμη της κοινωνιολογίας. Τα διαδικτυακά κοινωνικά δίκτυα είναι εξαιρετικά πλούσια σε περιεχόμενο όπως κείμενο, εικόνες και βίντεο αλλά και επιπλέον πληροφορία σχετικά με αυτό το περιεχόμενο, όπως το πότε και πού δημιουργήθηκε. Ο μεγάλος όγκος πληροφορίας που είναι διαθέσιμος στο περιβάλλον τέτοιων δικτύων επιτρέπει την ανάλυση των κοινωνικών δικτύων με τρόπο επιστημονικά και στατιστικά εμπεριστατωμένο και έχει δώσει μεγάλη ώθηση στους τομείς της ανακάλυψης γνώσης και εξόρυξης δεδομένων. Στην παρούσα υποενότητα θα εξετάσουμε τις σημαντικότερες ερευνητικές περιοχές και προσεγγίσεις που συσχετίζονται με την ανάλυση δεδομένων κοινωνικών δικτύων. Σε πολλές περιπτώσεις, τα συμπεράσματα που εξάγονται για τις διαδικτυακές κοινότητες μπορούν να εφαρμοστούν με ασφάλεια και στην περίπτωση των συμβατικών κοινωνικών δικτύων.

Στο πλαίσιο της ανάλυσης των κοινωνικών δικτύων, τα δεδομένα που αναλύονται ανήκουν σε δύο βασικές κατηγορίες:

- **Δεδομένα γράφου:** Στην ανάλυση γράφου, μελετάμε τη συνδεσμολογία του δικτύου ώστε να προσδιορίσουμε τους σημαντικούς κόμβους, κοινότητες, συνδέσμους και εξελισσόμενες περιοχές του δικτύου.
- **Περιεχόμενο χρηστών:** Στην ανάλυση περιεχομένου, επεξεργαζόμαστε το περιεχόμενο που παράγεται από τους χρήστες του κοινωνικού δικτύου, όπως αναρτήσεις κειμένου, μετα-δεδομένα φωτογραφιών και σχόλια χρηστών.

Μια σημαντική παράμετρος στο πλαίσιο των αλγορίθμων ανάλυσης κοινωνικών δικτύων έγκειται στο διαχωρισμό μεταξύ δυναμικής και στατικής ανάλυσης. Οι αλγόριθμοι στατικής ανάλυσης θεωρούν ότι η δομή ενός κοινωνικού δικτύου αλλάζει αργά στο χρόνο και μελετούν το δίκτυο αναδρομικά χρησιμοποιώντας συγκεκριμένα στιγμιότυπα. Αντιθέτως η δυναμική ανάλυση αφορά δίκτυα η δομή των οποίων αλλάζει συχνά, ως αποτέλεσμα τακτικών αλληλεπιδράσεων μεταξύ των μελών τους και παρουσιάζει αυξημένο ερευνητικό ενδιαφέρον τα τελευταία χρόνια λόγω της δυναμικής φύσης των διαδικτυακών κοινοτήτων. Τα δυναμικά κοινωνικά δίκτυα μπορούν να μοντελοποιηθούν ως δυναμικοί γράφοι οι ακμές των οποίων αλλάζουν συνεχώς στο χρόνο. Η επεξεργασία τέτοιων δυναμικών γράφων εμπεριέχει σημαντικές προκλήσεις εξαιτίας του μεγάλου όγκων συνδέσεων μεταξύ των κόμβων που πρέπει να παρακολουθούνται ταυτόχρονα.

Όπως αναφέρθηκε προηγουμένως, βασικό αντικείμενο της παρούσας διατριβής είναι η μελέτη της διάχυσης της επιρροής σε ένα κοινωνικό δίκτυο. Ένα κοινωνικό δίκτυο είναι μια δομή βασισμένη στις αλληλεπιδράσεις μεταξύ των ατόμων που την αποτελούν και ως τέτοια επιτρέπει τη διάχυση της πληροφορίας μεταξύ των ατόμων. Η ανάλυση της διάχυσης της επιρροής είναι συνεπώς μία από τις ερευνητικές περιοχές που ευνοήθηκε σημαντικά από τη διάδοση των κοινωνικών δικτύων, ειδικά σε σχέση με τον εντοπισμό των πιο επιδραστικών μελών του δικτύου, αυτών δηλαδή

που είναι ικανοί να επηρεάσουν τους υπόλοιπους χρήστες με τις πράξεις τους. Η μελέτη του προβλήματος αυτού αναλύεται εκτενώς στα επόμενα κεφάλαια της διατριβής. Στη συνέχεια θα εξετάσουμε συνοπτικά κάποιες άλλες ερευνητικές περιοχές στις οποίες έχει εφαρμογή η ανάλυση των κοινωνικών δικτύων.

Το πιο διαδεδομένο ίσως πρόβλημα σχετικό με την ανάλυση κοινωνικών δικτύων είναι ο εντοπισμός κοινοτήτων σε ένα δίκτυο, η εύρεση δηλαδή ομάδων ανθρώπων με κοινά χαρακτηριστικά οι οποίες αναφέρονται ως κοινότητες [8] [9] [10]. Το πρόβλημα αυτό εξετάζεται με μελέτη του γράφου τόσο σε στατικό όσο και σε δυναμικό πλαίσιο, ενώ σημαντική βελτίωση στην επίδοση των σχετικών μεθόδων επιτυγχάνεται εμπλουτίζοντας τη γραφική ανάλυση με πληροφορία σχετική με το περιεχόμενο των χρηστών. Στην παρούσα εργασία, όπου εξετάζουμε την επιρροή στο πλαίσιο μιας θεματικής κοινότητας, για τον ορισμό της κοινότητας βασιζόμαστε αποκλειστικά στο περιεχόμενα και συγκεκριμένα στη χρήση ενός κοινού hashtag από τα μέλη της κοινότητας, μια πρακτική αρκετά διαδεδομένη στη βιβλιογραφία.

Οι περισσότερες ερευνητικές προσεγγίσεις, συμπεριλαμβανομένων και των παραπάνω, χρησιμοποιούν τους συνδέσμους μεταξύ των χρηστών για να εξάγουν πληροφορία σχετικά με το κοινωνικό δίκτυο, όπως για παράδειγμα τις υπάρχουσες κοινότητες. Στα περισσότερα δίκτυα, ωστόσο, οι σύνδεσμοι είναι δυναμικοί και μπορούν να αλλάξουν σημαντικά στο χρόνο. Για παράδειγμα, στις διαδικτυακές κοινότητες, νέες φιλίες μεταξύ χρηστών δημιουργούνται σε πολύ συχνή βάση. Ένα ενδιαφέρον ζήτημα που προκύπτει, συνεπώς, είναι η πρόβλεψη μελλοντικών συνδέσμων χρησιμοποιώντας είτε τη δομή του δικτύου είτε τα χαρακτηριστικά των κόμβων. Το πρόβλημα αυτό ονομάζεται πρόβλεψη συνδέσμων και στη βιβλιογραφία έχουν προταθεί αρκετά μοντέλα για την επίλυσή του [11] [12].

Σημαντικό ερευνητικό ενδιαφέρον παρουσιάζει επίσης η ταξινόμηση των κόμβων ενός κοινωνικού δικτύου σε θεματικές κατηγορίες [13] [14]. Ένα συνηθισμένο πρόβλημα για πολλές εφαρμογές είναι, αν θεωρήσουμε κάποιους κόμβους γνωστούς και ταξινομημένους σε θεματικές κατηγορίες, να γίνει κατηγοριοποίηση και των υπόλοιπων κόμβων χρησιμοποιώντας τα χαρακτηριστικά και τη δομή του δικτύου. Ανάλογα με την εφαρμογή, οι κατηγορίες αυτές μπορεί να αφορούν δημογραφικές πληροφορίες, πολιτικές και κοινωνικές πεποιθήσεις, ενδιαφέροντα και ασχολίες ή άλλου είδους χαρακτηριστικά που αφορούν διαφορετικές πλευρές της προσωπικότητας ή της συμπεριφοράς ενός ατόμου. Για παράδειγμα σε μια εφαρμογή μάρκετινγκ, είναι συχνό φαινόμενο να γνωρίζουμε το ενδιαφέρον κάποιων χρηστών-κόμβων για ένα συγκεκριμένο προϊόν και να θέλουμε να εντοπίσουμε και άλλους χρήστες που πιθανόν να ενδιαφέρονται για το ίδιο ή κάποιο παραπλήσιο προϊόν. Καθώς η δημιουργία συνδέσμου μεταξύ δύο χρηστών συχνά προϋποθέτει την ύπαρξη κοινών στοιχείων μεταξύ τους (όπως ενδιαφέροντα ή δημογραφικά χαρακτηριστικά), η γνώση της δομής του δικτύου είναι πολύ σημαντική για τέτοιου είδους προβλέψεις και μπορεί να χρησιμοποιηθεί για την επέκταση της κατηγοριοποίησης από λίγους κόμβους στο υπόλοιπο δίκτυο. Το περιεχόμενο και τα χαρακτηριστικά των κόμβων, μπορούν επίσης να χρησιμοποιηθούν για να βελτιώσουν την ακρίβεια των προβλέψεων.

Η εξόρυξη πληροφορίας από κείμενο είναι ακόμη ένα διαδεδομένο ερευνητικό θέμα που σχετίζεται με τη μελέτη των κοινωνικών δικτύων. Τα κοινωνικά δίκτυα περιέχουν σημαντικό όγκο κειμένου σε διάφορες μορφές, όπως αναρτήσεις, σχόλια, συνδέσμους σε άρθρα, επικεφαλίδες φωτογραφιών. Η ανάλυση του περιεχομένου αυτού μπορεί να ενισχύσει σημαντικά την ποιότητα των συμπερασμάτων που μπορούμε να εξάγουμε για ένα κοινωνικό δίκτυο. Η εξόρυξη κειμένου βρίσκει

εφαρμογή σε τομείς όπως η ανίχνευση γεγονότων που χρησιμοποιεί τεχνικές ανάλυσης κειμένου για να εντοπίσει σημαντικά γεγονότα στο πλαίσιο ενός κοινωνικού δικτύου και η ανάλυση συναισθήματος, που αξιολογεί το συναίσθημα ενός ατόμου ή ομάδας χρηστών ως θετικό ή αρνητικό με βάση τις σχετικές αναρτήσεις του/τους. Η ανάλυση συναισθήματος χρησιμοποιείται και στην παρούσα εργασία καθώς η συναισθηματική πόλωση των χρηστών μιας κοινότητας βάσει του περιεχομένου τους υπολογίζεται και χρησιμοποιείται ως μετρική ομοιότητας μεταξύ της δραστηριότητας δύο κοινοτήτων.

Τέλος, σημαντικός αριθμός ερευνητικών εργασιών επικεντρώνεται στην μοντελοποίηση των κοινωνικών δικτύων, με στόχο τη δημιουργία συνθετικών γράφων εξελισσόμενων στο χρόνο και με ιδιότητες που προσομοιάζουν αυτές των πραγματικών δικτύων [15]. Η δημιουργία μηχανισμών που επιτρέπουν την παραγωγή ρεαλιστικών μοντέλων δικτύων είναι σημαντική για την αξιολόγηση και επαλήθευση υποθέσεων και προσομοιώσεων σε περιπτώσεις όπου οι πραγματικοί γράφοι είναι δύσκολο να βρεθούν ή να χρησιμοποιηθούν. Οι βασικές προκλήσεις που συσχετίζονται με τη μοντελοποίηση της δομής ενός δικτύου αφορούν την επιλογή του κατάλληλου μοντέλου παραγωγής συνθετικών γράφων και την αξιολόγηση του μοντέλου ως προς την ομοιότητα του παραγόμενου γράφου με το πραγματικό δίκτυο.

Η σελίδα αυτή είναι σκόπιμα λευκή.



# 3

## *Επιρροή στα κοινωνικά δίκτυα*

Η διάχυση της επιρροής στα κοινωνικά δίκτυα του πραγματικού κόσμου έχει απασχολήσει πολλές έρευνες τις τελευταίες δεκαετίες – για μια λεπτομερή πληροφόρηση δείτε το [16]. Όπως αναφέρθηκε και προηγουμένως, η ανάλυση της επιρροής είναι μία από τις ερευνητικές περιοχές που ευνοήθηκε σημαντικά από τη διάδοση των διαδικτυακών κοινωνικών δικτύων. Η καταγραφή της δραστηριότητας των χρηστών σε αυτά τα συστήματα επέτρεψε στους ερευνητές να μελετήσουν τη διάχυση της κοινωνικής επιρροής σε μια άνευ προηγουμένου κλίμακα. Στο παρόν κεφάλαιο, αρχικά εξετάζουμε τις βασικές θεωρίες επιρροής για το Twitter, οργανώνοντας τις σε μία δισδιάστατη ταξινόμια. Στη συνέχεια εξετάζουμε το ερευνητικό πλαίσιο γύρω από την αξιολόγηση των θεωριών επιρροής και αναλύουμε περαιτέρω τρία επιπλέον θέματα που σχετίζονται άμεσα με την ανάπτυξη και χρήση των θεωριών επιρροής.

### **3.1 Θεωρίες επιρροής**

Όπως εξηγήσαμε προηγουμένως, η πλατφόρμα του Twitter προσφέρεται για τη μελέτη της διάχυσης της επιρροής σε μεγάλη κλίμακα, με αποτέλεσμα πλήθος θεωριών επιρροής να έχουν αναπτυχθεί ειδικά για το συγκεκριμένο δίκτυο. Στην παρούσα ενότητα εξετάζουμε τις πιο σημαντικές, ώστε να επιλέξουμε αυτές που θα

συμπεριλάβουμε στην ανάλυσή μας. Για να διευκολύνουμε την κατανόηση, εισάγουμε μια δισδιάστατη ταξινόμια που τις κατηγοριοποιεί αναφορικά με το εύρος και τον τύπο μετρικής που χρησιμοποιούν για την αξιολόγηση της επιρροής.

Το εύρος μιας θεωρίας επιρροής εξετάζει την περιοχή εφαρμογής της, δηλαδή το κομμάτι του κοινωνικού δικτύου από το οποίο προέρχονται τόσο οι υποψήφιοι επιδραστικοί χρήστες όσο και οι ενδείξεις σχετικά με τη δραστηριότητα τους. Αυτή η διάσταση διαχωρίζει τις θεωρίες επιρροής στις ακόλουθες κατηγορίες:

- Οι θεωρίες καθολικής επιρροής λαμβάνουν υπόψη τους τη δραστηριότητα ολόκληρου του κοινωνικού δικτύου με στόχο να εντοπίσουν τους συνολικά πιο επιδραστικούς χρήστες. Για παράδειγμα, θεωρείστε μια θεωρία που αναγνωρίζει ως επιδραστικούς τους χρήστες με το μεγαλύτερο αριθμό ακόλουθων σε όλο το δίκτυο του Twitter (καθολικός βαθμός εισόδου).
- Οι θεωρίες τοπικής επιρροής λαμβάνουν υπόψη αποκλειστικά τη δραστηριότητα μιας θεματικής κοινότητας (δηλαδή μιας κοινότητας χρηστών με ενδιαφέρον για ένα συγκεκριμένο θέμα) με το σκοπό να εντοπίσουν τα πιο επιδραστικά μέλη της. Για παράδειγμα, θεωρείστε μια θεωρία που επιλέγει ως επιδραστικούς τους χρήστες με το μεγαλύτερο αριθμό ακολούθων μέσα στη θεματική κοινότητα (τοπικός βαθμός εισόδου).
- Οι θεωρίες μεικτής επιρροής αποτελούν υβρίδιο των παραπάνω καθώς λαμβάνουν υπόψη ενδείξεις από ολόκληρο το κοινωνικό δίκτυο για να εντοπίσουν τους επιδραστικούς χρήστες μιας συγκεκριμένης κοινότητας. Για παράδειγμα, θεωρείστε μια θεωρία που ορίζει ως επιδραστικά τα μέλη μιας θεματικής κοινότητας που έχουν το μεγαλύτερο καθολικό βαθμό εισόδου.

Αναφορικά με το εύρος, η ανάλυσή μας επικεντρώνεται στις θεωρίες τοπικής επιρροής του Twitter. Αυτές οι θεωρίες έχουν καλύτερη επίδοση στις προβλέψεις

τους και τη λειτουργία τους, καθώς λαμβάνουν υπόψη ενδείξεις από τη δραστηριότητα και την τοπολογία της συγκεκριμένης κοινότητας.

Η μετρική μιας θεωρίας επιρροής εξετάζει το είδος πληροφορίας που λαμβάνει υπόψη για να εκτιμήσει την επιρροή που ασκεί ένας συγκεκριμένος χρήστης. Η διάσταση αυτή διαχωρίζει τις θεωρίες επιρροής στις ακόλουθες κατηγορίες:

- Οι θεωρίες επιρροής βάσει γράφου υπολογίζουν την επιρροή βάσει της θέσης ενός κόμβου στον κοινωνικό γράφο. Στην κατηγορία αυτή ανήκουν γραφικά κριτήρια όπως ο βαθμός εισόδου και η κεντρικότητα του κόμβου.
- Οι θεωρίες επιρροής βάσει περιεχομένου εκτιμούν την επιρροή ενός χρήστη με αποκλειστικό κριτήριο το περιεχόμενο που παράγει. Δεδομένου ότι δεν υπάρχει αντικειμενικός τρόπος άμεσης αξιολόγησης του περιεχομένου χρηστών, στην πράξη χρησιμοποιούνται έμμεσοι τρόποι υπολογισμού. Στο πλαίσιο του Twitter, κριτήρια επιρροής αυτής της κατηγορίας είναι ο αριθμός των αναδημοσιεύσεων (retweets) και ο αριθμός των αναφορών (mentions) που έχει λάβει ο συγκεκριμένος χρήστης. Η έμμεση παραδοχή στην περίπτωση αυτή είναι ότι όσο υψηλότερη είναι η ποιότητα των μηνυμάτων που δημοσιεύει ένας χρήστης, τόσο συχνότερα αναπαράγεται το περιεχόμενό του ή αναφέρονται οι άλλοι χρήστες σε αυτόν.
- Οι ολιστικές θεωρίες επιρροής λειτουργούν υβριδικά, θεωρώντας τόσο το περιεχόμενο που παράγουν οι χρήστες όσο και τη θέση τους στο γράφο.

Σημειώστε ότι οι θεωρίες που ανήκουν στις δύο πρώτες κατηγορίες μπορούν να διαχωριστούν περαιτέρω με βάση την πολυπλοκότητά τους σε ατομικές και σύνθετες θεωρίες επιρροής: οι πρώτες λαμβάνουν υπόψη μια μοναδική ένδειξη (μία μετρική) ενώ οι δεύτερες εξετάζουν ένα συνδυασμό πολλαπλών μετρικών. Οι ολιστικές θεωρίες επιρροής είναι εξορισμού σύνθετες. Η ανάλυση μας επικεντρώνεται σε

ατομικές μετρικές βάσει περιεχομένου και βάσει γράφου, εφόσον οι ενδείξεις που λαμβάνουν υπόψη προέρχονται από την πληροφορία που περιέχεται σε συγκεκριμένες θεματικές κατηγορίες.

Τονίζουμε σε αυτό το σημείο ότι ο όρος θεωρία επιρροής χρησιμοποιείται καταχρηστικά σε αυτήν την εργασία, καθώς δεν αναφέρεται σε πραγματικές θεωρίες που αναλύουν ενδελεχώς την συμπεριφορά των χρηστών ενός κοινωνικού δικτύου, εξηγώντας γιατί κάποιοι από αυτούς μιμούνται τη συμπεριφορά άλλων. Αντί για αυτό, ο όρος αναφέρεται σε μεθόδους ταξινόμησης επιρροής, οι οποίες αναθέτουν σε κάθε χρήστη ένα βαθμό ανάλογο της εκτίμησης της επιρροής που ασκεί στα άλλα μέλη της κοινότητας. Κατά σύμβαση, οι μέθοδοι αυτές αναφέρονται ως θεωρίες επιρροής στη σχετική βιβλιογραφία και η σύμβαση αυτή ακολουθείται και εδώ.

Στον Πίνακα 3-1 απεικονίζεται η δισδιάστατη ταξινόμια που προτείνουμε, σε συνδυασμό με τις θεωρίες επιρροής που έχουμε κατηγοριοποιήσει βάσει αυτής. Παρατηρούμε ότι κάποιοι τύποι θεωριών επιρροής δεν έχουν ακόμα διερευνηθεί στη βιβλιογραφία, τουλάχιστον όσον αφορά το Twitter. Για παράδειγμα, δεν υπάρχει θεωρία μεικτής επιρροής που να εξετάζει το περιεχόμενο ως κριτήριο επιρροής. Αντίθετα, οι περισσότερες θεωρίες για το Twitter χρησιμοποιούν ολιστικές μετρικές για την αξιολόγηση της επιρροής και έχουν τοπικό εύρος. Στη συνέχεια εξετάζουμε την εσωτερική λειτουργία κάθε θεωρίας με σκοπό να αποφασίσουμε ποιες είναι κατάλληλες για την ανάλυσή μας.

		Μετρική				Ολιστική
		Βάσει Περιεχομένου		Βάσει Γράφου		
		ατομική	σύνθετη	ατομική	σύνθετη	
Ε ύ ρ ο ς	Καθολική	-	-	-	-	Bakshy et al. [1] Petrovic et al. [17]
	Τοπική	Cha et al. [4]	-	-	Li et al. [19]	Liu et al. [20] Purohit et al. [18]
	Μεικτή	-	-	Cha et al. [4]	-	Weng et al. [5]

**Πίνακας 3-1: Δισδιάστατη ταξινόμια των θεωριών επιρροής του Twitter**

Στο [4], οι συγγραφείς εξετάζουν ένα σύνολο θεωριών ατομικής επιρροής που χρησιμοποιούν την τοπική και καθολική δραστηριότητα των χρηστών. Η εργασία τους εξετάζει την επιρροή των συνηθισμένων χρηστών σε αντίθεση με αυτήν των διακεκριμένων χρηστών του Twitter όπως οι διασημότητες. Οι θεωρίες τοπικής επιρροής που εξετάζουν είναι οι Θεωρίες Βαθμού Εισόδου, Αναφορών και Αναδημοσιεύσεων που παρουσιάζονται αναλυτικά στην Ενότητα 4.3.

Στο [1], οι συγγραφείς θεωρούν την επιρροή ως την ικανότητα ενός χρήστη να δημοσιεύει URL που διαχέονται μέσω των ακολούθων του σε ολόκληρο το γράφο του Twitter. Αυτή είναι μια καθολική θεωρία που χρησιμοποιεί ένα δενδρικό μοντέλο παλινδρόμησης για την πρόβλεψη της επιρροής ενός χρήστη με βάση το μέσο μέγεθος των αλληλουχιών αναμετάδοσης που προκαλεί. Ως μεταβλητές του μοντέλου, οι συγγραφείς χρησιμοποιούν ένα συνδυασμό μετρικών περιεχομένου και γράφου: τον αριθμό των ακολούθων, τον αριθμό των φίλων, τον αριθμό των tweets την ημερομηνία εγγραφής και τον αριθμό αναδημοσιεύσεων τόσο από τους άμεσους

γείτονες (προηγούμενη τοπική επιρροή) όσο και από οποιονδήποτε χρήστη (προηγούμενη καθολική επιρροή). Τα αποτελέσματα των πειραματικών μελετών τους δείχνουν ότι η προηγούμενη τοπική επιρροή και ο αριθμός των ακολούθων είναι οι πιο αξιόπιστοι παράγοντες για τον καθορισμό της επιρροής. Οι συγγραφείς συμπεραίνουν επίσης ότι οι καμπάνιες προώθησης θα πρέπει να θεωρούν κάθε άτομο ως επιδραστικό και να στοχεύουν πολλούς συνηθισμένους χρήστες ώστε να επιτύχουν υψηλά ποσοστά αποδοχής. Την ιδέα αυτή την εξετάζουμε στο πλαίσιο των θεματικών κοινοτήτων μέσω της Θεωρίας Τυχαίας Επιρροής που παρουσιάζεται στην Ενότητα 4.3. Σημειώνουμε, επίσης, ότι το κριτήριο επιρροής που βασίζεται στο μέγεθος των αλληλουχιών αναδημοσίευσης είναι παρόμοιο με τη Θεωρία Αναδημοσιεύσεων όταν περιορίζεται στα όρια μιας θεματικής κοινότητας.

Στο [17], οι συγγραφείς χρησιμοποιούν ένα πλήθος μετρικών περιεχομένου και γράφου προκειμένου να καθορίσουν την πιθανότητα αναδημοσίευσης ενός tweet. Ανακάλυψαν ότι τα καθολικά κοινωνικά χαρακτηριστικά του συγγραφέα (και ειδικά ο βαθμός εισόδου) επιτυγχάνουν μεγαλύτερη ακρίβεια πρόβλεψης από τα χαρακτηριστικά που αφορούν το ίδιο το tweet. Δυστυχώς τα περισσότερα από τα χαρακτηριστικά είναι καθολικά και άρα ασύμβατα με την ανάλυση μας. Μόνο ο βαθμός εισόδου χρησιμοποιείται από τη Θεωρία Βαθμού Εισόδου, προσαρμοσμένος ωστόσο στο τοπικό πλαίσιο των θεματικών κοινοτήτων.

Στο [18], οι συγγραφείς εξετάζουν την τοπική επιρροή στο πλαίσιο εταιρικών σελίδων, λαμβάνοντας υπόψη στοιχεία σχετικά με τις αλληλεπιδράσεις των χρηστών και τα προφίλ τους. Στην πρώτη κατηγορία ανήκουν ο αριθμός των αναδημοσιεύσεων, ο αριθμός των απαντήσεων και ο αριθμός των αναφορών, ενώ στη δεύτερη ο αριθμός των ακολούθων και ο αριθμός των tweets πάνω στο θέμα. Όλες αυτές οι μετρικές συμπεριλαμβάνονται και στην ανάλυση μας: για παράδειγμα, ο

αριθμός των tweets πάνω στο θέμα είναι η βάση της Θεωρίας Δημοσιεύσεων και ο αριθμός των απαντήσεων σχετίζεται άμεσα με τη Θεωρία Αναφορών. Επιπροσθέτως, οι συγγραφείς εφαρμόζουν διαδομένους αλγόριθμους ανάλυσης συνδέσμων, όπως τα HITS και PageRank, στον κοινωνικό γράφο του Twitter αλλά και στο έμμεσο δίκτυο που σχηματίζεται από τα retweets. Τα ευρήματά τους δείχνουν ότι οι αλγόριθμοι αυτοί επιτυγχάνουν καλύτερη επίδοση από τα χαρακτηριστικά προφίλ του χρήστη, ειδικά για το πραγματικό δίκτυο του Twitter. Ωστόσο οι καθολικές ενδείξεις που χρησιμοποιούν είναι ασύμβατες με τη θεματο-κεντρική ανάλυσή μας. Κυρίως, όμως, οι γραφικές ενδείξεις είναι ακατάλληλες στο πλαίσιο των θεματικών κοινοτήτων του Twitter: τα μέλη τους συμμετέχουν σε αυτές μέσω του περιεχομένου που δημοσιεύουν και, συνεπώς, δεν είναι απαραίτητα συνδεδεμένα μεταξύ τους.

Μια γραφική προσέγγιση επιχειρείται, επίσης, από τους Li et al. [19]. Ο αλγόριθμος τους ποσοτικοποιεί την επιρροή κάθε μέλους της κοινότητας χρησιμοποιώντας τον υποκείμενο προσημασμένο γράφο, που περιέχει τις θετικές και αρνητικές σχέσεις μεταξύ όλων των ατόμων. Δεδομένου ότι το Twitter δεν είναι ένα ρητά προσημασμένο δίκτυο, δε μπορούμε να συμπεριλάβουμε την προσέγγιση αυτή στην ανάλυσή μας.

Στο [20], οι συγγραφείς εισάγουν ένα πιθανοτικό μοντέλο για την εξόρυξη άμεσης και έμμεσης επιρροής μεταξύ των κόμβων ετερογενών δικτύων. Το μοντέλο τους συνδυάζει μετρικές περιεχομένου και γράφου ώστε να προβλέψει την πιθανότητα ένας χρήστης να αναδημοσιεύσει τη δημοσίευση ενός φίλου του. Το μοντέλο αξιολογήθηκε πειραματικά πάνω στο Twitter και σε δύο ακόμα δίκτυα, με τα αποτελέσματα να συνιστούν σημαντική βελτίωση στην ακρίβεια πρόβλεψης. Η μέθοδος αυτή, ωστόσο, λειτουργεί στο επίπεδο ατομικών συνδέσμων μεταξύ χρηστών και συνεπώς ο στόχος της διαφέρει από το δικό μας που είναι η αξιολόγηση της

αποτελεσματικότητας των θεωριών τοπικής επιρροής στην ανίχνευση υψηλά επιδραστικών μελών της κοινότητας.

Τέλος οι Weng et al. [5] εισήγαγαν μια καινοτόμο θεωρία μεικτής επιρροής που βασίζεται σε μια τροποποιημένη εκδοχή του PageRank της Google που ονομάζεται TwitterRank. Για να μετρήσει την τοπική επιρροή των ατόμων, υπολογίζει την θεματική ομοιότητα μεταξύ δύο χρηστών με βάση την καθολική πληροφορία για όλες τις άλλες θεματικές κατηγορίες στις οποίες συμμετέχουν. Μαζί με ορισμένες μετρικές περιεχομένου, λαμβάνει υπόψη και τη δομή των συνδέσμων ολόκληρου του κοινωνικού γράφου του Twitter. Όπως εξηγήσαμε προηγουμένως, ωστόσο, τέτοιου είδους μετρικές είναι ασύμβατες με την ανάλυσή μας, καθώς αυτή συμπεριλαμβάνει αποκλειστικά θεωρίες επιρροής βασισμένες σε πληροφορία που περιέχεται στα όρια μιας συγκεκριμένης θεματικής κοινότητας.

