

**ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ**

**ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ**



**ΠΡΟΓΡΑΜΜΑ ΜΕΤΑΠΤΥΧΙΑΚΩΝ ΣΠΟΥΔΩΝ**

**ΣΤΑ ΠΡΟΗΓΜΕΝΑ ΣΥΣΤΗΜΑΤΑ ΠΛΗΡΟΦΟΡΙΚΗΣ**

**ΕΙΣΗΓΗΤΙΚΑ ΣΥΣΤΗΜΑΤΑ ΚΑΙ ΠΡΟΣΤΑΣΙΑ**

**ΠΡΟΣΩΠΙΚΩΝ ΔΕΔΟΜΕΝΩΝ**

**Δήμητρα Καλέμου**

Διπλωματική Εργασία υποβληθείσα στο Τμήμα Πληροφορικής του Πανεπιστημίου Πειραιώς ως μέρος των απαιτήσεων για την απόκτηση Μεταπτυχιακού Διπλώματος στο Προηγμένα Συστήματα Πληροφορικής

Πειραιάς, Απρίλιος, 2015

Πανεπιστήμιο Πειραιώς

**UNIVERSITY OF PIRAEUS**  
**DEPARTMENT OF INFORMATICS**



**MASTER PROGRAM IN**  
**IN ADVANCED SYSTEMS**

**RECOMMENDATION SYSTEM AND PRIVACY**  
**POLICIES**

**By**  
**[Dimitra Kalemou]**

Master Thesis submitted to the Department of Informatics of the University of Piraeus in partial fulfillment  
of the requirements for the degree of Master of Advanced Systems

**Piraeus, Greece, April 2015**

Πανεπιστήμιο Πειραιώς



## ***Περίληψη***

Σκοπός της διπλωματικής εργασίας είναι η μελέτη Εισηγητικών Συστημάτων Διαδικτύου που βασίζονται στη μοντελοποίηση των προτιμήσεων χρήστη και στη διήθηση της πληροφορίας με έμφαση στην προστασία των προσωπικών δεδομένων των χρηστών.

Τα τελευταία χρόνια υπάρχει έντονος προβληματισμός για τον τρόπο με τον οποίο μπορεί να λυθεί το πρόβλημα της επιτυχούς αναζήτησης πληροφοριών στις δεξαμενές γνώσης του Διαδικτύου. Τα εισηγητικά συστήματα παρέχουν μια λύση στο πρόβλημα με τις εξατομικευμένες συστάσεις τους. Για να επιτευχθεί η εξατομίκευση, η επιτυχής δημιουργία, η συντήρηση και η εκμετάλλευση του προφίλ κάθε χρήστη αποδεικνύονται ότι είναι ιδιαίτερα κρίσιμες διαστάσεις. Για το λόγο αυτό κρίθηκε χρήσιμη η παρουσίαση μιας ταξινομησης και σύντομης ανάλυσης αυτών των διαστάσεων.

Στη διπλωματική εργασία παρουσιάζονται οι μέθοδοι του φιλτραρίσματος περιεχομένου και του συνεργατικού φιλτραρίσματος καθώς και υβριδικά εισηγητικά συστήματα. Η μελέτη δίνει μεγάλη σημασία στην προστασία των χριστών και στα προσωπικά δεδομένα και ολοκληρώνεται με την παρουσίαση ενός μοντέλου εισηγητικού συστήματος που δίνει έμφαση στην ανωνυμία των χρηστών .

### **Λέξεις Κλειδιά**

Εισηγητικό Σύστημα, Διήθηση Πληροφοριών, Συνεργατικό Φιλτράρισμα, Φιλτράρισμα Περιεχομένου, Προστασία Προσωπικών Δεδομένων .

Πανεπιστήμιο Πειραιώς

## *Abstract*

The aim of the thesis is the study of Internet Recommendation Systems that are based on user preferences modeling and on information filtering, and the implementation of an application, as well.

Last years there is a lot of speculation among the members of the community of artificial intelligence concerning the way AI can help with the problem of successful information search in the reservoirs of knowledge of Internet. Recommendation systems provide a solution in this problem giving individualized recommendations. The successful generation, maintenance and exploitation of single user profiles are proved particularly critical dimensions in order to achieve individualization. That is why the presentation of a classification and short analysis of these dimensions is considered to be useful.

In the thesis are presented the methods of content-based and collaborative filtering and hybrid recommendation systems. The study provides important in the protection of the users and the personal data and it is completed with the presentation of a Recommendation System model that emphasizes in the user's anonymity

### Key Words

Recommender, Recommendation System, Information Filtering, Content-Based Filtering, Collaborative Filtering, k-Anonymity, Microaggregation

Πανεπιστήμιο Πειραιώς

## Πίνακας Περιεχομένων

Abstract .....	vii
Πίνακας Σχημάτων.....	x
Πίνακας Πινάκων .....	xi
<b>1. Εισαγωγή.....</b>	<b>13</b>
1.0 Γενικά .....	13
1.1. Ορισμός του Εισηγητικού Συστήματος.....	14
1.2. Δομή της Διπλωματικής Εργασίας.....	16
<b>2. Διατήρηση της προστασίας των προσωπικών δεδομένων στο συνεργατικό φιλτράρισμα.....</b>	<b>17</b>
2.0 Εισαγωγή.....	17
2.1 Η Απορρήτου Διατήρηση του Συνεργατικού φιλτραρίσματος.....	22
2.1.1 Η Κεντρική Μέθοδος PPCF .....	24
2.1.2 Η Αποκεντρωμένη Μέθοδος PPCF .....	26
2.1.3 Η Υβριδική Μέθοδος PPCF .....	32
2.1.3.1 Υβριδικές τεχνικές συνεργατικού φιλτραρίσματος .....	32
2.1.3.2 Υβριδικά CF και Content- Based χαρακτηριστικά.....	33
2.1.3.3 Υβριδικά CF και άλλα εισηγητικά συστήματα .....	36
2.1.3.4 Υβριδικά συστήματα σε σχέση με τα απλά CF.....	37
2.2 Νέες τεχνικές και θέματα που προκύπτουν.....	38
2.3 Επιπλέον προβλήματα της PPCF .....	40
<b>3. Βασικές έννοιες.....</b>	<b>41</b>
3.0 Εισαγωγή.....	41
3.1 Η Κανονική Κατανομή(Gauss).....	41
3.2 Η MDAV.....	43
3.3 Η SVD.....	46
3.3.1 Εφαρμογή της SVD στο συνεργατικό φιλτράρισμα .....	46
3.3.2 Singular Value Decomposition (SVD).....	47
<b>4. PPCF με k-Ανωνυμία μέσω μικροσυσσωμάτωσης.....</b>	<b>49</b>
4.0 Εισαγωγή.....	49

4.1 Η Κ – ανωνυμία με μικροσυσσωμάτωση.....	52
4.2 Το Υπόβαθρο .....	52
4.2.1 Η Προστασία των Προσωπικών Δεδομένων και η Διατήρηση του Συνεργατικού Φιλτραρίσματος .....	52
4.2.2 Έλεγχος αποκάλυψης στατιστικών στοιχείων και μικροσυσσωμάτωση .....	55
4.3 Η Προτεινόμενη Μέθοδος.....	56
4.3.1 Τα Πειραματικά Αποτελέσματα.....	59
4.3.2 Η Προστασία της ιδιωτικής ζωής .....	59
4.3.3 Η Προγνωστική ακρίβεια.....	63
4.4 Σύγκριση και συζήτηση .....	65
<b>5. Συμπεράσματα.....</b>	<b>67</b>
5.0 Σύνοψη και Συμπεράσματα.....	67
5.1 Μελλοντικές επεκτάσεις .....	67
<b>Βιβλιογραφία .....</b>	<b>69</b>

## **Πίνακας Σχημάτων**

Σχήμα 3.1 Η εξόδος του αλγορίθμου MDAV για το παράδειγμα των παιχνιδιών .....	45
Σχήμα 4.2 Σχήμα Κεντρικής PPCF. Ο χρήστης κάνει μια αίτηση σε ένα στοιχείο στο διακομιστή, ο οποίος απαντά με μια εξατομικευμένη πρόβλεψη. ....	56
Σχήμα 4.3 Σχήμα ομαδοποίησης MDAV. Τα βήματα ροής από δεξιά (αρχικό σύνολο δεδομένων) προς τα αριστερά (ασαφές σύνολο δεδομένων).....	57
Σχήμα 4.4 Προτεινόμενη βήμα προς βήμα μέθοδος .....	58
Σχήμα 4.5 Οι SEE τιμές της μεθόδου μας για Movielens των 100k δεδομένων .....	61
Σχήμα 4.6 Οι τιμές DR της μεθόδου μας, για Movielens των 100k δεδομένων.....	62
Σχήμα 4.7 τιμές SSE της μεθόδου GNA για Movielens των 100k δεδομένων.....	62
Σχήμα 4.8 Οι τιμές DR της μεθόδου GNA για Movielens των 100k δεδομένων .....	63
Σχήμα 4.9 Σχέση μεταξύ της SSE και DR για τις μεθόδους που αναλύθηκαν και αφορούν την Movielens των 100k βάση δεδομένων .....	66

## ***Πίνακας Πινάκων***

Πίνακας 2.1:Είναι παράδειγμα δεδομένων μήτρας όπου κάθε σειρά αντιστοιχεί σε ένα χρήστη $U_i$ και κάθε στήλη αντιστοιχεί σε ένα στοιχείο $I_i$ . Κάθε κελί αποθηκεύει μια ψηφοφορία στο πλαίσιο ενός εύρους τιμών από το 1 έως το 5. Τα κενά κελία δεν έχουν αξιολογηθεί από ένα συγκεκριμένο χρήστη. <i>ΣΗΜΕΙΩΣΗ</i> η στήλη με αριθμό 3, με γκριζό χρώμα, δεν είχε ακόμη εκτιμηθεί.[72].....	19
Πίνακας 2.2 : Το Περιεχόμενο του ενισχυμένου CF και οι παραλλαγές του (α) τα δεδομένα περιεχομένου και αρχικά ελάχιστα στοιχεία για την αξιολόγηση των δεδομένων b) ψευδοαξιολογημένα δεδομένα που γेमίζουν από προγνωστικό περιεχόμενο (γ) προβλέψεις από (σταθμισμένη) Pearson CF για τα ψευδοαξιολογημένα δεδομένα. ....	35
Πίνακας 4.3 Αποτελέσματα του MDAV βασισμένα στο PPCF. Τα SSE αποτελέσματα εμφανίζονται στην κλίμακα $10^3$ .....	64
Πίνακας 4.4 Αποτελέσματα του GNA βασισμένα στο PPCF. Τα SSE αποτελέσματα εμφανίζονται στην κλίμακα $10^3$ .....	64
Πίνακας 4.5 Το MAE περιέχει τιμές, συγκρίνοντας τις μήτρες πρόβλεψης με την αυθεντική βάση δεδομένων .....	65

Πανεπιστήμιο Πειραιώς



# **1. Εισαγωγή**

## **1.0 Γενικά**

Η ραγδαία ανάπτυξη του Διαδικτύου σηματοδότησε την είσοδο σε μια νέα περίοδο πληροφορίας. Ο Παγκόσμιος Ιστός παρέχει ένα καινούργιο μέσο επικοινωνίας το οποίο ξεπερνά κατά πολύ τα παραδοσιακά μέσα επικοινωνίας όπως το ραδιόφωνο, το τηλέφωνο και η τηλεόραση. Ο Ιστός, έχει ακόμα σημαντικό αντίκτυπο τόσο στην ακαδημαϊκή έρευνα όσο και στην καθημερινή ζωή. Έχει φέρει επανάσταση στον τρόπο με τον οποίο η πληροφορία συλλέγεται, αποθηκεύεται, επεξεργάζεται, παρουσιάζεται, μοιράζεται και χρησιμοποιείται. Δεδομένα υπό τη μορφή κειμένου, εικόνας και αρχείων video βρίσκονται άφθονα στο Διαδίκτυο και είναι εύκολα προσβάσιμα από όποιον τα ζητήσει.

Το γεγονός όμως ότι οι πληροφορίες γενικά είναι εύκολα προσβάσιμες δε σημαίνει πως είναι και εύκολα αναζητήσιμες οι πληροφορίες που αναζητάμε. Η τεράστια ποσότητα πληροφοριών που υπάρχει στο διαδίκτυο είναι δυνατόν να οδηγήσει σε αχρήστευση πολύ σημαντικού μέρους των πληροφοριών εφόσον μένει ανοργάνωτη, καθώς είναι δύσκολο να εντοπισθεί και να χρησιμοποιηθεί από τον απλό χρήστη. Είναι πολύ σημαντικό η αναζήτηση να γίνεται με αποδοτικό τρόπο ώστε να είναι γρήγορη και οικονομική. Αυτό συνεπάγεται πως χρειάζεται ένας αυτόματος τρόπος αναζήτησης.

Η ταξινόμηση των κειμένων είναι μία πολύ χρήσιμη διαδικασία η οποία μας επιτρέπει την οργάνωση ενός σώματος κειμένων με βάση κάποιο κοινό χαρακτηριστικό, διευκολύνοντας την αναζήτηση στο σώμα αυτό. Συχνά η ταξινόμηση κειμένων γίνεται σε κατηγορίες που έχουν προεπιλεγεί από ένα σύστημα με αποτέλεσμα να απαιτείται η προσαρμογή του χρήστη στους κανόνες λειτουργίας του συστήματος. Αυτό δυσκολεύει την αναζήτηση μιας και δεν είναι άμεσα κατανοητό από την ανθρώπινη λογική, σε ποια κατηγορία του συστήματος ανήκει ένα κείμενο. Άλλη μία λύση στο συγκεκριμένο πρόβλημα είναι η κατάταξη ενός κειμένου σε περισσότερες από μία κατηγορίες. Αυτή τη λύση έχει υιοθετήσει, για παράδειγμα το Yahoo[12] στους καταλόγους του.

Άλλες λύσεις αποτελούν οι εξατομικευμένες μηχανές αναζήτησης, οι ευφυείς πράκτορες λογισμικού και τα εισηγητικά συστήματα και έχουν γίνει ευρέως αποδεκτά από

την κοινότητα των χρηστών. Υπάρχουν εισηγητικά συστήματα για την πρόταση ιστοσελίδων, netnews, εστιατορίων, ταινιών κ.ά. Αυτά τα συστήματα βασίζονται σε ένα συνδυασμό μοντελοποίησης των προτιμήσεων συγκεκριμένων χρηστών, στη δημιουργία προτύπων περιεχομένου και στη μοντελοποίηση κοινωνικών προτύπων για να εισηγηθούν επιτυχημένων προτάσεων. [110]

### **1.1. Ορισμός του Εισηγητικού Συστήματος**

Θα ήταν λοιπόν χρήσιμο να οριστεί το εισηγητικό σύστημα (recommender) για να γίνει κατανοητό πως αποτελεί μια πολύπλοκη οντότητα, η επιτυχία της οποίας εξαρτάται από την αποτελεσματική συνεργασία των δομικών συστατικών της.

*Εισηγητικό είναι ένα σύστημα που αναλύει τις προτιμήσεις ενός συγκεκριμένου χρήστη και αναγνωρίζει ένα σύνολο στοιχείων που του προτείνει ως ενδιαφέροντα.*

Από τον ορισμό του συστήματος γίνεται κατανοητό πως είναι πολύ σημαντικό να μπορούν να αναγνωρισθούν, να αναλυθούν και να αξιοποιηθούν αποδοτικά τα ιδιαίτερα χαρακτηριστικά κάθε χρήστη (προφίλ). Τα μοντέλα και οι τεχνικές που ακολουθούνται κάθε φορά εξαρτώνται από τις ιδιαίτερες ανάγκες του κάθε συστήματος. Γενικά, μπορούν να αναγνωρισθούν τα ακόλουθα δομικά στοιχεία :

- Αναπαράσταση Προφίλ
- Αρχικό Προφίλ
- Τεχνικές εκμάθησης προφίλ
- Ανατροφοδότηση σχετικότητας

- Τεχνικές προσαρμογής προφίλ
- Μέθοδοι φιλτραρίσματος πληροφοριών
- Ταίριασμα στοιχείων- προφίλ χρήστη
- Ταίριασμα παραμέτρων χρήστη

Η παραγωγή και συντήρηση ενός ακριβούς προφίλ αποτελούν ένα βασικό τμήμα του συστήματος. Η επιλογή της κατάλληλης αναπαράστασης προηγείται κάθε άλλης ενέργειας, αφού οι άλλες τεχνικές βασίζονται σε αυτή. Άλλωστε, το σύστημα δεν μπορεί να λειτουργήσει αν δεν υπάρχουν παράμετροι χρήστη. Επιπλέον, το σύστημα χρειάζεται να γνωρίζει όσο το δυνατόν περισσότερα για το χρήστη ώστε να είναι σε θέση να του παρέχει ικανοποιητικά αποτελέσματα από την αρχή. Για αυτό το λόγο είναι συνήθως απαραίτητες και τεχνικές που δημιουργούν ένα αρχικό προφίλ.

Η συνεχής αλληλεπίδραση του χρήστη με το σύστημα είναι απαραίτητη όχι μόνο για την αναγνώριση των χαρακτηριστικών γνωρισμάτων του χρήστη και την αξιολόγηση των προτάσεων αλλά και για την προσαρμογή του προφίλ αφού είναι δυνατόν οι προτιμήσεις ενός ατόμου να αλλάζουν με την πάροδο του χρόνου. Οι πληροφορίες αυτές αποτελούν την ανατροφοδότηση σχετικότητας και μπορεί να εξαχθούν είτε με άμεσο τρόπο όπως, ερωτήσεις προς το χρήστη, είτε με έμμεσο, όπως με το χρόνο που αυτός διαθέτει σε μια ιστοσελίδα.

Για να αξιοποιηθούν οι παραπάνω πληροφορίες χρησιμοποιούνται τεχνικές εκμάθησης προφίλ, οι οποίες αναγνωρίζουν τις σχετικές πληροφορίες και τις χρησιμοποιούν για την προσαρμογή των προφίλ, ανάλογα βέβαια με την αναπαράσταση.

Εφόσον λοιπόν έχουν καθοριστεί τα παραπάνω για το προφίλ, ακολουθεί η εκμετάλλευσή του προκειμένου να προταθούν αντικείμενα ή υπηρεσίες. Οι αποφάσεις για τις προτάσεις που θα γίνουν λαμβάνονται σύμφωνα με τις υπάρχουσες πληροφορίες, και για την ακρίβεια, σύμφωνα με το χρήσιμο τμήμα αυτών των πληροφοριών. Υπάρχει δηλαδή η ανάγκη διήθησης των πληροφοριών ώστε η ποιότητα των συστάσεων να είναι καλύτερη αλλά και να επιτύχουμε οικονομία χώρου και χρόνου. Οι τεχνικές που χρησιμοποιούνται συνήθως είναι η δημογραφική, η συνεργατική και η ανάλυση περιεχομένου. Η δημογραφική αντιστοιχεί αντικείμενα προτάσεων σε τύπους χρηστών. Η συνεργατική λαμβάνει υπόψη την

ανατροφοδότηση άλλων χρηστών και η μέθοδος που βασίζεται στο περιεχόμενο αναγνωρίζει τη σχέση ενός συγκεκριμένου χρήστη με το περιεχόμενο των αντικειμένων. Κάθε μέθοδος έχει πλεονεκτήματα και μειονεκτήματα. Συνήθως οι μέθοδοι συνδυάζονται για να εξουδετερωθούν τα μειονεκτήματά τους.

Για να την αντιστοιχίσει χρηστών-αντικειμένων και το ταίριασμα χρηστών-χρηστών χρησιμοποιούνται διάφορες μέθοδοι. Οι πιο σημαντικές είναι η ομοιότητα συνημιτόνου, η συσχέτιση Pearson, η κατηγοριοποίηση και οι αλγόριθμοι πλησιέστερων γειτόνων.[110]

## **1.2. Δομή της Διπλωματικής Εργασίας**

Στο επόμενο κεφάλαιο θα παρουσιάσουμε μια προσπάθεια ταξινόμησης των συστημάτων recommender σύμφωνα με το είδος και τις τεχνικές που εφαρμόζονται στα δομικά συστατικά τους. Παρουσιάζονται αναλυτικά οι κεντρικές αποκεντρωμένες και υβριδικές μέθοδοι και με τη χρήση ποιών μεθόδων εφαρμόζονται στη βιβλιογραφία. . Γίνεται αναφορά στα πλεονεκτήματα αλλά και τα μειονεκτήματα της κάθε μεθόδου

Στο κεφάλαιο 3 γίνεται ιδιαίτερη αναφορά στις βασικές μεθόδους πάνω στις οποίες στηρίζονται περεταίρω θέματα και αποτελούν το υπόβαθρο ώστε να κατανοηθούν οι περεταίρω έννοιες. Παρουσιάζεται η Γκαουσιανή (GAUSSIAN) μέθοδος , η MDAV και η SVD.

Στο κεφάλαιο 4 ειδικά παρουσιάζεται μια εφαρμογή που έχει προστεθεί στην έρευνα και στη βιβλιογραφία και αφορά την προστασία των προσωπικών δεδομένων των χρηστών που τροφοδοτούν τα εισηγητικά συστήματα.

Τέλος, αναφέρονται τα συμπεράσματα που προκύπτουν από τη μελέτη των εισηγητικών συστημάτων γενικά, αλλά και ειδικά στην περίπτωση μας που δίνουμε έμφαση στην προστασία των δεδομένων, καθώς και ιδέες για μελλοντική εργασία.

## **2. Διατήρηση της προστασίας των προσωπικών δεδομένων στο συνεργατικό φιλτράρισμα**

### **2.0 Εισαγωγή**

Τα αυτόματα συστήματα συστάσεων έχουν γίνει ακρογωνιαίος λίθος του ηλεκτρονικού εμπορίου, ειδικά μετά την μεγάλη αποδοχή της συμμετοχής και της αλληλεπίδρασης των χρηστών του Διαδικτύου στο Web 2.0. Το συνεργατικό φιλτράρισμα (CF) είναι ένα σύστημα συστάσεων που γίνεται ολοένα και πιο σημαντικό για τη βιομηχανία λόγω της ανάπτυξης του Διαδικτύου, στο οποίο έχει καταστήσει πολύ πιο δύσκολο να εξάγουμε αποτελεσματικά χρήσιμες πληροφορίες.

Στην εργασία αυτή, θα παρουσιάσουμε μια ταξινόμηση των διαφόρων CF και θα συζητήσουμε τις πιο σχετικές με την Προστασία Προσωπικών Δεδομένων και τη Διατήρηση τους με συνεργατικό φιλτράρισμα (PPCF), μεθόδους που αναφέρονται και στη βιβλιογραφία. Επιπλέον, για να καταλάβουν τις εγγενείς προκλήσεις του PPCF, εμείς διεξάγουμε μια ανασκόπηση των τρεχουσών τάσεων και τα σημαντικά μειονεκτήματα των συστημάτων συστάσεων του τύπου αυτού, προτείνοντας διάφορες στρατηγικές για να ξεπεραστούν κάποιες ελλείψεις.

Τα εισηγητικά συστήματα [1], προέρχονται από την ανακάλυψη της γνώσης σε βάσεις δεδομένων (KDD) [2], [3]. Σήμερα, η μεγαλύτερη πηγή των συστάσεων είναι το Διαδίκτυο. Από τη μία πλευρά, το Διαδίκτυο παρέχει έναν πλούτο πληροφοριών σχετικά με μια τεράστια ποικιλία προϊόντων και υπηρεσιών που μπορεί να είναι χρήσιμη για τους δυνητικούς αγοραστές. Από την άλλη πλευρά, μια τέτοια ποσότητα των πληροφοριών μπορεί να γίνει ένα πρόβλημα και όχι η λύση, διότι μπορεί να παρεμποδίσει την λήψη αποφάσεων.

Το συνεργατικό φιλτράρισμα (CF) [4], [5] είναι ένα συνιστών σύστημα το οποίο έχει ως στόχο να κάνει προτάσεις για τα στοιχεία (π.χ. βιβλία, μουσική, ή ταινίες) με βάση τις προτιμήσεις των χρηστών που έχουν ήδη αποκτήσει ή / και βαθμολογήσει αυτά τα στοιχεία. Η CF εφαρμόζεται, προκειμένου να παρέχετε αυτόματη σύσταση σε ένα ψηφιακό περιβάλλον. Είναι ενεργά παρούσα σε όλο το δίκτυο:

- Εταιρείες ηλεκτρονικού εμπορίου (π.χ. Amazon, eBay, Barnes & Noble) μπορούν να επωφεληθούν από την CF, προκειμένου να λάβουν αποτελεσματικά οφέλη.

- Η CF είναι μέρος της έννοιας του Web 2.0, η οποία ορίζεται ως το νέο τρόπο να χρησιμοποιούμε το δίκτυο. Αυτή η νέα τάση δίνει μεγάλη σημασία στην ενεργό συμμετοχή των χρηστών στην υποδομή (π.χ. blogs, κοινωνικά δίκτυα και Πληροφορίες & πύλες υπηρεσία).
- Χρησιμοποιείται για την ανάλυση των προτιμήσεων των πληροφοριών. Για παράδειγμα, οι ιστοσελίδες που επισκέπτεται ο χρήστης μπορεί να παρακολουθούνται και να χρησιμοποιηθεί αυτή η πληροφορία ώστε να προτείνονται παρόμοιες ιστοσελίδες με άλλες που αφορούν χρήστες με παρόμοια ενδιαφέροντα.
- Χρησιμοποιείται ευρέως στο μουσικό και οπτικοακουστικό πλαίσιο, με υπηρεσίες όπως το last.fm, το Spotify, MyStrands, Netflix και Moviefinder όπου η CF μπορεί να ωφελήσει με τη χρήση σε μορφή συμβουλών και παρέχει επίσης μια σαφή έρευνα αγοράς για τις εταιρείες.