### **3.2 Αξιολόγηση θεωριών τοπικής επιρροής**

Σκοπός της παρούσας διατριβής είναι ο ορισμός ενός πλαισίου που επιτρέπει την αξιολόγηση θεωριών τοπικής επιρροής σε κοινωνικά δίκτυα σε μεγάλη κλίμακα. Από όσο γνωρίζουμε, καμία από τις ως τώρα ερευνητικές δουλειές δεν έχει προτείνει μια επίσημη μεθοδολογία για τη σύγκριση της επίδοσης των θεωριών επιρροής σε μεγάλη κλίμακα. Η μόνη εργασία που είναι σχετική με το πλαίσιο αξιολόγησης που προτείνουμε είναι ένα μοντέλο που προτείνεται στο [21] για την αξιολόγηση των θεωριών καθολικής επιρροής αναφορικά με την ανοχή τους σε κακόβουλους χρήστες. Βάσει αυτού του μοντέλου, μια επιτυχημένη θεωρία θα πρέπει να τοποθετεί σε χαμηλές θέσεις κατάταξης λογαριασμούς spamming ή προώθησης ενώ οι υψηλές θέσεις κατάταξης θα πρέπει να ανατίθενται σε έγκυρους πιστοποιημένους λογαριασμούς. Ο συγγραφέας πραγματοποίησε μια συγκριτική ανάλυση αρκετών καθολικών, γραφικών θεωριών επιρροής (δείτε την Ενότητα 3 για τη σχετική



κατηγοριοποίηση) μεταξύ των οποίων οι PageRank, HITS, TwitterRank [5] και κάποιες μετατροπές τους. Τα αποτελέσματα δείχνουν ότι οι περισσότεροι αλγόριθμοι επιτυγχάνουν παρόμοια επίδοση από άποψη ανοχής σε κακόβουλους χρήστες, με την εξαίρεση του TwitterRank που σημείωσε σημαντικά χαμηλότερη επίδοση από τους άλλους.

### **3.3 Πίεση από τον κοινωνικό περίγυρο**

Ένα σημαντικό πρόβλημα στην ανάλυση κοινωνικών δικτύων είναι η πρόβλεψη της επίδρασης της κοινωνικής επιρροής στην συμπεριφορά συνδεδεμένων χρηστών. Ο στόχος είναι η πρόβλεψη των δράσεων συγκεκριμένων χρηστών με τη χρήση μιας θεωρίας επιρροής και ιστορικών δεδομένων σχετικά με τη δική τους δραστηριότητα και τη δραστηριότητα των φίλων τους.

Στο [22], οι συγγραφείς προτείνουν ορισμένα στατικά και χρονικά εξαρτώμενα μοντέλα που υπολογίζουν την πιθανότητα ένας χρήστης να μιμηθεί τους φίλους του στην συμμετοχή ενός συγκεκριμένου γκρουπ καθώς και τη χρονική στιγμή που αυτό θα συμβεί (μέσα σε στενά χρονικά όρια). Τα μοντέλα τους λαμβάνουν υπόψη και την τοπολογία του κοινωνικού δικτύου και το ιστορικό της δραστηριότητας του χρήστη. Δοκιμάστηκαν σε πραγματικά δεδομένα από το Flickr και επέτυχαν μεγάλη ακρίβεια.

Στο [23] οι συγγραφείς προτείνουν ένα Noise Tolerant Time-varying Factor Graph Model (NTT-FGM) για τη μοντελοποίηση και την πρόβλεψη των κοινωνικών δραστηριοτήτων. Για μεγαλύτερη ακρίβεια πρόβλεψης, το NTT-FGM συμπεριλαμβάνει τα χαρακτηριστικά του χρήστη σε συνδυασμό με την δομή του κοινωνικού δικτύου και το ιστορικό των δράσεων του χρήστη. Οι συγγραφείς εισάγουν επίσης την έννοια της λανθάνουσας κατάστασης, η οποία προσδιορίζει την πιθανότητα ένας χρήστης να εκτελέσει μια πράξη μια συγκεκριμένη χρονική στιγμή.

Το μοντέλο αξιολογήθηκε πάνω σε πραγματικά δεδομένα από τρία ετερογενή κοινωνικά δίκτυα και βρέθηκε να έχει σταθερά καλύτερη επίδοση από τις μεθόδους βάσης.

Στο [24], οι συγγραφείς αξιολογούν δυο ορισμούς της επιρροής σε σχέση με την πιθανότητα υιοθέτησης μιας συμπεριφοράς ως αποτέλεσμα πίεσης από τον κοινωνικό περίγυρο (peer pressure). Ο πρώτος ορισμός βασίζεται σε στιγμιαίες παρατηρήσεις του δικτύου τη χρονική στιγμή αμέσως πριν το άτομο αποφασίσει αν θα υιοθετήσει μια συμπεριφορά, ενώ ο δεύτερος ορισμός βασίζεται σε λεπτομερείς χρονικές δυναμικές. Η σχέση μεταξύ αυτών των δύο τύπων επιρροής εξετάστηκε με χρήση δεδομένων από τη Wikipedia, με τα αποτελέσματα να επιβεβαιώνουν την υπόθεση ύπαρξης άμεσης σύνδεσης μεταξύ τους.

### **3.4 Εφαρμογές των θεωριών επιρροής**

Η έννοια της κοινωνικής επιρροής αποτελεί τη βάση πολλών επιχειρησιακών εφαρμογών. Στόχος τους είναι η εφαρμογή των θεωριών επιρροής σε εκστρατείες διαφήμισης και προώθησης ώστε να μεγιστοποιηθεί η αποτελεσματικότητά τους επιτυγχάνοντας παρόμοια επίπεδα υιοθέτησης με μικρότερο κόστος. Το πρόβλημα αυτό ορίζεται ως πρόβλημα μεγιστοποίησης επιρροής και περιλαμβάνει τον εντοπισμό ενός μικρού υποσυνόλου χρηστών που θα μπορούσαν να μεγιστοποιήσουν τη διάχυση της επιρροής στο κοινωνικό δίκτυο [25]. Στο πλαίσιο του viral marketing, το πρόβλημα αυτό μεταφράζεται στη μεγιστοποίηση της κατανάλωσης ενός προϊόντος στοχεύοντας σε ένα μικρό αριθμό υψηλά επιδραστικών χρηστών και δίνοντάς τους κίνητρα για να το υιοθετήσουν. Ο μικρός αριθμός χρηστών κρατάει το κόστος χαμηλό, ενώ η υψηλή επιδραστική τους ικανότητα εξασφαλίζει μεγάλη διάχυση του προϊόντος με προώθηση από στόμα σε στόμα.

Στο [26], οι συγγραφείς θεωρούν την αγορά ως ένα κοινωνικό δίκτυο που μπορεί να μοντελοποιηθεί ως τυχαίο πεδίο Markov. Εισάγουν επίσης την έννοια της αξίας δικτύου πελάτη που υπολογίζει τα αναμενόμενα κέρδη από πωλήσεις σε άλλα άτομα που επηρεάστηκαν από ένα συγκεκριμένο πελάτη. Στο [27], οι συγγραφείς θεωρούν το πρόβλημα της επιλογής των πιο επιδραστικών κόμβων και το μελετούν στο πλαίσιο των πιο διαδεδομένων μοντέλων στην ανάλυση κοινωνικών δικτύων. Αποδεικνύουν ότι, υπό αυτές τις συνθήκες, αποτελεί ένα NP-hard πρόβλημα βελτιστοποίησης και παρέχουν διασφαλίσεις προσέγγισης για αποδοτικούς αλγόριθμους. Η εργασία αυτή αποτελεί τη βάση πολλών ακόμα πρόσφατων μελετών πάνω στο θέμα.

Από μια άλλη ερευνητική σκοπιά, στο [4] οι συγγραφείς εξετάζουν τη δυναμική της επιρροής στο Twitter και ανακαλύπτουν ότι οι πιο επιδραστικοί χρήστες λαμβάνουν δυσανάλογα περισσότερες αναφορές από τους συνηθισμένους χρήστες. Συμπεραίνουν, συνεπώς, ότι η διάχυση της πληροφορίας μπορεί να μεγιστοποιηθεί στοχεύοντας ένα μικρό αριθμό ηγετικών φυσιογνωμιών. Τα ευρήματα αυτά είναι συμβατά με την παραδοσιακή οπτική στη διάχυση επιρροής [16], αλλά έρχονται σε αντίθεση με τις μοντέρνες θεωρίες που δίνουν έμφαση στο ρόλο των διαπροσωπικών σχέσεων ανάμεσα σε συνηθισμένους χρήστες [1]. Οι τελευταίες ισχυρίζονται ουσιαστικά ότι οι καμπάνιες προώθησης, για να είναι επιτυχείς, θα πρέπει να στοχεύουν σε ένα μεγάλο αριθμό συνηθισμένων χρηστών.

Τέλος, σχετική με τη μεγιστοποίηση της επιρροής είναι και η διάχυση της καινοτομίας. Στο [28], οι συγγραφείς προτείνουν μια νέα μέθοδο αντιμετώπισης αυτού του προβλήματος μέσω της θεωρίας διάχυσης της θερμότητας στη φυσική. Αναπτύσσουν τρία μοντέλα διάχυσης σε συνδυασμό με εξειδικευμένους αλγόριθμους για την επιλογή των καταλληλότερων ατόμων για την αποστολή προωθητικών

δειγμάτων. Τα αποτελέσματα επιβεβαιώνουν ότι το πλαίσιο αυτό έχει καλύτερη επίδοση σε σχέση με προηγούμενες εργασίες και είναι ιδιαίτερα αποτελεσματικό για την παρακολούθηση της διάχυσης αρνητικής πληροφορίας. Στο [29], οι συγγραφείς εξετάζουν τη μακροσκοπική διάχυση σε σχέση με την υιοθέτηση συμπεριφοράς, χρησιμοποιώντας χαρακτηριστικά τοπολογίας όπως η διάχυση βαθμού, καθώς και την επιρροή από στόμα σε στόμα που προέρχεται από γειτονικούς κόμβους.

### **3.5 Ομοφιλία ή κοινωνική επιρροή;**

Οι προαναφερθείσες εργασίες θεωρούν την κοινωνική επιρροή ως το μοναδικό παράγοντα που επηρεάζει τα άτομα να αλλάξουν τη συμπεριφορά τους ώστε να συμμορφωθούν με τη δραστηριότητα των γειτόνων τους. Η κοινωνική επιρροή, όμως, δεν είναι η μοναδική αιτία κοινωνικής συσχέτισης. Το φαινόμενο αυτό μπορεί να είναι αποτέλεσμα ομοφιλίας, ενός φαινομένου που εξετάζεται αναλυτικά στο [30]. Η εργασία αυτή ισχυρίζεται ότι κοινωνικοί δεσμοί σχηματίζονται με υψηλή πιθανότητα μεταξύ ατόμων με κοινά αρκετά από τα θεμελιώδη στοιχεία της προσωπικότητας τους, όπως πεποιθήσεις, ενδιαφέροντα και δημογραφικά χαρακτηριστικά (π.χ. ηλικία, φυλή και τοποθεσία). Σημαντικός είναι, συνεπώς, ο διαχωρισμός της κοινωνικής επιρροής από την ομοφιλία και άλλες αδιόρατες συγγεόμενες μεταβλητές που μπορούν να προκαλέσουν στατιστική συσχέτιση μεταξύ φίλων σε ένα κοινωνικό δίκτυο. Το έργο αυτό απασχολεί πολλές ερευνητικές εργασίες [31], [32] [33].

Άλλες εργασίες εξετάζουν τη δυναμική της κοινωνικής επιρροής σε σχέση με την ομοφιλία. Στο [34], οι συγγραφείς προτείνουν ένα μοντέλο που συνδυάζει αυτούς τους δύο παράγοντες κοινωνικής συσχέτισης και εξηγούν πώς η ισορροπία τους μπορεί να ελεγχθεί μέσω μιας μοναδικής παραμέτρου. Οι συγγραφείς μελετούν, επίσης, πώς η δυναμική που σχετίζεται από κάθε παράγοντα επηρεάζει την κατάτμηση του δικτύου και συμπεραίνουν ότι η ομοφιλία έχει ως αποτέλεσμα ένα

μεγάλο αριθμό μικρών συστάδων ενώ η κοινωνική επιρροή δημιουργεί μεγάλες και συνεκτικές συστάδες. Στο [35] οι συγγραφείς εξετάζουν πώς η ομοφιλία και η επιρροή επηρεάζουν τη μοντελοποίηση δυναμικών δικτύων και εισάγουν τυποποιημένους ορισμούς και για τους δύο παράγοντες. Αναπτύσσουν επίσης μετρικές για την εκτίμηση της συσχέτισης μεταξύ συνδέσμων και χαρακτηριστικών πάνω από διαφορετικές στρατηγικές χρήσης ιστορικών δεδομένων του δικτύου. Τα αποτελέσματά τους δείχνουν ότι η σημασία των επιμέρους χαρακτηριστικών στη δημιουργία συνδέσμων αλλάζουν με το χρόνο.

Η σελίδα αυτή είναι σκόπιμα λευκή.

# 4

## *Ορισμός Πλαισίου Αξιολόγησης*

Στο κεφάλαιο αυτό ορίζουμε το προτεινόμενο πλαίσιο αξιολόγησης καθώς και τις βασικές έννοιες που σχετίζονται μαζί του. Αρχικά εξετάζουμε τα βασικά χαρακτηριστικά του Twitter που αποτελούν βάση των τοπικών θεωριών επιρροής, στη συνέχεια διατυπώνουμε τις θεωρίες που συμμετέχουν στην αξιολόγηση και τέλος ορίζουμε τις πέντε συνθήκες που αποτελούν το πλαίσιο μας.

### **4.1 Παρουσίαση του Twitter**

Το Twitter αποτελεί μία από τις πιο δημοφιλείς υπηρεσίες micro-blogging, με μια βάση περισσότερων από 500 εκατομμύρια χρηστών που δημοσιεύουν πάνω από 340 εκατομμύρια σύντομα μηνύματα τη μέρα. Είναι επίσης ένα από τα πιο δημοφιλή κοινωνικά δίκτυα του διαδικτύου μεταξύ των ερευνητών που μελετούν τη διάχυση της επιρροής [1] [4] [5]. Σημαντικό μέρος της δημοφιλίας του Twitter οφείλεται στα θεμελιώδη χαρακτηριστικά του:

1. Στους χρήστες επιτρέπεται η δημοσίευση μόνο σύντομων μηνυμάτων έως 140 χαρακτήρων που ονομάζονται *tweets*. Αυτό ενεργοποιεί τους χρήστες ώστε να επιστρατεύσουν το ταλέντο τους στη δημιουργία πρωτότυπων,

αυτόνομων και έξυπνων μηνυμάτων που απαιτούν από τους αναγνώστες ελάχιστη προσοχή και χρόνο. Συνεπώς τα tweets μπορούν εύκολα να απομνημονευτούν και να αναπαραχθούν από άλλους χρήστες, όπως συμβαίνει και με τα προωθητικά σλόγκαν.

2. Οι λογαριασμοί στο Twitter είναι από προεπιλογή δημόσιοι, ενθαρρύνοντας έτσι την αλληλεπίδραση μεταξύ των χρηστών. Οποιοσδήποτε χρήστης μπορεί ελεύθερα να ακολουθήσει (να εγγραφεί δηλαδή) τους λογαριασμούς άλλων ώστε να λαμβάνει τα πιο πρόσφατα tweets τους. Ακολουθώντας ένα συγκεκριμένο χρήστη  $u$  ο εγγραφόμενος χρήστης άμεσα δηλώνει ότι ο  $u$  τον ενδιαφέρει ιδιαίτερα είτε λόγω κοινών χαρακτηριστικών (π.χ. ένα χόμπι) είτε λόγω ποιότητας περιεχομένου (π.χ. υπηρεσίες ειδήσεων)

Σημαντικά επίσης για την επιτυχία του Twitter ήταν πρότυπα χρήσης που καθιερώθηκαν από τους χρήστες:

1. Τα tweets μπορούν εύκολα να ενταχθούν σε θεματικές κατηγορίες που έχουν οριστεί ελεύθερα από άλλους χρήστες. Αυτό γίνεται με την προσθήκη – συνήθως στο τέλος του μηνύματος - ενός ή περισσότερων *hashtags*. Αυτά είναι ειδικές επισημάνσεις που αποτελούνται από το χαρακτήρα #, ακολουθούμενο από μια ή περισσότερες λέξεις ή αλφαριθμητικά συνενωμένα (π.χ., *#twitter*). Τα *hashtags* μπορούν συνεπώς να χρησιμοποιηθούν για τον εντοπισμό ομάδων ανθρώπων που ενδιαφέρονται για το ίδιο θέμα.
2. Το Twitter μπορεί να χρησιμοποιηθεί ως πλατφόρμα συζήτησης μεταξύ των μελών του. Προσθέτοντας την επισήμανση “@user” στα tweets του, ένας χρήστης  $u$  μπορεί άμεσα να απευθυνθεί στο χρήστη *user* ο οποίος στη



συνέχεια μπορεί να απαντήσει με τον ίδιο τρόπο (@u). Αυτός ο τρόπος αλληλεπίδρασης ονομάζεται αναφορά (*mention*) και επιτρέπει τον εντοπισμό χρηστών που εμπλέκονται σε διμερείς συζητήσεις.

3. Οι χρήστες μπορούν να μοιραστούν με τους ακολούθους τους tweets που βρίσκουν ενδιαφέροντα αλλά έχουν γραφτεί από άλλους χρήστες. Αυτό γίνεται με την αναδημοσίευση του αρχικού tweet με την ειδική επισήμανση “*RT @user*” ώστε να δείξουν αναγνώριση στον αρχικό χρήστη – στην περίπτωσή μας τον *user*. Η πρακτική αυτή ονομάζεται *retweet* και επιτρέπει την παρακολούθηση της διάχυσης ενός συγκεκριμένου tweet και την εκτίμηση της επίδρασης του. Όσο μεγαλύτερη είναι η αλληλουχία από *retweets* που προκαλεί ένα μήνυμα, τόσο πιο επιδραστικό είναι [1].

Αυτές οι πρακτικές του Twitter αποτελούν βασικό στοιχείο της ανάλυσης των προτύπων αλληλεπίδρασης μεταξύ των μελών συγκεκριμένων κοινοτήτων. Χάριν συντομίας, ορίζουμε ως *επισημασμένα* τα μηνύματα τα οποία περιέχουν μια αναφορά ή αναδημοσίευση. Ορίζουμε επίσης ως *πολωμένα* tweets τα μηνύματα τα οποία εκφράζουν είτε θετικό είτε αρνητικό συναίσθημα περιλαμβάνοντας το αντίστοιχο εικονίδιο (*emoticon*). Ένα θετικό tweet περιέχει μια από τις εξής ακολουθίες χαρακτήρων “:)", “:-)”, “:)”, “:D” ή “=)”, ενώ ένα αρνητικό μια από τις ακόλουθες: “:(”, “:-(", ή “:(” [36]. Τα tweets που περιέχουν και αρνητικά και θετικά *emoticons* θεωρούνται διφορούμενα και αποκλείονται από την ανάλυσή μας.

Σε αυτό το σημείο αξίζει να αναφέρουμε ότι στη βιβλιογραφία έχουν προταθεί πιο προχωρημένες μέθοδοι για τον καλύτερο αυτόματο προσδιορισμό του συναισθήματος ενός μηνύματος [37]. Η εφαρμογή, ωστόσο, ενός πιο εξελιγμένου σχήματος ταξινόμησης για την συναισθηματική ανάλυση είναι εκτός του πεδίου της παρούσας

εργασίας. Ο λόγος γι' αυτό είναι ότι τέτοιου είδους τεχνικές είναι συνήθως εξαρτημένες από συγκεκριμένη γλώσσα και εφαρμογή, καθώς πρέπει να ξεπεράσουν τις έμφυτες προκλήσεις του περιεχομένου των χρηστών του Twitter [38]:

- Υψηλά επίπεδα θορύβου με τη μορφή ορθογραφικών σφαλμάτων και ελλιπούς ή λανθασμένης πληροφορίας
- Αραιή πληροφορία, καθώς οι περιορισμοί στο μέγεθος ελαχιστοποιούν την πληροφορία που περιέχουν τα επιμέρους μηνύματα
- Πολυγλωσσία, καθώς ένα μόνο μήνυμα μπορεί να περιέχει λέξεις σε πολλαπλές γλώσσες. Ακόμα κι όταν μια μόνο γλώσσα χρησιμοποιείται, δεν υπάρχουν μετα-πληροφορίες που να την υποδεικνύουν ή να επιτρέπουν την αυτόματη αναγνώριση της.
- Εξελισσόμενο, μη τυποποιημένο λεξιλόγιο, με τη μορφή αργκό και διαλέκτων που χρησιμοποιούνται συχνά στην ανεπίσημη επικοινωνία μεταξύ χρηστών κοινωνικών δικτύων.

Ελλείπει, λοιπόν, μιας γενικής προσέγγισης στην κατηγοριοποίηση συναισθήματος στο Twitter, χρησιμοποιούμε emoticons για την αναγνώριση του συναισθήματος των tweets. Η προσέγγιση αυτή υπερνικά τις παραπάνω προκλήσεις και είναι αρκετά γενική ώστε να μπορεί να εφαρμοστεί σε κάθε θεματική κοινότητα. Είναι επίσης αρκετά αξιόπιστη για την επίδειξη της λειτουργίας του πλαισίου αξιολόγησης που προτείνουμε, καθώς συχνά χρησιμοποιείται στη βιβλιογραφία ως βάση σύγκρισης για το συναισθηματικά πολωμένο περιεχόμενο [39].

## **4.2 Βασικοί ορισμοί**

Για να διευκολύνουμε την κατανόηση των τύπων που θα ακολουθήσουν, καθώς και της ανάλυσης που θα παρουσιαστεί στις επόμενες ενότητες, στον Πίνακα 4-1 συνοψίζουμε τις βασικές έννοιες που θα χρησιμοποιηθούν.

Symbol	Description
$\langle x, y \rangle$	κατευθυνόμενη ακμή που δείχνει ότι ο χρήστης $x$ ακολουθεί το χρήστη $y$
$\langle m, u, t \rangle$	τριάδα που δείχνει ότι το μήνυμα $m$ δημοσιεύτηκε από το χρήστη $u$ τη χρονική στιγμή $t$
$G_T, G_h, G_p$	ο κοινωνικός γράφος του Twitter, της θεματικής κοινότητας $h$ , διακεκριμένης ομάδας $p$
$G_p@k$	η διακεκριμένη ομάδα που περιέχει τους $k$ κορυφαίους χρήστες με τη μεγαλύτερη επιρροή
$V_T, V_h, V_p$	οι χρήστες/κόμβοι του Twitter, της θεματικής κοινότητας $h$ , διακεκριμένης ομάδας $p$
$ V_T ,  V_h ,  V_p $	το μέγεθος (αριθμός χρηστών) του Twitter, της θεματικής κοινότητας $h$ , διακεκριμένης ομάδας $p$
$E_T, E_h, E_p$	οι κοινωνικοί σύνδεσμοι (ακμές) του Twitter, της θεματικής κοινότητας $h$ , διακεκριμένης ομάδας $p$
$N(h), M(G_p), M(u)$	τα κοινωνικά μετα-δεδομένα (tweets) της θεματικής κοινότητας $h$ , διακεκριμένης ομάδας $p$ , του χρήστη $u$
$ASM(G_l)$	η τιμή της μετρικής περίληψης δραστηριότητας για την (υπο-) κοινότητα $G_l$
$PR(G_l)$	η τιμή της αναλογίας πόλωσης για την (υπο-) κοινότητα $G_l$
$M(G_l, ASM, t_i, t_j)$	τα κοινωνικά μετα-δεδομένα της $G_l$ που δημοσιεύτηκαν κατά τη χρονική περίοδο $[t_i, t_j]$ και είναι σχετικά με την $ASM$
$Neg(G_l), Pos(G_l)$	τα αρνητικά/θετικά κοινωνικά μετα-δεδομένα που παράχθηκαν στην (υπο-) κοινότητα $G_l$
$CAC(G_l, ASM, t_1, t_2)$	η τιμή της συγκεντρωτικής καμπύλης δραστηριότητας αναφορικά με την $ASM$ για τη $G_l$ κατά το διάστημα $[t_1, t_2]$
$F_{rnd}$	η θεωρία τυχαίας επιρροής
$F_{ind}$	η θεωρία επιρροής βαθμού εισόδου
$F_{mnt}$	η θεωρία επιρροής αναφορών
$F_{rtw}$	η θεωρία επιρροής αναδημοσιεύσεων
$F_{twt}$	η θεωρία επιρροής όγκου δεδομένων

**Πίνακας 4-1.** Σύνοψη των βασικών εννοιών που χρησιμοποιούνται στους ορισμούς

Τα μέλη του Twitter και οι συνδέσεις μεταξύ τους, τυπικά μοντελοποιούνται μέσω ενός γράφου. Ο κάθε χρήστης αναπαρίσταται ως ένας κόμβος και κάθε σχέση μεταξύ δύο χρηστών ως μια ακμή που συνδέει τους αντίστοιχους κόμβους. Κάθε ακμή  $\langle x,$

$y$  είναι κατευθυνόμενη, με κατεύθυνση από τον *ακόλουθο*  $x$  προς τον *ακολουθούμενο*  $y$ . Η αναπαράσταση αυτή επιτυγχάνει να αναπαραστήσει την τοπολογία του υποκείμενου κοινωνικού δικτύου αλλά δεν περιλαμβάνει την δραστηριότητά του. Για τη μοντελοποίηση της τελευταίας, χρησιμοποιούμε επίσης την έννοια των κοινωνικών μετα-δεδομένων, δηλαδή το σύνολο των τριάδων που καταγράφουν όλο το περιεχόμενο που παράγεται από τους χρήστες στο χρόνο. Μια τριάδα αναπαρίσταται ως  $\langle m, u, t \rangle$ , όπου  $m$  είναι ένα μήνυμα που δημοσιεύεται από το χρήστη  $u$  τη χρονική στιγμή  $t$ .

Συνδυάζοντας τις δύο αναπαραστάσεις, ορίζουμε τον κοινωνικό γράφο του Twitter ως εξής:

**Ορισμός 4.1.** *Ο κοινωνικός γράφος του Twitter είναι ένας γράφος  $G_T = \{V_T, E_T, M\}$ , όπου  $V_T$  το σύνολο των κόμβων που αναπαριστούν τους χρήστες,  $E_T$  το σύνολο των κατευθυνόμενων ακμών που αναπαριστούν τους κοινωνικούς συνδέσμους μεταξύ τους, και  $M$  μια συνάρτηση που συσχετίζει τον κάθε χρήστη  $u \in V_T$  με τα μηνύματα που έχει δημοσιεύσει  $M(u) = \{\langle m, u, t \rangle\}$ .*

Η  $M(u)$  εμπεριέχει ολόκληρη τη συγγραφική δραστηριότητα του χρήστη  $u$ , που ονομάζεται *μετα-δεδομένα χρήστη*. Είναι απαραίτητο να επισημανθεί ο διαχωρισμός σε σχέση με τα *μετα-δεδομένα του θέματος*  $h$ , τα οποία ορίζονται από μια συνάρτηση  $N(h)$  και περιλαμβάνουν όλα τα μηνύματα που περιέχουν ένα συγκεκριμένο hashtag  $h$  μαζί με τους συγγραφείς τους και τη χρονική στιγμή που δημοσιεύτηκαν. Πιο επίσημα:

$$N(h) = \{\langle m, u, t \rangle : \text{hashtag}(m, h) = \text{true}\},$$

Όπου η  $\text{hashtag}(m, h)$  είναι μια boolean συνάρτηση που επιστρέφει *true* αν το μήνυμα  $m$  περιέχει το hashtag  $h$  και *false* διαφορετικά.

Γενικά, το σύνολο όλων των tweets που έχουν κοινό ένα συγκεκριμένο hashtag ορίζουν μια *θεματική κοινότητα* η οποία περιλαμβάνει και τους συγγραφείς τους. Η πρακτική αυτή υιοθετείται από πολλές σχετικές εργασίες στη βιβλιογραφία: στο [40], οι συγγραφείς χρησιμοποιούν τα hashtags ως υποκατάστατα των θεμάτων, στο [41] η χρήση σχετικών hashtags χρησιμοποιείται για να εκτιμηθεί η συνάφεια ενός χρήστη με ένα συγκεκριμένο θέμα, ενώ στο [42] οι συγγραφείς εξετάζουν τη δυναμική θεματικών κοινοτήτων σχετικών με την πολιτική όπως αυτές ορίζονται από τα hashtags. Ακολουθώντας αυτή τη σύμβαση, ορίζουμε τη *θεματική κοινότητα* ως τα μετα-δεδομένα ενός θέματος  $h$  μαζί με τις κοινωνικές σχέσεις μεταξύ των χρηστών που τα παρήγαγαν. Πιο επίσημα:

**Ορισμός 4.2.** Δεδομένου του κοινωνικού γράφου του Twitter  $G_T$  και μιας θεματικής κατηγορίας  $h$ , η *θεματική κοινότητα* που αντιστοιχεί στο  $h$  είναι ένας υπο-γράφος  $G_h = \{V_h, E_h, N\}$ , όπου  $N(h)$  είναι το σύνολο των μετα-δεδομένων του θέματος,  $V_h$  είναι το σύνολο όλων των χρηστών στο  $V_T$  που έχουν δημοσιεύσει τουλάχιστον ένα μήνυμα με την ετικέτα  $h$  ( $V_h = \{u \in V_T: M(u) \cap N(h) \neq \emptyset\}$ ) και το  $E_h$  περιλαμβάνει όλες τις ακμές του  $E_T$  και οι δύο ακμές των οποίων περιέχονται στο  $V_h$  (i.e.,  $E_h = \{\langle u_1, u_2 \rangle \in E_T: u_1 \in V_h \wedge u_2 \in V_h\} \subseteq E_T$ ). Ο αριθμός των χρηστών που συμμετέχουν στην κοινότητα αποτελεί το *μέγεθος της κοινότητας* και αναπαρίσταται ως  $|V_h|$ .