Η Μέθοδος CF χρησιμοποιείται σε μεγάλες βάσεις δεδομένων που περιέχουν πληροφορίες σχετικά με τις τιμές / αγορές αγαθών των χρηστών. Αυτά τα δεδομένα μοντελοποιούνται ως μήτρες που αποτελούνται από τους χρήστες  $n$  και  $m$  τα είδη και κάθε κελί  $(n_i, m_j)$  αποθηκεύει την αξιολόγηση του  $i$  χρήστη για το στοιχείο  $j$ . Αυτά αποθηκεύουν τις αξιολογήσεις και εκπροσωπούνται από ένα εύρος τιμών (π.χ. μεταξύ 0 και 10) ή απλώς με ένα δυαδικό ψηφίο, το οποίο μπορεί να είναι θετικό ή αρνητικό (ή αγοράζονται και δεν αγοράζονται). Φαίνεται στον **Πίνακα 2.1** ένα μικρό παράδειγμα ενός πίνακα παιχνιδιού CF. Υπάρχουν πολλά παραδείγματα δεδομένων CF που αναφέρονται [6] στη βιβλιογραφία: Eachmovie, MovieLens, Jester. Αυτά χρησιμοποιούνται συχνά ως σημείο αναφοράς για την αξιολόγηση της αποτελεσματικότητας, της ποιότητας και της αξιοπιστίας της CF μεθόδου [7].

**Πίνακας 2.1:** Είναι παράδειγμα δεδομένων μήτρας όπου κάθε σειρά αντιστοιχεί σε ένα χρήστη  $U_i$  και κάθε στήλη αντιστοιχεί σε ένα στοιχείο  $I_i$ . Κάθε κελί αποθηκεύει μια ψηφοφορία στο πλαίσιο ενός εύρους τιμών από το 1 έως το 5. Τα κενά κελία δεν έχουν αξιολογηθεί από ένα συγκεκριμένο χρήστη. **ΣΗΜΕΙΩΣΗ** η στήλη με αριθμό 3, με γκριζό χρώμα, δεν είχε ακόμη εκτιμηθεί.[72]

	$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$
$U_a$	2	4			1	2
$U_b$		3		2		
$U_c$	3	2			2	2
$U_d$		5		1		1

Οι συστάσεις που παρέχονται βασίζονται σε μεθόδους CF με βάση την υπόθεση ότι παρόμοιοι χρήστες θα ενδιαφερθούν για τα ίδια προϊόντα. Με αυτόν τον τρόπο, τα αντικείμενα που αγοράζονται από ένα χρήστη  $u_a$  μπορούν να συνιστανται σε άλλο χρήστη  $u_b$ , αν  $u_a$  και  $u_b$  είναι παρόμοιοι χρήστες. Επίσης, υπάρχουν προσεγγίσεις που κάνουν συστάσεις με βάση στην ομοιότητα μεταξύ των στοιχείων. Στο πλαίσιο αυτό, πολλοί παρόμοιοι χρήστες ή αντικείμενα συνθέτουν μια γειτονιά.

Η τελευταία τεχνολογία στην CF, προτείνει μια ταξινόμηση των μεθόδων σε τρεις κύριες κατηγορίες, ανάλογα με τα δεδομένα που χρησιμοποιούν για να κάνουν τη σύσταση: η μέθοδος βασισμένη σε μνήμη (που χρησιμοποιούν την πλήρη μήτρα με όλες τις ειδικότητες), η μέθοδος με βάση το μοντέλο (που χρησιμοποιούν στατιστικά μοντέλα και τις λειτουργίες του πίνακα δεδομένων, αλλά όχι την πλήρη μήτρα δεδομένων) και η υβριδική μεθόδους (η οποία συνδυάζει τις προηγούμενες μεθόδους με περιεχόμενο που βασίζεται [6]στη μέθοδο σύστασης).

Στη μέθοδο CF που βασίζεται στη μνήμη, οι συστάσεις διατυπώνονται σε δύο βήματα: αναζήτηση γειτονιάς και την πρόβλεψη σύστασης. Δεδομένου ενός χρήστη  $u_a$ , η συσχέτιση και η απόσταση είναι λειτουργίες που χρησιμοποιούνται για να υπολογίσουν τη γειτονιά του. Οι πιο χρησιμοποιούμενες λειτουργίες συσχέτισης είναι η Pearson Correlation [5], ομοιότητα συνημίτονου [8] και η Ευκλείδεια απόσταση. Άλλες γνωστές μέθοδοι είναι οι μέσες τετραγωνικές διαφορές [9], η περιορισμένη Pearson [10] και τη σταθμισμένη Pearson [11]. Η ομοιότητα μεταξύ των χρηστών μπορεί επίσης να υπολογιστεί με ένα πολύ πιο αποτελεσματικό τρόπο, σύμφωνα με την συμπεριφορά τους όταν ψηφίζουν. Τα παραδείγματα



αυτών φαίνονται στο [12], στο οποίο υπολογίζονται οι τάσεις των χρηστών, ή [13], στο οποίο οι αντιστοιχίες των χρηστών υπολογίζονται με σεβασμό την ιδιωτική ζωή. Γενικά, για τον υπολογισμό των ομοιοτήτων μεταξύ των χρηστών, πρέπει να υπάρχει ένας επαρκής αριθμός κοινών βαθμολογούμενων αντικειμένων. Οι γειτονιές μπορούν να υπολογιστούν τόσο με στατικό τρόπο (τον προσδιορισμό του κοντινότερου γείτονα της  $u_a$ ) ή με ένα συγκεντρωτικό / δυναμικό τρόπο (Να καθοριστεί ο εγγύτερος γείτονας της  $u_a$  χρησιμοποιώντας το κέντρο βάρους της νέας ομάδας για να συγκεντρώσει τον επόμενο γείτονα). Από τη στιγμή που έχουμε σχηματίσει την γειτονιά του χρήστη  $u_a$ , μπορούμε να κάνουμε μια σύσταση χρησιμοποιώντας τις μεθόδους που προτείνονται στο [5], [14], ομοιότητα σύντηξης [15] ή σταθμισμένη πρόβλεψη με πλειοψηφία [16]. Αυτές οι μέθοδοι μπορούν να χρησιμοποιηθούν για να προβλέψουν την ψηφοφορία ή να συστήσει στοιχεία *top-N* για τον  $u_a$ .

Η μέθοδος με βάση το μοντέλο CF μπορεί να επεξεργαστεί ένα μοντέλο από την πλήρη μήτρα στο οποίο θα διατυπώνει συστάσεις. Η έμφαση δίνεται σε αυτές τις μεθόδους για τους περιορισμούς της memory based CF όσον αφορά την επεκτασιμότητα, την πολυπλοκότητα του υπολογισμού και τη σπανιότητα. Μερικές καλές και γνωστές μέθοδοι για τη μείωση της διάστασης της μήτρας είναι η Singular Value Decomposition (SVD) και η Principal Component Analysis (PCA). Ωστόσο, η χρήση της διάστασης των μεθόδων μείωσης θα μπορούσε να επηρεάσει την ποιότητα των συστάσεων εφόσον μειώνουν το εύρος των δεδομένων. Υπάρχει μια τεράστια ποικιλία από μοντέλα με βάση τις μεθόδους CF: εκτός από οι προαναφερόμενες μέθοδοι μείωσης διάστασης (SVD [2], υπάρχει RSVD, Βελτιωμένη RSVD, NSVD2 και SVD ++ [12], [17]), λανθάνουσα σημασιολογική μεθόδους [18], γραμμική παλινδρόμηση [19], οι μέθοδοι ομαδοποίησης [6], [20], Bayesian δικτυακά μοντέλα [6].

Οι σύνθετες μέθοδοι CF συνδυάζουν τη μέθοδο που βασίζεται σε μνήμη και την model based μέθοδο κατά τρόπο που να διατηρήσει τα πλεονεκτήματα των αλγορίθμων, ενώ, ελαχιστοποιεί τα μειονεκτήματα και τις ελλείψεις. Παραδείγματα αυτών είναι η διάγνωση της προσωπικότητας [22] και το Πιθανοτικό μοντέλο που βασίζεται σε μνήμη [23]. Οι Υβριδικοί μέθοδοι μπορούν επίσης να ληφθούν με τον συνδυασμό του μοντέλου της μνήμης βάσης και του μοντέλου με βάση τις μεθόδους, με περιεχόμενο που βασίζονται στα συστήματα



συστάσεων. Μερικά γνωστά παραδείγματα [6] είναι: Filterbots, Contentboosted, Fab και Ripper.

Ανεξάρτητα από τη μέθοδο CF που χρησιμοποιείται, υπάρχουν διάφορες εγγενείς περιορισμοί σε αυτού του είδους τα σύστημα συστάσεων. Μερικοί από τους πιο σημαντικούς περιορισμούς είναι η ιδιαίτερα χαμηλή επεκτασιμότητα, οι επιθέσεις, η συνωνυμία, η δωροδοκία, το αντίγραφο του προφίλ, και η έλλειψη της ιδιωτικής ζωής [6], [24].

Η χρήση των συστημάτων εμπιστοσύνης αυξήθηκε στο Διαδίκτυο κατά τη διάρκεια των τελευταίων ετών. Μια δήλωση εμπιστοσύνης ορίζεται ως η ρητή άποψη που εκφράζεται από ένα χρήστη σε έναν άλλο χρήστη όσον αφορά την αντιληπτή ποιότητα ορισμένων χαρακτηριστικών αυτού του χρήστη [24]. Η έννοια της εμπιστοσύνης χρησιμοποιείται ευρέως, για παράδειγμα, σε μηχανές αναζήτησης όπως η Google, η οποία χρησιμοποιεί την παγκόσμια μέτρηση εμπιστοσύνης [25], όπως η Pagerank [26] και, επίσης, στο ηλεκτρονικό εμπόριο (eBay) όπου οι χρήστες εκφράζουν το επίπεδο ικανοποίησής τους από την αγορά ενός προϊόντος .

Οι Δηλώσεις εμπιστοσύνης που εκφράζονται από κάθε χρήστη μπορούν να αθροιστούν για να παραχθεί μια κοινότητα ή μια γειτονιά [27], όπως φαίνεται, για παράδειγμα, σε κοινωνικά δίκτυα. Στις αξιολογήσεις όπου οι χρήστες είναι λίγοι, σε πολλές περιπτώσεις, η εύρεση παρόμοιων χρηστών γίνεται αδύνατη. Με την εφαρμογή ενός trust που βασίζεται σε ευρετικά δίκτυα εμπιστοσύνης, μας δίνει τη δυνατότητα να βρούμε τους χρήστες που είναι πιο αξιόπιστοι για την ενεργό χρήστη. Μπορούμε να συνδυάσουμε τις πληροφορίες που παρέχονται από τα δίκτυα εμπιστοσύνης με CF μήτρες για να σχηματίσουμε ένα έμπιστο σύστημα συστάσεων [24], το οποίο απευθύνεται σε προβλήματα όπως μια αραιή αναπαράσταση δεδομένων, εκκίνηση με ψυχρούς χρήστες και πλαστές ταυτότητες, πιο αποτελεσματικά. Παρά τα οφέλη, τα συστήματα εμπιστοσύνης έχουν επίσης προβλήματα με αμφιλεγόμενους χρήστες, που μπορεί να ελαχιστοποιηθούν με τη χρήση τοπικών μετρήσεων εμπιστοσύνης [25].

Η παρούσα εργασία παρουσιάζει συνοπτική state-of-the-art των υφιστάμενων μεθόδων CF και κατάταξη των πιο σχετικών PPCF. Επίσης παρουσιάζει τα τρέχοντα μειονεκτήματα των μεθόδων αυτών και παρουσιάζει νέες βελτιώσεις.

## **2.1 Η Απορρήτου Διατήρηση του Συνεργατικού Φιλτραρίσματος**

Όπως αναφέρθηκε προηγουμένως, η διαδεδομένη χρήση του CF στο Διαδίκτυο προσφέρει μεγάλες ευκαιρίες και οφέλη τόσο για επιχειρήσεις όσο και για τους χρήστες, αλλά υπάρχει ένα σημαντικό μειονέκτημα: η έλλειψη της ιδιωτικής ζωής των χρηστών. Η σημασία της προστασίας της ιδιωτικής ζωής σε συστήματα CF επιτείνεται από τον ρυθμό αύξησης με τον οποίο πληροφορίες του κάθε χρήστη συλλέγονται και αποθηκεύονται. Η απρόσεκτη διαχείριση των προσωπικών πληροφοριών, εκτός από το να είναι παράνομη, θα μπορούσε να οδηγήσει σε σοβαρές συνέπειες για τους χρήστες, των οποίων οι πληροφορίες αποθηκεύονται, όπως καθώς και για τις επιχειρήσεις.[72]

Ένα από τα κύρια προβλήματα στην CF είναι ότι οι πελάτες οι οποίοι πιστεύουν ότι οι προτιμήσεις / προφίλ τους μπορεί να εκτεθεί, τους κάνουν να μην θέλουν να δώσουν την εκτίμησή τους για ένα συγκεκριμένο θέμα ή, αν δοθεί, δεν το κάνουν σωστά ή το κάνουν παραμορφωμένα[28]. Αυτή η συμπεριφορά των χρηστών, που προέρχεται από την έλλειψη σεβασμού του αισθήματος της ιδιωτικής ζωής, συνεπάγεται μείωση τόσο στον αριθμό των αξιολογήσεων καθώς και την ποιότητα τους. Ένα άλλο μειονέκτημα είναι ότι οι εταιρείες μπορούν να αποκτήσουν τα δεδομένα των προτιμήσεων των πολλών χρηστών σε μια δεδομένη αγορά, έτσι έχουν ένα μεγάλο πλεονέκτημα έναντι των νέων ανταγωνιστών, εφόσον αποφασίσουν να επεκταθούν σε άλλες αγορές. Ένα άλλο μειονέκτημα που αφορά το χρήστη είναι ότι υπάρχουν μεγάλα μονοπώλια στο Διαδίκτυο (Google, Amazon) έτσι τα δεδομένα τους μπορούν να μεταφερθούν μεταξύ των διαφόρων φορέων, τα οποία διαχειρίζονται από μεγάλες εταιρείες, εν αγνοία του χρήστη ή χωρίς τη συγκατάθεση του.