Τα μέλη μιας κοινότητας διαφέρουν στο επίπεδο επιρροής που εξασκούν στους υπόλοιπους χρήστες. Κάποιοι χρήστες είναι κυρίως παθητικοί, ενώ άλλοι διαπρέπουν σε κάποιο είδος δραστηριότητας της κοινότητας, επηρεάζοντας την συμπεριφορά των άλλων μελών και δημιουργώντας σχετικές τάσεις. Ονομάζουμε *επιδραστικούς* ή *διακεκριμένους χρήστες* τα μέλη μιας κοινότητας που έχουν κατακτήσει διακεκριμένη

θέση μέσα σε αυτή. Στο σύνολό τους, οι επιδραστικοί χρήστες μιας κοινότητας αποτελούν μια υπο-κοινότητα που ονομάζεται *διακεκριμένη ομάδα*. Επίσημα, η ομάδα αυτή ορίζεται ως εξής:

**Ορισμός 4.3.** Μέσα σε μια κοινότητα, η *διακεκριμένη ομάδα* είναι ένας υπο-γράφος  $G_p = \{V_p, E_p, M\}$ , όπου  $V_p$  είναι το σύνολο των διακεκριμένων χρηστών που περιλαμβάνονται στο  $V_h$ ,  $E_p$  το σύνολο των ακμών μεταξύ τους στο  $E_h$  (i.e.,  $E_p = \{ \langle u_1, u_2 \rangle \in E_h : u_1 \in V_p \wedge u_2 \in V_p \}$ ) και  $M$  η συνάρτηση που επιστρέφει τα μετα-δεδομένα ενός διακεκριμένου χρήστη (η ίδια συνάρτηση όπως στον ορισμό 3.1). Το μέγεθος  $|V_p|$  της διακεκριμένης ομάδας ονομάζεται *διακεκριμένο μέγεθος*.

Για λόγους απλότητας, τα συνολικά μετα-δεδομένα μιας διακεκριμένης ομάδας  $G_p$  συμβολίζονται ως  $M(G_p)$ ; δηλαδή,  $M(G_p) = \bigcup_{u \in V_p} M(u)$ . Σημειώνουμε επίσης ότι επεκτείνουμε το συμβολισμό αυτό για κάθε (υπο-) κοινότητα  $G_i$ , δηλαδή, για κάθε υπο-γράφο του κοινωνικού γράφου του Twitter που περιλαμβάνει ένα συγκεκριμένο σύνολο χρηστών μαζί με τις κοινωνικές επαφές τους και τα μηνύματα που έχουν γράψει. Στη συνέχεια, χρησιμοποιούμε τον όρο (υπο-) κοινότητα για να αναφερόμαστε συνδυαστικά σε μια θεματική κοινότητα και τις διακεκριμένες ομάδες της.

### 4.3 Θεωρίες τοπικής επιρροής

Μια θεωρία τοπικής επιρροής είναι ένα τυποποιημένο μοντέλο που στόχο έχει την ανίχνευση επιδραστικών χρηστών μέσα σε μια θεματική κοινότητα. Για το λόγο αυτό, συσχετίζει κάθε χρήστη με ένα *βαθμό επιρροής*, μια ταξινομημένη τιμή, δηλαδή, που υποδηλώνει το ατομικό επίπεδο επιρροής. Για παράδειγμα, θεωρείστε ένα βαθμό

επιρροής που ισούται με τον τοπικό βαθμό εισόδου του χρήστη, δηλαδή τον αριθμό των ακολούθων που έχει μέσα στην κοινότητα. Αφού τοποθετήσει τους χρήστες σε μια κλίμακα βαθμού επιρροής, μια θεωρία επιρροής τους ταξινομεί σε φθίνουσα σειρά, και αυτοί που τοποθετούνται στις  $k$  κορυφαίες θέσεις αποτελούν την **διακεκριμένη ομάδα -  $k$**  – με συμβολισμό  $G_p@k$ . Πιο επίσημα, μια θεωρία τοπικής επιρροής ορίζεται ως εξής:

**Ορισμός 4.4.** Για μια θεματική κοινότητα  $G_h$ , μια **θεωρία τοπικής επιρροής** είναι μια αντιστοίχιση των μελών της  $V_h$  προς ένα πλήρως ταξινομημένο σύνολο  $O_h$ , το οποίο αναθέτει σε κάθε χρήστη  $u \in V_h$  ένα βαθμό επιρροής μέσω συνάρτησης  $F: V_h \times G_h \rightarrow R$  και τους ταξινομεί σε φθίνουσα σειρά σύμφωνα με τα εξής:  $F(u_i, G_h) \leq F(u_j, G_h) \leftrightarrow o(u_i) \geq o(u_j)$ , όπου ο  $o(u_i)$  υποδηλώνει τη θέση κατάταξης του χρήστη  $u_i$  στο  $O_h$ . Όσο χαμηλότερη είναι η θέση κατάταξης  $o(u_i)$  για το χρήστη  $u_i$ , τόσο ψηλότερη είναι η επιρροή του, με τον  $o(u_i) = 1$  να αντιστοιχεί στον πιο επιδραστικό χρήστη.

Σε σχέση με τον ορισμό αυτό, δύο σημεία αξίζει να αποσαφηνιστούν:

- Ο βαθμός επιρροής που ανατίθεται σε ένα συγκεκριμένο χρήστη μπορεί να εξαρτάται όχι μόνο από τη δική του δραστηριότητα (δηλαδή τα μεταδεδωμένα χρήστη), αλλά και από πληροφορίες που προέρχονται από ολόκληρη την κοινότητα (π.χ. ο τοπικός βαθμός εισόδου).
- Κάθε ένδειξη που προέρχεται εκτός των ορίων της θεματικής κοινότητας δεν λαμβάνεται υπόψη. Για παράδειγμα, ένας εξωτερικός χρήστης  $u_j \notin G_h$  που ακολουθεί ένα μέλος της κοινότητας  $u_i \in G_h$  δεν συνυπολογίζεται στον υπολογισμό του τοπικού βαθμού εισόδου του  $u_i$ , έως ότου ο  $u_j$  συμμετάσχει κι αυτός στην ίδια κοινότητα.

Βάσει του ορισμού 4.4, τυποποιούμε τώρα τις θεωρίες τοπικής επιρροής που επιλέχθηκαν στην Ενότητα 3.1 για την αξιολόγηση μεγάλης κλίμακας που σκοπεύουμε να πραγματοποιήσουμε:

1. Η **Θεωρία Τυχαίας Επιρροής** ( $F_{rnd}$ ) υποστηρίζει ότι κάθε μέλος μιας θεματικής κοινότητας είναι επιδραστικό. Όλοι οι χρήστες, συνεπώς, έχουν το ίδιο επίπεδο επιρροής και τους αντιστοιχεί ο ίδιος βαθμός επιρροής:

$$F_{rnd}(u, G_h) = c \forall u \in V_h.$$

2. Η **Θεωρία Βαθμού Εισόδου** ( $F_{ind}$ ) βασίζεται στο συλλογισμό ότι σε όσο μεγαλύτερο κοινό απευθύνεται κάποιος, τόσο μεγαλύτερες πιθανότητες έχει να υιοθετηθούν οι απόψεις του και, συνεπώς, τόσο μεγαλύτερη είναι η επιρροή του. Συνεπώς, ο αντίστοιχος βαθμός επιρροής είναι ανάλογος της δημοφιλίας του χρήστη όπως υποδηλώνεται από τον αριθμό των ακολούθων που έχει μέσα στη θεματική κοινότητα (δηλαδή τον τοπικό βαθμό εισόδου του κόμβου του). Συγκεκριμένα:

$$F_{ind}(u, G_h) = |\{ \langle x, u \rangle \in E_h \}|.$$

3. Η **Θεωρία Αναφορών** ( $F_{mnt}$ ) σχετίζει την επιρροή ενός χρήστη με την ικανότητα του να εμπλέκεται σε συζητήσεις με τα υπόλοιπα μέλη της κοινότητας. Όσο πιο κοινωνικός είναι ένας χρήστης, τόσο περισσότεροι είναι οι άνθρωποι με τους οποίους έρχεται σε επαφή και μεγαλύτερη η επιρροή του. Με βάση τα παραπάνω, ποσοτικά ο βαθμός επιρροής ενός χρήστη ισούται με τον αριθμό των εσωτερικών αναφορών σε αυτόν.

$$F_{mnt}(u_i, G_h) = |\{ \langle m, u_j, t \rangle \in N(h) \wedge e(m, u_i) = true \wedge u_i \neq u_j \}|,$$

όπου  $e(m, u_i)$  μια boolean συνάρτηση που επιστρέφει *true* αν το tweet  $m$  αναφέρει το χρήστη  $u_i$  και *false* αλλιώς.



4. Η **Θεωρία Αναδημοσιεύσεων** ( $F_{rtw}$ ) θεωρεί ότι η επιρροή ενός χρήστη εξαρτάται από την ποιότητα του περιεχομένου που δημοσιεύει. Όσο πιο ενδιαφέροντα είναι τα μηνύματά του, τόσο περισσότεροι χρήστες θα τα διαβάσουν και τόσο μεγαλύτερη θα είναι η επιρροή που ασκεί. Η μετρική αυτή είναι υποκειμενική και μπορεί πρακτικά να εξαχθεί από τη *συχνότητα εσωτερικών αναδημοσιεύσεων*, η οποία δείχνει πόσες φορές τα μηνύματά του έχουν αναπαραχθεί από τα υπόλοιπα μέλη της θεματικής κοινότητας. Πιο επίσημα:

$$F_{rtw}(u_i, G_h) = |\{ \langle m, u_j, t \rangle \in N(h) \wedge r(m, u_i) = true \wedge u_i \neq u_j \}|,$$

όπου  $r(m, u_i)$  μια boolean συνάρτηση που επιστρέφει *true* αν ένα tweet  $m$  είναι αναδημοσίευση (retweet) ενός μηνύματος που αρχικά δημοσιεύτηκε από τον  $u_i$  και *false* αλλιώς.

5. Η **Θεωρία Όγκου Δεδομένων** ( $F_{twt}$ ) θεωρεί τον όγκο περιεχομένου που παράγεται από ένα συγκεκριμένο χρήστη ως σημαντική ένδειξη του επιπέδου επιρροής του. Όσο πιο παραγωγικός είναι ένας χρήστης, αυξάνεται η πιθανότητα ένα μέλος της κοινότητας να διαβάζει τα μηνύματά του και αυξάνεται και η επιρροή του. Στην περίπτωση αυτή, λοιπόν, ο βαθμός επιρροής ενός χρήστη  $u$  ισούται με τον αριθμό των tweets έχει δημοσιεύσει πάνω σε αυτό το θέμα:

$$F_{twt}(u, G_h) = |M(u) \cap N(h)|.$$

Έχοντας ορίσει τις θεωρίες επιρροής που θα αναλύσουμε, μπορούμε τώρα να δημιουργήσουμε ένα καθορισμένο πλαίσιο για να αξιολογήσουμε την ακρίβεια πρόβλεψής τους.

## 4.4 Πλαίσιο αξιολόγησης

Η συμπεριφορά των διακεκριμένων χρηστών είναι πολύ σημαντική για μια θεματική κοινότητα, καθώς τις πράξεις και τις απόψεις τους μιμείται και αναπαράγει ένα σημαντικό μέρος των μελών της κοινότητας – ακόμα και αυτοί στους οποίους δεν έχουν άμεση πρόσβαση. Το φαινόμενο αυτό είναι τόσο έντονο ώστε κάποιες πλευρές της συνολικής δραστηριότητας μιας κοινότητας να μπορούν να θεωρηθούν ως αντανάκλαση της αντίστοιχης συμπεριφοράς των επιδραστικών χρηστών. Για να εκτιμήσουμε ποσοτικά το βαθμό στον οποίο συμβαίνει αυτό, πρέπει να συγκρίνουμε τη συνολική δραστηριότητα των διακεκριμένων χρηστών με αυτήν της κοινότητας, βάσει μιας αντικειμενικά μετρήσιμης μετρικής. Μα τέτοια μετρική παίρνει ως είσοδο (ένα μέρος από) τον κοινωνικό γράφο και αντιστοιχίζει τη σχετική δραστηριότητα σε μια αριθμητική τιμή. Την ονομάζουμε *μετρική περίληψης δραστηριότητας* και την ορίζουμε επίσης ως εξής:

**Ορισμός 4.5.** Δεδομένης μιας (υπο-)κοινότητας  $G_1$ , μια *μετρική περίληψης δραστηριότητας* (Activity Summary Metric - ASM) είναι μια συνάρτηση που αντιστοιχεί την  $G_1$  στο χώρο των πραγματικών αριθμών, βάσει ενός συγκεκριμένου υποσυνόλου της πληροφορίας που εμπεριέχει.

Αυτός είναι ένας γενικός ορισμός που μπορεί να συμπεριλάβει πολλές μετρικές. Για παράδειγμα, μπορούμε να θεωρήσουμε μια συνάρτηση που υπολογίζει τη συχνότητα των εμφανίσεων συγκεκριμένων όρων σχετικών με το θέμα στα tweets των μελών της κοινότητας. Δεδομένου, ωστόσο, ότι τέτοιες μετρικές εξαρτώνται και από τη γλώσσα και από το θέμα, ο θόρυβος και η πολυγλωσσία του περιεχομένου του Twitter τις καθιστούν ασύμβατες με το πλαίσιο αξιολόγησής μας, το οποίο δε θέτει προϋποθέσεις σχετικά με την ποιότητα και τα χαρακτηριστικά των θεματικών

κοινοτήτων που χρησιμοποιεί ως δεδομένα αναφοράς. Δεδομένου επίσης ότι ο στόχος της ανάλυσής μας είναι να δείξει την λειτουργία του πλαισίου, θεωρούμε ως *ASM* αποκλειστικά μία μετρική που βασίζεται σε πολωμένο περιεχόμενο, όπως ορίζεται στην Ενότητα 3.1. Η μετρική ποσοτικοποιεί το αθροιστικό συναίσθημα που εκφράζεται από τα tweets μιας κοινότητας ή μέρους των χρηστών της. Την ονομάζουμε Βαθμό Πόλωσης (*Polarity Ratio PR*) και το ορίζουμε ως εξής:

**Ορισμός 4.6.** Δεδομένης μιας (υπο-)κοινότητας  $G_l$ , ο **βαθμός πόλωσης**  $PR(G_l)$  ορίζεται ως εξής:

$$PR(G_l) = \begin{cases} \frac{|Pos(G_l)| + 1}{|Neg(G_l)| + 1} - 1, & \text{αν } |Neg(G_l)| < |Pos(G_l)| \\ -\frac{(|Neg(G_l)| + 1)}{|Pos(G_l)| + 1} + 1, & \text{αν } |Pos(G_l)| < |Neg(G_l)| \end{cases}$$

όπου  $Neg(G_l) \subseteq M(G_l)$  και  $Pos(G_l) \subseteq M(G_l)$  τα υποσύνολα των κοινωνικών μετα-δεδομένων της  $G_l$  που αντιστοιχούν σε αρνητικά και θετικά tweets, αντίστοιχα – όπως ορίζονται στην Ενότητα 3.1 βάσει των εικονιδίων έκφρασης συναισθήματος (δηλαδή με ένα τρόπο ανεξάρτητο από γλώσσα και με ανοχή σε θόρυβο). Τα  $|Neg(G_l)|$  και  $|Pos(G_l)|$  συμβολίζουν το μέγεθός τους, με  $|Neg(G_l)| + |Pos(G_l)| \leq |M(G_l)|$ .

Η  $PR(G_l)$  λαμβάνει τιμές στο διάστημα  $(-|M(G_l)|, +|M(G_l)|)$ , με τις θετικές τιμές να υποδηλώνουν την επικράτηση των θετικών tweets και αντίστροφα. Πιο λεπτομερώς, μια θετική τιμή  $n$  υποδηλώνει ότι τα θετικά tweets ξεπερνούν τα αρνητικά κατά ένα παράγοντα (περίπου)  $n + 1$ ; το αντίθετο ισχύει για τις αρνητικές τιμές. Τα ουδέτερα συναισθήματα αντιστοιχούν σε ισορροπία μεταξύ αρνητικών και θετικών tweets ( $|Neg(G_l)| \approx |Pos(G_l)|$ ), δίνοντας τιμές πολύ κοντά στο 0.

Βασιζόμενοι στην έννοια της  $ASM$  και την αρχικοποίηση της μέσω της  $PR$ , εισάγουμε τώρα τον επίσημο ορισμό του πλαισίου αξιολόγησης μεγάλης κλίμακας για τις θεωρίες επιρροής. Συγκεκριμένα, θεωρούμε ότι οι ακόλουθες συνθήκες αποτελούν ισχυρές ενδείξεις πραγματικής επιρροής όταν ικανοποιούνται συνδυαστικά:

- i. **Συνθήκη συσχέτισης.** Η τιμή μιας  $ASM$  επιτρέπει την αξιολόγηση της ομοιότητας της συμπεριφοράς των διακεκριμένων χρηστών με αυτήν της υπόλοιπης κοινότητας. Υψηλή συσχέτιση μεταξύ  $ASM(G_p)$  και  $ASM(G_h)$  σημαίνει ότι η συμπεριφορά της διακεκριμένης ομάδας είναι αντιπροσωπευτική ολόκληρης της θεματικής κοινότητας. Στην πραγματικότητα, όσο πιο κοντά είναι οι τιμές τους κατά τη διάρκεια του χρόνου, τόσο μεγαλύτερη η επιρροή των διακεκριμένων χρηστών πάνω στην κοινότητα. Επίσημα, η συνθήκη αυτή εκφράζεται ως εξής:

$$ASM(G_p) \approx ASM(G_h).$$

Για παράδειγμα, αν θεωρήσουμε το βαθμό πόλωσης: οι πραγματικά επιδραστικοί χρήστες θα πρέπει να είναι σε θέση να καθορίσουν τη συνολική γνώμη μιας κοινότητας με αδιαμφισβήτητο τρόπο, και, συνεπώς, οι  $PR(G_p)$  και  $PR(G_h)$  θα πρέπει να έχουν το ίδιο πρόσημο και παρόμοιες τιμές μεγεθών.

- ii. **Συνθήκη μεγέθους.** Το σχετικό μέγεθος της διακεκριμένης ομάδας ( $|V_p|$ ) παίζει σημαντικό ρόλο στη συσχέτιση μεταξύ  $ASM(G_p)$  και  $ASM(G_h)$ : όσο περισσότερους χρήστες θεωρήσουμε ως επιδραστικούς, τόσο πλησιάζει η συμπεριφορά της διακεκριμένης ομάδας σε αυτή της κοινότητας. Η απόλυτη προσέγγιση ( $ASM(G_p) = ASM(G_h)$ ) ουσιαστικά αντιστοιχεί στην ακραία περίπτωση κατά την οποία η  $G_p$  περιλαμβάνει ολόκληρη την κοινότητα  $G_h$ . Στην πράξη ωστόσο, όσο μικρότερη είναι η διακεκριμένη ομάδα, τόσο πιο

αποτελεσματική από άποψη κόστους είναι η αντίστοιχη εφαρμογή (π.χ., viral marketing). Είναι θεμελιώδες, συνεπώς, να θεωρήσουμε διακεκριμένες ομάδες με μέγεθος πολύ μικρότερο της συνολικής κοινότητας. Η συνθήκη αυτή εκφράζεται ως εξής:

$$|V_p| \ll |V_h|.$$

- iii. **Συνθήκες όγκου.** Ο όγκος του περιεχομένου που παράγεται από μια διακεκριμένη ομάδα είναι κριτικός παράγοντας στον καθορισμό της πραγματικής της επιρροής πάνω στην κοινότητα. Για παράδειγμα, στην περίπτωση της *PR*, ας θεωρήσουμε μια κοινότητα στην οποία τα πολωμένα tweets προέρχονται σχεδόν αποκλειστικά από μια μικρή ομάδα χρηστών, ενώ οι υπόλοιποι χρήστες είναι σχετικά φειδωλοί, συνεισφέροντας λιγότερο από 10% του πολωμένου περιεχομένου. Αν και στην περίπτωση αυτή δε μπορούμε να μιλήσουμε για πραγματική σχέση επιρροής, ο βαθμός πόλωσης αναπόφευκτα εμφανίζει υψηλή συσχέτιση μεταξύ  $PR(G_p)$  και  $PR(G_h)$ . Είναι σαφές, λοιπόν, ότι μια διακεκριμένη ομάδα θα πρέπει να ευθύνεται για ένα μικρό μόνο ποσοστό του περιεχομένου της κοινότητας που σχετίζεται με την επιλεγμένη *ASM*. Ακολουθώντας αντίστοιχο συλλογισμό, μια διακεκριμένη ομάδα δε μπορεί να αποτελείται από υπερδραστήριους χρήστες που κατακλύζουν τους υπόλοιπους χρήστες με το περιεχόμενό τους. Αντίθετα, οι υποψήφιοι διακεκριμένοι χρήστες θα πρέπει να ευθύνονται για ένα μικρό ποσοστό του συνολικού περιεχομένου της κοινότητας. Πιο συγκεκριμένα, το πλαίσιο μας απαιτεί από τους πραγματικά επιδραστικούς χρήστες να πληρούν τις ακόλουθες δύο συνθήκες σχετικά με τον όγκο του περιεχομένου τους:

- Η **συνθήκη όγκου συνολικού περιεχομένου** απαιτεί οι διακεκριμένοι χρήστες να είναι υπεύθυνοι για ένα μικρό ποσοστό

των μετα-δεδομένων ολόκληρης της κοινότητας, δηλαδή,  $|M(G_p, t_1, t_2)| \ll |M(G_h, t_1, t_2)|$ , όπου το  $M(G_l, t_i, t_j)$  υποδηλώνει το σύνολο των μετα-δεδομένων που δημοσιεύτηκαν από τους χρήστες της  $G_l$  κατά το χρονικό διάστημα  $[t_i, t_j]$  και το  $|M(G_l, t_i, t_j)|$  συμβολίζει το μέγεθός του.

- Η **συνθήκη όγκου περιεχομένου μετρικής** απαιτεί οι διακεκριμένοι χρήστες να παράγουν δεδομένα σχετικά με την επιλεγμένη  $ASM$ , αλλά αυτά να αποτελούν ένα μικρό ποσοστό των tweets της κοινότητας που σχετίζονται με αυτή την  $ASM$ , i.e.,  $0 < |M(G_p, ASM, t_1, t_2)| \ll |M(G_h, ASM, t_1, t_2)|$ , όπου το  $M(G_l, ASM, t_i, t_j)$  υποδηλώνει το σύνολο των μεταδεδομένων που σχετίζονται με την  $ASM$  τα οποία δημοσιεύτηκαν από τους χρήστες της  $G_l$  κατά το χρονικό διάστημα  $[t_i, t_j]$  και το  $|M_l(G_l, ASM, t_i, t_j)|$  συμβολίζει το μέγεθός του. Στην περίπτωση του  $PR$ , η συνθήκη όγκου περιεχομένου μετρικής ορίζει ένα ανώτατο όριο για το ποσοστό των πολωμένων tweets που δημοσιεύουν οι διακεκριμένοι χρήστες

- iv. **Συνθήκη χρόνου.** Άλλη μία σημαντική πλευρά της πραγματικής επιρροής είναι η χρονική σχέση μεταξύ  $ASM(G_p)$  και  $ASM(G_h)$ , καθώς το περιεχόμενο που αναφέρεται στην  $ASM(G_p)$  πρέπει να προηγείται χρονικά του περιεχομένου που αναφέρεται στην  $ASM(G_h)$ . Για την κατανόηση αυτής της συνθήκης, ας θεωρήσουμε μια μικρή ομάδα χρηστών που απλά ακολουθεί το γενικό συναίσθημα της κοινότητας και είναι φειδωλή στη δημοσίευση πολωμένου περιεχομένου. Η ομάδα αυτή ικανοποιεί όλες τις προηγούμενες

συνθήκες και μπορεί, συνεπώς, να θεωρηθεί λανθασμένα ως διακεκριμένη ομάδα, αν και στην πραγματικότητα δεν ασκεί καμία επιρροή στην κοινότητα.

Δεδομένης μιας *ASM*, υπολογίζουμε το σχετικό ρυθμό παραγωγής περιεχομένου μεταξύ δύο ομάδων για ένα συγκεκριμένο χρονικό διάστημα  $[t_1, t_2]$  μέσω της καμπύλης συγκεντρωτικής δραστηριότητας, μιας μετρικής που ορίζεται ως εξής:

**Ορισμός 4.7.** Δεδομένης μιας *ASM* και μιας (υπο-)κοινότητας  $G_l$  ενεργής κατά το χρονικό διάστημα  $[t_1, t_2]$ , η καμπύλη συγκεντρωτικής δραστηριότητας (*Cumulative Activity Curve - CAC*) ποσοτικοποιεί το ρυθμό παραγωγής του περιεχομένου της  $G_l$  που αναφέρεται στην *ASM* και ορίζεται ως εξής:

$$CAC(G_l, ASM, t_1, t_2) = \int_{t_1}^{t_2} \frac{M_l(G_l, ASM, t_1, t)}{M_l(G_l, ASM, t_1, t_2)} dt.$$

Η μονότονα αυξητική συνάρτηση του  $M(G_l, ASM, t_i, t_j)$  είναι κατάλληλη για την αποτύπωση του φαινομένου των πρώιμων οπαδών (early adopter). Η *CAC* λαμβάνει τιμές στο διάστημα  $[0,1]$ , με τις μεγαλύτερες τιμές να υποδηλώνουν ταχύτερη σύγκλιση στον τελικό όγκο περιεχομένου. Στις ρυθμίσεις μας που βασίζονται στο *PR*, η *CAC* εκφράζει πόσο γρήγορα το πολωμένο περιεχόμενο μιας κοινότητας  $G_h$  (ή μιας διακεκριμένη ομάδα της  $G_p$ ) παράχθηκε μέσα στο χρονικό διάστημα  $[t_1, t_2]$ .

Για να κατανοήσουμε το σκεπτικό πίσω από αυτή τη μετρική, ας θεωρήσουμε το δισδιάστατο χώρο του Πίνακα 4-2: ο οριζόντιος άξονας ( $x$ ) αντιστοιχεί στη χρονική περίοδο μεταξύ  $t_1$  και  $t_2$ , και ο κάθετος ( $y$ ) στο ποσοστό των πολωμένων tweets. Στην πράξη, η *CAC* υπολογίζει το εμβαδό κάτω από την καμπύλη που σχηματίζεται αν αντιστοιχίσουμε στο χώρο αυτό την τιμή της αναλογίας  $\frac{M_l(G_l, ASM, t_1, t)}{M_l(G_l, ASM, t_1, t_2)}$  μεταξύ των χρονικών στιγμών  $t_1$  και  $t_2$ . Όσο πιο γρήγορα η καμπύλη συγκλίνει στη γραμμή

$y = 1$  τόσο υψηλότερη είναι η τιμή της  $CAC$  και τόσο υψηλότερος ο ρυθμός παραγωγής.

Τα ενδεικτικά δεδομένα που παρουσιάζονται στην Πίνακας 4-2 δείχνουν τις τρεις γενικές μορφές της  $CAC$ : Η  $G_1$  αναπαριστά μια κοινότητα που έχει παράξει το περισσότερο πολωμένο περιεχόμενο της πριν το μέσο της χρονικής περιόδου και παίρνει υψηλή τιμή  $CAC$ :  $CAC(G_1, PR, t_1, t_2) = 0.70$ . Αντίθετα, η  $G_3$  υποδηλώνει μια ομάδα χρηστών που δημοσίευσαν τη μεγάλη πλειοψηφία των πολωμένων tweets προς το τέλος του εξεταζόμενου χρονικού διαστήματος, λαμβάνοντας χαμηλή τιμή  $CAC$ :  $CAC(G_3, PR, t_1, t_2) = 0.37$ . Στη μέση των δύο αυτών άκρων, βρίσκεται η  $G_2$ , το πολωμένο περιεχόμενο της οποίας κατανέμεται ισομερώς κατά το διάστημα της εξεταζόμενης περιόδου, έχοντας ως αποτέλεσμα μια ισορροπημένη τιμή για την  $CAC$ :  $CAC(G_2, PR, t_1, t_2) = 0.50$ .

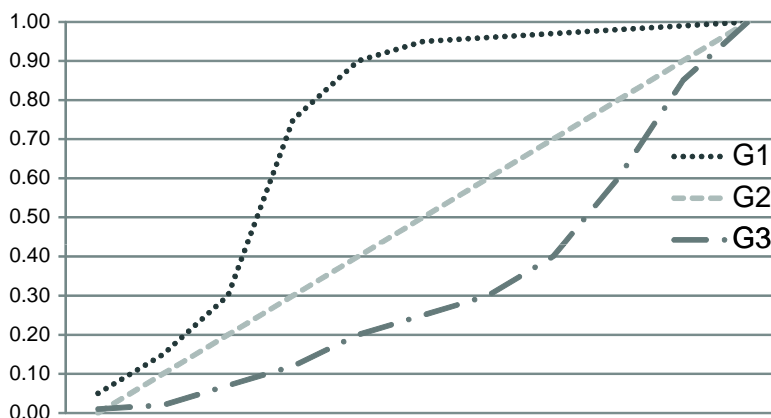
Συνολικά, η συνθήκη χρόνου για ύπαρξη πραγματικής επιρροής, εκφράζεται ως εξής:

$$CAC(G_p, ASM, t_1, t_2) \gg CAC(G_h, ASM, t_1, t_2).$$

Επισημαίνουμε ότι η  $CAC$  είναι εξορισμού μια κανονικοποιημένη μετρική. Συνεπώς, είναι κατάλληλη για να συγκρίνουμε τη δραστηριότητα δύο ομάδων χρηστών (π.χ.,  $G_h$  και  $G_p$ ) που ήταν ενεργές κατά τη διάρκεια της ίδιας χρονικής περιόδου *σε ισότιμη βάση*: η ομάδα με την πιο πρώιμη παραγωγή περιεχομένου παίρνει υψηλότερη τιμή  $CAC$ , ανεξάρτητα από το σχετικό μέγεθος των αντίστοιχων μεταδεδομένων.

Συνδυάζοντας τις παραπάνω πέντε συνθήκες, μπορούμε να τυποποιήσουμε τη διαδικασία της αξιολόγησης της επίδοσης μιας θεωρίας επιρροής ως εξής:





**Πίνακας 4-2: Αναπαράσταση της CAC με ενδεικτικά δεδομένα από τρεις κοινότητες  $G_1$ ,  $G_2$  και  $G_3$**

**ΠΡΟΒΛΗΜΑ 1.** Έστω μια ASM και μια θεματική κοινότητα  $G_h$  που ήταν ενεργή κατά τη διάρκεια ενός χρονικού διαστήματος  $[t_1, t_2]$ . Μία αποδοτική θεωρία τοπικής επιρροής ορίζει ως διακεκριμένη ομάδα ένα υποσύνολο  $G_p \subset G_h$  που ικανοποιεί όλες τις ακόλουθες συνθήκες:

$$(1) |V_p| \ll |V_h|,$$

$$(2) |M(G_p, t_1, t_2)| \ll |M(G_h, t_1, t_2)|,$$

$$(3) 0 < |M(G_p, ASM, t_1, t_2)| \ll |M(G_h, ASM, t_1, t_2)|,$$

$$(4) ASM(G_p) \approx ASM(G_h), \text{ και}$$

$$(5) CAC(G_p, ASM, t_1, t_2) \gg CAC(G_h, ASM, t_1, t_2).$$

Όσο περισσότερο ισχύουν οι συνθήκες αυτές, τόσο πιο αποδοτική είναι η θεωρία επιρροής.

Επισημαίνουμε ότι ο παραπάνω ορισμός είναι αρκετά γενικός ώστε να μπορεί να δεχθεί ένα ευρύ φάσμα μετρικών περίληψης δραστηριότητας. Για τους λόγους που αναφέρθηκαν παραπάνω, στην ανάλυση που ακολουθεί χρησιμοποιούμε ως μετρική αποκλειστικά το βαθμό πόλωσης. Επισημαίνουμε, επίσης, ότι οι τελευταίες δύο συνθήκες θα εξεταστούν μόνο σε περίπτωση που οι άλλες τρεις προϋποθέσεις

ισχύουν για τις διακεκριμένες ομάδες υπό εξέταση. Ουσιαστικά, οι πρώτες τρεις συνθήκες φιλτράρουν τις διακεκριμένες ομάδες που είναι ασύμβατες με το πλαίσιο αξιολόγησής μας: δε μπορούμε να εξάγουμε ασφαλή συμπεράσματα για διακεκριμένους χρήστες που αποτελούν ένα μεγάλο μέρος της βάσης χρηστών της κοινότητας ή για χρήστες που είναι υπερβολικά παραγωγικοί και παράγουν υπερβολικά μεγάλο μέρος του συνολικού ή του σχετικού με τη μετρική περιεχομένου της κοινότητας. Το ίδιο ισχύει για διακεκριμένες ομάδες που δεν έχουν δημοσιεύσει περιεχόμενο σχετικό με την επιλεγμένη μετρική. Στην επόμενη ενότητα, εφαρμόζουμε το πλαίσιό μας στις παραπάνω θεωρίες επιρροής ακολουθώντας αυτή τη σειρά εκτέλεσης.

Η σελίδα αυτή είναι σκόπιμα λευκή.

# 5

## Αξιολόγηση θεωριών τοπικής επιρροής

Στην παρούσα ενότητα θέτουμε το πλαίσιο αξιολόγησης σε εφαρμογή, συγκρίνοντας τις θεωρίες επιρροής της προηγούμενης ενότητας με βάση το Πρόβλημα 1. Αρχικά, εξηγούμε τη διαδικασία δημιουργίας του συνόλου πραγματικών δεδομένων αναφοράς μεγάλης κλίμακας που χρησιμοποιήσαμε. Στη συνέχεια, εξετάζουμε κάθε συνθήκη ξεχωριστά, με τη σειρά κατά την οποία θα πρέπει να εφαρμοστεί. Η μόνη εξαίρεση είναι η Συνθήκη Μεγέθους: αντί αυτής, θεωρούμε το καθορισμένο σύνολο διακριμένων μεγεθών  $k \in \{10, 20, 50, 100\}$ . Επιλέξαμε το συγκεκριμένο σύνολο τιμών  $k$  για τους εξής λόγους: πρώτον, επιτρέπει τη σύγκριση των θεωριών επιρροής σε ισότιμη βάση κάτω από διαφορετικές συνθήκες. Δεύτερον, κάθε τιμή  $k_i$  αποτελεί κατάλληλο μέγεθος ομάδας για αρκετές εφαρμογές ενώ ικανοποιεί τη συνθήκη  $k_i \ll |V_h|$ . Τρίτον, επιλέγοντας μεγέθη μεγαλύτερης ακρίβειας (π.χ.,  $k \in \{10, 20, 30, 40, \dots\}$ ) καταλήγουμε σε αμελητέες διακυμάνσεις μεταξύ των μεγεθών, ενώ μεγέθη μικρότερης ακρίβειας (π.χ.,  $k \in \{10, 100, 1000, \dots\}$ ) οδηγούν σε απώλεια σημαντικής πληροφορίας. Αντιθέτως, το επιλεγμένο φάσμα μεγεθών έχει ως αποτέλεσμα ισορροπημένες μεταβολές στην επίδοση που επιτρέπει τον προσδιορισμό

του βέλτιστου διακεκριμένου μεγέθους, εκείνου δηλαδή που δίνει την καλύτερη επίδοση για το Πρόβλημα 1 στις περισσότερες κοινότητες αναφοράς.

## **5.1 Twitter Dataset**

Στην πειραματική μας ανάλυση, χρησιμοποιήσαμε το σύνολο δεδομένων του Twitter που χρησιμοποιήθηκε στο [43]. Περιέχει πάνω από 475 εκατομμύρια tweets που δημοσιεύτηκαν από περισσότερους από 17 εκατομμύρια διακριτούς χρήστες σε μια χρονική περίοδο 7 μηνών – από την αρχή του Ιουνίου του 2009 μέχρι το τέλος του Δεκεμβρίου του 2009. Συνολικά, καταγράφει περίπου το 20%-30% της συνολικής δραστηριότητας του Twitter κατά τη διάρκεια της περιόδου αυτής, αποτελεί συνεπώς ένα αντιπροσωπευτικό δείγμα της δραστηριότητας του Twitter για την περίοδο αυτή, επαρκές για την εξαγωγή ασφαλών συμπερασμάτων για τις σχέσεις επιρροής μεταξύ των μελών του. Δεν περιέχει, όμως καμία πληροφορία για τις κοινωνικές σχέσεις μεταξύ των μελών του.

Για να ξεπεράσουμε αυτή την έλλειψη, λάβαμε υπόψη το γράφο που σχηματίζεται από τους κατευθυνόμενους συνδέσμους μεταξύ των χρηστών και συγκεκριμένα την αποτύπωση του γράφου του Twitter που χρησιμοποιήθηκε στο [44], η οποία τοποθετείται ημερολογιακά τον Αύγουστο του 2009 και συνεπώς συμπίπτει χρονικά με τα δεδομένα που εξετάζουμε. Μπορέσαμε, συνεπώς, να συνενώσουμε τα δύο σύνολα δεδομένων μέσω των ονομάτων χρηστών (λογαριασμού Twitter) που έχουν από κοινού – κατά μέσο όρο, περίπου τα μισά από τα μέλη κάθε κοινότητας συμπεριλαμβάνονταν στον αποτυπωμένο γράφο. Δυστυχώς, η αποτύπωση του γράφου δε εμπεριέχει πληροφορία για τη χρονική διάσταση των συνδέσεων μεταξύ των χρηστών, τη χρονική στιγμή δηλαδή που σχηματίστηκαν και την εξέλιξη τους μέσα στη χρονική περίοδο που εξετάζουμε. Κατά συνέπεια, τα συμπεράσματα που

εξάγουμε σχετικά με τη γραφική δομή του Twitter αποτελούν προσέγγιση των πραγματικών φαινομένων.

Όπως εξηγήσαμε στην Ενότητα 3.2, οι θεματικές κοινότητες ορίζονται μέσω hashtags. Η χρήση hashtags, όμως, δε συνοδεύεται απαραίτητα από αλληλεπίδραση με άλλα μέλη της θεματικής κοινότητας – ο συγγραφέας μπορεί να μην δει καν άλλα tweets με το ίδιο hashtag. Για να εξασφαλίσουμε ένα ελάχιστο επίπεδο κοινωνικής αλληλεπίδρασης μέσα σε μια θεματική κοινότητα, θεωρούμε ως έγκυρες για την ανάλυσή μας κοινότητες που περιέχουν:

- τουλάχιστον 500 εσωτερικές αναφορές (αναφορές, δηλαδή, σε ένα μέλος της κοινότητας από ένα άλλο μέλος), και
- τουλάχιστον 500 εσωτερικές αναδημοσιεύσεις (αναδημοσιεύσεις δηλαδή μηνυμάτων ενός μέλους της κοινότητας από άλλα μέλη)

Επιβάλλουμε επίσης τους παρακάτω περιορισμούς ώστε να εξασφαλίσουμε ότι οι κοινότητες που εξετάζουμε είναι αρκετά μεγάλες ώστε να εξάγουμε ασφαλή συμπεράσματα<sup>1</sup>:

- περιέχει πάνω από 5,000 διακριτά tweets,
- περιλαμβάνει πάνω από 500 χρήστες,
- περιέχει πάνω από 500 πολωμένα tweets.

---

<sup>1</sup>Επισημαίνουμε ότι όλα τα όρια τέθηκαν σε στρογγυλοποιημένες τιμές που προσεγγίζουν τη μέση τιμή πάνω από όλες τις θεματικές κοινότητες.

	Min.	Mean	Median	Max.	Σύνολο
Χρήστες	502	8,851	5,901	88,156	663,855
Χρήστες στο γράφο	30.89%	48.67%	48.94%	69.98%	-
Tweets	7,771	80,148	39,603	1,191,345	6,011,101
Tweets ανά χρήστη	3.26	8.35	6.60	30.33	-
Διάρκεια σε ημέρες	92	195	203	207	-
Αρνητικά Tweets	38	668	273	12,902	50,080
Θετικά Tweets	28	917	669	8,501	68,751
Εσωτερικές αναφορές	548	8,419	3,438	172,432	631,415
Εσωτερικές αναδημοσιεύσεις	570	18,998	5,661	471,588	1,424,835

**Πίνακας 5-1. Τεχνικά χαρακτηριστικά των 75 θεματικών κοινοτήτων που αποτελούν το σύνολο δεδομένων αναφοράς για την πειραματική μελέτη.**

Το αρχικό σύνολο δεδομένων περιέχει πάνω από 49 εκατομμύρια tweets επισημασμένα με τουλάχιστον ένα hashtag, που αντιστοιχούν σε περίπου 3 εκατομμύρια διακριτά θέματα συνολικά. Ανάμεσά τους, 728 κοινότητες ικανοποιούν όλες τις παραπάνω συνθήκες. Για να επιλέξουμε τις πιο ενεργητικές, τις αξιολογήσαμε σε φθίνουσα σειρά με βάση το μέσο αριθμό tweets που δημοσιεύει το κάθε μέλος. Όσο πιο ψηλό είναι αυτό το νούμερό, τόσο πιο δραστήριοι αναμένεται να είναι οι επιμέρους χρήστες. Επιλέξαμε στη συνέχεια τα 100 πρώτα σε κατάταξη hashtags και μελετήσαμε προσεκτικά το περιεχόμενό τους για να εξασφαλίζουμε ότι ορίζουν πραγματικές θεματικές κοινότητες.

Στην προσπάθεια αυτή, βασιστήκαμε σε προηγούμενες εργασίες που αξιολογούν τα hashtags αναφορικά με την ικανότητά τους να ορίσουν πραγματικές θεματικές κατηγορίες. Στο [45], οι συγγραφείς ξεχωρίζουν τα hashtags σε αυτά που ορίζουν πραγματικά γεγονότα και στα λεγόμενα memes, τα οποία είναι θέματα συζήτησης ή διαδεδомένοι τρόποι επισήμανσης που δε σχετίζονται με πραγματικά γεγονότα. Παρόμοια προσέγγιση ακολουθείται και από στο [46]. Με βάση αυτή την κατηγοριοποίηση, αφαιρούμε τα memes από την ανάλυσή μας, καθώς είναι απίθανο

να ορίζουν πραγματικές θεματικές κοινότητες και είναι, συνεπώς, ακατάλληλα για την αξιολόγηση που επιχειρούμε. Παραδείγματα memes που αφαιρέσαμε είναι τα hashtags *#ff*, που χρησιμοποιείται για την πρόταση χρηστών που αξίζει κάποιος να ακολουθήσει, και *#job*, που συνοδεύει αγγελίες εύρεσης εργασίας. Εκτός από τα memes, αφαιρέθηκαν τα hashtags που είναι υπερβολικά γενικά για να θεωρήσουμε ότι ορίσουν μια θεματική κοινότητα, όπως *#web*, καθώς και διαφορετικά hashtags που μπορεί να αναφέρονται σε δύο ή και περισσότερα θέματα. Ένα παράδειγμα για το τελευταίο είναι το *#gr88*, που αναφέρεται στην πράσινη επανάσταση (καθώς το περσικό ημερολογιακό έτος είναι το 1388) αλλά χρησιμοποιείται επίσης ως αναφορά σε διαδικτυακό καζίνο (*gr88.com*). Οι 75 θεματικές κοινότητες που απέμειναν μετά την αφαίρεση, σχηματίζουν τα δεδομένα αναφοράς μας.

Τα τεχνικά χαρακτηριστικά τους παρουσιάζονται στον Πίνακα 5-1. Κατά μέσο όρο, κάθε κοινότητα του συνόλου δεδομένων ήταν ενεργή για πάνω από 6 μήνες και περιείχε σχεδόν 9,000 διακριτά μέλη που δημοσίευσαν ο καθένας πάνω από 8 tweets. Από άποψη περιεχομένου, κάθε κοινότητα αντιστοιχεί σε 80,000 tweets κατά μέσο όρο, το 1/10 των οποίων αποτελούν συζητήσεις μεταξύ των μελών της κοινότητας. Οι εσωτερικές αναδημοσιεύσεις καλύπτουν το 1/4 του συνολικού περιεχομένου, ενώ 650 tweets επισημαίνονται ως αρνητικά και άλλα 900 ως θετικά. Επισημαίνουμε, ωστόσο, ότι η υψηλότερη συχνότητα των θετικών tweets δεν υπονοεί ότι υπερέχουν σε όλες τις θεματικές κατηγορίες. Στην πραγματικότητα, τα αρνητικά tweets υπερέχουν αριθμητικά των θετικών στο 1/3 όλων των κοινοτήτων.

## 5.2 Συνθήκες όγκου

Η ενότητα αυτή ερευνά την απόδοση των επιλεγμένων θεωριών επιρροής σε σχέση με τις δύο συνθήκες όγκου. Ξεκινάμε με τη συνθήκη όγκου δεδομένων που ορίζονται από τη μετρική, η οποία μετράει το ποσοστό του πολωμένου περιεχομένου ανά



θεματική κατηγορία το οποίο προέρχεται από τους διακεκριμένους χρήστες που ορίζονται από κάθε θεωρία επιρροής. Στην πραγματικότητα, ο στόχος μας είναι να διακρίνουμε τις διακεκριμένες ομάδες σε τρεις βασικές κατηγορίες:

		Min.	Median	Max.	Mean±SD	Ανενεργές	Υπερδραστήριες
<b><math>G_p@10</math></b>	<b><math>F_{mnt}</math></b>	0.04	7.33	61.95	10.54±11.92	0	9
	<b><math>F_{rtw}</math></b>	0.00	6.70	62.30	10.54±11.86	2	9
	<b><math>F_{tw}</math></b>	0.66	18.61	82.11	22.70±18.00	0	25
	<b><math>F_{ind}</math></b>	0.00	0.00	0.27	0.02±0.05	15	0
	<b><math>F_{rnd}</math></b>	0.01	0.14	3.84	0.33±0.54	0	0
<b><math>G_p@20</math></b>	<b><math>F_{mnt}</math></b>	0.06	11.18	63.54	14.99±14.37	0	14
	<b><math>F_{rtw}</math></b>	0.06	12.44	63.54	15.11±13.01	0	12
	<b><math>F_{tw}</math></b>	0.85	26.95	82.40	30.95±19.49	0	41
	<b><math>F_{ind}</math></b>	0.00	0.01	0.43	0.04±0.08	5	0
	<b><math>F_{rnd}</math></b>	0.01	0.36	3.93	0.58±0.64	0	0
<b><math>G_p@50</math></b>	<b><math>F_{mnt}</math></b>	0.34	19.26	71.89	22.04±16.86	0	25
	<b><math>F_{rtw}</math></b>	0.23	19.60	65.31	22.35±15.25	0	28
	<b><math>F_{tw}</math></b>	4.05	40.68	83.60	43.36±20.06	0	60
	<b><math>F_{ind}</math></b>	0.00	0.04	0.86	0.09±0.15	0	0
	<b><math>F_{rnd}</math></b>	0.03	0.96	10.37	1.61±1.80	0	0
<b><math>G_p@100</math></b>	<b><math>F_{mnt}</math></b>	0.40	26.82	78.49	29.50±19.40	0	39
	<b><math>F_{rtw}</math></b>	1.95	25.43	66.55	27.47±15.92	0	41
	<b><math>F_{tw}</math></b>	15.59	52.16	89.98	52.77±19.73	0	68
	<b><math>F_{ind}</math></b>	0.00	0.07	1.15	0.15±0.21	0	0
	<b><math>F_{rnd}</math></b>	0.12	2.00	19.08	3.01±3.01	0	0

**Πίνακας 5-2: Ποσοστό του πολωμένου περιεχομένου που παράγεται από τους διακεκριμένους χρήστες για όλες τις θεωρίες επιρροής**

- Οι αδρανείς διακεκριμένες ομάδες δε δημοσιεύουν ούτε θετικά ούτε αρνητικά tweets. Συνεπώς δεν ικανοποιούν το πρώτο μέρος της συνθήκης όγκου περιεχομένου μετρικής (δηλ.,  $0 < |M(G_p, ASM, t_1, t_2)|$ ) και δε μπορούν να ασκήσουν πραγματική επιρροή.

- Οι υπερδραστήριες διακεκριμένες ομάδες μονοπωλούν το πολωμένο περιεχόμενο μιας κοινότητας, παράγοντας δυσανάλογα μεγάλο όγκο θετικών και αρνητικών tweets. Η συνθήκη όγκου περιεχομένου μετρικής δεν ικανοποιείται για αυτές, καθώς η επιρροή τους εμπεριέχει μεγάλο κόστος. Στη μελέτη αυτή, ορίσαμε το δεύτερο μέρος της συνθήκης ως εξής:

$$4 \cdot |M(G_p, ASM, t_1, t_2)| \leq |M(G_h, ASM, t_1, t_2)|.$$

Αυτό σημαίνει ότι θεωρούμε μια διακεκριμένη ομάδα ως υπερδραστήρια αν παράγει πάνω από το 25% του συνολικού πολωμένου περιεχομένου. Στις περιπτώσεις αυτές, το εναπομείναν πολωμένο περιεχόμενο μπορεί να είναι ανεπαρκές για την εξαγωγή ασφαλών συμπερασμάτων για το συναίσθημα των υπολοίπων μελών της κοινότητας.

- Οι έγκυρες διακεκριμένες ομάδες δεν είναι ούτε αδρανείς ούτε υπερδραστήριες, ικανοποιώντας συνεπώς τη συνθήκη όγκου περιεχομένου μετρικής.

Τα αποτελέσματα της ανάλυσης αυτής για όλες τις θεωρίες επιρροής παρουσιάζονται στον Πίνακα 5-2. Επισημαίνουμε ότι προκειμένου να εξάγουμε ασφαλή συμπεράσματα για το  $F_{rnd}$ , το εφαρμόσαμε 1,000 φορές για κάθε μέγεθος διακεκριμένης ομάδας  $k$ , χρησιμοποιώντας σε κάθε επανάληψη ένα τυχαίο δείγμα  $k$  χρηστών από κάθε κοινότητα.

Παρατηρούμε ότι οι αδρανείς διακεκριμένες ομάδες είναι αρκετά σπάνιες, καθώς εμφανίζονται μόνο σε δώδεκα θεματικές κατηγορίες για  $k = 10$ . Αυξάνοντας το  $k$ , όμως, αυξάνεται και η συμμετοχή των επιδραστικών χρηστών σε πολωμένο περιεχόμενο. Αυτό έχει ως αποτέλεσμα οι αδρανείς διακεκριμένες ομάδες να περιορίζονται σε μόνο 5 για  $k = 20$  και να εξαφανίζονται εντελώς για μεγαλύτερα μεγέθη. Κατά τον ίδιο τρόπο, η αύξηση του  $k$  αυξάνει και το ποσοστό των υπερδραστήριων ομάδων. Στην περίπτωση των  $F_{mnt}$  και  $F_{rwt}$ , παραμένουν

μειοψηφία για μεγέθη έως 50, αλλά γίνονται η πλειοψηφία για  $k = 100$ . Η σχέση μεταξύ μέσης τιμής και διαμέσου φανερώνει ότι οι κατανομές των θεωριών αυτών έχουν θετική προδιάθεση (προς τα μικρότερα ποσοστά πολωμένου περιεχομένου). Για το  $k = 100$ , ωστόσο, η μέση τιμή και η διάμεσος ξεπερνούν κατά ελάχιστα το όριο του 25%, διαχωρίζοντας ομοιόμορφα τις επιμέρους θεματικές κοινότητες σε έγκυρες και υπερδραστήριες. Αντιθέτως, η πλειοψηφία των διακεκριμένων ομάδων που ορίζονται από την  $F_{twt}$  είναι υπερδραστήριες ήδη από το  $k = 20$  και για μεγαλύτερα μεγέθη, το ποσοστό των έγκυρων διακεκριμένων ομάδων γίνεται λιγότερο από το 1/5. Αυτό είναι σε ένα βαθμό αναμενόμενο, καθώς οι διακεκριμένοι χρήστες της θεωρίας είναι εξορισμού παραγωγικοί με συνέπεια να παράγουν και μεγάλο όγκο πολωμένου περιεχομένου.

Συνεπείς εξαιρέσεις σε αυτά τα πρότυπα αποτελούν οι  $F_{ind}$  και  $F_{rnd}$ : καμία από τις διακεκριμένες ομάδες που ορίζουν σε όλες τις θεωρίες επιρροής και τα μεγέθη διακεκριμένων ομάδων δεν είναι υπερδραστήρια. Συνεπώς ικανοποιούν τη συνθήκη όγκου περιεχομένου μετρικής σε όλες τις περιπτώσεις – με εξαίρεση τις λίγες αδρανείς ομάδες που ορίζει η  $F_{ind}$  για μικρές διακεκριμένες ομάδες.

Συνεχίζουμε την ανάλυσή μας με τη συνθήκη όγκου συνολικού περιεχομένου, η οποία υπολογίζει το ποσοστό όλων των tweets που δημοσιεύονται από τους διακεκριμένους χρήστες. Στην περίπτωση αυτή, στόχος μας είναι να διαχωρίσουμε τις διακεκριμένες ομάδες σε υπερδραστήριες και έγκυρες<sup>2</sup>. Ως υπερδραστήριες θεωρούμε τις διακεκριμένες ομάδες που έχουν δημοσιεύσει πάνω από 25% του συνολικού περιεχομένου της κοινότητας, ενώ οι υπόλοιπες διακεκριμένες ομάδες

---

<sup>2</sup> Επισημαίνουμε ότι δεν έχει νόημα να ορίσουμε αδρανείς διακεκριμένες ομάδες αναφορικά με τη συνθήκη συνολικού περιεχομένου. Εξορισμού, κάθε χρήστης που συμμετέχει σε μία κοινότητα έχει δημοσιεύσει τουλάχιστον ένα σχετικό tweet και συνεπώς κάθε διακεκριμένη ομάδα έχει παράξει ελάχιστο αριθμό tweets ίσο με το μέγεθός της, ανεαρτήτως της θεωρίας επιρροής που την ορίζει.

θεωρούνται έγκυρες. Προφανώς, από τις δύο αυτές κατηγορίες, μόνο οι έγκυρες ομάδες επιδραστικών χρηστών θεωρούνται συμβατές με το πλαίσιο μας.

		Min.	Median	Max.	Mean±SD	Υπερδραστήριες
$G_p@10$	$F_{mnt}$	0	6.41	45.68	9.38±9.26	5
	$F_{rtw}$	0.63	6.85	50.01	9.84±9.58	5
	$F_{twt}$	4.16	20.25	55.54	22.96±11.90	25
	$F_{ind}$	0.00	0.01	0.35	0.02±0.05	0
	$F_{rnd}$	0.01	0.15	2.18	0.28±0.32	0
$G_p@20$	$F_{mnt}$	0.89	10.32	56.89	12.88±11.18	7
	$F_{rtw}$	2.65	10.91	55.14	13.76±10.74	8
	$F_{twt}$	6.88	27.84	65.57	30.50±13.62	42
	$F_{ind}$	0.00	0.01	0.43	0.04±0.08	0
	$F_{rnd}$	0.02	0.35	3.91	0.56±0.60	0
$G_p@50$	$F_{mnt}$	2.41	16.91	70.00	19.59±13.98	17
	$F_{rtw}$	4.66	19.39	63.38	20.85±12.48	21
	$F_{twt}$	12.93	40.07	81.71	42.18±15.35	65
	$F_{ind}$	0.00	0.04	0.86	0.09±0.15	0
	$F_{rnd}$	0.05	0.86	10.28	1.42±1.53	0
$G_p@100$	$F_{mnt}$	4.80	22.95	79.45	26.46±16.18	34
	$F_{rtw}$	6.45	25.68	72.14	26.40±13.09	40
	$F_{twt}$	19.65	49.21	90.96	51.56±15.61	74
	$F_{ind}$	0.00	0.07	0.98	0.12±0.14	0
	$F_{rnd}$	0.13	1.73	19.37	2.89±2.96	0

**Πίνακας 5-3: Ποσοστό συνολικού περιεχομένου που παράγεται από τις διακεκριμένες ομάδες για όλες τις θεωρίες επιρροής.**

Τα αποτελέσματα τις ανάλυσής μας παρουσιάζονται στον Πίνακα 5-3. Επισημαίνουμε ότι η απόδοση της  $F_{rnd}$  υπολογίστηκε χρησιμοποιώντας τα ίδια 1,000 δείγματα  $k$  χρηστών από κάθε θεματική κοινότητα όπως στον Πίνακα 5-2. Παρατηρούμε ότι όλες οι θεωρίες εμφανίζουν παρόμοια συμπεριφορά με τη συνθήκη όγκου περιεχομένου μετρικής: όλες οι διακεκριμένες ομάδες που ορίζονται από τις  $F_{ind}$  και  $F_{rnd}$  είναι έγκυρες για όλες τις θεωρίες επιρροής και τα μεγέθη, ενώ η  $F_{twt}$

κυριαρχείται από υπερδραστήριους χρήστες ήδη για  $k = 20$ . Για τις  $F_{mnt}$  και  $F_{rwt}$ , το ποσοστό των υπερδραστήριων ομάδων αυξάνεται για μεγαλύτερα μεγέθη διακεκριμένων ομάδων, αλλά περιορίζεται σε λιγότερο από 1/3 για  $k$  έως 50. Για  $k = 100$ , όμως, οι διακεκριμένες ομάδες τους κατανέμονται ισομερώς μεταξύ έγκυρων και υπερδραστήριων.

Συνολικά μπορούμε να συμπεράνουμε ότι η  $F_{twt}$  είναι ασύμβατη με το πλαίσιο αξιολόγησής μας: ορίζει υπερδραστήριες ομάδες επιδραστικών χρηστών ακόμα και για μικρά μεγέθη ομάδων, παραβιάζοντας και τις δύο συνθήκες όγκου. Για το λόγο αυτό, την αφαιρούμε από τις αναλύσεις συσχέτισης και χρόνου. Αντίθετα, οι  $F_{rnd}$  και  $F_{ind}$  ικανοποιούν με συνέπεια και τις δύο συνθήκες όγκου για όλες τις θεωρίες και τα μεγέθη. Οι  $F_{mnt}$  και  $F_{rwt}$  παρουσιάζουν παρόμοια συμπεριφορά, ορίζοντας έγκυρες διακεκριμένες ομάδες και για τις δύο συνθήκες όγκου για τις περισσότερες θεματικές κοινότητες. Για να εξάγουμε ασφαλή συμπεράσματα για την επίδοσή τους, αποκλείουμε από τις περαιτέρω αναλύσεις μας τις θεματικές κοινότητες που αντιστοιχούν σε υπερδραστήριες ομάδες για οποιαδήποτε θεωρία επιρροής και για οποιοδήποτε μέγεθος ομάδας. Από τις 75 θεματικές ομάδες που θεωρήσαμε στη μελέτη μας, οι 30 ικανοποιούν με συνέπεια και τις δύο συνθήκες όγκου για τις επιλεγμένες θεωρίες επιρροής. Στη συνέχεια εξετάζουμε αποκλειστικά αυτές τις κοινότητες. Επισημαίνουμε, ωστόσο, ότι η  $F_{ind}$  ορίζει αδρανείς διακεκριμένες ομάδες για κάποιες από αυτές τις κοινότητες: υπάρχουν 10 τέτοιες κοινότητες για  $k = 10$ , ενώ για  $k = 20$  περιορίζονται στις 4. Στις ακόλουθες αναλύσεις, οι κοινότητες αυτές δε λαμβάνονται υπόψη για τον υπολογισμό της απόδοσης της  $F_{ind}$ .