Είναι ενδιαφέρον ότι, αν και η μεθόδους CF προκαλεί συσκότιση (ασάφεια) της προστασίας της ιδιωτικής ζωής και / ή να αποκρύπτει τις πληροφορίες του προφίλ του χρήστη, όμως η δημιουργία ομάδων παρόμοιων χρηστών, το οποίο είναι ένα πολύ κοινό γεγονός στο δίκτυο, μπορεί να γίνει ένα δίκικο μαχαίρι. Πρώτον, οι χρήστες μπορούν εύκολα να βρουν αξιόπιστες συστάσεις σχετικά με αντικείμενα από κοινότητες σε ένα συγκεκριμένο πλαίσιο. Δεύτερον, μπορεί να υποστεί ένα ομοιομορφικό πρόβλημα στο δίκτυο, έτσι ώστε οι συστάσεις εκτός του πλαισίου της κοινότητας θα δώσουν αποτελέσματα με νόημα, ακριβώς λόγω της ομοιογένειας της ομάδας. Προκειμένου να επιλύσει τα ζητήματα προστασίας της ιδιωτικής ζωής που θέτει η συστηματική συλλογή ιδιωτικών πληροφοριών, η οποία είναι

απαραίτητη για τη σωστή χρήση του CF, η παρούσα εργασία εστιάζει στη Διατήρηση των Προσωπικών Δεδομένων με Συνεργατική Μέθοδο φιλτραρίσματος (PPCF).

Σε μια δυναμική αγορά όπως το Internet, οι επιχειρήσεις ενδιαφέρονται να συνεργαστούν για να επιτύχουν καλύτερες συστάσεις για τους πελάτες τους. Λόγω της ιδιωτικής ζωής και τις ανησυχίες των επιχειρήσεων, τα δεδομένα δεν θα πρέπει να γνωστοποιούνται μεταξύ των εταιρειών. Σε αυτό το πλαίσιο, τα δεδομένα μπορεί να διαμοιράζονται μεταξύ των διαφόρων μερών με διαφορετικά τρόπους:

**Κάθετη στεγανοποίηση (VP)**, στην οποία οι εταιρείες κατέχουν συνεχές σύνολα στοιχείων, αλλά για τους ίδιους χρήστες. Αυτή η κατάσταση μπορεί να βρεθεί από τρίτους στην ίδια εταιρεία, αλλά η VP είναι συνήθως πιο κατάλληλη για τη λήψη πληροφοριών που αφορά πρόσωπα που διέρχονται από μεγάλες ποσότητες πληροφοριών από διαφορετικό είδος των βάσεων δεδομένων.

**Οριζόντια στεγανοποίηση (HP)**, στην οποία διαφορετικά μέρη κατέχουν ξένα σύνολα των χρηστών με τις απόψεις των ίδιων ειδών. Για τις Παγκόσμιες κοινότητες ή τις επιχειρήσεις ηλεκτρονικού εμπορίου είναι κατάλληλο αυτό το είδος του μοντέλου ομαδοποίησης των δεδομένων. Για παράδειγμα, πολλές εταιρείες στον ίδιο τομέα της αγοράς μπορούν να συνεργάζονται να επιτύχουν καλύτερες συστάσεις και προβλέψεις για τους χρήστες, και να αυξηθούν τα οφέλη χωρίς απώλειες.

**Αυθαίρετη στεγανοποίηση (AP)** [30], κατά την οποία δεν υπάρχει πρότυπο του πώς τα δεδομένα διανέμονται. Εάν ολόκληρη η σειρά ορίζεται από ένα  $m \times n$  πίνακα χρήστη-αντικείμενο, το τμήμα A κατέχει ένα υποσύνολο των χρηστών  $m_a = \leq m$ , ενώ ένα άλλο μέρος B κατέχει το υπόλοιπο  $m_b = m - m_a$  και το ίδιο εφαρμόζεται για τα στοιχεία. Παρατηρήστε ότι η VP και HP είναι ειδικές περιπτώσεις του AP. Αυτό είναι ένα πιο ρεαλιστικό σύστημα στο ηλεκτρονικό εμπόριο στο οποίο οι εταιρείες προσφέρουν διάφορα προϊόντα μέσα στο ίδιο πλαίσιο και οι χρήστες μπορεί να ενδιαφέρονται για ένα απροσδιόριστο αριθμό διαφορετικών εταιρειών.

Όπως παρουσιάζεται στην βιβλιογραφία, υπάρχουν διάφοροι τρόποι για να προστατευτεί η προστασία της ιδιωτικής ζωής στις βάσεις δεδομένων. Για την ανωνυμία των προφίλ σε

μεγάλες βάσεις δεδομένων μπορούμε να χρησιμοποιήσουμε μεθόδους για την παροχή της  $k$  ανωνυμίας [31]. Ενδιαφέρων παρουσιάζουν οι ιδιότητες, του δημόσιου κλειδιού για ομομορφική κρυπτογράφηση [32], [33], οι ασφαλείς υπολογισμοί πολύ-τμημάτων [34] και της κρυπτογράφησης πρωτόκολλων που χρησιμοποιούνται στο Internet. Άλλες μέθοδοι προσθέτουν θόρυβο ώστε τα δεδομένα να νοθεύσουν τις αξίες τους, κατά τρόπο που επηρεάζει το δυνατόν λιγότερο τα στατιστικά χαρακτηριστικά της μήτρας, όπως οι μέσες αξιολογήσεις των χρηστών για τα αντικείμενα.

Στις επόμενες υποενότητες, προτείνουμε μια ταξινόμηση του PPCF σε κεντρικές μεθόδους, σε αποκεντρωμένες μεθόδους και σε υβριδικές μεθόδους, και θα συνοψίσουμε τα πιο σημαντικά από αυτά σε κάθε κατηγορία.

### **2.1.1 Η Κεντρική Μέθοδος PPCF**

Μια μέθοδος PPCF γίνεται κεντρική εάν χρησιμοποιεί ένα τρίτο μέρος για να κάνει ενδιάμεσους υπολογισμούς μεταξύ των χρηστών ή οντοτήτων. Μία μέθοδος θεωρείται επίσης συγκεντρωτική αν έχουν αποθηκευτεί οι βαθμολογίες σε έναν server όπου οι συστάσεις και προβλέψεις υπολογίζονται εκεί. Περιπτώσεις στις οποίες μοιράστηκαν τα δεδομένα δεν θεωρείται κεντρική μέθοδος, καθώς τα δεδομένα κατανέμονται σε διαφορετικά μέρη. Συνήθως, η κεντρική μέθοδος PPCF προσφέρει υψηλότερη απόδοση από τις αποκεντρωμένες, διότι η ομοιότητα και πρόβλεψη των υπολογισμών οδηγεί στην αποφυγή των πολλών επικοινωνιακών εξόδων.[72]

Αρκετές προτάσεις της κεντρικής μεθόδου PPCF μπορούν να βρεθούν στη βιβλιογραφία. Ένα σαφές παράδειγμα είναι η μέθοδος που εισάγεται από το Polat et al. στο [35], το οποίο δείχνει μια τεχνική για να επιτευχθεί μια καλή ισορροπία μεταξύ της προστασίας της ιδιωτικής ζωής και της ποιότητας των συστάσεων. Λόγω των ανησυχιών παραβίασης της ιδιωτικής ζωής, ο κεντρικός server δεν θα πρέπει να αποθηκεύσει τα δεδομένα πραγματικών χρηστών. Για να αποφύγετε την αποκάλυψη του προφίλ των πραγματικών χρηστών, οι χρήστες στρεβλώνουν τα δεδομένα τους με την προσθήκη τυχαίων διανυσμάτων, που δημιουργούνται μετά από μια κατανομή Gauss, πριν από την αποστολή τους στον διακομιστή. Για να θολώσουν τα δεδομένα, ο διακομιστής αποφασίζει ένα εύρος  $[-x, x]$ , γνωστό από τους χρήστες, για να περικόψει τις τυχαίες τιμές που παράγονται από την κατανομή Gauss. Αργότερα, κάθε  $u_i$  υπολογίζει το z-scores  $z_{ij}$  από τα στοιχεία τα οποία ο  $u_i$

έχει βαθμολογήσει. Τέλος, κάθε  $u_i$  δημιουργεί  $n_i$  ομοιόμορφους τυχαίους αριθμούς  $r_{ij}$  στην περιοχή  $[-x, x]$ , όπου  $n_i$  είναι ο αριθμός της βαθμολογίας του προϊόντος από το χρήστη. Μετά από αυτό, κάθε  $u_i$  συγκεντρώνει τους τυχαίους αριθμούς στις αξιολογήσεις z-score και παράγει τα συγκεκαλυμμένο z-score  $z_{ij} = z_{ij} + r_{ij}$ . Στη συνέχεια, οι χρήστες στέλνουν το συγκεκαλυμμένο z-score στο διακομιστή. Μετά την απόκτηση των συγκεκαλυμμένων z-score  $z_{ij}$  των διαφόρων χρηστών, ο διακομιστής είναι σε θέση να στείλει τις συγκεντρωτικές πληροφορίες στο δραστικό χρήστη ο οποίος πρέπει να υπολογίσει τις προβλέψεις τοπικά.

Υπάρχει τρόπος να βελτιώσουμε την ταχύτητα υπολογισμού και την αποδοτικότητα όπως φαίνεται στο [35], Polat et al. όπου προτείνεται μία μέθοδος στην [36], η οποία χρησιμοποιεί το SVD για τη μείωση της διάστασης του πρωτοτύπου της μήτρας. Ο στόχος της είναι να κάνει προβλέψεις χρησιμοποιώντας μια νέα μήτρα με βελτιωμένα χαρακτηριστικά, όπως μειωμένη διασπορά των δεδομένων. Η μέθοδος λειτουργεί ακριβώς όπως εκείνη που είχε προταθεί παλαιότερα στο [35], αλλά στην περίπτωση αυτή ο διακομιστής υπολογίζει εκ των προτέρων το SVD της μήτρας, που δημιουργήθηκε με τα ασαφή (συσκοτισμένα) z-score που αποστέλλονται από τους χρήστες, για να υπολογίσει μετά τις προβλέψεις.

Μετά την γραμμή της συσκοτίσης των δεδομένων, θα βρούμε τη μέθοδο NeNDS που προτείνει το Parameswaran et al. [38]. Η πρόταση αυτή χρησιμοποιεί ένα διακομιστή CF όπου διαφορετικές πηγές συνδυάζουν τα δεδομένα τους για να αποκτήσουν επαρκείς πληροφορίες και μπορεί να ασχοληθεί καλύτερα με μια αραιή αναπαράσταση δεδομένων. Το πρόγραμμα υποθέτει ότι οι φορείς έχουν τρεις τύπους βάσεων δεδομένων: Ο χρήστης-info, σημείο-info και βαθμός-info. Σε αυτήν την περίπτωση, οι βάσεις δεδομένων είναι ασαφείς και αποστέλλονται στον κεντρικό server, ο οποίος εκτελεί μια συνάθροιση δεδομένων και επιστρέφει τις βάσεις δεδομένων με νέες τιμές για κάθε πηγή. Για να διατηρηθούν οι ιδιότητες του κάθε περιεχομένου της βάσης δεδομένων, κάθε πεδίο αντιμετωπίζεται ξεχωριστά και παραλλαγές παρόμοιων στοιχείων που εκτελούνται για να συσκοτίσει τα δεδομένα. Το όφελος είναι μια ασαφή βάση δεδομένων χωρίς να επηρεάζεται η αξία ολόκληρων υποομάδων. Δεδομένου ότι η προαναφερόμενη συσκοτίση θα μπορούσε να είναι αδύναμη όσο αφορά τους όρους της ιδιωτικής ζωής, αφού η ασαφή βάση δεδομένων έχει ληφθεί, η συγγραφείς χρησιμοποιούν γεωμετρικούς μετασχηματισμούς (κλιμάκωση, περιστροφή και μετάφραση), επειδή διατηρούν τόσο τις υποομάδες της βάση δεδομένων όσο και τις αποστάσεις μεταξύ των στοιχείων.