### **5.3 Συνθήκη συσχέτισης**

Στόχος αυτής της ενότητας είναι να διερευνήσει σε ποιο βαθμό το συνολικό συναίσθημα των διακεκριμένων ομάδων ταυτίζεται με αυτό των υπολοίπων μελών

της κοινότητας. Όπως εξηγήσαμε παραπάνω, ποσοτικοποιούμε τη σχέση αυτή μέσω του βαθμού πόλωσης των αντίστοιχων ομάδων χρηστών. Στην πραγματικότητα, στόχος μας είναι να υπολογίσουμε τη συσχέτιση μεταξύ  $PR(G_h)$  και  $PR(G_p)$  σε όλες τις ομάδες των δεδομένων αναφοράς για όλες τις θεωρίες επιρροής και τα μεγέθη διακεκριμένων ομάδων.

Προς αυτή την κατεύθυνση, ορίσαμε αρχικά δύο μεταβλητές για κάθε θεωρία επιρροής  $F_i$  και μέγεθος διακεκριμένης ομάδας  $k_i$ : η  $X_{i,l}$  εκφράζει το βαθμό πόλωσης των αντίστοιχων διακεκριμένων ομάδων για τις 30 θεματικές κοινότητες που εξετάσαμε σε αυτή τη μελέτη, ενώ η  $Y_{i,l}$  αντιπροσωπεύει το βαθμό πόλωσης των υπόλοιπων μελών της κοινότητας. Στη συνέχεια, για κάθε ζεύγος μεταβλητών  $X_{i,l}$  και  $Y_{i,l}$  υπολογίσαμε το συντελεστή συσχέτισης Pearson<sup>3</sup> ( $\rho_{X_{i,l},Y_{i,l}}$ ): όσο ψηλότερη είναι η τιμή του  $\rho_{X_{i,l},Y_{i,l}}$  τόσο μεγαλύτερη είναι η ισχύς της συνθήκης συσχέτισης για την αντίστοιχη θεωρία επιρροής.

Τα αποτελέσματα της ανάλυσης μας παρουσιάζονται στον Πίνακα 5-4. Επισημαίνουμε ότι για την  $F_{rnd}$  και για κάθε μέγεθος διακεκριμένης ομάδας  $k$ , θεωρήσαμε τα ίδια 1,000 τυχαία δείγματα από  $k$  χρήστες για κάθε κοινότητα όπως και στην προηγούμενη ενότητα.

---

<sup>3</sup> Ο συντελεστής συσχέτισης Pearson (Pearson correlation coefficient)  $\rho_{X,Y}$  είναι μια καθιερωμένη μετρική που εκφράζει τη γραμμική εξάρτηση μεταξύ δύο μεταβλητών  $X$  και  $Y$ . Λαμβάνει τιμές στο διάστημα  $[-1,1]$  και οι υψηλότερες απόλυτες τιμές αντιστοιχούν σε ισχυρότερη συσχέτιση μεταξύ  $X$  και  $Y$ . Ουσιαστικά η τιμή  $|\rho_{X,Y}| = 1$  υποδεικνύει πλήρως γραμμική σχέση της μορφής  $X = \alpha \cdot Y + \beta$ , με  $\alpha, \beta \in R$  και  $0 < \alpha$  αν  $\rho_{X,Y} = 1$ , ή  $\alpha < 0$  αν  $\rho_{X,Y} = -1$ .

	$G_p@10$	$G_p@20$	$G_p@50$	$G_p@100$
$F_{mnt}$	0.09	0.40	0.55	0.66
$F_{rtw}$	0.24	0.30	0.64	0.79
$F_{ind}$	0.26	0.49	0.63	0.66
$F_{rnd}$	$0.19 \pm 0.14$	$0.26 \pm 0.15$	$0.35 \pm 0.16$	$0.42 \pm 0.13$

**Πίνακας 5-4. Συσχέτιση Pearson μεταξύ του συνολικού συναισθήματος των διακεκριμένων ομάδων και των υπόλοιπων μελών της κοινότητας.**

Παρατηρούμε ότι υπάρχει θετική συσχέτιση για κάθε συνδυασμό θεωρίας επιρροής και μεγέθους διακεκριμένης ομάδας, που αυξάνει αναλογικά με την αύξηση του  $k$ : όσο αυξάνεται το μέγεθος των διακεκριμένων ομάδων, αυξάνεται και το ποσοστό πολωμένου περιεχομένου που παράγουν και συνεπώς το ομαδικό τους συναίσθημα πλησιάζει σε αυτό ολόκληρης της κοινότητας (περισσότερες λεπτομέρειες σχετικά υπάρχουν στην επόμενη ενότητα). Μπορούμε να συμπεράνουμε, συνεπώς, ότι οι μικρές διακεκριμένες ομάδες μπορούν να προβλέψουν με μικρότερη ακρίβεια τη συνολική διάθεση των υπολοίπων χρηστών.

Συγκρίνοντας τις επιμέρους θεωρίες επιρροής για τα τέσσερα μεγέθη ομάδας, παρατηρούμε ότι ο συντελεστής συσχέτισης για την  $F_{rnd}$  είναι κατά 1/3 χαμηλότερος από ό,τι για τις άλλες θεωρίες επιρροής για τα περισσότερα μεγέθη. Μοναδική εξαίρεση για  $k = 10$ , όπου η  $F_{mnt}$  επιδεικνύει ακόμα χαμηλότερη συσχέτιση. Οι  $F_{rtw}$  και  $F_{ind}$  λαμβάνουν αντίστοιχες τιμές στις περισσότερες περιπτώσεις, με την  $F_{mnt}$  να ακολουθεί με μικρή απόσταση. Η συνολικά υψηλότερη τιμή αντιστοιχεί στην  $F_{rtw}$  και είναι πολύ κοντά στο 1, υποδηλώνοντας μια σχεδόν γραμμική σχέση ανάμεσα στο ομαδικό συναίσθημα των διακεκριμένων χρηστών και αυτό της υπόλοιπης κοινότητας.

Στην επόμενη ενότητα, εξετάζουμε αν αυτές οι συσχετίσεις προέρχονται από μηνύματα των διακεκριμένων χρηστών που προηγούνται χρονικά αυτών των υπολοίπων μελών.

## 5.4 Συνθήκη χρόνου

Στην ενότητα αυτή, αξιολογούμε τη συνθήκη χρόνου χρησιμοποιώντας τις 30 έγκυρες θεματικές κοινότητες για όλες τις θεωρίες επιρροής και τα μεγέθη διακεκριμένων ομάδων. Υπενθυμίζουμε ότι η συνθήκη αυτή εκφράζει το σχετικό ρυθμό παραγωγής πολωμένου περιεχομένου μεταξύ των διακεκριμένων ομάδων και των υπόλοιπων μελών της κοινότητας και ορίζεται ως εξής:

$$CAC(G_p, ASM, t_1, t_2) \gg CAC(G_h, ASM, t_1, t_2),$$

όπου το  $CAC(G_p, ASM, t_1, t_2)$  αντιστοιχεί στην τιμή CAC για τη διακεκριμένη ομάδα και το  $CAC(G_h, ASM, t_1, t_2)$  στην τιμή του CAC για την υπόλοιπη κοινότητα. Για να ποσοτικοποιήσουμε αυτή τη γενική σχέση, ορίζουμε μια νέα μετρική που ονομάζεται *απόκλιση CAC* ( $\Delta CAC$ ), και εκφράζει τη διαφορά ανάμεσα στις δύο μετρικές:

$$\Delta CAC = CAC(G_p, ASM, t_1, t_2) - CAC(G_h, ASM, t_1, t_2).$$

Προφανώς, η συνθήκη χρόνου ικανοποιείται μόνο για θετική απόκλιση CAC ( $\Delta CAC > 0$ ), ενώ οι υψηλότερες τιμές είναι ενδεικτικές του φαινομένου πρώιμης υιοθέτησης (early adopter) για τους επιδραστικούς χρήστες.

Η κατανομή του  $\Delta CAC$  ανά συνδυασμό θεωρίας επιρροής και μεγέθους διακεκριμένης ομάδας παρουσιάζεται στον Πίνακα 5-5. Επισημαίνουμε ότι η πιο δεξιά στήλη με το όνομα *Θετικές* εκφράζει τον αριθμό των κοινοτήτων για τις οποίες το  $\Delta CAC$  παίρνει θετικές τιμές.

Για την  $F_{rnd}$  χρησιμοποιήσαμε πάλι τα ίδια 1,000 τυχαία δείγματα ανά θεματική κοινότητα και μέγεθος. Παρατηρούμε ότι η διάμεσος και η μέση τιμή  $\Delta CAC$  παίρνουν



αρνητικές τιμές για κάθε περίπτωση. Το ίδιο ισχύει και για τη μέγιστη τιμή για  $k = 10$  και  $k = 20$ . Κατά συνέπεια, για τα μεγέθη αυτά, η στήλη *Θετικές* είναι μηδενική, κάτι που σημαίνει ότι τη πολωμένη δραστηριότητα των επιδραστικών χρηστών που ορίζονται από την  $F_{rnd}$  έπεται χρονικά αυτής των συνηθισμένων χρηστών για όλες τις θεματικές κατηγορίες. Κατά συνέπεια, οι χρήστες αυτοί δεν ασκούν πραγματικοί επιρροή στους υπόλοιπους. Στην πραγματικότητα, κοιτάζοντας τη στήλη *Θετικές* διαπιστώνουμε ότι η  $F_{rnd}$  ορίζει πραγματικές διακεκριμένες ομάδες μόνο για μεγάλα μεγέθη ομάδας και για ελάχιστες κοινότητες. Ακόμα και σε αυτές τις περιπτώσεις, όμως, οι τιμές του  $\Delta CAC$  είναι πολύ χαμηλές, το οποίο σημαίνει ότι η χρονική διαφορά είναι πολύ περιορισμένη.

		Min	Median	Max.	Mean $\pm$ SD	Θετικές
$G_p@10$	$F_{mnt}$	-0.079	0.017	0.078	$0.012 \pm 0.029$	23
	$F_{rtw}$	-0.079	0.013	0.043	$0.001 \pm 0.024$	21
	$F_{ind}$	-0.042	0.021	0.085	$0.017 \pm 0.028$	16
	$F_{rnd}$	-0.529	-0.395	-0.076	$-0.372 \pm 0.123$	0
$G_p@20$	$F_{mnt}$	-0.079	0.015	0.051	$0.011 \pm 0.025$	24
	$F_{rtw}$	-0.079	0.010	0.040	$0.004 \pm 0.027$	20
	$F_{ind}$	-0.042	0.018	0.058	$0.017 \pm 0.028$	20
	$F_{rnd}$	-0.441	-0.243	-0.012	$-0.225 \pm 0.120$	0
$G_p@50$	$F_{mnt}$	-0.034	0.010	0.038	$0.009 \pm 0.018$	22
	$F_{rtw}$	-0.041	0.010	0.029	$0.004 \pm 0.020$	22
	$F_{ind}$	-0.072	0.017	0.074	$0.015 \pm 0.028$	22
	$F_{rnd}$	-0.258	-0.051	0.002	$-0.068 \pm 0.069$	3
$G_p@100$	$F_{mnt}$	-0.037	0.008	0.031	$0.007 \pm 0.017$	23
	$F_{rtw}$	-0.047	0.002	0.027	$-0.001 \pm 0.017$	19
	$F_{ind}$	-0.020	0.016	0.057	$0.016 \pm 0.020$	23
	$F_{rnd}$	-0.085	-0.003	0.007	$-0.014 \pm 0.024$	7

**Πίνακας 5-5: Τιμές  $\Delta CAC$  για τις 30 επιλεγμένες κοινότητες**

Όσον αφορά τις υπόλοιπες θεωρίες επιρροής, η τελευταία στήλη δείχνει ότι ορίζουν με συνέπεια πραγματικές διακεκριμένες ομάδες, για τουλάχιστον τα 2/3 από τις κοινότητες που εξετάζουμε. Η μόνη εξαίρεση είναι η  $F_{ind}$  για  $k=10$ , η οποία ορίζει

πραγματικά επιδραστικούς χρήστες για 16 μόνο θεματικές κατηγορίες. Αυτό οφείλεται, ωστόσο, κατά κύριο λόγο στην αδράνεια των διακεκριμένων χρηστών που παρατηρείται σε 10 από τις κοινότητες αυτές. Παρατηρούμε, επίσης, ότι καθώς αυξάνεται το μέγεθος της διακεκριμένης ομάδας, η μέγιστη τιμή για τις  $F_{mnt}$ ,  $F_{rtw}$  και  $F_{ind}$  μειώνεται ενώ η ελάχιστη τιμή αυξάνεται. Αυτό σημαίνει ότι όσο μεγαλύτερες είναι οι διακεκριμένες ομάδες, τόσο μικρότερη είναι η μεταβολή της ΔCAC μεταξύ των 30 θεματικών ομάδων και άρα η συμπεριφορά τους γίνεται πιο συνεπής. Αξίζει επίσης να επισημάνουμε τη μονότονη μείωση της μέσης τιμής και της διαμέσου της ΔCAC καθώς αυξάνεται το μέγεθος της διακεκριμένης ομάδας. Η συμπεριφορά αυτή δείχνει ότι όσο μεγαλύτερες είναι οι διακεκριμένες ομάδες που ορίζονται από τις  $F_{mnt}$ ,  $F_{rtw}$  και  $F_{ind}$ , τόσο μειώνεται η απόδοσή τους στην πρόβλεψη του συνολικού συναισθήματος της κοινότητας. Αξίζει τέλος να αναφέρουμε ότι η  $F_{ind}$  έχει σταθερά καλύτερη απόδοση από τις  $F_{mnt}$  και  $F_{rtw}$  για όλα τα μεγέθη διακεκριμένης ομάδας καθώς επιτυγχάνει υψηλότερες τιμές για όλες τις στατιστικές μετρικές.

## 5.5 Συζήτηση

Η παραπάνω ανάλυση προσφέρει σημαντικά συμπεράσματα για την απόδοση των επιλεγμένων θεωριών επιρροής.

Καταρχήν φανερώνει ότι η  $F_{ind}$  αποτυγχάνει να αποτυπώσει πραγματική άσκηση επιρροής. Ο βασικός λόγος είναι ότι αγνοεί τη δραστηριότητα των επιμέρους χρηστών και ορίζει τους επιδραστικούς χρήστες χωρίς να λαμβάνει υπόψη οποιαδήποτε ποιοτική πληροφορία. Η συμπεριφορά της σε σχέση με τη συνθήκη χρόνου αποδεικνύει ότι θα είναι λάθος να αποδώσουμε τις θετικές συσχετίσεις που παρατηρούνται στον Πίνακα 5-4 στο φαινόμενο των πρώιμων οπαδών (early adopter). Αντίθετα, πρόκειται απλά για ένδειξη του ότι οι διακεκριμένοι χρήστες που

ορίζονται από την  $F_{rnd}$  ακολουθούν την πολωμένη δραστηριότητα του περίγυρού τους.

Δεύτερον, η ανάλυσή μας επιβεβαιώνει ότι οι  $F_{mnt}$ ,  $F_{rtw}$  και  $F_{ind}$  επιδεικνύουν πραγματική επιρροή για τις περισσότερες συνθήκες. Ουσιαστικά, οι συνθήκες του πλαισίου μας αποκαλύπτουν τα εξής πρότυπα συμπεριφορών: όσο αυξάνεται το μέγεθος μιας διακεκριμένης ομάδας, τόσο αυξάνεται η ποσότητα του περιεχομένου που παράγουν και τόσο πλησιάζει το πολωμένο περιεχόμενο που εκφράζουν στο συνολικό συναίσθημα της κοινότητας. Ταυτόχρονα όμως, παρατηρείται καθυστέρηση στο χρόνο πολωμένης δραστηριότητας, καθώς οι αντίστοιχες διακεκριμένες ομάδες είναι λιγότερο συνεπείς χρονικά στην πρόβλεψη της διάθεσης της κοινότητας.

Τρίτον, η ανάλυση μας τονίζει τη σχετική απόδοση των τριών αυτών θεωριών επιρροής. Οι επιδραστικοί χρήστες που ορίζει η  $F_{rtw}$  καταφέρνουν σε μεγάλο βαθμό να προσεγγίσουν το συνολικό συναίσθημα της κοινότητας, ενώ αυτοί της  $F_{ind}$  επιτυγχάνουν προβλέψεις με καλύτερη χρονική απόδοση. Και στις δύο περιπτώσεις, η  $F_{mnt}$  ακολουθεί με μικρή απόσταση. Αξίζει να επισημάνουμε ότι η  $F_{ind}$  καταφέρνει να έχει ανταγωνιστική επίδοση παρά το χαμηλότερο (κατά δύο τάξεις μεγέθους) όγκο πολωμένου και συνολικού περιεχομένου που παράγουν οι επιδραστικοί χρήστες της σε σχέση με αυτούς των  $F_{mnt}$  και  $F_{rtw}$ . Η επίδοση αυτή αποδίδεται στο μεγάλο κοινό των επιδραστικών χρηστών της οι οποίοι, αν και παράγουν ένα αμελητέο αριθμό πολωμένων μηνυμάτων, είναι ικανοί να επηρεάσουν μεγάλο μέρος της κοινότητας.

Τέλος η ανάλυσή μας δείχνει ότι η  $F_{twt}$  είναι ασύμβατη με το πλαίσιο αξιολόγησής μας: για τις περισσότερες θεματικές κοινότητες, η τεράστια δραστηριότητα των επιδραστικών χρηστών που ορίζει έχει σαν αποτέλεσμα να παραβιάζονται και οι δύο συνθήκες όγκου, ακόμα και για μικρά μεγέθη διακεκριμένων ομάδων.

Η σελίδα αυτή είναι σκόπιμα λευκή.

# 6

## *Εσωτερική ανάλυση θεωριών*

### *Θεματικής επιρροής*

Στην παρούσα ενότητα, εισάγουμε μια μεθοδολογία για την εσωτερική ανάλυση των τοπικών θεωριών επιρροής, που μας επιτρέπει δηλαδή να εξετάσουμε τη δυναμική μεταξύ των χρηστών στις διακεκριμένες ομάδες που ορίζουν. Η μεθοδολογία εξετάζει τρεις επιμέρους παραμέτρους των σχέσεων μεταξύ των διακεκριμένων χρηστών:

- την ομοφιλία τους, όπως συνεπάγεται από τους αμοιβαίους συνδέσμους μεταξύ τους,
- την επικοινωνία τους, με βάση τη συχνότητα των μεταξύ τους αλληλεπιδράσεων, και
- την πολυπλευρικότητά τους, με κριτήριο την επικάλυψη μεταξύ των διακεκριμένων ομάδων διαφορετικών θεωριών.

Ο απώτερος σκοπός της μεθοδολογίας αυτής είναι διπλός: (i) να εξηγήσει τη σχετική απόδοση των θεωριών επιρροής την οποία παρατηρήσαμε στην προηγούμενη ενότητα και (ii) να εκτιμήσει το βέλτιστο μέγεθος διακεκριμένης ομάδας ώστε να επιτυγχάνεται υψηλή απόδοση για τις περισσότερες θεωρίες επιρροής.

## 6.1 Ομοφιλία

Οι χρήστες του Twitter είναι ελεύθεροι να δημοσιεύσουν μηνύματα πάνω σε οποιοδήποτε θέμα, ανά πάσα στιγμή και με τον τρόπο αυτό σχηματίζουν κοινότητες αυθόρμητα και στιγμιαία. Αυτό έχει ως αποτέλεσμα τα μέλη μιας κοινότητας να μην συνδέονται απαραίτητα με σαφείς ακμές πάνω στον κοινωνικό γράφο. Στην πράξη, οι σύνδεσμοι που ανήκουν στην  $E_h$  το πιθανότερο είναι να αντιπροσωπεύουν σχέσεις που σχηματίστηκαν εκτός της κοινότητας  $G_h$ , κάτι που υποδηλώνει ανθρώπους που έχουν περισσότερα κοινά μεταξύ τους από το ενδιαφέρον τους στο θέμα  $h$ . Επιπλέον, οι αμοιβαίοι σύνδεσμοι<sup>4</sup> υποδηλώνουν χρήστες με ακόμα υψηλότερα επίπεδα ομοφιλίας, δηλαδή με παρόμοιο προσωπικό υπόβαθρο σε ουσιαστικά θέματα από κοινωνικής άποψης [5].

Στόχος της παρούσας ενότητας είναι να εξετάσει τα επίπεδα ομοφιλίας που μοιράζονται οι επιδραστικοί χρήστες κάθε θεωρίας επιρροής. Μπορούμε να εξάγουμε την πληροφορία αυτή από την τοπολογία της  $E_p$ : πυκνά συνδεδεμένοι υπογράφοι υποδεικνύουν σημαντικές ομοιότητες μεταξύ των χρηστών και αντίστροφα. Βάσει των παραπάνω, για όλες τις θεωρίες επιρροής και τα μεγέθη διακεκριμένων ομάδων εκτιμήσαμε το βαθμό αμοιβαιότητας των αντίστοιχων διακεκριμένων ομάδων υπολογίζοντας το ποσοστό των ζευγών διακεκριμένων χρηστών που συνδέονται αμοιβαία. Συγκεκριμένα, ο βαθμός αμοιβαιότητας μιας διακεκριμένης ομάδας  $G_p$  ορίζεται ως εξής:

$$\begin{aligned} \text{reciprocity}(G_p) &= \frac{\{ \langle x, y \rangle \in E_p : \langle y, x \rangle \in E_p \} / 2}{|G_p| \cdot (|G_p| - 1) / 2} \cdot 100\% \\ &= \frac{\{ \langle x, y \rangle \in E_p : \langle y, x \rangle \in E_p \}}{|G_p| \cdot (|G_p| - 1)} \cdot 100\%, \end{aligned}$$

---

<sup>4</sup>Στο Twitter, δύο χρήστες είναι αμοιβαία συνδεδεμένοι αν ακολουθούν ο ένας τον άλλον.

όπου ο αριθμητής υπολογίζει τον αριθμό των υπάρχοντων αμοιβαίων συνδέσμων, ενώ ο παρονομαστής ισούται με το συνολικό αριθμό των δυνατών αμοιβαίων συνδέσμων μέσα στην  $G_p$  (δηλαδή όλους τους πιθανούς συνδυασμούς ζευγών διακεκριμένων χρηστών). Ο βαθμός αμοιβαιότητας λαμβάνει τιμές στο διάστημα  $[0, 100]$ , με υψηλές τιμές να αντιστοιχούν σε υψηλά ποσοστά αμοιβαίων συνδέσμων. Τα αποτελέσματα της ανάλυσης μας συνοψίζονται στον Πίνακα 6-1.

		Min.	Median	Max.	Mean±SD
$G_p@10$	$F_{mnt}$	0.00	19.44	75.56	23.40±18.17
	$F_{rtw}$	0.00	22.22	71.11	24.29±18.44
	$F_{twt}$	0.00	8.33	80.56	14.77±18.01
	$F_{ind}$	0.00	0.00	17.86	0.44±2.32
	$F_{rnd}$	0.00	5.00	32.14	6.33±5.75
$G_p@20$	$F_{mnt}$	0.65	14.29	66.84	19.49±15.33
	$F_{rtw}$	0.00	15.03	63.68	19.11±15.74
	$F_{twt}$	0.00	8.09	64.74	12.70±13.93
	$F_{ind}$	0.00	0.00	8.09	0.39±1.25
	$F_{rnd}$	0.18	5.04	34.25	6.58±5.83
$G_p@50$	$F_{mnt}$	0.87	9.41	50.32	13.03±10.69
	$F_{rtw}$	0.44	10.18	41.75	12.78±10.32
	$F_{twt}$	0.00	6.78	57.82	9.29±10.11
	$F_{ind}$	0.00	0.13	8.42	0.61±1.31
	$F_{rnd}$	0.65	5.08	29.28	6.43±5.26
$G_p@100$	$F_{mnt}$	0.56	7.13	39.76	9.43±8.31
	$F_{rtw}$	0.42	5.93	33.82	8.79±8.10
	$F_{twt}$	0.51	4.76	44.17	6.62±7.10
	$F_{ind}$	0.00	0.18	9.28	0.60±1.28
	$F_{rnd}$	0.49	5.26	31.32	6.33±5.33

**Πίνακας 6-1. Βαθμός αμοιβαιότητας διακεκριμένων ομάδων για όλες τις θεωρίες επιρροής και τα μεγέθη ομάδων.**

Επισημαίνουμε ότι για την  $F_{rnd}$  και για κάθε μέγεθος διακεκριμένης ομάδας, επαναλάβαμε τις μετρήσεις πάνω στα 1,000 τυχαία δείγματα που χρησιμοποιήσαμε

και στην Ενότητα 4. Παρατηρούμε ότι η συμπεριφορά τους παραμένει σχετικά σταθερή για όλες τις τιμές του  $k$ : η τιμή διαμέσου του βαθμού αμοιβαιότητας βρίσκεται κοντά στο 5% και η μέση τιμή κοντά στο 6% για όλες τις περιπτώσεις. Συνεπώς οι διακεκριμένοι χρήστες που ορίζει η  $F_{rnd}$  παρουσιάζουν αμελητέα επίπεδα ομοφιλίας.

Σχετικά με την  $F_{ind}$ , παρατηρούμε ότι οι διακεκριμένες ομάδες που ορίζει παρουσιάζουν σημαντικά χαμηλότερα επίπεδα αμοιβαιότητας από αυτές της  $F_{rnd}$ , κάτι που αποτελεί σαφή ένδειξη ότι οι πιο δημοφιλείς χρήστες μιας θεματικής κοινότητας έχουν τυπικά ελάχιστα κοινό υπόβαθρο.

Αντιθέτως, οι τιμές του βαθμού αμοιβαιότητας για τις  $F_{mnt}$ ,  $F_{rtw}$ , and  $F_{twt}$  είναι σημαντικά υψηλότερες από αυτές της  $F_{rnd}$ . Οι μέγιστες τιμές αντιστοιχούν στις  $F_{mnt}$  και  $F_{rtw}$ , με την  $F_{twt}$  να ακολουθεί με μικρή διαφορά. Σε όλες τις περιπτώσεις, η αμοιβαιότητα λαμβάνει τα υψηλότερα επίπεδα για  $k = 10$  και μειώνεται αναλογικά με την αύξηση του μεγέθους της διακεκριμένης ομάδας. Αυτό σημαίνει ότι η τετραγωνική αύξηση του παρονομαστή (τα ζεύγη των διακεκριμένων ομάδων) είναι μεγαλύτερη από την αύξηση του αριθμητή (αμοιβαία συνδεδεμένοι διακεκριμένοι χρήστες). Για  $k = 100$ , ο βαθμός αμοιβαιότητας των τριών θεωριών πλησιάζει αρκετά τις τιμές της  $F_{rnd}$ , ειδικά στις τιμές της διαμέσου. Αυτό υποδεικνύει ότι τα μέλη των μεγάλων διακεκριμένων ομάδων παρουσιάζουν κατά μέσο όρο τα ίδια επίπεδα ομοφιλίας με οποιοδήποτε τυχαίο ζεύγος χρηστών. Μπορούμε συνεπώς να θεωρήσουμε το μέγεθος διακεκριμένης ομάδας  $k = 50$  ως το κριτικό όριο πάνω από το οποίο η έννοια της διακεκριμένης ομάδας εκφυλίζεται.

Συνολικά, μπορούμε να συμπεράνουμε ότι οι διακεκριμένες ομάδες των  $F_{mnt}$ ,  $F_{rtw}$ , and  $F_{twt}$  παρουσιάζουν υψηλά επίπεδα ομοφιλίας για μεγέθη έως 50, σχηματίζοντας



διακεκριμένες ομάδες με πυκνά συνδεδεμένους γράφους. Αυτό δεν ισχύει ούτε για την  $F_{rnd}$  ούτε για την  $F_{ind}$  – ανεξαρτήτως μεγέθους διακεκριμένης ομάδας.

## 6.2 Επικοινωνία

Η ενότητα αυτή επιχειρεί να μετρήσει την επικοινωνία μεταξύ των διακεκριμένων χρηστών, το επίπεδο, δηλαδή, στο οποίο γνωρίζονται και διατηρούν κοινωνικές σχέσεις. Ως μέσο αλληλεπίδρασης θεωρούμε την αναφορά από ένα χρήστη σε έναν άλλον. Στο πλαίσιο του Twitter, αυτό γίνεται αποκλειστικά μέσω των επισημασμένων tweets δηλαδή των αναφορών (mentions) και των αναδημοσιεύσεων (retweets). Μπορούμε να θεωρήσουμε ότι όσο περισσότερο ένας χρήστης του Twitter ‘αναφέρεται’ σε έναν άλλον, τόσο μεγαλύτερη είναι η επικοινωνία μεταξύ τους. Στόχος μας, συνεπώς, είναι να εξετάσουμε αν οι διακεκριμένοι χρήστες είναι πιο πιθανό να ‘αναφέρονται’ στους υπόλοιπους επιδραστικούς χρήστες από ό,τι στα υπόλοιπα μέλη της ίδιας κοινότητας. Για να ποσοτικοποιήσουμε αυτή την ιδέα χρησιμοποιούμε τις παρακάτω μετρικές:

**Ορισμός 6.1.** Έστω μια θεματική κοινότητα,  $G_h$ , και η διακεκριμένη ομάδα της,  $G_p$ . Η πιθανότητα εσωτερικής αναφοράς ( $P_{im}$ ) εκφράζει την πιθανότητα ένας διακεκριμένος χρήστης της  $G_p$  να αναφερθεί σε ένα άλλο μέλος της διακεκριμένης ομάδας. Ορίζεται ως εξής:

$$P_{im}(G_p, G_h) = \frac{\text{mentions}(G_p, G_p)}{\text{mentions}(G_p, G_h)} \cdot 100\%,$$

όπου  $\text{mentions}(x, y)$  μια συνάρτηση που επιστρέφει τον αριθμό των αναφορών από χρήστες της ομάδας  $x$  σε χρήστες της ομάδας  $y$ .