Στην περίπτωση του [39], Shorki et al. προτείνει να μοιραστούν οι βαθμολογίες μεταξύ παρόμοιων χρηστών να συσκοτιστούν τα πραγματικά δεδομένα, έτσι ώστε τα προφίλ των χρηστών να γίνουν εν μέρει μικτά. Με τον τρόπο αυτό, τροποποιημένα προφίλ μπορούν να αποσταλούν σε έναν κεντρικό server, διατηρώντας την ιδιωτικότητα των χρηστών που συμμετέχουν.

### **2.1.2 Η Αποκεντρωμένη Μέθοδος PPCF**

Όλες αυτές οι μέθοδοι οι οποίες χρησιμοποιούν τα μέλη της ένα καταναμημένο δίκτυο, προκειμένου να εκτελέσουν ενδιάμεσους υπολογισμούς και προβλέψεις για τα στοιχεία βαθμολογίας, μπορούν να χαρακτηριστούν ως αποκεντρωμένες μέθοδοι PPCF. Τα μέλη αυτά θα πρέπει να θεωρούνται στις περισσότερες περιπτώσεις ως χρήστες. Η χρήση αποκεντρωμένων συστημάτων γενικά εξασφαλίζει ότι οι πληροφορίες που εκτίθενται είναι πολύ λιγότερες από ό, τι στην περίπτωση των κεντρικών συστημάτων, αλλά αυτό συνεπάγεται τη χρήση δαπανηρών πρωτόκολλων και πιο σύνθετους υπολογισμούς. Συνήθως, στην αποκεντρωμένη Μεθόδους PPCF, οι χρήστες αποθηκεύουν τις δικές τους αξιολογήσεις. Είναι ευρέως γνωστό ότι οι PPCF, οι οποίες περιλαμβάνουν διάφορα μέρη μοιράζονται τα δεδομένα τους για να εκτελέσουν την CF με περισσότερες παραπομπές, έτσι θεωρούνται επίσης αποκεντρωμένες μέθοδοι.[72]

Αρκετές προσεγγίσεις με τμηματοποιημένο καλάθι αγοράς μιας βάσης δεδομένων, έχουν προταθεί στην βιβλιογραφία. Αυτά τα είδη των βάσεων δεδομένων είναι κατάλληλα για να κάνει top-N συστάσεις με υψηλή ακρίβεια και χαμηλό υπολογιστικό κόστος λόγω της δυαδικής βαθμολογίας που έχει ως περιεχόμενο. Ένα παράδειγμα έδειξε από Polat et al. στο [40], στο οποίο μια μέθοδος για να κάνει συστάσεις κάθετα και οριζόντια προτείνεται να τμηματοποιήσει το καλάθι αγοράς των δεδομένων της. Σε αυτά τα είδη των δεδομένων, οι αξιολογήσεις αποθηκεύονται ως  $r_{ij} = 0$  εάν το στοιχείο  $ij$  δεν έχει αγοραστεί για το χρήστη  $u_i$  διαφορετικά  $r_{ij} = 1$ . Οι συγγραφείς υποστηρίζουν ότι η αγορά ενός προϊόντος δεν ευχαριστεί αναγκαστικά το χρήστη, όπως ο ίδιος μπορεί να είναι απογοητευμένος γι 'αυτό, αλλά εξακολουθεί να έχει υψηλό δείκτη στην προτίμηση του. Τέλος πάντων, είναι καλύτερα να δείξει τους χρήστες τις προτιμήσεις του ως αν τους άρεσε το στοιχείο τότε είναι (1) ή αν όχι τότε είναι (0). Υποθέτοντας ότι τα δεδομένα καταναμήθηκαν μεταξύ των μερών A και B, η



HPD και η VPD προτείνουν τη χρήση πρωτοκόλλων που σέβονται την ιδιωτική ζωή και τυχαίες παραλλαγές για να αποφευχθεί η αποκάλυψη των ιδιωτικών πληροφοριών μεταξύ των μερών. Οι συγγραφείς σημειώνουν ότι ένας χρήστης είναι ένας δυνητικός καταναλωτής των προϊόντων που έχουν γείτονες οι οποίοι είναι αγοραστές, αλλά επίσης μπορεί να ενδιαφέρονται για τα στοιχεία τα οποία δεν έχουν αγοραστεί από ανόμιους χρήστες. Ως εκ τούτου, πριν από την επιλογή των στοιχείων top-N ώστε να κάνουμε την πρόταση στον ενεργό χρήστη  $u_a$ , οι τιμές των πιο ανόμιων χρηστών του  $u_a$  να έχουν αλλάξει από το 0 έως το 1. Τέλος, τα N στοιχεία που περιέχουν περισσότερα 1s μπορούν να συσταθούν. Οι συγγραφείς επίσης παρουσιάζουν στο [40] ότι το συστήματά τους μπορεί να επεκταθεί και σε συστήματα πολλαπλών τμημάτων.

Το Polat et al. [41] πρότεινε τη χρήση μιας naive Bayesian ταξινόμησης (NBC) για την εκτέλεση των συστάσεων PPCF στην APD. Η προτεινόμενη μέθοδος έχει διαφοροποιημένα offline και online στάδια για να αυξηθεί η συνολική απόδοση. Χρησιμοποιούν ασφαλή πρωτόκολλα στο στάδιο offline για την κατασκευή του μοντέλου πρόβλεψης NBC. Αυτό το μοντέλο κατασκευής πρόβλεψης έχει δύο σαφή βήματα. Το πρώτο περιλαμβάνει την εκτίμηση των πιθανοτήτων, η οποία διεξάγεται με τη χρήση ασφαλών πρωτοκόλλων που βασίζονται στο Paillier [32] σύστημα ομομορφικής κρυπτογράφησης του δημοσίου κλειδιού. Το δεύτερο, το οποίο περιλαμβάνει τον εκ των προτέρων υπολογισμό, χρησιμοποιεί τυχαίοποιημένη τεχνική διαταραχής για τον υπολογισμό της πιθανότητας της ύπαρξης 1 ή 0 σε ένα στοιχείο στόχο  $q$ . Το online στάδιο συνεπάγεται την εκτίμηση της σύστασης, η οποία πραγματοποιείται από τους απευθείας σε σύνδεση με πρωτόκολλο εκτίμησης συστάσεων.

Μια παρόμοια προσέγγιση NBC στα δυαδικά δεδομένα έχει επίσης προτείνει Kaleli et al. [42], αλλά χρησιμοποιώντας ένα peer-to-peer (P2P) δίκτυο. Οι συγγραφείς υποστηρίζουν ότι η κεντρική αποθήκευση δημιουργεί πολλούς κινδύνους για τους χρήστες, διότι μια ενιαία οντότητα ελέγχει τα δεδομένα των χρηστών. Για να αποφύγουν το ζήτημα της ιδιωτικής ζωής, προτείνουν ένα δίκτυο P2P, το οποίο οι χρήστες (που είναι ενεργεί ως συνάδελφοι) χρησιμοποιούν για να επικοινωνούν και να ανταλλάσσουν δεδομένα μεταξύ τους, προκειμένου να προβεί σε συστάσεις.

Επιπλέον, για την NBC όσο αφορά την ιδιωτική ζωή, έχει διεξαχθεί αξιολόγηση σε αριθμητικά πλαίσια στο [43]. Είναι ενδιαφέρον να τονίσουμε ότι οι δυαδικές μέθοδοι είναι εφαρμόσιμες σε οποιαδήποτε βάση δεδομένων αξιολόγησης CF, αν οι εκτιμήσεις

γενικευτούν σε θετικές ή αρνητικές αξιολογήσεις, έχει όμως σαν αποτέλεσμα την απώλεια πληροφοριών.

Υπάρχουν, επίσης, προσεγγίσεις με συστήματα εμπιστοσύνης που διατηρούν την ιδιωτική ζωή όπως στο [48]. Εδώ, ο Dokoochaki και οι συνεργάτες του εκφράζουν το πρόβλημα της αποκάλυψης των δεδομένων του δικτύου εμπιστοσύνης, διότι αυτές οι πληροφορίες μπορεί να αποκαλύψουν τις συμπεριφορές των χρηστών. Για να αποφευχθεί αυτό, θα παρουσιαστεί ένα πρόγραμμα το οποίο θεωρείται αποκεντρωμένο επειδή οι υπολογισμοί γίνονται σε τοπικό επίπεδο εμπιστοσύνης από κάθε χρήστη.

Κατ' αρχάς, κάθε χρήστης υπολογίζει το z-score του και μεταμορφώνει τα δεδομένα του ακολουθώντας ένα πρωτόκολλο απόκρυψης. Μετά από αυτό, κάθε χρήστης υπολογίζει με ένα καταναμημένο τρόπο την ιδιωτική εκτίμηση εμπιστοσύνης μεταξύ άλλων χρηστών. Τέλος, για να γίνει η παραγωγή συστάσεων, οι χρήστες παίρνουν τις πληροφορίες που παρέχονται από την ιδιωτική εκτίμηση εμπιστοσύνης και τα καμουφλαρισμένα προφίλ. Δεδομένου ότι τα z-scores χρησιμοποιούνται, οι χρήστες θα πρέπει να απομαλοποιήσουν τα αποτελέσματα της πρόβλεψης για να πάρουν τις πραγματικές αξίες.

Στα συστήματα εμπιστοσύνης, επίσης προτείνονται η VPD για την διατήρηση του απόρρητου της ιδιωτικής ζωής. Το σύστημα που προτείνεται στο [49] αποτελείται τόσο σε offline και όσο και σε online στάδιο. Πρώτα απ' όλα, μια offline προσέγγιση χρησιμοποιείται για τον υπολογισμό επιμέρους εμπιστοσύνης μεταξύ των χρηστών με προστασία της ιδιωτικής ζωής. Αφού οι τιμές αυτές υπολογίζονται, ένας αλγόριθμος με βάση την απόσταση της ιδιωτικής διαλογής (DPSA) χρησιμοποιείται από τα συμβαλλόμενα μέρη για τον προσδιορισμό της γειτονιάς του κάθε χρήστη, που σχηματίζεται από  $k$  πιο έμπιστους χρήστες. Κατά τη διάρκεια της διαδικασίας offline, τα συμβαλλόμενα μέρη χρησιμοποιούν το πρωτόκολλο DPSA  $n$  φορές για να καθορίσουν τις γειτονιές για όλους τους χρήστες. Μόλις υπολογίζονται αυτές οι γειτονιές, χρησιμοποιείται ένα συστημένο ιδιωτικό πρωτόκολλο, βάσει του ομομορφικού συστήματος κρυπτογράφησης Paillier, για την παροχή on-line προβλέψεων.

Στο [50], οι Hsieh et al. χρησιμοποιούν τη γνωστή Pearson Correlation μέθοδο για τον υπολογισμό των ομοιοτήτων μεταξύ των χρηστών στο HPD. Για τον υπολογισμό της συσχέτισης του Pearson σε HPD, η προτίμηση του χρήστη  $U_i$  στο στοιχείο  $i_j$  αποκαλύπτεται



επειδή η μέση τιμή του  $U_i$  πρέπει να κατανέμεται μεταξύ των φορέων. Για να αποφευχθεί αυτή η αποκάλυψη, οι συγγραφείς προτείνουν ένα ασφαλές πρωτόκολλο για τον υπολογισμό της συσχέτισης του Pearson μεταξύ των χρηστών χρησιμοποιώντας την ομομορφική κρυπτογράφηση ElGamal [33]. Μόλις καθοριστούν οι γειτονιές, ο υπολογισμός των συστάσεων είναι πολύ απλός. Σε μια μεταγενέστερη προσέγγιση, καθώς ο αριθμητής της συσχέτισης του Pearson είναι ένα εσωτερικό γινόμενο, οι Zhan et al. [51] χρησιμοποιούν ένα ασφαλές πρωτόκολλο που υπολογίζει το εσωτερικό γινόμενο μεταξύ δύο διανυσμάτων  $X_\alpha$  και  $X_\beta$ . Οι συγγραφείς δείχνουν ότι το πρωτόκολλο βαθμιδωτού ασφαλούς προϊόντος είναι περισσότερο αποδοτικό από ό,τι η κρυπτογράφηση ElGamal [33], επειδή η ύπαρξη της εμπιστοσύνης, που είναι απαραίτητη στο πρωτόκολλο για την παραγωγή τυχαίων αριθμών, μπορεί να παρέχει αυτούς τους αριθμούς πριν την έναρξη του πρωτοκόλλου. Αυτό μπορεί να γίνει γνωρίζοντας πόσα βαθμιδωτά προϊόντα πρέπει να υπολογιστούν μεταξύ των φορέων που θέλουν να εκτελέσουν PPCF. Με τον τρόπο αυτό, αυτοί οι φορείς μπορούν να μειώσουν δραστικά την ανταλλαγή πληροφοριών σε ένα βήμα και όχι σε  $n$ .

Στο [52], ο J. Canny προτείνει το σύστημα peer-to-peer με τη χρήση κρυπτογραφημένων δεδομένων. Σε αυτό το σύστημα, οι χρήστες έχουν το δικό τους φορέα αξιολόγησης και μόνο τα συγκεντρωτικά στοιχεία εκτίθενται. Οι χρήστες μπορούν να αποφασίσουν αν πρέπει να μοιράζονται τα στοιχεία τους για την μέθοδο CF ή όχι, έτσι ώστε αυτό το σύστημα να μπορεί να ενθαρρύνει το σχηματισμό των ομάδων από τη χρήση των ομοιοτήτων για να υπολογίσει τις συστάσεις. Για να αποκτήσει το σύνολο των αξιολογήσεων με το οποίο θα προβεί σε συστάσεις, ο συγγραφέας χρησιμοποιεί μια επαναληπτική μέθοδο SVD, η οποία υπολογίζεται χρησιμοποιώντας την μέθοδο των συζυγών διευθύνσεων (6) [53]. Οι επαναληπτικές προσθήκες, που είναι αναγκαίες για να υπολογιστεί ο πίνακας με τις μειωμένες διαστάσεις, γίνονται μέσω ομομορφισμών χρησιμοποιώντας το ElGamal [33] κρυπτογραφικό σύστημα του δημόσιου κλειδιού. Παρά την επικοινωνία και τις υπολογιστικές δαπάνες του πρωτοκόλλου, ο συγγραφέας δείχνει ότι η μέθοδος έχει καλή επίδοση σε όρους υπολογιστικής ταχύτητας και ποιότητας των συστάσεων. Αργότερα, στο [54], ο Canny προτείνει τη χρήση της παραγοντικής ανάλυσης Μέγιστης Προσδοκίας [55], η οποία ασχολείται με ελάχιστα στοιχεία για να υπολογίσει τον πίνακα μειωμένων διαστάσεων. Στην περίπτωση αυτή, ο συγγραφέας δείχνει μια αύξηση στην ταχύτητα και καλή ποιότητα των συστάσεων σε σχέση με το σύστημα peer-to-peer με τη χρήση κρυπτογραφημένων δεδομένων.

Οι Berkvosky et al. [56] προτείνουν επίσης τη χρήση του peer-to-peer αποκεντρωμένου συστήματος για να εκτελέσουν την PPCF. Σε αυτήν την περίπτωση, κάθε χρήστης έχει το δικό του διάνυσμα βαθμολογίας και δεν υπάρχει ανάγκη για κεντρικό server. Σε περίπτωση που ένας ενεργός χρήστης  $U_a$  στέλνει ένα αίτημα, μεταξύ άλλων, κατά τη διάρκεια μιας αξιολόγησης ενός στοιχείου  $i_a$ , κάθε χρήστης αποφασίζει είτε να παρέχει στοιχεία για την CF είτε όχι. Εάν ο  $U_b$  χρήστης αποφασίσει να συνεργαστεί, υπολογίζει την ομοιότητα μεταξύ  $U_b$  και  $U_a$  χρησιμοποιώντας το μέτρο ομοιότητας συνημίτονου [6] και στέλνει το αποτέλεσμα στον  $U_a$ , με την αποτίμηση του στοιχείου επί του οποίου έγινε το αίτημα. Μόλις συλλέγονται τα δεδομένα, οι  $k$  πλησιέστεροι γείτονες του  $U_a$  έχουν επιλεγεί για τον καθορισμό της τιμής πρόβλεψης του  $i_a$  για τον  $U_a$ , και ο σταθμισμένος μέσος όρος υπολογίζεται με βάση την εγγύτητά τους. Αυτή η προσέγγιση αυξάνει την προστασία της ιδιωτικής ζωής των χρηστών καθώς και τα προφίλ τους, είναι αποθηκευμένα με ένα καταναμημένο τρόπο και η ομοιότητα των υπολογισμών γίνονται σε τοπικό επίπεδο. Οι χρήστες απλά εκθέτουν μια βαθμολογία εφόσον συνεργάζονται και όχι το πλήρες προφίλ τους. Αν και η έκθεση του προφίλ των χρηστών ελαχιστοποιείται, οι ‘‘ανέντιμοι’’ χρήστες θα μπορούσαν να κάνουν αιτήματα πολλαπλών επιθέσεων για την ανοικοδόμηση του προφίλ των χρηστών, οπότε το μέτρο αυτό δεν είναι αρκετό. Για να αποκρύψουν τις πραγματικές βαθμολογίες του χρήστη, οι αποτιμήσεις γίνονται ασαφείς ακολουθώντας την κατανομή Gauss.

Στο [57], οι Tada et al. προτείνουν μια μέθοδο PPCF με βάση την ομοιότητα μεταξύ των στοιχείων. Υποστηρίζουν ότι οι ομοιότητες μεταξύ των αντικειμένων μπορούν να δημοσιοποιηθούν χωρίς διαρροές της προστασία της ιδιωτικής ζωής, διότι τα προφίλ των χρηστών δεν είναι εκτεθειμένα, ούτε οι ομοιότητες μεταξύ τους. Οι συγγραφείς εκμεταλλεύονται τις ομομορφικές ιδιότητες του κρυπτογραφικού συστήματος του δημοσίου κλειδιού του Paillier και την ταχύτερη διαδικασία αποκρυπτογράφησης (σε σχέση με το άρθρο του ElGamal), για να υπολογίσουν την απλή συνημιτονιακή συσχέτιση μεταξύ των στοιχείων [6] με την ιδιωτική ζωή. Στο [58], οι Kikuchi et al. προτείνουν ένα πρόγραμμα το οποίο είναι παρόμοιο με εκείνο που ήδη δείχτηκε στο [57]. Εδώ, για να μειωθεί η επιβάρυνση των προτεινόμενων πρωτόκολλων, οι συγγραφείς παρουσιάζουν συστάδες των στοιχείων και συστάδες των χρηστών. Και στις δύο περιπτώσεις, τα αποτελέσματα δείχνουν μια μείωση τόσο στα έξοδα επικοινωνίας όσο και υπολογιστικά σε σχέση με τα αποτελέσματα που προέκυψαν με το αρχικό σύνολο δεδομένων.

Οι Basu et al. [59] προτείνουν τη χρήση της σταθμισμένης Slope One PPCF στο cloud computing, η οποία χρησιμοποιεί τα στοιχεία απόκλισης για να κάνει προβλέψεις. Σε αυτό το σύστημα, οι χρήστες αποθηκεύουν τις δικές τους αξιολογήσεις. Για να αποσυνδέσει τους χρήστες από τις βαθμολογίες τους, εφόσον αυτές έχουν υποβληθεί, οι συγγραφείς προτείνουν τη χρήση γνωστών συστημάτων ανωνυμίας. Το προτεινόμενο σύστημα έχει δύο διαφοροποιημένα μέρη. Στο πρώτο μέρος, οι πληθάρημοι και οι αποκλίσεις μεταξύ των στοιχείων υπολογίζονται και αποθηκεύονται σε ένα πίνακα πληθικότητας και σε ένα πίνακα αποκλίσεων, αντίστοιχα. Αν και εφόσον οι χρήστες δεν έχουν σχέση με τις πληροφορίες που παρέχουν και με οποιαδήποτε συμπαγής αξιολόγηση αποστέλλεται, οι πράξεις αυτές δεν αποκαλύπτουν καμία πληροφορία. Μόλις ληφθούν οι παραπάνω πίνακες, οι χρήστες πρέπει να στείλουν το διάλυμα της βαθμολογίας τους στο Cloud, προκειμένου να γίνουν προβλέψεις. Για να αποτραπεί η αποκάλυψη των δεδομένων, τα προφίλ τους αποστέλλονται κρυπτογραφημένα χρησιμοποιώντας μια τροποποιημένη έκδοση του κρυπτογραφικού συστήματος του δημόσιου κλειδιού. Αυτή η τροποποίηση είναι απαραίτητη για να λειτουργούν με αρνητικούς αριθμούς διότι οι αποκλίσεις, οι οποίες συμμετέχουν στη διαδικασία πρόβλεψης δεν μπορεί να είναι θετικές. Το Cloud λαμβάνει το προφίλ του χρήστη κρυπτογραφημένο με το δημόσιο κλειδί και παράγει συμπεριφορές με ομοιορφικές λειτουργίες με τα δεδομένα που έχει και τις πληροφορίες για το προφίλ του χρήστη για τη λήψη των κρυπτογραφημένων συστάσεων. Στη συνέχεια, στέλνει τις κρυπτογραφημένες συστάσεις πίσω στο χρήστη, ο οποίος θα τις αποκρυπτογραφήσει με το προσωπικό του κλειδί. Το σύστημα που προτείνεται έχει ένα σημαντικό ζήτημα προστασίας προσωπικών δεδομένων, το οποίο είναι ότι ο χρήστης αποκαλύπτει τον αριθμό των στοιχείων που έχει βαθμολογήσει στη φάση πρόβλεψης. Ωστόσο, αυτό το πρόβλημα μπορεί να επιλυθεί εάν οι χρήστες διενεργούν τις προβλέψεις σε τοπικό επίπεδο, παραλαμβάνοντας τις απαιτούμενες πληροφορίες από το Cloud, ή επίσης αν κάνουν αιτήματα για περιττά αντικείμενα.

Από μια άλλη άποψη, μεγάλο μέρος της έρευνας στον τομέα της προστασίας προσωπικών δεδομένων και κρυπτογραφικών συστημάτων, προτείνει τη χρήση των mixservers και των αξιόπιστων φορέων για την εκτέλεση ορισμένων υπολογισμών μεταξύ των διαφόρων τμημάτων, χωρίς να εκθέτουν τόσο τα ενδιάμεσα δεδομένα όσο και τις μαθηματικές εργασίες που εκτελούνται. Η PPCF μπορεί να πραγματοποιηθεί με τη χρήση αξιόπιστων φορέων, όπως προτείνεται στο Sharemind [60]. Το Sharemind είναι ένα εικονικό περιβάλλον για κατανεμημένους υπολογισμούς, το οποίο χρησιμοποιεί μια σειρά από φορείς

για την εκτέλεση υπολογισμών διαφύλαξης της ιδιωτικής ζωής από ένα σύνολο δεδομένων. Τα δεδομένα διανέμονται χρησιμοποιώντας ένα πρόσθετο μυστικό σύστημα καταμερισμού επί του δακτυλίου  $Z_2$  [32], το οποίο διαφέρει από το κλασικό σύστημα Shamir [61]. Το πρωτόκολλο προστασίας προσωπικών δεδομένων μπορεί να επιτευχθεί μόνο αν η αναλογία των κακόβουλων χρηστών στο σύνολο των χρηστών είναι λιγότερο από  $1/3$  και οι χρήστες δεν συνωμοτούν. Αυτή η μέθοδος μπορεί να εκτελέσει επιπρόσθετες, πολλαπλασιαστικές και συγκριτικές εργασίες με ασφαλή τρόπο, αρκετές ενέργειες για τους περισσότερους αλγόριθμους στατιστικής ανάλυσης και εξόρυξης δεδομένων. Παρά την παρεχόμενη προστασία της ιδιωτικής ζωής, η μέθοδος αυτή είναι αργή σε πραγματικό χρόνο τόσο λόγω του μυστικού συστήματος καταμερισμού όσο και του εγγενούς πρωτοκόλλου Sharemind.

### **2.1.3 Η Υβριδική Μέθοδος PPCF**

#### **2.1.3.1 Υβριδικές τεχνικές συνεργατικού φιλτραρίσματος**

Τα Υβριδικά συστήματα συνδυάζονται με άλλες τεχνικές σύστασης (συνήθως με τα συστήματα που βασίζονται στο περιεχόμενο) για προβλέψεις ή συστάσεις. [6]

Με βάση το περιεχόμενο τα εισηγητικά συστήματα κάνουν συστάσεις αναλύοντας το περιεχόμενο του κειμένου πληροφοριών, όπως έγγραφα, διευθύνσεις, τα μηνύματα από τις ειδήσεις, αρχεία καταγραφής ιστού, στοιχεία από περιγραφές, και τα προφίλ για τους χρήστες όπως τα γούστα, τις προτιμήσεις, και τις ανάγκες, και η εύρεση κανονικότητας στο περιεχόμενο [74]. Πολλά στοιχεία συμβάλλουν στη σημασία του περιεχομένου του κειμένου, όπως παρατήρηση των χαρακτηριστικών περιήγησης των λέξεων ή σελίδων (Π.χ. ο όρος συχνότητα και αντίστροφη συχνότητα εγγράφου), και ομοιότητα μεταξύ των στοιχείων που ένας χρήστης άρεσε στο παρελθόν [75]. Ένα σύστημα με βάση το περιεχόμενο συστάσεων στη συνέχεια χρησιμοποιεί ευρετικές μεθόδους ή αλγόριθμοι ταξινόμησης για να διατυπώνει συστάσεις [76]. Οι Τεχνικές με βάση το περιεχόμενο έχουν το πρόβλημα εκκίνησης, τα οποία πρέπει να διαθέτουν επαρκείς πληροφορίες για να οικοδομήσει μια αξιόπιστη ταξινομητή. Επίσης, αυτές περιορίζονται από τα χαρακτηριστικά που συνδέονται ρητά με τα αντικείμενα που προτείνονται (μερικές φορές αυτά τα χαρακτηριστικά είναι δύσκολο να τα εξάγει), ενώ το συνεργατικό φιλτράρισμα μπορεί να προβαίνει σε συστάσεις χωρίς περιγραφικά στοιχεία. Επίσης, οι τεχνικές που βασίζονται στο περιεχόμενο έχουν το πρόβλημα του overspecialization, δηλαδή, μπορούν να συστήσουν μόνο τα στοιχεία που θα

έχουν ιδιαίτερα υψηλό σκορ κατά το προφίλ ενός χρήστη ή την ιστορία του / της αξιολόγησης [82, 77].