**Ορισμός 6.2.** Έστω μια θεματική κοινότητα,  $G_h$ , και η διακεκριμένη ομάδα της,  $G_p$ . Η **πιθανότητα εσωτερικής αναδημοσίευσης** ( $P_{ir}$ ) εκφράζει την πιθανότητα ένας διακεκριμένος χρήστης της  $G_p$  να αναδημοσιεύσει ένα μήνυμα που αρχικά δημοσίευσε ένα άλλο μέλος της διακεκριμένης ομάδας. Ορίζεται ως εξής:

$$P_{ir}(G_p, G_h) = \frac{\text{retweets}(G_p, G_p)}{\text{retweets}(G_p, G_h)} \cdot 100\%,$$

όπου  $\text{retweets}(x, y)$  μια συνάρτηση που επιστρέφει τον αριθμό των μηνυμάτων που δημοσιεύτηκαν αρχικά από χρήστες της ομάδας  $y$  και αναδημοσιεύτηκαν από χρήστες της ομάδας  $x$ .

Και οι δύο μετρικές ορίζονται στο διάστημα  $[0,100]$ , με τις υψηλότερες τιμές να αναλογούν σε συχνότερες αλληλεπιδράσεις μεταξύ επιδραστικών χρηστών. Συνεπώς, όσο πιο υψηλές είναι αυτές οι πιθανότητες, τόσο μεγαλύτερη η επικοινωνία μεταξύ των επιδραστικών χρηστών.

Υπολογίσαμε τις πιθανότητες αυτές για όλες τις θεματικές κοινότητες για τις γνωστές θεωρίες επιρροής. Για την  $F_{rnd}$ , εξετάσαμε τις ίδιες 1,000 ομάδες  $k$  τυχαίων χρηστών ανά θεματική κατηγορία. Επισημαίνουμε ότι, εξορισμού, υπάρχει ισχυρή συσχέτιση μεταξύ  $P_{im}$  και  $F_{mnt}$  όπως και μεταξύ  $P_{ir}$  και  $F_{rtw}$ : οι επιδραστικοί χρήστες με υψηλό ποσοστό αναφορών (αναδημοσιεύσεων) είναι αναμενόμενο να λαμβάνουν συχνά αναφορές (αναδημοσιεύσεις) και από τους υπόλοιπους επιδραστικούς χρήστες, οδηγώντας μας συνεπώς σε ταυτόσημη παρατήρηση. Παρόλα αυτά, συμπεριλαμβάνουμε τους συνδυασμούς αυτούς στην ανάλυσή μας ως μια ένδειξη των υψηλότερων δυνατών πιθανοτήτων.

Τα αποτελέσματα της ανάλυσής μας παρουσιάζονται στον Πίνακα 6-2. Βλέποντας τα αποτελέσματα, παρατηρούμε ότι η  $F_{rnd}$  παρουσιάζει πολύ χαμηλές τιμές και για τις δύο πιθανότητες για όλα τα μεγέθη διακεκριμένης ομάδας: η διάμεσος είναι ίση ή σχεδόν ίση με 0, ενώ η μέση τιμή αυξάνει αναλογικά με το  $k$ , παραμένει όμως χαμηλότερη από 3% σε όλες τις περιπτώσεις. Η συμπεριφορά αυτή φανερώνει αμελητέα επίπεδα αλληλεπίδρασης μεταξύ τυχαία επιλεγμένων διακεκριμένων χρηστών, καθώς οι τελευταίοι είναι πιθανότερο να αναφερθούν σε μέλη της κοινότητας εκτός της διακεκριμένης ομάδας στην οποία ανήκουν.

Στο άλλο άκρο βρίσκονται οι  $F_{mnt}$  και  $F_{rtw}$ , οι οποίες εμφανίζουν σημαντικά υψηλότερες τιμές και για τις δύο πιθανότητες για όλα τα μεγέθη διακεκριμένης ομάδας. Σε όλες τις περιπτώσεις, η διάμεσος είναι σχεδόν ίση με τη μέση τιμή, το οποίο είναι ενδεικτικό κανονικής κατανομής στις 75 κοινότητες. Για την  $P_{im}$ , η  $F_{rtw}$  παίρνει τιμές 10% με 20% χαμηλότερες από τις μέγιστες (δηλαδή αυτές της ταυτολογίας της  $F_{mnt}$ ). Παρόμοια συμπεριφορά έχει και η  $P_{ir}$ , όπου η  $F_{mnt}$  πλησιάζει στο 10% τις τιμές της ταυτολογίας της  $F_{rtw}$ . Σε κάθε περίπτωση, όσο αυξάνεται το μέγεθος της διακεκριμένης ομάδας, τόσο αυξάνονται και οι αντίστοιχες πιθανότητες. Συγκεκριμένα, για  $k = 10$  και  $k = 20$ , η διάμεσος και η μέση πιθανότητα παίρνουν τιμές χαμηλότερες από 50%, συνιστώντας έτσι μέτρια επίπεδα επικοινωνίας μεταξύ των διακεκριμένων χρηστών. Τα μεγαλύτερα μεγέθη διακεκριμένης ομάδας, ωστόσο, παρουσιάζουν σημαντικά υψηλότερες τιμές διαμέσου και μέσης πιθανότητας, που ξεπερνούν σημαντικά το 50%. Αυτό υποδηλώνει αρκετά ισχυρή επικοινωνία, καθώς οι διακεκριμένοι χρήστες είναι πολύ πιθανότερο να αναφερθούν σε έναν άλλο επιδραστικό χρήστη παρά σε κάποιο άλλο μέλος της κοινότητας.

Ανάμεσα στα δύο αυτά άκρα βρίσκονται οι  $F_{twt}$  και  $F_{ind}$ , που παρουσιάζουν παρόμοια συμπεριφορά και για τις δύο πιθανότητες – ανεξάρτητα από το μέγεθος διακεκριμένης ομάδας. Για  $k = 10$ , η συμπεριφορά τους πλησιάζει αυτήν της  $F_{rnd}$  καθώς διάμεσος και μέση τιμή είναι χαμηλότερες από το ήμισυ της διαφοράς μεταξύ  $F_{rnd}$  και  $F_{mnt}/F_{rtw}$ . Για μεγαλύτερες τιμές  $k$ , όμως, οι δεσμοί επικοινωνίας μεταξύ των διακεκριμένων χρηστών ενισχύονται σημαντικά, γεφυρώνοντας σε ένα βαθμό τη διαφορά τους με τις  $F_{mnt}$  και  $F_{rtw}$ . Παρόλα αυτά, οι τιμές τους τόσο για την  $P_{im}$  όσο και για την  $P_{ir}$  σπάνια υπερβαίνουν το 50%, το οποίο σημαίνει ότι οι επιδραστικοί χρήστες που ορίζουν είναι πιθανότερο να αλληλεπιδρούν με την υπόλοιπη κοινότητα παρά μεταξύ τους. Μπορούμε, συνεπώς, να συμπεράνουμε ότι, κατά μέσο όρο, οι δύο αυτές θεωρίες παρουσιάζουν πρότυπα επικοινωνίας μέτριου βαθμού για μεγάλα μεγέθη διακεκριμένων ομάδων.

Για να κατανοήσουμε τις διακυμάνσεις των  $P_{im}$  και  $P_{ir}$  ανά κοινότητα, εξετάζουμε πώς επηρεάζονται από το μέγεθος της κοινότητας. Γενικά, περιμένουμε ότι όσο μεγαλύτερος είναι ο αριθμός των μελών μιας κοινότητας, τόσο μικρότερη είναι η πιθανότητα να υπάρχει αλληλεπίδραση μεταξύ δύο διακεκριμένων χρηστών. Για να εξετάσουμε αυτή την υπόθεση, υπολογίσαμε το βαθμό συσχέτισης Pearson μεταξύ του μεγέθους της κοινότητας και των πιθανοτήτων  $P_{im}$  και  $P_{ir}$  για όλες τις θεωρίες επιρροής και τα μεγέθη επιδραστικής ομάδας. Οι τιμές των αποτελεσμάτων είναι, όπως αναμενόταν, σαφώς αρνητικές: ο βαθμός συσχέτισης διακυμαίνεται μεταξύ -0.20 και -0.45, με τις χαμηλότερες τιμές να αντιστοιχούν σε μεγαλύτερες διακεκριμένες ομάδες. Το φαινόμενο αυτό παρατηρείται σε όλες τις θεωρίες επιρροής, τόσο για την  $P_{im}$  όσο και για την  $P_{ir}$ , και είναι ιδιαίτερα έντονο για τις  $F_{ind}$  και  $F_{rnd}$ .

	$P_{im}$				$P_{ir}$			
	Min.	Med.	Max.	Mean±SD	Min.	Med.	Max.	Mean±SD
<b><math>G_p@10</math></b>								
$F_{mnt}$	0.00	32.97	88.00	35.63±21.75	0.00	24.36	72.49	29.11±18.62
$F_{rtw}$	0.00	24.29	87.37	26.39±19.10	<b>0.00</b>	<b>31.56</b>	<b>73.18</b>	<b>34.08±20.47</b>
$F_{twt}$	0.00	8.05	91.13	15.03±18.43	0.00	12.08	80.39	16.81±16.38
$F_{ind}$	0.00	6.45	100.00	15.56±21.17	0.00	11.11	100.00	18.69±19.90
$F_{rnd}$	0.00	0.00	2.05	0.18±0.37	0.00	0.02	1.20	0.16±0.28
<b><math>G_p@20</math></b>								
$F_{mnt}$	<b>0.00</b>	<b>49.17</b>	<b>92.37</b>	<b>48.72±21.63</b>	0.00	41.76	87.49	40.79±21.46
$F_{rtw}$	0.00	36.09	92.51	35.74±20.11	<b>0.00</b>	<b>48.51</b>	<b>85.46</b>	<b>48.66±19.77</b>
$F_{twt}$	0.00	17.24	91.44	21.09±19.75	0.00	24.05	70.67	27.12±18.23
$F_{ind}$	0.00	15.09	90.91	21.57±21.47	0.00	22.03	76.47	25.45±18.75
$F_{rnd}$	0.00	0.19	2.51	0.41±0.52	0.00	0.27	1.99	0.40±0.44
<b><math>G_p@50</math></b>								
$F_{mnt}$	<b>22.14</b>	<b>66.43</b>	<b>97.76</b>	<b>64.94±18.69</b>	15.30	58.48	91.59	56.04±19.52
$F_{rtw}$	0.00	53.38	94.98	50.61±21.23	<b>25.32</b>	<b>68.66</b>	<b>96.27</b>	<b>66.65±17.11</b>
$F_{twt}$	0.00	28.62	95.24	33.49±23.61	0.00	39.29	87.01	40.44±22.07
$F_{ind}$	0.61	29.17	87.06	33.38±17.70	0.00	38.69	85.71	38.26±20.76
$F_{rnd}$	0.00	0.78	9.40	1.24±1.50	0.00	0.91	8.40	1.30±1.39
<b><math>G_p@100</math></b>								
$F_{mnt}$	<b>32.98</b>	<b>77.11</b>	<b>99.49</b>	<b>75.26±16.36</b>	22.05	68.18	96.67	65.71±17.75
$F_{rtw}$	0.16	66.29	98.52	61.42±20.93	<b>47.85</b>	<b>77.51</b>	<b>100.00</b>	<b>77.96±14.61</b>
$F_{twt}$	0.17	38.45	98.25	42.26±24.36	3.23	52.70	93.93	51.22±21.71
$F_{ind}$	10.74	41.18	92.58	42.31±18.05	11.46	47.27	92.05	48.65±19.13
$F_{rnd}$	0.11	1.57	18.36	2.34±2.66	0.12	1.39	17.49	2.42±2.66

**Πίνακας 6-2. Πιθανότητα εσωτερικής αναφοράς (αριστερά) και πιθανότητα εσωτερικής αναδημοσίευσης (δεξιά) για όλες τις θεματικές κατηγορίες σε όλες τις θεωρίες επιρροής και τα μεγέθη επιδραστικής ομάδας. Οι ταυτολογίες εμφανίζονται με έντονα γράμματα.**

Συνολικά, μπορούμε να συμπεράνουμε ότι οι επιδραστικοί χρήστες των  $F_{mnt}$  και  $F_{rtw}$  σχηματίζουν στενά συνδεδεμένες υπο-ομάδες με συχνές αναφορές μεταξύ τους, ειδικά για μεγέθη διακεκριμένης ομάδας μεγαλύτερα από  $k = 20$ . Αντίθετα, οι διακεκριμένοι χρήστες της  $F_{rnd}$  δεν εμφανίζουν ουσιαστικά κανένα δεσμό μεταξύ τους παρά μόνο για μεγάλες διακεκριμένες ομάδες. Τέλος οι  $F_{twt}$  και  $F_{ind}$  ορίζουν διακεκριμένες ομάδες με μέτρια επικοινωνία που αυξάνεται αναλογικά με το  $k$ .

### 6.3 Πολυπλευρικότητα

Στην ενότητα αυτή, εξετάζουμε αν διαφορετικές θεωρίες επιρροής θεωρούν τους ίδιους χρήστες ως επιδραστικούς στο πλαίσιο των επιμέρους θεματικών κοινοτήτων. Δεδομένου ότι κάθε θεωρία αντιπροσωπεύει μια διαφορετική πλευρά μιας κοινότητας, οι επιδραστικοί χρήστες που υπάρχουν σε δύο από αυτές αντιστοιχούν σε ιδιαίτερα επιδραστικά μέλη με πολλαπλή δραστηριότητα ικανή να επηρεάσει τους υπόλοιπους χρήστες με διαφορετικούς τρόπους. Για να μελετήσουμε αυτό το θέμα, εξετάζουμε την *επικάλυψη μεταξύ θεωριών*, το ποσοστό δηλαδή των κοινών χρηστών μεταξύ των διακεκριμένων ομάδων *ίσου* μεγέθους και της *ίδιας* κοινότητας όπως ορίζονται από *διαφορετικές* θεωρίες επιρροής.

Για τη μέτρηση της επικάλυψης χρησιμοποιούμε το *συντελεστή ομοιότητας Jaccard* (*Jaccard similarity coefficient*). Για δύο διακεκριμένες ομάδες ίσου μεγέθους  $k$ ,  $G_1@k$  και  $G_2@k$ , η ομοιότητα Jaccard  $J(G_1@k, G_2@k)$ , ορίζεται ως το μέγεθος της τομής τους προς το μέγεθος της ένωσης τους:

$$J(G_1@k, G_2@k) = \frac{G_1@k \cap G_2@k}{G_1@k \cup G_2@k} \cdot 100\%$$

Ουσιαστικά, η ομοιότητα Jaccard εκφράζει το ποσοστό των χρηστών που είναι κοινοί ανάμεσα στις δύο ομάδες. Παίρνει τιμές στο διάστημα  $[0,100]$  και οι υψηλότερες τιμές αντιστοιχούν μεγαλύτερη ομοιότητα (δηλαδή επικάλυψη).

Επισημαίνουμε ότι στην ανάλυσή μας δεν συμπεριλάβαμε την  $F_{rnd}$ , καθώς δεν προσφέρει χρήσιμα συμπεράσματα, αφού η επικάλυψή της με οποιαδήποτε άλλη θεωρία είναι εξορισμού τυχαία. Για τις υπόλοιπες θεωρίες επιρροής, θεωρήσαμε τους έξι πιθανούς συνδυασμούς ανά δύο, υπολογίζοντας την επικάλυψη μεταξύ των διακεκριμένων ομάδων τους για τα τέσσερα συνήθη μεγέθη και για όλες τις

κοινότητες. Συνολικά έχουμε 24 διαφορετικούς συνδυασμούς μεγεθών και ζευγών θεωριών επιρροής.

Τα αποτελέσματα της ανάλυσης μας παρουσιάζονται στον Πίνακα 6-3. Παρατηρούμε ότι η ομοιότητα της  $F_{ind}$  με όλες τις άλλες θεωρίες είναι αμελητέα, ανεξαρτήτως του μεγέθους της διακεκριμένης ομάδας. Η ελάχιστη τιμή είναι 0 σε όλες τις περιπτώσεις, ενώ, με δύο μόνο εξαιρέσεις, το ίδιο ισχύει και για τη διάμεσο. Αυτό σημαίνει ότι υπάρχει επικάλυψη μεταξύ θεωριών μόνο για ελάχιστες κοινότητες σε κάθε περίπτωση. Ακόμα και για τις κοινότητες αυτές, η επικάλυψη των διακεκριμένων ομάδων παραμένει αρκετά χαμηλή, με μέγιστη τιμή 10% και μέση τιμή μικρότερη από 1.5%.

Ένα σημαντικό ερώτημα που προκύπτει από τα παραπάνω είναι αν οι μη μηδενικές ομοιότητες μεταξύ της  $F_{ind}$  και των υπόλοιπων θεωριών είναι παράπλευρη συνέπεια του μεγέθους των αντίστοιχων κοινοτήτων: όσο μικρότερη είναι μια κοινότητα, τόσο μεγαλύτερη η πιθανότητα δύο θεωρίες να θεωρούν τους ίδιους χρήστες ως επιδραστικούς. Για να απαντήσουμε το ερώτημα αυτό, συγκρίναμε τα μεγέθη των κοινοτήτων στις οποίες υπάρχει επικάλυψη με το μέσο μέγεθος και διαπιστώσαμε ότι το πρώτο είναι μικρότερο από το δεύτερο σε όλες τις περιπτώσεις. Συγκεκριμένα, υπάρχουν 7 διακεκριμένα θέματα που παρουσιάζουν επικάλυψη μεταξύ της  $F_{ind}$  και κάποιας άλλης θεωρίας για  $k = 10$ . Το μέσο μέγεθός τους είναι 5,300 χρήστες, με ελάχιστη τιμή 2,600 και μέγιστη 10,500, είναι συνεπώς σημαντικά μικρότερες από το μέσο μέγεθος κοινότητας των 13,000 χρηστών. Συνεπώς, οι χρήστες με υψηλό βαθμό εισόδου παίζουν ενεργό, πολύπλευρο ρόλο μόνο σε αρκετά μικρές κοινότητες. Συνολικά μπορούμε να συμπεράνουμε ότι η  $F_{ind}$  είναι εντελώς ανεξάρτητη από τις άλλες θεωρίες. Αυτό σημαίνει ότι οι χρήστες με πολλούς ακόλουθους σπάνια είναι

και πολύ παραγωγικοί, όπως και σπάνια αναφέρονται ή αναδημοσιεύονται από τους υπόλοιπους χρήστες.

Θεωρία 1	$F_{ind}$	$F_{ind}$	$F_{ind}$	$F_{mnt}$	$F_{mnt}$	$F_{rtw}$
Θεωρία 2	$F_{mnt}$	$F_{rtw}$	$F_{twt}$	$F_{rtw}$	$F_{twt}$	$F_{twt}$
<b><math>G_p@10</math></b>						
<b>Minimum.</b>	0.00%	0.00%	0.00%	5.00%	0.00%	0.00%
<b>Median</b>	0.00%	0.00%	0.00%	25.00%	10.00%	15.00%
<b>Maximum.</b>	10.00%	0.00%	10.00%	45.00%	35.00%	35.00%
<b>Mean±SD</b>	0.40%±1.69	0.35%±1.78	0.20%±1.21	26.55%±9.37	11.45%±8.51	12.95%±7.92
<b><math>G_p@20</math></b>						
<b>Minimum</b>	0.00%	0.00%	0.00%	10.00%	0.00%	0.00%
<b>Median</b>	0.00%	0.00%	0.00%	27.50%	12.50%	12.50%
<b>Maximum</b>	5.00%	5.00%	7.50%	45.00%	32.50%	30.00%
<b>Mean±SD</b>	0.45%±1.20	0.35%±1.12	0.28%±1.00	26.08%±8.16	12.20%±7.70	13.93%±7.00
<b><math>G_p@50</math></b>						
<b>Minimum</b>	0.00%	0.00%	0.00%	11.00%	2.00%	2.00%
<b>Median</b>	0.00%	0.00%	0.00%	29.00%	14.00%	16.00%
<b>Maximum</b>	6.00%	6.00%	6.00%	43.00%	32.00%	30.00%
<b>Mean±SD</b>	0.91%±1.35	0.82%±1.18	0.68%±1.12	27.65%±7.37	14.05%±6.96	15.91%±5.85
<b><math>G_p@100</math></b>						
<b>Minimum</b>	0.00%	0.00%	0.00%	11.00%	3.50%	5.00%
<b>Median</b>	0.00%	1.00%	0.50%	29.50%	15.50%	18.00%
<b>Maximum</b>	10.00%	7.50%	7.50%	45.00%	33.00%	29.50%
<b>Mean±SD</b>	0.91%±1.35	1.42%±1.57	1.25%±1.46	27.73%±6.97	15.94%±6.92	17.69%±5.50

**Πίνακας 6-3. Ομοιότητες Jaccard μεταξύ των διακεκριμένων ομάδων που ορίζονται από διαφορετικές θεωρίες επιρροής πάνω στην ίδια θεματική κατηγορία για τα συνήθη μεγέθη διακεκριμένης ομάδας.**

Σχετικά με την επικάλυψη μεταξύ των άλλων τριών θεωριών, ( $F_{mnt}$ ,  $F_{rtw}$  και  $F_{twt}$ ), είναι αξιοσημείωτο ότι οι 8 από τις 12 περιπτώσεις έχουν ελάχιστη τιμή μεγαλύτερη από 0. Αυτό σημαίνει ότι οι θεωρίες αυτές έχουν τουλάχιστον ένα κοινό επιδραστικό χρήστη για κάθε θέμα ανεξαρτήτως μεγέθους ομάδας. Είσοι αξιοσημείωτες είναι οι σχεδόν ίσες τιμές που παίρνουν η μέση τιμή και η διάμεσος, κάτι που αποτελεί ισχυρή ένδειξη κανονικής κατανομής στις επιμέρους κοινότητες. Η ένδειξη αυτή

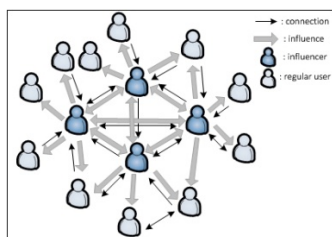


ενισχύεται από το γεγονός ότι στις περισσότερες περιπτώσεις η διάμεσος είναι ίση με τη μέση τιμή της ελάχιστης και της μέγιστης τιμής.

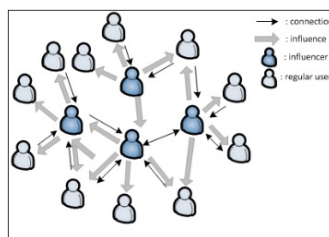
Σημειώνουμε, ωστόσο, ότι η τιμή της επικάλυψης παρουσιάζει σημαντικές διακυμάνσεις για κάθε συνδυασμό θεωριών. Μεταξύ  $F_{mnt}$  και  $F_{rtw}$  είναι αρκετά υψηλή, με μέση τιμή πάνω από 25%. Αυτό σημαίνει ότι το 1/4 των διακεκριμένων χρηστών που εμπλέκονται σε συζητήσεις με τα άλλα μέλη της κοινότητας, παράγουν επίσης περιεχόμενο υψηλού επιπέδου το οποίο αναπαράγεται προς ένα πιο ευρύ κοινό και αντίστροφα. Αυτό είναι συμβατό με τη διαισθητική αντίληψη ότι και τα δύο κριτήρια είναι ενδεικτικά της ικανότητας ενός χρήστη να προκαλέσει απευθείας αντίδραση από άλλους χρήστες. Από την άλλη πλευρά, οι ομοιότητες μεταξύ  $F_{twt}$  και  $F_{mnt}$  όπως και μεταξύ  $F_{twt}$  και  $F_{rtw}$  παραμένουν σχετικά χαμηλές, διακυμαίνονται δηλαδή από 10% έως 20%. Αυτό δείχνει ότι περίπου το 1/6 των  $k$  πιο παραγωγικών χρηστών συχνά αναδημοσιεύονται ή συζητιούνται, με το ποσοστό αυτό να αυξάνεται καθώς αυξάνεται και το  $k$ . Τα παραπάνω φαινόμενα επικάλυψης μεταξύ των  $F_{mnt}$ ,  $F_{rtw}$  και  $F_{twt}$  παραμένουν σχετικά σταθερά για όλα τα μεγέθη διακεκριμένης ομάδας και είναι συμβατά με προηγούμενα ευρήματα που εξετάζουν τη σχέση μεταξύ των τριών θεωριών σε καθολικό (μη τοπικό) επίπεδο (για ολόκληρο δηλαδή τον κοινωνικό γράφο του Twitter και τη συνολική δραστηριότητα των χρηστών τους). [4]

## 6.4 Συζήτηση

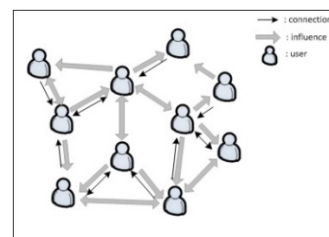
Τα αποτελέσματα της ανάλυσης μας δείχνουν ότι οι θεωρίες επιρροής μπορούν να χωριστούν σε τρεις βασικές κατηγορίες.



**Εικόνα 6-1. Εσωτερική δυναμική των διακεκριμένων ομάδων για  $F_{mnt}$  και  $F_{rtw}$ .**



**Εικόνα 6-2. Εσωτερική δυναμική των διακεκριμένων ομάδων για  $F_{ind}$  και  $F_{twt}$ .**



**Εικόνα 6-3. Εσωτερική δυναμική των διακεκριμένων ομάδων για  $F_{rnd}$ .**

Η πρώτη περιλαμβάνει τους επιδραστικούς χρήστες που ορίζουν διακεκριμένες ομάδες με ισχυρούς δεσμούς μεταξύ των μελών τους. Οι επιδραστικοί χρήστες γνωρίζονται μεταξύ τους και σχηματίζουν στενά συνδεδεμένες ομάδες με οργανωμένη, ομοιογενή συμπεριφορά που εξασκεί σημαντική επιρροή στην υπόλοιπη κοινότητα. Αυτή η εσωτερική δυναμική διευκολύνει τη συνεργασία μεταξύ των επιδραστικών ατόμων, επιτρέποντάς τους να μεγιστοποιούν την επιρροή τους και να επιτυγχάνουν υψηλή επίδοση αναφορικά με το προτεινόμενο πλαίσιο αξιολόγησης. Η δυναμική αυτή φαίνεται στην Εικόνα 6-1. Από τις θεωρίες που εξετάσαμε στην Ενότητα 5, οι  $F_{mnt}$  και  $F_{rtw}$  ανήκουν σε αυτή την κατηγορία. Οι επιδραστικοί χρήστες τους παρουσιάζουν υψηλά επίπεδα ομοφιλίας (κυρίως για μεγέθη διακεκριμένης ομάδας  $k \leq 50$ ) και υψηλά επίπεδα επικοινωνίας (κυρίως για μεγέθη διακεκριμένης ομάδας  $k \geq 50$ ). Σημαντικό μέρος των χρηστών αυτών εμφανίζουν ιδιαίτερα πολλαπλή δραστηριότητα, συμβάλλοντας σημαντικά στο περιεχόμενο και τις συζητήσεις της θεματικής κοινότητας καθώς και στην προώθηση των προσωπικών απόψεων. Τα φαινόμενα αυτά είναι πιο έντονα σε διακεκριμένες ομάδες με 50 μέλη, παρέχοντας έτσι μια αξιόπιστη εκτίμηση ως προς το βέλτιστο μέγεθος διακεκριμένης ομάδας.

Η δεύτερη κατηγορία αποτελείται από θεωρίες επιρροής που αναγνωρίζουν ως διακεκριμένους χρήστες εξέχοντα άτομα που όμως δε σχηματίζουν ομάδα με άλλα επιδραστικά μέλη. Ο λόγος είναι ότι παρουσιάζουν χαμηλά επίπεδα ομοφιλίας, αλληλεπιδρούν μεταξύ τους σε περιορισμένο βαθμό και η δραστηριότητά τους μέσα στην κοινότητα είναι συνήθως μονομερής, παρουσιάζοντας μικρή πολυπλευρικότητα. Η εσωτερική δυναμική αυτή παρουσιάζεται στην Εικόνα 6-2, όπου φαίνεται ότι οι επιδραστικοί χρήστες που ορίζονται από αυτές τις θεωρίες σπάνια δρουν συλλογικά. Στην κατηγορία αυτή ανήκουν οι  $F_{ind}$  και  $F_{twt}$ , από τις οποίες η δεύτερη είναι ασύμβατη με το προτεινόμενο πλαίσιο αξιολόγησης, σε αντίθεση με την πρώτη που παρουσίασε σημαντικά υψηλές επιδόσεις. Οι επιδραστικοί χρήστες της είναι φειδωλοί στο περιεχόμενο που συνεισφέρουν στην κοινότητα, είναι όμως ικανοί να προβλέψουν με ακρίβεια και έγκαιρα ορισμένες πλευρές της δραστηριότητας της κοινότητας.