Άλλα εισηγητικά συστήματα περιλαμβάνουν δημογραφικά-based εισηγητικά συστήματα, τα οποία χρησιμοποιούν πληροφορίες του προφίλ του χρήστη όπως το φύλο, τον ταχυδρομικό κώδικα, το επάγγελμα, και ούτω καθεξής [78] με βάση τη χρησιμότητα των συστημάτων συστάσεων και συστάσεων που βασίζεται στη γνώση, απαιτούν γνώση σχετικά με πώς ένα συγκεκριμένο αντικείμενο ικανοποιεί τις ανάγκες των χρηστών [79, 80]. Δεν θα συζητήσουμε αυτά τα συστήματα λεπτομερώς σε αυτή την εργασία .

Ελπίζοντας να αποφύγει τους περιορισμούς του κάθε συνιστών συστήματος και τη βελτίωση των επιδόσεων, τα υβριδικά CF συστήματα συστάσεων έχουν ένα συνδυασμό με την προσθήκη χαρακτηριστικών με βάση το περιεχόμενο για τα μοντέλα CF συνδυάζοντας με CF με βάση το περιεχόμενο ή άλλα συστήματα, ή με συνδυασμό διαφορετικών αλγορίθμων ΤΣ [82, 81].

### ***2.1.3.2 Υβριδικά CF και Content- Based χαρακτηριστικά.***

Το περιεχόμενο του ενισχυμένου αλγόριθμου CF χρησιμοποιεί Naïve Bayes ταξινομητή ως περιεχόμενο, τότε συμπληρώνει τις άγνωστες τιμές του πίνακα αξιολόγησης με τις προβλέψεις για να σχηματίσει μια μήτρα βαθμολογίας, στην οποία οι παρατηρούμενες αξιολογήσεις κρατήθηκαν ανέπαφες και οι βαθμολογίες λείπουν αντικαθίστανται από τις προβλέψεις του προγνωστικού περιεχόμενου. Στη συνέχεια, κάνει προβλέψεις στη μήτρα με τις αξιολογήσεις που προέκυψαν χρησιμοποιώντας ένα σταθμισμένο αλγόριθμο συσχέτισης Pearson-based CF, ο οποίος δίνει μεγαλύτερο βάρος για ένα στοιχείο το οποίο περισσότεροι χρήστες βαθμολογούσαν, και δίνει ένα μεγαλύτερο βάρος για τον ενεργό χρήστη [16] (δείτε μια εικόνα στον Πίνακα 5). Το περιεχόμενο του ενισχυμένου αλγόριθμου CF έχει βελτιωμένη απόδοση πρόβλεψης πάνω σε ορισμένες καθαρά με βάση το περιεχόμενο συστάσεις. Ξεπερνά επίσης το πρόβλημα εκκίνησης με ψυχρό (αδιάφορος) χρήστη και αντιμετωπίζει το πρόβλημα της αραιότητας.

Στο Ansari et al. [83] προτείνουν ένα Bayesian μοντέλο προτίμησης που ενσωματώνει στατιστικά διάφορους τύπους πληροφοριών που είναι χρήσιμες για τη διατύπωση συστάσεων, όπως τις προτιμήσεις του χρήστη, ο χρήστης και τα χαρακτηριστικά του στοιχείου. Χρησιμοποιούν Markov αλυσίδα Monte Carlo (MCMC) μέθοδο[84], για δειγματοληψία, το

οποίο περιλαμβάνει δειγματοληψία με βάση την εκτίμηση των παραμέτρων από την πλήρη συνθήκη κατανομής των παραμέτρων. Πέτυχαν καλύτερη απόδοση από το καθαρό συνεργατικό φιλτράρισμα.

Ο συνιστών Fab, που προτάθηκε από τους Balabanović και Shoham [81], διατηρεί προφίλ χρηστών για τις σελίδες που τους ενδιαφέρουν στο διαδίκτυο με τη χρήση τεχνικών που βασίζονται σε περιεχόμενο, και χρησιμοποιεί τεχνικές CF για τον εντοπισμό προφίλ με παρόμοια γούστα. Στη συνέχεια, μπορεί να συστήσει έγγραφα σε προφίλ χρηστών. Στο Sarwar et al. [85] τίθεται σε εφαρμογή μια σειρά από «filterbots» που βασίζεται στη γνώση ως τεχνητοί χρήστες που χρησιμοποιούν ορισμένα κριτήρια. Ένα απλό παράδειγμα filterbot είναι genrebot, η οποία βασίζεται τη γνώμη της αποκλειστικά στο είδος του στοιχείου, για παράδειγμα, ένας "jazzbot" θα δώσει ένα πλήρη σήμα σε ένα CD απλά επειδή είναι στην τζαζ κατηγορία, ενώ θα δίνουν χαμηλή βαθμολογία για οποιοδήποτε άλλο CD σε μία βάση δεδομένων. Στο Mooney και Roy [86] χρησιμοποιούν την πρόβλεψη από το σύστημα CF ως είσοδο σε ένα περιεχόμενο που βασίζεται σε συστάσεις. Το Condiff et al. [77] προτείνει ένα Bayesian μοντέλο που αναμιγνύει τα αποτελέσματα και ενσωματώνει τις αξιολογήσεις του χρήστη, το χρήστη, και τα χαρακτηριστικά του στοιχείου σε ένα ενιαίο ενοποιημένο πλαίσιο. Το CF σύστημα Ripper, προτείνει το Basu et al. [87], χρησιμοποιεί και τις αξιολογήσεις των χρηστών και τα χαρακτηριστικά των περιεχομένων για να παράγουν συστάσεις.



**Πίνακας 2.2** : Το Περιεχόμενο του ενισχυμένου CF και οι παραλλαγές του (α) τα δεδομένα περιεχομένου και αρχικά ελάχιστα στοιχεία για την αξιολόγηση των δεδομένων b) ψευδοαξιολογημένα δεδομένα που γεμίζουν από προγνωστικό περιεχόμενο (γ) προβλέψεις από (σταθμισμένη) Pearson CF για τα ψευδοαξιολογημένα δεδομένα.

(A)

Content information					Rating matrix				
	Age	Sex	Career	zip	I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	I <sub>4</sub>	I <sub>5</sub>
U <sub>1</sub>	32	F	Writer	22904			4		
U <sub>2</sub>	27	M	Student	10022	2		4	3	
U <sub>3</sub>	24	M	Engineer	60402		1			
U <sub>4</sub>	50	F	Other	60804		3	3	3	3
U <sub>5</sub>	28	M	Educator	85251	1				

(B)

Pseudo rating data				
I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	I <sub>4</sub>	I <sub>5</sub>
2	3	4	3	2
2	2	4	3	2
3	1	3	4	3
3	3	3	3	3
1	2	4	1	2

(C)

Pearson-CF prediction				
$I_1$	$I_2$	$I_3$	$I_4$	$I_5$
2	3	4	2	3
3	4	2	2	3
3	3	2	3	3
3	3	3	3	3
1	3	1	2	2

### 2.1.3.3 Υβριδικά CF και άλλα εισηγητικά συστήματα

Ένα σταθμισμένο υβριδικό σύστημα συστάσεων συνδυάζει διαφορετικές τεχνικές σύστασης από τα βάρη τους, τα οποία υπολογίζονται από τα αποτελέσματα του συνόλου των διαθέσιμων τεχνικών που υπάρχουν στο σύστημα [79] σύστασης. Ο συνδυασμός μπορεί να είναι γραμμική, τα βάρη μπορεί να είναι ρυθμιζόμενα [88], και η πλειοψηφία να είναι σταθμισμένη ψηφοφορία [74, 90], ή να είναι σταθμισμένη η μέση ψηφοφορία [89] και να μπορεί να χρησιμοποιηθεί. Για παράδειγμα, το PTango σύστημα [88] δίνει αρχικά CF και με βάση το περιεχόμενο σύστασης ισούται με το βάρος, αλλά σταδιακά προσαρμόζει τη στάθμιση όπως επιβεβαιώνονται οι προβλέψεις σχετικά με τις αξιολογήσεις των χρηστών ή τις ανεπιβεβαίωτες προβλέψεις. Η στρατηγική του συστήματος P-Tango είναι παρόμοια στην ενίσχυση [91].

Μια υβριδική μεταγωγή συστάσεων αλλάζει μεταξύ των τεχνικών συστάσεων που χρησιμοποιούν ορισμένα κριτήρια, όπως τα επίπεδα εμπιστοσύνης για τις τεχνικές σύστασης. Όταν το CF σύστημα δεν μπορεί να κάνει μια σύσταση με επαρκή εμπιστοσύνη, τότε ένα



άλλο σύστημα συστάσεων, όπως ένα σύστημα με βάση το περιεχόμενο επιχειρείται να εφαρμοστεί. Η υβριδική μεταγωγή συστάσεων επίσης, εισάγει την πολυπλοκότητα της παραμετροποίησης για τα κριτήρια μεταγωγής [79].

Άλλα υβριδικά συστήματα σε αυτήν την κατηγορία περιλαμβάνουν μικτά υβριδικά συστήματα συστάσεων [92], αλληλουχικά υβριδικά συστήματα [79] και ούτω καθεξής.

Πολλοί ερευνητές συγκρίνουν την απόδοση του υβριδικού συστήματος με το απλό CF και με βάση το περιεχόμενο και βρέθηκε ότι τα υβριδικά συστήματα μπορούν να κάνουν πιο ακριβείς συστάσεις, ιδιαίτερα για τον νέο χρήστη και νέες καταστάσεις όπου ένας κανονικός αλγόριθμος CF δεν μπορεί να κάνει ικανοποιητικές συστάσεις. Ωστόσο, τα υβριδικά συστήματα στηρίζονται σε εξωτερικές πληροφορίες που δεν είναι συνήθως διαθέσιμες, και γενικά έχουν αυξημένη πολυπλοκότητα στην εφαρμογή τους [74, 79, 93].

#### ***2.1.3.4 Υβριδικά συστήματα σε σχέση με τα απλά CF***

Ο δύο κύριες κατηγορίες CF, που βασίζεται σε μνήμη και στο μοντέλο CF μπορούν να συνδυαστούν για να σχηματίσουν υβριδικές προσεγγίσεις. Οι παραστάσεις σύστασης από αυτούς τους αλγορίθμους είναι γενικά καλύτερες από κάποια με βάση τη μνήμη CF και ένα αλγόριθμο που βασίζεται στο μοντέλο [23, 94].

Η Πιθανοτική με βάση τη μνήμη για το συνεργατικό φιλτράρισμα (PMCF) συνδυάζει τεχνικές βάσει τη μνήμη και με βάση το μοντέλο [23]. Χρησιμοποιούν ένα μοντέλο μίγμα που κατασκευάστηκε με βάση ένα σύνολο των αποθηκευμένων προφίλ του χρήστη και τη χρήση της κατανομής των αξιολογήσεων του χρήστη που μπορεί να κάνει προβλέψεις. Για την αντιμετώπιση των νέων προβλημάτων του χρήστη, μια ενεργή επέκταση μάθησης στο PMCF σύστημα μπορεί να χρησιμοποιηθεί για την αναζήτηση ενός ενεργού χρήστη για πρόσθετες πληροφορίες όταν δεν υπάρχουν επαρκείς πληροφορίες. Για να μειώσει το χρόνο υπολογισμού, η PMCF επιλέγει ένα μικρό υποσύνολο, που ονομάζεται χώρος προφίλ από το σύνολο της βάσης δεδομένων των αξιολογήσεων των χρηστών και δίνει τις προβλέψεις από το μικρό χώρο προφίλ αντί ολόκληρης της βάσης δεδομένων. Η PMCF έχει μεγαλύτερη ακρίβεια από ό,τι η Pearson CF και το CF μοντέλο με βάση τη χρήση Naive Bayes.

Η Personality diagnosis (PD) είναι ένα αντιπροσωπευτικό υβριδικό CF, μια προσέγγιση που συνδυάζει το CF που βασίζεται στη μνήμη και τη CF που βασίζεται στο μοντέλο και διατηρεί κάποια πλεονεκτήματα και των δύο αλγορίθμων [94]. Στην PD, ο ενεργός χρήστης θεωρητικά δημιουργείται επιλέγοντας έναν στην τύχη από άλλους χρήστες και προσθέτοντας θόρυβο Gauss τις αξιολογήσεις του. Δεδομένων των γνωστών βαθμολογιών του ενεργού χρήστη, μπορούμε να υπολογίσουμε την πιθανότητα ότι αυτός ή αυτή είναι ο ίδιος "τύπος προσωπικότητας" με τους υπόλοιπους χρήστες, και ποια η πιθανότητα αυτός ή αυτή να ήθελα τα νέα στοιχεία. Η PD μπορεί επίσης να θεωρηθεί ως μια μέθοδος ομαδοποίησης με ακριβώς ένα χρήστη ανά σύμπλεγμα.