Τέλος, η τρίτη κατηγορία περιλαμβάνει τις θεωρίες επιρροής οι οποίες ορίζουν διακεκριμένους χρήστες που σχηματίζουν ομάδα μόνο κατ'εμφημισμόν. Δε συσχετίζονται ο ένας με τον άλλον καθώς είναι πιθανό να είναι άγνωστοι μεταξύ τους ή να έχουν διαφορετικό κοινωνικό υπόβαθρο (έλλειψη ομοφιλίας). Εμφανίζονται επίσης αρκετά αδιάφοροι για τη δραστηριότητα της κοινότητας, αφού σπάνια ξεχωρίζουν για τη συμμετοχή τους σε κάποια πλευρά της. Είναι αναμενόμενο, λοιπόν, να μην συγχρονίζονται ώστε να επηρεάσουν συλλογικά τους υπόλοιπους χρήστες, παρουσιάζοντας έτσι χαμηλές επιδόσεις με βάση το προτεινόμενο πλαίσιο αξιολόγησης ανεξαρτήτως μεγέθους διακεκριμένης ομάδας. Στην ανάλυσή μας, παράδειγμα τέτοιας θεωρίας είναι η  $F_{rnd}$ . Η δυναμική των διακεκριμένων χρηστών της παρουσιάζεται στην Εικόνα 6-3..

Ένα ακόμη σημείο που αξίζει να εξετάσουμε είναι αν τα ευρήματά μας για το Twitter (δηλαδή οι επιδραστικοί χρήστες που επηρεάζουν το συναίσθημα μιας θεματικής κοινότητας) αντιστοιχούν σε ουσιαστική επιρροή στον πραγματικό κόσμο. Σχετικές έρευνες ([47], [48]) ισχυρίζονται ότι το συναίσθημα που εκφράζεται στα μέσα κοινωνικής δικτύωσης σχετίζεται άμεσα με τις απόψεις των ανθρώπων στον πραγματικό κόσμο, σε βαθμό μάλιστα που η ανάλυση συναισθήματος στα κοινωνικά δίκτυα θα μπορούσε να χρησιμοποιηθεί ως υποκατάστατο των δημοσκοπήσεων. Από την άποψη αυτή, φαίνεται ότι οι επιδραστικοί χρήστες των κοινωνικών δικτύων μπορούν να έχουν επίδραση στις απόψεις του συνολικού πληθυσμού.

$F_{ind}$	$F_{mnt}$	$F_{rtw}$	$F_{twt}$
cnnbrk	t	oxfordgirl	dominiquerdr
persiankiwi	iranbaan	t	tehranweekly
mashable	oxfordgirl	lissnup	perry1949
google	onlymehdi	persianbanoo	lissnup
mousavi1388	tehranweekly	manic77	khoshkeleodoc
nprpolitics	lissnup	iranbaan	edwand
neilhimsel	dadashiii	onlymehdi	mwolda
timoreilly	k	k	t4tx
oxfordgirl	manic77	iran88	fardinzamani
anamariex	cnn	parsa4	persia_news
nansen	persiankiwi	madyar	artemis_ia
zaibatsu	gr	gr	dadashiii
nprnews	persianbanoo	tehranweekly	iranlaya
iamdiddy	dominiquerdr	persiankiwi	atlsafa
brooksbayne	madyar	dadashiii	uncoolbobby
leolaporte	iran88	fahimn	persia_max_news
dcagle	Lotfan	lotfan	joannemichele
huffingtonpost	fahimn	sheydaj	sannri
cnn	austinheap	sheydajahanbin	iranwwp
nytimeskristof	patrickaltoft	dominiquerdr	persia_news_2

**Πίνακας 6-4.** Η  $G_p@20$  που ορίζεται από την κάθε θεωρία επιρροής για τη θεματική κοινότητα #iranelection.

Προκειμένου να κάνουμε πιο εμφανή τη σύγκριση μεταξύ πραγματικού και δικτυακού κόσμου και να κατανοήσουμε καλύτερα τα αποτελέσματα της ανάλυσης μας, εξετάζουμε τους διακεκριμένους χρήστες που ορίζονται από την κάθε θεωρία επιρροής σε ό,τι αφορά την ιδιότητα και την ταυτότητά τους στον πραγματικό κόσμο

(όταν αυτό είναι δυνατό). Για το σκοπό αυτό επιλέξαμε δύο θεματικές κοινότητες: την #iranelection, η οποία είναι σχετική με τις ιρανικές προεδρικές εκλογές της 12ης Ιουνίου 2009, και τη #noh8, που αφορά την καμπάνια NOH8 η οποία προωθεί την ισότητα σε θέματα φύλου και ανθρωπίνων δικαιωμάτων. Οι διακεκριμένοι χρήστες για τις δύο θεματικές κατηγορίες και για  $k = 20$  παρουσιάζονται στους Πίνακες 6-7 και 6-8 αντίστοιχα.

Ξεκινώντας με την #iranelection, παρατηρούμε ότι η διακεκριμένη ομάδα της  $F_{rtw}$  αποτελείται σχεδόν εξολοκλήρου από δημοσιογράφους, bloggers και ακτιβιστές, η πλειοψηφία των οποίων είναι ιρανοί. Η διακεκριμένη ομάδα που ορίζεται από την  $F_{ind}$  περιλαμβάνει εξίσου υψηλό ποσοστό ειδησεογραφικών πρακτορείων και δημοσιογράφων/ bloggers. Οι επιδραστικοί χρήστες της  $F_{mnt}$  είναι παρόμοιοι με αυτούς της  $F_{rtw}$ , με την προσθήκη μερικών τυχαίων χρηστών και κάποιων διασημοτήτων που δε φαίνεται να σχετίζονται με το θέμα. Η διακεκριμένη ομάδα της  $F_{twt}$  φαίνεται να είναι η λιγότερο έγκυρη, καθώς περιέχει λίγους μόνο λογαριασμούς που σχετίζονται με το Ιράν και το Πράσινο Κίνημα. Ενδιαφέρον παρουσιάζει το γεγονός ότι δύο λογαριασμοί επιλέγονται ως επιδραστικοί από όλες τις θεωρίες – εκτός της  $F_{twt}$ : οι *oxfordgirl* και *persiankiwi*. Και οι δύο αντιστοιχούν σε συνηθισμένους χρήστες από το Ιράν οι οποίοι, αν και αρχικά άγνωστοι, απέκτησαν φήμη ως αξιόπιστες πηγές ενημέρωσης μέσω της αφοσίωσής τους και της ενεργής συμμετοχής τους στην κοινότητα. Αξίζει επίσης να σημειώσουμε ότι ένα μεγάλο μέρος των επιδραστικών χρηστών αυτής της κοινότητας (συμπεριλαμβανομένων των *oxfordgirl* και *persiankiwi*) εμφανίζονται επίσης στις διακεκριμένες ομάδες σχετικών θεμάτων όπως το #neda.

$F_{ind}$	$F_{mnt}$	$F_{rtw}$	$F_{tw}$
drdrew	noh8campaign	bouska	greenbean55
bouska	bouska	biolawyer	noh8campaign
emmyrossum	ste_vee	kylechristian	angelbenton
hollymadison123	greenbean55	ryanleejohnson	johnprather
queerunity	lukasrossi	greenbean55	geisha_boy
lovebscott	hellomattwalker	kristicansler	khajiitchick
torianddean	drdrew	noh8campaign	matthewlush
bluecrystalsky	scoutmasterson	scoutmasterson	alchey
gaycivilrights	looneypyrodude	jeffrago	kylechristian
egheitasean	rosemcgowan	bear54	unlivedlife
lukasrossi	jeffparshley	matthewlush	andrewxanarchy
onemoreslesbian	chimasimone	willfalls	crystallewis60
ontd_fluffy	dawnrichard	mabergel	sthprknt
yezbok	madalyngrimm	ladygagast	aristokatcy
rufuscoolkitty	matthewlush	rakefet27	kausinkonfusion
artemisrex	shannamoakler	egheitasean	ste_vee
jaysays	chris_gorham	tiffanyrinehart	jamesfirestein
tiffanyrinehart	hollymadison123	afishe	prettyynikki
aramina	kimzolciak	ladyspeaker	scoutmasterson
stephjonesmusic	kylechristian	Etejeday	Hasthepotential

**Πίνακας 6-5. Η  $G_p@20$  που ορίζεται από την κάθε θεωρία επιρροής για τη θεματική κοινότητα #noh8.**

Για το #noh8, παρατηρούμε ότι οι  $F_{ind}$ ,  $F_{mnt}$  και  $F_{rtw}$  περιέχουν ένα ισοσταθμισμένο μείγμα από διασημότητες και τυχαίους χρήστες μαζί με κάποιους λογαριασμούς που ανήκουν σε οργανισμούς ανθρωπίνων δικαιωμάτων που αγωνίζονται εναντίον των διακρίσεων. Εξετάζοντας πιο αναλυτικά τους λογαριασμούς αυτούς, όμως, προκύπτει ότι οι περισσότεροι τυχαίοι χρήστες είναι ακτιβιστές ή χρήστες με ισχυρή συμμετοχή σε θέματα φυλετικής ισότητας. Και σε αυτήν την περίπτωση, η  $F_{tw}$  παράγει τα λιγότερο σχετικά αποτελέσματα, ενώ περισσότερο από τα δύο τρίτα των λογαριασμών είναι πλέον απενεργοποιημένοι, κάτι που αποτελεί ισχυρή ένδειξη ότι παρουσίασαν συμπεριφορές που σχετίζονται με spamming ή διαφήμιση. Αξίζει επίσης να επισημάνουμε ότι για καμία θεωρία επιρροής δεν υπάρχει επικάλυψη μεταξύ των διακεκριμένων ομάδων των δύο θεματικών κοινοτήτων.

Τα παραπάνω ευρήματα είναι συνεπή με τα αποτελέσματα των προηγούμενων αναλύσεων, που δείχνουν ότι οι  $F_{ind}$  και  $F_{rtw}$  είναι οι πιο αξιόπιστες μέθοδοι για τον

εντοπισμό των επιδραστικών χρηστών, ενώ η  $F_{twt}$  βρίσκεται στο άλλο άκρο. Συμπεραίνουμε, λοιπόν, ότι οι  $F_{ind}$  και  $F_{rtw}$  φανερωθούν διαφορετικούς, αλλά εξίσου σημαντικούς τύπους επιρροής, με την πρώτη να θεωρεί τη δημοφιλία των χρηστών ως αξιόπιστο κριτήριο της επιρροής τους και τη δεύτερη να βασίζεται στην ικανότητά τους να παράξουν ενδιαφέρον περιεχόμενο. Τέλος, είναι σημαντικό να τονίσουμε ότι η επιρροή τους φαίνεται να έχει ισχύ εντός συγκεκριμένου περιεχομένου, καθώς παρουσιάζεται αμελητέα ή ανύπαρκτη επικάλυψη μεταξύ των διακεκριμένων ομάδων διαφορετικών κοινοτήτων, εκτός αν τα αντίστοιχα θέματα σχετίζονται νοηματικά.

Η σελίδα αυτή είναι σκόπιμα λευκή.



# 7

## *Ανάσυρση δεδομένων από πολλαπλά κοινωνικά δίκτυα*

Στο παρόν κεφάλαιο, παρουσιάζουμε την υλοποίηση μιας εφαρμογής που επιτρέπει το συνδυασμό δεδομένων και λειτουργικότητας από πολλαπλές πλατφόρμες κοινωνικής δικτύωσης με ενιαίο τρόπο. Συγκεκριμένα, παρέχει μια διεπαφή που αποτελεί κοινό σημείο πρόσβασης για τις υποκείμενες πλατφόρμες, εκθέτοντας ένα σύνολο μεθόδων που εμπεριέχουν τη λειτουργικότητά τους. Βασικό συστατικό του εργαλείου ανάσυρσης δεδομένων είναι το μοντέλο δεδομένων το οποίο ορίζει ένα σύνολο οντοτήτων που αναπαριστούν θεμελιώδεις έννοιες των κοινωνικών δικτύων και τις μεταξύ τους συσχετίσεις.

### ***7.1 Εργασία με δεδομένα από πολλαπλά κοινωνικά δίκτυα***

Όπως αναφέρθηκε σε προηγούμενα κεφάλαια, η μεγάλη αύξηση της δημοτικότητας διαδικτυακών κοινοτήτων, όπως οι πλατφόρμες κοινωνικής δικτύωσης, είχε ως αποτέλεσμα την ύπαρξη ενός μεγάλου όγκου περιεχομένου που παράγεται και ανανεώνεται συνεχώς από τα μέλη αυτών των κοινοτήτων. Το περιεχόμενο αυτό δεν

αποτελείται μόνο από δεδομένα που μοιράζονται άμεσα οι ίδιοι οι χρήστες, όπως αναρτήσεις κειμένου, φωτογραφίες και βίντεο, αλλά και από ένα σημαντικό όγκο κοινωνικής πληροφορίας που εξάγεται έμμεσα από τη δραστηριότητα των χρηστών, όπως για παράδειγμα τα ενδιαφέροντα και οι προτιμήσεις τους αλλά και οι σχέσεις τους με τους άλλους χρήστες.

Αν και μέρος του περιεχομένου παραμένει ιδιωτικό, κάθε άλλο παρά αμελητέο είναι το ποσοστό του περιεχομένου αυτού που οι χρήστες καθιστούν δημόσια διαθέσιμο. Εξίσου σημαντικό είναι το γεγονός ότι τα ίδια τα κοινωνικά δίκτυα, αντί να περιορίζουν την προσβασιμότητα τρίτων στη λειτουργικότητα και τα δεδομένων τους, τα διαθέτουν ελεύθερα ως θεμελιώδες στοιχείο των παροχών τους στους τελικούς χρήστες και ενθαρρύνουν τη δημιουργία εφαρμογών βασισμένων σε αυτά. Αυτή τη στιγμή οι πιο δημοφιλείς υπηρεσίες κοινωνικής δικτύωσης και πολυμέσων όπως τα Twitter, Facebook και Youtube παρέχουν πρόσβαση σε μέρος της λειτουργικότητάς τους μέσω προγραμματιστικών διεπαφών. Χρησιμοποιώντας αυτές τις διεπαφές, κάθε χρήστης ή εξωτερική εφαρμογή μπορεί να αποκτήσει πρόσβαση στο περιεχόμενο και τις λειτουργίες των υπηρεσιών αυτών.

Τόσο ο επιχειρησιακός όσο και ο ερευνητικός κόσμος έχει από καιρό αναγνωρίσει το μεγάλο περιθώριο αξιοποίησης του κοινωνικού γράφου και του περιεχομένου χρηστών που βρίσκονται στα κοινωνικά δίκτυα και αναζητά τρόπους εκμετάλλευσής τους. Για τις επιχειρήσεις που διαθέτουν τα κατάλληλα εργαλεία διαχείρισής τους, τα δεδομένα των κοινωνικών δικτύων είναι πιθανές πηγές κέρδους, καθώς μπορούν να αποδειχτούν πολύτιμα εφόδια για καμπάνιες εξατομικευμένης διαφήμισης και viral marketing. Όσον αφορά την έρευνα, όπως είδαμε και προηγουμένως, η διάδοση των κοινωνικών δικτύων αναζωπύρωσε το ενδιαφέρον γύρω από πολλές ερευνητικές περιοχές, συμπεριλαμβανομένης και της μελέτης επιρροής που εξετάσαμε αναλυτικά

στα προηγούμενα κεφάλαια, καθώς ο μεγάλος όγκος δεδομένων που παράγουν οι χρήστες και οι σαφώς επισημασμένες μεταξύ τους σχέσεις επιτρέπουν την ανάλυση των δεδομένων και τη μελέτη της κοινωνικής δυναμικής σε πολύ μεγαλύτερη κλίμακα.

Αν και το περιεχόμενο των κοινωνικών δικτύων είναι άπλετο σε όγκο και εύκολα προσβάσιμο, η αξιοποίησή του συνεχίζει να παρουσιάζει σημαντικές προκλήσεις. Παρά τις ομοιότητες σε έννοιες και βασικές λειτουργίες, η αναπαράσταση των δεδομένων μεταξύ διαφορετικών κοινωνικών δικτύων είναι σημαντικά ετερογενής. Επιπλέον, κάθε δίκτυο διαθέτει τις δικές του προγραμματιστικές διεπαφές (API) μέσα από τις οποίες κάνει διαθέσιμη τη λειτουργικότητά του. Αυτό σημαίνει ότι, ελλείψει ενός μη εμπορικού εργαλείου που επιτρέπει την ταυτόχρονη πρόσβαση σε πολλαπλά δίκτυα από μια κοινή διεπαφή, ένας χρήστης που θέλει να συνδυάσει δεδομένα από δύο ή περισσότερα δίκτυα θα πρέπει να καλέσει ξεχωριστά τις διεπαφές τους και να μετατρέψει τα δεδομένα σε μια κοινή μορφοποίηση πριν τα επεξεργαστεί.

Στο παρόν κεφάλαιο, παρουσιάζουμε την υλοποίηση μιας εφαρμογής που επιτρέπει το συνδυασμό δεδομένων και λειτουργικότητας από πολλαπλές πλατφόρμες κοινωνικής δικτύωσης με ενιαίο τρόπο. Συγκεκριμένα, παρέχει μια διεπαφή που αποτελεί κοινό σημείο πρόσβασης για τις υποκείμενες πλατφόρμες, εκθέτοντας ένα σύνολο μεθόδων που εμπεριέχουν τη λειτουργικότητά τους. Για κάθε υποστηριζόμενη πλατφόρμα έχει υλοποιηθεί ένας προσαρμογέας, ενώ η εφαρμογή μπορεί να επεκταθεί ώστε να υποστηρίζει περισσότερα κοινωνικά δίκτυα υλοποιώντας τον αντίστοιχο προσαρμογέα. Τη στιγμή αυτή, η εφαρμογή παρέχει υποστήριξη για τις εξής εφτά πλατφόρμες: Facebook, Twitter, Flickr, Dailymotion, YouTube, Google+ και Instagram. Η εφαρμογή φέρνει τα δεδομένα από τα κοινωνικά δίκτυα σε πραγματικό χρόνο και δεν πραγματοποιεί καμία αποθήκευση, μόνιμη είτε

προσωρινή, των δεδομένων. Ασχολείται κυρίως με δημόσια διαθέσιμο περιεχόμενο, με την εξαίρεση δύο μεθόδων που απαιτούν πιστοποίηση εκ μέρους του χρήστη.

## **7.2 Σχετικές ερευνητικές εργασίες και υπάρχουσες εφαρμογές**

Όπως εξηγήθηκε προηγουμένως, το πλήθος και η ποικιλία των προγραμματιστικών διεπαφών των κοινωνικών δικτύων καθιστούν αναγκαία την ύπαρξη μιας μετα-διεπαφής που θα δρα ως σημείο συγκέντρωσης και θα παρέχει ενιαία πρόσβαση σε όλο το φάσμα περιεχομένου παραγόμενου από τους χρήστες. Στην αγορά υπάρχουν κάποια εμπορικά προϊόντα που ικανοποιούν την ανάγκη αυτή. Καταρχήν η GNIP [49], μια εταιρία που αγοράστηκε από το Twitter τον Απρίλιο του 2014, παρέχει πρόσβαση σε πολλές πηγές κοινωνικών δεδομένων, τόσο σε πραγματικό χρόνο όσο και πρόσβαση σε ιστορικά δεδομένα. Παρόμοια προσέγγιση ακολουθούν και οι εφαρμογές HootSuite [50] και DataSift [51], οι οποίες επιτρέπουν στους χρήστες τους να διαχειρίζονται με κοινό τρόπο πολλά κοινωνικά δίκτυα ενώ ταυτόχρονα προσφέρουν υπηρεσίες ανάλυσης δεδομένων.

Μια αρκετά διαδεδομένη προσέγγιση όσον αφορά την ενιαία διαχείριση λογαριασμών σε κοινωνικά δίκτυα είναι οι ιστοσελίδες συλλέκτες: σελίδες που παρέχουν μια κοινή διεπαφή από την οποία ο χρήστης έχει πρόσβαση στους λογαριασμούς του στα διάφορα κοινωνικά δίκτυα. Η βασική διαφορά των σελίδων αυτών από τις εφαρμογές που αναφέρθηκαν στην προηγούμενη παράγραφο είναι ότι δεν στοχεύουν στην εννοιολογική ταυτοποίηση των βασικών εννοιών των κοινωνικών δικτύων, αντίθετα περικλείουν τα δεδομένα που επιστρέφουν οι διαφορετικές πλατφόρμες κοινωνικής δικτύωσης σε μια κοινή διεπαφή χωρίς να τα αναλύουν περαιτέρω ή να τα συνδυάζουν μεταξύ τους. Αντίθετα, στόχος της εφαρμογής που παρουσιάζουμε εδώ είναι, αφού αναγνωριστούν οι κοινές έννοιες που

υπάρχουν στις υποκείμενες πλατφόρμες, να μετασχηματιστούν σε ένα κοινό μοντέλο δεδομένων ώστε να μπορούν να συνδυαστούν και να αναλυθούν περαιτέρω με ενιαίο τρόπο.

Μια σχετική ερευνητική εργασία με σκοπό τη δημιουργία μιας οντολογίας για την αναπαράσταση των κοινωνικών δικτύων είναι η [52], στην οποία οι συγγραφείς χρησιμοποιούν ως αφετηρία για τη δουλειά τους, όχι τα υπάρχοντα κοινωνικά δίκτυα αλλά μια θεωρητική δομή κοινωνικού δικτύου. Η προσαρμογή του μοντέλου αυτού για τις δημοφιλείς πλατφόρμες κοινωνικής δικτύωσης θα απαιτούσε σημαντική προσπάθεια, καθώς οι έννοιές τους έχουν αναπτυχθεί αυτόνομα. Στην προσέγγισή μας προτιμήσαμε τη χρήση ενός μοντέλου δεδομένων που προέκυψε από την ανάλυση πραγματικών κοινωνικών δικτύων. Βασιζόμενοι στο πρότυπο OpenSocial [53], προχωρήσαμε στη δημιουργία ενός νέου μοντέλου δεδομένων που επιτρέπει την αναπαράσταση των οντοτήτων των υποκείμενων δικτύων με ενιαίο τρόπο. Η προσέγγιση αυτή περιγράφεται πιο αναλυτικά στην ακόλουθη ενότητα.

### **7.3 Μοντέλο δεδομένων**

Στο πλαίσιο της ανάπτυξης ενός εργαλείου που επιτρέπει την αλληλεπίδραση με πολλαπλά κοινωνικά δίκτυα, βασικό βήμα ήταν η ανάπτυξη ενός μοντέλου δεδομένων για την ενιαία αναπαράσταση των εννοιών που βρίσκονται στα υποκείμενα δίκτυα. Παρά τις ομοιότητες σε έννοιες και βασικές λειτουργίες, η αναπαράσταση των δεδομένων μεταξύ διαφορετικών κοινωνικών δικτύων είναι σημαντικά ετερογενής. Κάθε κοινωνικό δίκτυο χρησιμοποιεί το δικό του μοντέλο για την αναπαράσταση των οντοτήτων (χρήστες, αναρτήσεις, κλπ) και των μεταξύ τους σχέσεων. Εξετάζοντας τα μοντέλα δεδομένων των πιο δημοφιλών κοινωνικών δικτύων ανακαλύψαμε οντότητες που αν και διαφέρουν λόγω της διαφορετικής

λειτουργίας του κάθε δικτύου, είναι ωστόσο αρκετά κοντινές και διαθέτουν αρκετά κοινά χαρακτηριστικά ώστε να μπορούν να ομαδοποιηθούν κάτω από την ίδια κατηγορία. Για παράδειγμα, μια φωτογραφία στο Instagram και ένα βίντεο στο YouTube μπορούν και τα δύο να χαρακτηριστούν “Δημοσίευση” ενώ διαθέτουν παρόμοιες ιδιότητες όπως τίτλος, ημερομηνία ανάρτησης, αριθμός σχολίων κλπ. Βασισμένοι στα παραπάνω, εντοπίσαμε όλα τα αντικείμενα που συνδέονται εννοιολογικά με αυτό τον τρόπο και ορίσαμε ένα σύνολο κατηγοριών που δύναται να αναπαραστήσει τα δεδομένα που προέρχονται από οποιοδήποτε από τα υποκείμενα κοινωνικά δίκτυα.

Το μοντέλο δεδομένων που προκύπτει ορίζει ένα σύνολο οντοτήτων που αναπαριστούν έννοιες θεμελιώδεις στο πλαίσιο των κοινωνικών δικτύων. Το μοντέλο δεδομένων μας αποτελείται από της ακόλουθες έννοιες:

- **Άτομο:** Το προφίλ ενός χρήστη. Περιέχει το λογαριασμό του χρήστη και τις προσωπικές πληροφορίες που περιέχονται στο προφίλ του.
- **Δημοσίευση:** Μια ανάρτηση που δημοσιεύτηκε από ένα χρήστη σε ένα κοινωνικό δίκτυο. Η δημοσίευση μπορεί να αφορά ανάρτηση βίντεο, εικόνας ή κειμένου.
- **Ενέργεια:** Πληροφορίες σχετικά με μια ενέργεια που πραγματοποίησε ο χρήστης σε ένα κοινωνικό δίκτυο.
- **Σχόλιο:** Ένα σχόλιο που ανάρτησε ένας χρήστης πάνω σε μια δημοσίευση σε ένα κοινωνικό δίκτυο.
- **Κοινωνικό Δίκτυο:** Υποστηριζόμενο κοινωνικό δίκτυο, ένα από τα Flickr, Facebook, Twitter, Youtube, Dailymotion, Google+ και Instagram.

- **Ταυτότητα:** Αναγνωριστικό ενός αντικείμενου σε ένα συγκεκριμένο κοινωνικό δίκτυο.
- **Διεύθυνση:** Πεδία που μπορούν να χρησιμοποιηθούν για να ορίσουν μια τοποθεσία.
- **Όνομα:** Πεδία που ορίζουν το πλήρες όνομα ενός χρήστη.

Κάθε ένα από τα παραπάνω αντικείμενα αντιστοιχεί σε μία ή περισσότερες ενότητες σε κάθε κοινωνικό δίκτυο, ενώ κάποια αντικείμενα εμφανίζονται μόνο σε ένα υποσύνολο των υποστηριζόμενων δικτύων (για παράδειγμα μόνο τα Google+, Facebook, YouTube και Dailymotion έχουν έννοια αντίστοιχη της Ενέργειας).

Κύρια Αντικείμενα				
Άτομο	Δημοσίευση	Ενέργεια	Σχόλιο	Κοινωνικό Δίκτυο
Ταυτότητα	Ταυτότητα	Ταυτότητα	Ταυτότητα	FLICKR
Κοινωνικό δίκτυο	Κοινωνικό δίκτυο	Κοινωνικό δίκτυο	Κοινωνικό δίκτυο	FACEBOOK
Περιγραφή	Ημερομηνία	Ημερομηνία	Ημερομηνία	TWITTER
Διευθύνσεις	Τίτλος	Τίτλος	Περιγραφή	YOUTUBE
Ημερομηνία γέννησης	thumbnailUrl	Περιγραφή	Χρήστης	DAILYMOTION
Τρέχουσα τοποθεσία	Περιγραφή	Τοποθεσία	Όνομα χρήστη	GOOGLEP
Όνομα χρήστη	Διάρκεια	Χρήστης	# θετικών ψήφων	INSTAGRAM
Email	Τοποθεσία	Τύπος αντικείμενου		
Φύλο	Γλώσσα	Δημοσιεύσεις		
Όνομα	Άδεια	Σχετιζόμενα Άτομα		
Φωτογραφίες	Μέγεθος αρχείου	Σχετικές Ενέργειες		
URL προφίλ	Βαθμολογία			
Μέλος Από	# βαθμολογιών			
thumbnailUrl	# θετικών ψήφων			
utcOffset	# αρνητικών ψήφων			
# φίλων	# σχολίων			
# ακολούθων	# προβολών			
# χρηστών που ακολουθεί	# αναδημοσιεύσεων			
	# προτιμώμενων			
	ετικέτες			
	Σχετιζόμενα Ατόμα			
	Τύπος			
	URL			
	Χρήστης			
	Σχόλια			

**Πίνακας 7-1: Κύρια αντικείμενα του μοντέλου δεδομένων της εφαρμογής ανάσχυσης δεδομένων από ετερογενή κοινωνικά δίκτυα.**

Τα αντικείμενα χωρίζονται σε κύρια (Άτομο, Δημοσίευση, Ενέργεια, Σχόλιο, Κοινωνικό Δίκτυο) και δευτερεύοντα (Ταυτότητα, Διεύθυνση, Όνομα). Τα πρώτα αναπαριστούν θεμελιώδεις έννοιες ενός κοινωνικού δικτύου ενώ τα δεύτερα αντιστοιχούν σε σύνθετα πεδία μέσα στα κύρια αντικείμενα. Εκτός από τα παραπάνω,

ορίσαμε έναν αριθμό αντικειμένων-φίλτρων, καθένα από τα οποία περιέχει ένα σύνολο από εννοιολογικά σχετιζόμενες παραμέτρους που μπορούν να χρησιμοποιηθούν για αναζήτηση. Τα αντικείμενα αυτά τα χρησιμοποιούμε ως παραμέτρους εισόδου σε μεθόδους αναζήτησης. Τα κύρια αντικείμενα του μοντέλου δεδομένων παρουσιάζονται στην Πίνακα 7-1.

Όπως αναφέρθηκε και παραπάνω, ως βάση για τη δημιουργία του μοντέλου δεδομένων χρησιμοποιήσαμε τις προδιαγραφές που ορίζει το Opensocial. Οι προδιαγραφές αυτές στόχευαν στον ορισμό ενός μοντέλου για τη δημιουργία νέων εφαρμογών που μπορούν να τρέξουν εντός διαφόρων κοινωνικών δικτύων. Για το λόγο αυτό, σχεδιάστηκαν έτσι ώστε να μπορούν να εξυπηρετήσουν όσο το δυνατόν περισσότερο τις πιθανές απαιτήσεις τέτοιου είδους εφαρμογών. Ο σχεδιασμός αυτός δεν είναι κατάλληλος, όμως, για ένα μοντέλο που ως στόχο έχει την αναπαράσταση δεδομένων που βρίσκονται στα υπάρχοντα κοινωνικά δίκτυα. Αυτό συμβαίνει γιατί πολλές από τις έννοιες που ορίζονται στις αρχικές προδιαγραφές είναι περιττές για το σκοπό αυτό. Για παράδειγμα, τα πεδία σχετικά με ένα άτομο όπως ορίζεται από το Opensocial, καλύπτουν πολλές διαφορετικές πλευρές της οντότητας ενός ανθρώπου (όπως για παράδειγμα σωματότυπο, οικογενειακή κατάσταση, θρησκευτικές πεποιθήσεις), τα οποία ωστόσο σε ελάχιστες περιπτώσεις υποστηρίζονται από τα ίδια τα κοινωνικά δίκτυα και σε ακόμα λιγότερες είναι διαθέσιμα μέσω των API τους. Τέτοια πεδία έχουν αφαιρεθεί από το μοντέλο δεδομένων μας, ενώ ταυτόχρονα έχουν προστεθεί πεδία που δε συμπεριλαμβάνονταν στις αρχικές προδιαγραφές αλλά περιέχουν σημαντική πληροφορία για την θέση του χρήστη στο γράφο (όπως ο αριθμός των σχολίων ή των αναδημοσιεύσεων μιας δημοσίευσης). Το μοντέλο δεδομένων που προέκυψε, αν και διατηρεί πολλές από τις ονομαστικές συμβάσεις και



βασικές έννοιες που ορίζει το OpenSocial, παρουσιάζει σημαντικές διαφορές από το αρχικό μοντέλο.

## **7.4 Λεπτομέρειες υλοποίησης**

Το εργαλείο που υλοποιήσαμε λειτουργεί ως ενιαία διεπαφή για μερικά από τα πιο δημοφιλή κοινωνικά δίκτυα, παρέχοντας κοινό σημείο πρόσβασης στις προγραμματιστικές διεπαφές τους. Από άποψη υλοποίησης, το εργαλείο που δημιουργήσαμε αποτελείται από ένα κεντρικό στοιχείο που λειτουργεί ως το μοναδικό σημείο αναφοράς για τις εφαρμογές – πελάτες και που είναι διαθέσιμο τόσο σε SOAP όσο και REST και ένα σύνολο από προσαρμογείς. Η εφαρμογή αναπτύχθηκε σε Java . Οι προσαρμογείς ουσιαστικά περικλείουν τη λειτουργία των μεμονωμένων προγραμματιστικών διεπαφών που παρέχουν τα κοινωνικά δίκτυα με τις οποίες επικοινωνούν μέσω REST.

Όπως είναι αναμενόμενο, η εφαρμογή εξαρτάται σε μεγάλο βαθμό από τις διεπαφές των υποκείμενων κοινωνικών δικτύων. Αυτό αφορά τόσο τη λειτουργία που προσφέρει το κάθε κοινωνικό δίκτυο, όσο και την επίδοση και τους περιορισμούς του. (π.χ. αριθμός κλήσεων ανά μονάδα χρόνου). Η εφαρμογή ανασύρει όλα τα δεδομένα σε πραγματικό χρόνο, χωρίς να αποθηκεύει δεδομένα χρηστών με οποιοδήποτε τρόπο. Όταν από τις αρχικές μεθόδους των κοινωνικών δικτύων απαιτείται πιστοποίηση χρήστη, οι παράμετροι πιστοποίησης δίνονται ως είσοδος από την εφαρμογή πελάτη.

Η διεπαφή παρέχει ένα σύνολο μεθόδων οι οποίες επιτρέπουν στο χρήστη να αλληλεπιδρά με τις διεπαφές των υποκείμενων κοινωνικών δικτύων. Κάποιες από τις μεθόδους αυτές επιχειρούν να ανασύρουν πληροφορίες σχετικά με συγκεκριμένα αντικείμενα ενώ άλλες αναζητούν δεδομένα που ανταποκρίνονται σε ορισμένα

κριτήρια. Κατά συνέπεια οι παράμετροι εισόδου είτε προσδιορίζουν την ταυτότητα των αντικειμένων που θέλουμε να ανασύρουμε, είτε λειτουργούν ως φίλτρα για την αναζήτηση. Το αντικείμενο που επιστρέφεται περιέχει τα αποτελέσματα των κλήσεων που ήταν επιτυχείς και λεπτομερή μηνύματα λάθους για αυτές που απέτυχαν (αν υπάρχουν).

Οι μέθοδοι που υποστηρίζει το εργαλείο εξόρυξης δεδομένων από πολλαπλά κοινωνικά δίκτυα μπορούν να ομαδοποιηθούν ως εξής:

- **Μέθοδοι αναζήτησης:** Πραγματοποιούν αναζήτηση ατόμου/δημοσίευσης /ενέργειας με βάση φίλτρα που καθορίζονται από το χρήστη (π.χ. αναζήτηση δημοσιεύσεων με βάση την τοποθεσία ή λέξεις κλειδιά).
- **Μέθοδοι ανάσυρσης με βάση την ταυτότητα:** Επιστρέφουν άτομα/ δημοσιεύσεις/ενέργειες/σχόλια με βάση την ταυτότητά τους.
- **Μέθοδοι ανάσυρσης συνδεδεμένων αντικειμένων:** Επιστρέφουν άτομα/ δημοσιεύσεις/ενέργειες/σχόλια που συνδέονται με ένα αντικείμενο (άτομο /δημοσίευση/ενέργεια) με βάση την ταυτότητα του αντικειμένου (π.χ. ανάσυρση των δημοσιεύσεων ή των φίλων ενός ατόμου).

Όπως αναφέρθηκε και προηγουμένως, η υλοποίηση μιας μεθόδου για το κάθε κοινωνικό δίκτυο εξαρτάται από τους περιορισμούς που επιβάλλουν τα ίδια τα κοινωνικά δίκτυα. Για παράδειγμα οι μέθοδοι που αναφέρονται σε μια δραστηριότητα (οντότητα Ενέργεια) υλοποιούνται μόνο από τα κοινωνικά δίκτυα που υποστηρίζουν την έννοια της δραστηριότητας (δηλαδή τα Google+, Facebook, YouTube and Dailymotion όπως είδαμε παραπάνω). Επιπλέον η λειτουργία της κάθε μεθόδου μπορεί να διαφέρει σημαντικά από το ένα κοινωνικό δίκτυο στο άλλο ανάλογα με τη λειτουργία που παρέχει το αντίστοιχο δίκτυο. Για παράδειγμα ανάλογα με την

υλοποίηση στο κάθε δίκτυο, η λίστα των ατόμων που επιστρέφει η μέθοδος ανάσυρσης ατόμων που συνδέονται με μία δημοσίευση μπορεί να περιλαμβάνει μόνο τον ιδιοκτήτη του αντικειμένου ή μια πλήρη λίστα με τα άτομα που το έχουν σχολιάσει, ψηφίσει θετικά ή αναφερθεί σε αυτό.

Από άποψη αρχιτεκτονικής, το εργαλείο εξόρυξης δεδομένων από πολλαπλά κοινωνικά δίκτυα αποτελείται από ένα κεντρικό στοιχείο λογισμικού, την υπηρεσία διαχείρισης, η οποία προσφέρει την προγραμματιστική διεπαφή στις εφαρμογές πελάτες και ένα σύνολο από προσαρμογείς. Η υπηρεσία διαχείρισης είναι υπεύθυνη για να αρχικοποιεί και να συντονίζει τους προσαρμογείς, καθώς επίσης και για να συγκεντρώνει τα αποτελέσματα και να τα επιστρέφει στην εφαρμογή/πελάτη. Οι προσαρμογείς είναι στοιχεία λογισμικού που ενσωματώνουν την επιμέρους λειτουργία των προγραμματιστικών διεπαφών των κοινωνικών δικτύων με τις οποίες επικοινωνούν χρησιμοποιώντας το πρωτόκολλο REST. Οι προσαρμογείς είναι μέρος του εργαλείου και έχουν υλοποιηθεί ως Java αντικείμενα που δημιουργούνται και αρχικοποιούνται από την υπηρεσία διαχείρισης σε κάθε κλήση μεθόδου.

Κάθε μέθοδος που υλοποιεί η υπηρεσία διαχείρισης και προσφέρεται μέσω της διεπαφής αντιστοιχεί σε μία μέθοδο ή σε συνδυασμό πολλών απλούστερων μεθόδων που υλοποιούν οι προσαρμογείς. Ο κάθε προσαρμογέας περιέχει μια υλοποίηση της μεθόδου προσαρμοσμένη για το συγκεκριμένο κοινωνικό δίκτυο στο οποίο αντιστοιχεί. Όλοι οι προσαρμογείς αποτελούν υλοποίηση μιας διεπαφής Java στην οποία ορίζονται οι μέθοδοι που πρέπει να υλοποιούν. Στην παρούσα του μορφή, το εργαλείο περιέχει προσαρμογείς για επτά δημοφιλή κοινωνικά δίκτυα, τα Twitter, Facebook, Flickr, Dailymotion, YouTube, Google+ και Instagram. Επιπλέον προσαρμογείς μπορούν να προστεθούν στην πλατφόρμα δημιουργώντας μια υλοποίηση της διεπαφής για το νέο κοινωνικό δίκτυο.

Οι προσαρμογείς των κοινωνικών δικτύων αρχικοποιούνται και συντονίζονται από την υπηρεσία διαχείρισης. Όταν γίνεται κλήση μίας μεθόδου, η υπηρεσία διαχείρισης επεξεργάζεται τη λίστα των παραμέτρων, επιλέγει τους προσαρμογείς που πρέπει να συμπεριληφθούν και δημιουργεί τα ορίσματα εισόδου που θα περαστούν στον κάθε προσαρμογέα. Για παράδειγμα, η μέθοδος ανάλυσης ατόμων με βάση τις ταυτότητές τους, παίρνει ως είσοδο μια λίστα από παραμέτρους τύπου Ταυτότητα, κάθε μια από τις οποίες περιέχει ένα αναγνωριστικό ταυτότητας και ένα κοινωνικό δίκτυο και επιστρέφει μια λίστα από αντικείμενα τύπου Άτομο.

Κατά την κλήση της μεθόδου, η υπηρεσία διαχείρισης επεξεργάζεται τη λίστα με τις ταυτότητες και της χωρίζει σε έναν αριθμό μικρότερων λιστών, ίσων με τον αριθμό των κοινωνικών δικτύων που περιέχονται στην αρχική λίστα. Η λίστα που θα περαστεί στον κάθε προσαρμογέα περιέχει μόνο τα αναγνωριστικά ταυτότητας που ανήκουν στο συγκεκριμένο κοινωνικό δίκτυο. Αφού δημιουργήσει τα ορίσματα εισόδου, η υπηρεσία διαχείρισης θα αρχικοποιήσει τη λίστα με τους προσαρμογείς και θα τους καλέσει ώστε να συγκεντρώσει τα αποτελέσματα και να τα επιστρέψει στο χρήστη. Τα αποτελέσματα που επιστρέφονται από τους προσαρμογείς συνδυάζονται σε ένα σύνθετο αντικείμενο που περιέχει επιπλέον τυχόν λάθη που συνέβησαν κατά τη διάρκεια της εκτέλεσης και επιστρέφονται στο χρήστη που έκανε την κλήση. Με τον τρόπο αυτό, τυχόν σφάλματα που συμβαίνουν σε ένα προσαρμογέα, δε θα επηρεάσουν τη λειτουργία των υπολοίπων προσαρμογέων.

Η σελίδα αυτή είναι σκόπιμα λευκή.

# 8

## Σύνοψη

Στην παρούσα διατριβή προτείνουμε ένα πλαίσιο αξιολόγησης της σχετικής επίδοσης θεωριών τοπικής επιρροής. Στόχος του πλαισίου είναι να εξετάσει αν οι θεωρίες αυτές εντοπίζουν ως επιδραστικούς χρήστες μικρές ομάδες ανθρώπων που είναι ικανοί να επηρεάσουν ορισμένα αντικειμενικά μετρήσιμα χαρακτηριστικά μιας κοινότητας με ελάχιστη προσπάθεια. Προκειμένου να θέσουμε σε εφαρμογή την ιδέα μας, ορίσαμε πέντε αναγκαίες και ικανές συνθήκες τις οποίες μια αποδοτική θεωρία επιρροής θα πρέπει να ικανοποιεί για ένα εύρος θεματικών κοινοτήτων. Οι τρεις πρώτες συνθήκες αποτελούν αναγκαίες προϋποθέσεις που χρησιμοποιούνται για να αποκλείσουμε θεωρίες που καταλήγουν σε πολυπληθείς ή υπερπαραγωγικές επιδραστικές ομάδες. Για τις άλλες δύο συνθήκες, εισαγάγαμε μετρικές αξιολόγησης που υπολογίζουν τη σημασιολογική και χρονολογική σχέση μεταξύ του μέσου συναισθήματος των επιδραστικών χρηστών και της υπόλοιπης κοινότητας. Και οι δύο βασίζονται στη χρήση εικονιδίων (emojis) για την έκφραση του συναισθήματος, μια μέθοδος που προτιμήθηκε από πιο προηγμένες τεχνικές ανάλυσης συναισθήματος, οι οποίες είναι επίσης συμβατές με το πλαίσιο μας.

Για να δείξουμε τη χρήση του πλαισίου μας, αξιολογήσαμε σε ένα μεγάλο σε όγκο σύνολο πραγματικών δεδομένων από το Twitter τη σχετική επίδοση πέντε

καθιερωμένων στη βιβλιογραφία θεωριών τοπικής επιρροής. Οι θεωρίες αυτές είναι η θεωρία τυχαίας επιρροής, η θεωρία επιρροής βαθμού εισόδου, η θεωρία επιρροής αναφορών, η θεωρία επιρροής αναδημοσιεύσεων και η θεωρία επιρροής όγκου δεδομένων. Το σύνολο δεδομένων αποτελείται από 75 θεματικές κατηγορίες ορισμένες με βάση τα hashtags, κάθε μία από τις οποίες περιέχει 80,000 tweets και σχεδόν 9,000 χρήστες κατά μέσο όρο. Προκειμένου να εξασφαλίσουμε συμβατότητα με το πλαίσιο μας, περιορίσαμε την ανάλυσή μας στις κοινότητες για τις οποίες μπορούμε να βγάλουμε ασφαλή συμπεράσματα, λαμβάνοντας έτσι υπόψη μέρος μόνο της δραστηριότητας της κοινότητας, τα μηνύματα αυτά δηλαδή που είναι πολωμένα και περιέχουν hashtags. Τα αποτελέσματα της ανάλυσής μας δείχνουν ότι η θεωρία επιρροής όγκου δεδομένων είναι ασύμβατη με το πλαίσιο μας ενώ η θεωρία τυχαίας επιρροής ορίζει ομάδες χωρίς πραγματική επιρροή στους υπόλοιπους χρήστες. Αντίθετα οι υπόλοιπες τρεις θεωρίες επιρροής επιτυγχάνουν υψηλή επίδοση για τις περισσότερες κοινότητες.

Προκειμένου να αναλύσουμε περαιτέρω τα αποτελέσματα των πειραμάτων μας, εισαγάγαμε μια νέα μεθοδολογία που εξετάζει την εσωτερική δυναμική των επιδραστικών ομάδων της κάθε θεωρίας. Τα αποτελέσματα της ανάλυσης δείχνουν ότι οι επιδραστικές θεωρίες μπορούν να ταξινομηθούν σε τρεις βασικές κατηγορίες. Η πρώτη περιέχει θεωρίες που ορίζουν επιδραστικές ομάδες με μεγάλη συνοχή, τα μέλη των οποίων δρουν ως ομάδα και παρουσιάζουν υψηλά επίπεδα ομοφιλίας και επικοινωνίας. Οι άλλες δύο κατηγορίες επιλέγουν ως επιδραστικούς χρήστες που δε λειτουργούν συντονισμένα. Διαφέρουν, ωστόσο, στην ποιότητα των επιδραστικών χρηστών που ορίζουν, αν δηλαδή πρόκειται για συνηθισμένα μέλη της κοινότητας ή για μέλη που ξεχωρίζουν σε κάποιο τομέα της δραστηριότητας της κοινότητας.

Στο πλαίσιο της εργασίας με δεδομένα από πολλαπλά κοινωνικά δίκτυα, αναπτύξαμε επίσης μια εφαρμογή που επιτρέπει το συνδυασμό δεδομένων και λειτουργικότητας από ετερογενείς πλατφόρμες κοινωνικής δικτύωσης από μια εννιάια διεπαφή. Βασικό συστατικό της εφαρμογής αυτής είναι το μοντέλο δεδομένων το οποίο ορίζει ένα σύνολο οντοτήτων που αναπαριστούν θεμελιώδεις έννοιες των κοινωνικών δικτύων και τις μεταξύ τους συσχετίσεις.

Στο μέλλον, σκοπεύουμε να επεκτείνουμε τα αποτελέσματα της έρευνας μας ώστε να περιλαμβάνει και άλλα κοινωνικά δίκτυα εκτός από το Twitter, χρησιμοποιώντας την εφαρμογή που περιγράψαμε. Σκοπεύουμε, επίσης, να εξετάσουμε εναλλακτικές μετρικές περίληψης δραστηριότητας εκτός από το βαθμό πόλωσης ώστε να συγκρίνουμε και να μελετήσουμε διαφορετικούς τύπους επιρροής. Σχεδιάζουμε τέλος να ερευνήσουμε τη δυνατότητα να συνδυάσουμε τις τρεις αποδοτικές θεωρίες (βαθμού εισόδου, αναφορών και αναδημοσιεύσεων) σε μία συνδυαστική θεωρία που επιτυγχάνει καλύτερη απόδοση από της επιμέρους θεωρίες που την αποτελούν.



## ***Βιβλιογραφικές Αναφορές***

- [1] E. Bakshy, J. Hofman, W. Mason και D. Watts, «Everyone's an influencer: quantifying influence on twitter,» σε *ourth ACM international conference on Web search and data mining (WSDM '11)*, Hong Kong, 2011.
- [2] E. Keller, B. Fay και J. Berry, «Leading the Conversation: Influencers' Impact on Word of Mouth and the Brand Conversation,» σε *WOMMA Measuring Word of Mouth: Current Thinking on Research and Measurement of Word of Mouth Marketing*, 2007.
- [3] D. Brown και N. Hayes, *Influencer marketing: who really influences your customers?*, Elsevier/Butterworth-Heinemann, 2008.
- [4] M. Cha, H. Haddadi, F. Benevenuto και K. Gummadi, «Measuring User Influence in Twitter: The Million Follower Fallacy,» σε *ICWSM*, Washington, 2010.
- [5] J. Weng, E.-P. Lim, J. Jiang και Q. He, «TwitterRank: finding topic-sensitive influential twitterers,» σε *Third ACM international conference on Web search and data mining (WSDM '10)*, New York, 2010.
- [6] S. Milgram, «The small-world problem,» *Psychology Today*, τόμ. 1, αρ. 1, pp. 61-67, 1967.
- [7] J. Kleinberg, «Complex Networks and Decentralized Search Algorithms,» *Proceedings of the International Congress of Mathematicians (ICM)*, 2006.
- [8] Y.-R. Lin, Y. Chi, S. Zhu, H. Sundaram και B. L. Tseng, «Facetnet: a framework for analyzing communities and their evolutions in dynamic networks.,» σε *Proceedings of the 17th international conference on World Wide Web (WWW '08)*, New York, USA, 2008.
- [9] Z. Zhao, S. Feng, Q. Wang, J. Z. Huang, G. J. Williams και J. Fan, «Topic oriented community detection through social objects and link analysis in social networks,» *Knowledge-Based Systems*, τόμ. 26, p. 164–173, 2012.

- [10] Y. Ding, «Community detection: Topological vs. topical,» *Journal of Informetrics*, τόμ. 5, αρ. 4, p. 498–514, 2011.
- [11] 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 '03)*, New York, NY, USA, 2003.
- [12] C. Wang, V. Satuluri και S. Parthasar, «Local Probabilistic Models for Link Prediction,» σε *Proceedings of the 2007 Seventh IEEE International Conference on Data Mining (ICDM '07)*, Washington, DC, USA, 2007.
- [13] G. Chaudhari, V. Avadhanula και S. Sarawagi, «A few good predictions: selective node labeling in a social network,» σε *In Proceedings of the 7th ACM international conference on Web search and data mining (WSDM '14)*, 2014.
- [14] S. A. Macskassy και F. Provost, «A Simple Relational Classifier,» σε *Proceedings of the Second Workshop on Multi-Relational Data Mining (MRDM-2003) at KDD-2003*, 2003.
- [15] J. Leskovec , D. Chakrabarti, J. Kleinberg, C. Faloutsos και Z. Ghahramani, «Kronecker Graphs: An Approach to Modeling Networks,» *The Journal of Machine Learning Research*, τόμ. 11, pp. 985-1042, 2010.
- [16] E. Katz, P. Lazarsfeld και E. Roper, *Personal influence: The part played by people in the flow of mass communications*, Transaction Publishers, 2005.
- [17] S. Petrovic, M. Osborne και V. Lavrenko, «RT to Win! Predicting Message Propagation in Twitter,» σε *ICWSM*, Barcelona, Spain, 2011.
- [18] H. Purohit, J. Ajmera, S. Joshi, A. Verma και A. Sheth, «Finding Influential Authors in Brand-Page Communities,» σε *ICWSM*, Dublin, Ireland, 2012.
- [19] H. Li, S. Bhowmick και A. Sun, «CASINO: towards conformity-aware social influence analysis in online social networks,» σε *CIKM*, Glasgow, UK, 2011.
- [20] L. Liu, J. Tang, J. Han, M. Jiang και S. Yang, «Mining topic-level influence in heterogeneous networks,» σε *CIKM*, Toronto, Canada, 2010.
- [21] D. Gayo-Avello, «Nepotistic relationships in Twitter and their impact on rank prestige algorithms,» *Information Processing & Management*, pp. 1250-1280, 2013.
- [22] A. Goyal, F. Bonchi και L. Lakshmanan, «Learning influence probabilities in social networks,» σε *WSDM*, New York, USA, 2010.
- [23] C. Tan, J. Tang, J. Sun, Q. Lin και F. Wang, «Social action tracking via noise tolerant time-varying factor graphs,» σε *16th ACM SIGKDD international conference on Knowledge*

- discovery and data mining (KDD '10)*, Washington, 2010.
- [24] D. Cosley, D. Huttenlocher, J. Kleinberg, X. Lan και S. Suri, «Sequential Influence Models in Social Networks,» σε *ICWSM*, Washington, 2010.
- [25] W. Chen, C. Wang και Y. Wang, «Scalable influence maximization for prevalent viral marketing in large-scale social networks,» σε *16th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '10)*, Washington, 2010.
- [26] P. Domingos και M. Richardson, «Mining the network value of customers,» σε *KDD*, 2001.
- [27] D. Kempe, J. Kleinberg και E. Tardos, «Maximizing the spread of influence through a social network,» σε *KDD*, Washington, DC, USA, 2003.
- [28] H. Ma, H. Yang, M. Lyu και I. King, «Mining social networks using heat diffusion processes for marketing candidates selection,» σε *CIKM*, Napa Valley, California, 2008.
- [29] M. D. Luu, E.-P. Lim, T.-A. Hoang και F. C. T. Chua, «Modeling Diffusion in Social Networks Using Network Properties,» σε *ICWSM*, Dublin, Ireland, 2012.
- [30] M. McPherson, L. Smith-Lovin και J. Cook, «Birds of a feather: Homophily in social networks,» *Annual review of sociology*, τόμ. 27, pp. 415-444, 2001.
- [31] A. Anagnostopoulos, R. Kumar και M. Mahdian, «Influence and correlation in social networks,» σε *14th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '08)*, Las Vegas, 2008.
- [32] S. Aral, L. Muchnik και A. Sundararajan, «Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks,» *Proceedings of the National Academy of Sciences*, 2009.
- [33] T. L. Fond και J. Neville, «Randomization tests for distinguishing social influence and homophily effects,» σε *WWW*, 2010.
- [34] P. Holme και M. E. J. Newman, «Nonequilibrium phase transition in the coevolution of networks and opinions,» *Phys. Rev. E*, 2006.
- [35] J. Scripps, P.-N. Tan και A.-H. Esfahanian, «Measuring the effects of preprocessing decisions and network forces in dynamic network analysis,» σε *15th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '09)*, Paris, France, 2009.
- [36] L. Barbosa και J. Feng, «Robust sentiment detection on Twitter from biased and noisy data,» σε *23rd International Conference on Computational Linguistics: Posters (COLING)*

- '10), Beijing, China, 2010.
- [37] M. Tsytsarau και T. Palpanas, «Survey on mining subjective data on the web,» *Data Mining and Knowledge Discovery*, τόμ. 24, αρ. 3, pp. 478-514, 2012.
- [38] G. Giannakopoulos, P. Mavridi, G. Paliouras, G. Papadakis και K. Tserpes, «Representation models for text classification: a comparative analysis over three web document types,» σε *2nd International Conference on Web Intelligence, Mining and Semantics (WIMS '12)*, Craiova, Romania, 2012.
- [39] A. Go, R. Bhayani και L. Huang, *Twitter Sentiment Classification using Distant Supervision*, 2010.
- [40] D.-A. Ernesto, D. Lucas, G. Zeno, S.-T. Lars και N. Wolfgang, «What is Happening Right Now ... That Interests Me? Online Topic Discovery and Recommendation in Twitter,» σε *21st ACM international conference on Information and knowledge management (CIKM '12)*, New York, 2012.
- [41] A. Pal και S. Counts, «Identifying Topical Authorities in Microblogs,» σε *fourth ACM international conference on Web search and data mining (WSDM '11)*, New York, 2011.
- [42] A. Bruns και J. Burgess, «The use of Twitter hashtags in the formation of ad hoc publics,» σε *6th European Consortium for Political Research General Conference*, Reykjavik, Iceland, 2011.
- [43] J. Yang και J. Leskovec, «Patterns of temporal variation in online media,» σε *Fourth ACM international conference on Web search and data mining (WSDM '11)*, Hong Kong, 2011.
- [44] H. Kwak, C. Lee, H. Park και S. Moon, «What is Twitter, a social network or a news media?,» σε *WWW*, Raleigh, North Carolina, USA, 2010.
- [45] A. Zubiaga, D. Spina, V. Fresno και R. Martínez, «Classifying trending topics: a typology of conversation triggers on Twitter.,» σε *CIKM*, Glasgow, 2011.
- [46] A. Cui, M. Zhang, Y. Liu, S. Ma και K. Zhang, «Discover breaking events with popular hashtags in twitter,» σε *21st ACM international conference on Information and knowledge management (CIKM '12)*, Maui, Hawaii, 2012.
- [47] B. O'Connor, R. Balasubramanyan, B. R. Routledge και N. A. Smith, *From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series*, 2010.
- [48] M. Abbasi, S. Chai, H. Liu και K. Sagoo, «Real-World Behavior Analysis through a Social Media Lens,» σε *In Proceedings of the 5th international conference on Social Computing, Behavioral-Cultural Modeling and Prediction (SBP'12)*, Maryland, 2012.

- [49] «The Source for Social Data - Gnip.,» [Ηλεκτρονικό]. Available: <http://gnip.com/>. [Πρόσβαση 18 7 2014].
- [50] «Social Media Management Dashboard - Hootsuite,» [Ηλεκτρονικό]. Available: <https://hootsuite.com/>. [Πρόσβαση 10 1 2015].
- [51] «DataSift | Powering the Social Economy.,» [Ηλεκτρονικό]. Available: <http://datasift.com/>. [Πρόσβαση 10 1 2015].
- [52] P. Mika, «Ontologies Are Us: A Unified Model of Social Networks and Semantics,» σε *The Semantic Web – ISWC 2005*, Y. Gil, E. Motta, V. R. Benjamins, and M. A. Musen, Eds. Springer Berlin Heidelberg., 2005, p. 522–536.
- [53] «OpenSocial,» 2015. [Ηλεκτρονικό]. Available: <http://www.w3.org/blog/2014/12/opensocial-foundation-moves-standards-work-to-w3c-social-web-activity/>. [Πρόσβαση 12 2 2015].



Η Δρ. **Μαγδαληνή Δ. Καρδάρá** γεννήθηκε στην Αθήνα το 1981. Αποφοίτησε από τη Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών του Εθνικού Μετσόβιου Πολυτεχνείου το 2004, ενώ το 2006 απέκτησε μεταπτυχιακό δίπλωμα στα Προηγμένα θέματα της Επιστήμης των Υπολογιστών (MSc in Advanced Computing) από το Imperial College London. Το 2015, ολοκλήρωσε τη διδακτορική της διατριβή στο Εθνικό Μετσόβιο Πολυτεχνείο με θέμα την ανάπτυξη τεχνικών ανάκτησης περιεχομένου και ανάλυσης διάχυσης της επιρροής στα κοινωνικά δίκτυα. Από το 2007 εργάζεται ως ερευνητικός συνεργάτης στο Εργαστήριο Κατανεμημένων Συστημάτων Διαχείρισης Γνώσης της Σχολής Ηλεκτρολόγων Μηχανικών και Μηχανικών Ηλεκτρονικών Υπολογιστών του Εθνικού Μετσόβιου Πολυτεχνείου και συμμετέχει σε ερευνητικά προγράμματα χρηματοδοτούμενα από την Ευρωπαϊκή Ένωση. Τα ερευνητικά της ενδιαφέροντα εστιάζονται στον τομέα της ανάλυσης δεδομένων από κοινωνικά δίκτυα.