## ***2.2 Νέες τεχνικές και θέματα που προκύπτουν***

Για να επιτευχθεί μεγαλύτερη ακρίβεια στις συστάσεις πρέπει να λάβουν περισσότερες πληροφορίες. Η χρήση των δικτύων εμπιστοσύνης [24], με προσεγγίσεις που έγιναν στις καταναμημένες βάσεις δεδομένων, και οι επιπλέον πληροφορίες που παρέχονται με βάση το περιεχόμενο των μεθόδων, βοηθούν ώστε να ασχοληθεί με τους «κρύους» χρήστες και τα προβλήματα σποραδικότητας. Υπό αυτή την έννοια, πολύ γνωστές μέθοδοι καταλογισμού [55] [62] θα μπορούσε επίσης να χρησιμοποιηθούν για να γεμίσουν μήτρες CF και στην άμβλυνση των προβλημάτων που οφείλονται στην έλλειψη πληροφοριών. Σύμφωνα με μια προηγούμενη μελέτη το σύνολο δεδομένων, μπορεί να επιλεγεί η πιο εξοπλισμένη μέθοδος καταλογισμού για να αποκτήσουν ακριβή και ρεαλιστικά αποτελέσματα. Επιπλέον, είναι επίσης σημαντικό να βρούμε και τον εντοπισμό των ακραίων τιμών [63] τόσο σε κεντρικό επίπεδο όσο και σε αποκεντρωμένα συστήματα, ώστε να αυξηθεί η ποιότητα στα αποθηκευμένα δεδομένα.[72]

Στην αρένα της ιδιωτικής ζωής, η τρέχουσα τάση είναι να χρησιμοποιείτε ασφαλή πρωτόκολλα υπολογισμού και τις ομοιορφικές ιδιότητες του δημοσίου κλειδιού για να προστατεύουν τα δεδομένα των χρηστών, όσο διεξάγετε η CF. Εκτός από την ακρίβεια και την προστασία της ιδιωτικής ζωής, οι μέθοδοι CF χρησιμοποιούνται ευρέως για την σύσταση προϊόντων για τους χρήστες σε πραγματικό χρόνο, υπολογίζουν και λαμβάνουν υπόψη τα γενικά έξοδα επικοινωνίας. Αυτό το γεγονός προσθέτει τη μεταβλητή χρόνου στη γνωστή ακρίβεια και στην εξισορρόπηση των προβλημάτων ιδιωτικότητας.

Αρκετές προσεγγίσεις [59], [60] προτείνουν να αποσυνδέσουμε τους χρήστες από τις αξιολογήσεις τους χρησιμοποιώντας αναγνωρισμένες τεχνικές όπως mixservers ή συστήματα ανωνυμίας τους. Αν και οι μέθοδοι δεν εγγυώνται την ασφάλεια για τους δικούς τους, η χρήση τους θα πρέπει να είναι γενικευμένη και να ενισχυθεί η προστασία της ιδιωτικής ζωής των ήδη γνωστών CF και μεθόδων PPCF.

Όπως αναφέρθηκε προηγουμένως, είναι δυνατόν να ταξινομηθεί η PPCF σε κεντρικές και αποκεντρωμένες μεθόδους. Τα κεντρικά συστήματα προσφέρουν πολλά οφέλη, αλλά και ζητήματα προστασίας της ιδιωτικής ζωής. Από τη μία πλευρά, η συνολική απόδοση είναι αυξημένη επειδή η επιπλέον επικοινωνία αποφεύγεται λόγω εξόδων μεταξύ των χρηστών και των τρίτων μερών. Αφετέρου στις κεντρικές μεθόδους, τα δεδομένα διαχειρίζονται ένα μέρος τα οποία έχει τον απόλυτο έλεγχο πάνω τους, με τα θέματα προστασίας της ιδιωτικής ζωής που συνεπάγεται ότι τα δεδομένα έχουν χαμηλό επίπεδο προστασίας. Δυστυχώς, κανείς δεν μπορεί να εγγυηθεί ότι οι επιχειρήσεις δεν θα πωλούν τις πληροφορίες σε άλλους σε περίπτωση πτώχευσης [52]. Παρ' όλα αυτά, μπορούμε να χρησιμοποιήσουμε τις SDC γνωστές μεθόδους [64] για την παροχή K-ανωνυμίας[31] και επίσης ενός δημόσιου κλειδιού κρυπτογράφησης των συστημάτων [32], [33] για να αποθηκεύσουμε τις ευαίσθητες και προσωπικές πληροφορίες με έναν ασφαλή τρόπο και ως εκ τούτου, εμποδίζει τη γνωστοποίηση τους.

Τα αποκεντρωμένα συστήματα έχουν συνήθως υψηλότερη προστασία των δεδομένων από τα συγκεντρωτικά. Αυτό συμβαίνει επειδή οι εμπλεκόμενοι χρήστες και τα τμήματα διαχειρίζονται τα δικά τους σύνολα δεδομένων, καθώς και τα δεδομένα που αποστέλλονται μεταξύ αυτών είναι συνήθως ελεγχόμενα από ασφαλή πρωτοκόλλα, για να αποφευχθεί η παράνομη αποκάλυψη των προσωπικών δεδομένων, η οποία θα μπορούσε επίσης να δώσει ένα οικονομικό όφελος στους ανταγωνιστές. Ωστόσο, η επιπλέον προστασία της ιδιωτικής ζωής που επιτυγχάνεται έχει ένα σημαντικό κόστος, καθώς αυξάνει η εναέρια επικοινωνία, ως εκ τούτου, επηρεάζεται την αποτελεσματικότητα των συστάσεων σε πραγματικό χρόνο.

Αν και τα θέματα του χρόνου μπορούν να αντιμετωπιστούν εν μέρει με την συνεχή αναβάθμιση των αξιολογήσεων της μήτρας σε ένα κεντρικό σύστημα, η λύση αυτή είναι δύσκολο να εφαρμοστεί στα συστήματα όπου οι χρήστες διαχειρίζονται τις δικές τους αξιολογήσεις και οποιοσδήποτε τρίτος έχει χρησιμοποιηθεί για την εκτέλεση ενδιάμεσων υπολογισμών, διότι στην περίπτωση αυτή, ένα οποιοδήποτε δεδομένο είναι αποθηκευμένο.

Επιπλέον, οι χρήστες οι οποίοι ελέγχουν τα δικά τους προφίλ τείνουν να σχηματίζουν κοινότητες με ήδη γνωστούς χρήστες ή με αυτούς που έχουν παρόμοιες συμπεριφορές. Όπως έχουμε ήδη σχολιάσει στο τμήμα I, το γεγονός αυτό μπορεί να συνεπάγεται σε ένα πολύ γνωστό homophily πρόβλημα [29], το οποίο θα μπορούσε να λυθεί εν μέρει με την προσθήκη διαφορετικών κοινοτήτων, προκειμένου να δημιουργηθεί μια μεγάλη κοινωνία από ομάδες, μέσω τεχνικών ομαδοποίησης. Με αυτό τον τρόπο, οι χρήστες αναζητούν συστάσεις ποιότητας για ένα δοσμένο πλαίσιο, έχοντας αρκετές πληροφορίες.

### **2.3 Επιπλέον προβλήματα της PPCF**

Το συνεργατικό φιλτράρισμα είναι ένα σύστημα συστάσεων που επιτρέπει την αυτόματη σύσταση προς ομάδες χρηστών σε πολλά πλαίσια. Στο κεφάλαιο αυτό, πραγματοποιήσαμε μια σύντομη επισκόπηση της state-of-the-art διαφορετικές οικογένειες CF και μεθόδους. Παρά τα μεγάλα πλεονεκτήματα του CF, έχουμε επισημάνει τον σημαντικό κίνδυνο που σχετίζεται με την προστασία της ιδιωτικής ζωής των χρηστών. Ως εκ τούτου, εκτελέσαμε μια ανάλυση των σημαντικότερων μεθόδων PPCF και τις διαφοροποιήσαμε κυρίως από το αν χρησιμοποιούν μία κεντρική ή κατανεμημένη αρχιτεκτονική για την προστασία των δεδομένων και τη δημιουργία συστάσεων. Έχει επίσης διεξαχθεί μια σύντομη ανάλυση από τα σημαντικότερα μειονεκτήματα της PPCF και έχουν συζητηθεί διάφοροι τρόποι για να ξεπεραστούν αυτά τα προβλήματα. [72]

Ένα μεγάλο μέρος των μεθόδων CF χρειάζεται μια μελέτη σε βάθος, δεδομένου ότι εξακολουθούν να υπάρχουν προκλήσεις που πρέπει να ξεπεραστούν. Η επαρκής προστασία της ιδιωτικής ζωής των χρηστών, διατηρώντας παράλληλα ένα υψηλό βαθμό ακρίβειας και την ποιότητα των συστάσεων, είναι ίσως ένα από τα πιο σημαντικά θέματα. Θα πρέπει να επικεντρώνεται σε θέματα όπως το πώς να επεξεργάζονται αραιά δεδομένα πιο αποτελεσματικά, ενώ ταυτόχρονα προστατεύεται και η ιδιωτική ζωή των χρηστών. Ένας πρόσθετος περιορισμός είναι ότι τα παραπάνω θα πρέπει να γίνεται με τη χαμηλότερη δυνατή επικοινωνία και τον χαμηλότερο δυνατό υπολογισμό.

### 3. Βασικές έννοιες

#### 3.0 Εισαγωγή

Θεωρητικά είναι δεδομένο ότι υπάρχει κάποιο υπόβαθρο σχετικά με τον έλεγχο αποκάλυψης στατιστικών στοιχείων και μικροσυσσωμάτωσης. Όμως παρακάτω γίνεται μια μικρή αναφορά σε βασικές έννοιες που χρησιμοποιούνται στο επόμενο κεφάλαιο και που είναι απαραίτητες για να κατανοήσουμε και να συγκρίνουμε τις μεθόδους μεταξύ τους.

#### 3.1 Η Κανονική Κατανομή (Gauss)

Η **κανονική κατανομή (normal distribution)** θεωρείται η σπουδαιότερη κατανομή της *Θεωρίας Πιθανοτήτων* και της *Στατιστικής*. Οι λόγοι που εξηγούν την εξέχουσα θέση της, είναι βασικά δύο:

- i) Πολλές τυχαίες μεταβλητές περιγράφονται ικανοποιητικά από την κανονική κατανομή ή περιγράφονται από κατανομές που μπορούν να προσεγγισθούν από την κανονική κατανομή.
- ii) Οι ιδιότητες της κανονικής κατανομής αξιοποιούνται στη Στατιστική Συμπερασμαματολογία. Ουσιαστικά, η κανονική κατανομή, αποτελεί το θεμέλιο της Στατιστικής Συμπερασμαματολογίας.

Το «μυστικό» που εξηγεί το μεγάλο εύρος εφαρμογών της *κανονικής κατανομής*, βρίσκεται σε ένα εκπληκτικά ισχυρό θεωρητικό αποτέλεσμα της *Θεωρίας Πιθανοτήτων* το οποίο επιβεβαιώνεται και πειραματικά. Πρόκειται για το **Κεντρικό Οριακό Θεώρημα (Central Limit Theorem)** τις βάσεις του οποίου έθεσαν δύο μεγάλοι Μαθηματικοί. Ο *Abraham De Moivre* το 1733 και, έναν αιώνα περίπου αργότερα, το 1812, ο *Laplace*. Σε αυτό το σημείο, δε θα διατυπώσουμε αυστηρά, ούτε θα αποδείξουμε, το *Κεντρικό Οριακό Θεώρημα*. Θα προσπαθήσουμε να εξηγήσουμε μόνο το νόημα και τη σημασία του.

Σύμφωνα με το *Κεντρικό Οριακό Θεώρημα*, το άθροισμα και -επομένως- η μέση τιμή, μεγάλου αριθμού ανεξάρτητων παρατηρήσεων, ακολουθεί κατά προσέγγιση *κανονική κατανομή*, ανεξαρτήτως από το ποια κατανομή ακολουθούν οι παρατηρήσεις. Πώς, όμως, αυτό το αποτέλεσμα ερμηνεύει τη μεγάλη εφαρμοσιμότητα της *κανονικής κατανομής*; Είναι απλό. Σε πολλά φαινόμενα και πειράματα, οι τιμές διαφόρων χαρακτηριστικών

(μεταβλητών), είναι αποτέλεσμα αθροιστικής επίδρασης πολλών ανεξάρτητων αιτίων-παραγόντων κανένα από τα οποία δεν υπερισχύει των άλλων.

Για παράδειγμα, ο χρόνος αναμονής σε μια ουρά, είναι αποτέλεσμα πολλών παραγόντων, όπως, η ημέρα της εβδομάδας, η ώρα της ημέρας, η αποτελεσματικότητα του υπαλλήλου, το είδος της συναλλαγής που διεκπεραιώνεται, κ.ά. Επίσης, το βάρος των ζώων μιας κτηνοτροφικής μονάδας, οφείλεται σύμφωνα με τους ειδικούς, σε πληθώρα παραγόντων όπως, η ατομικότητα του ζώου, η φυλή, το γένος, οι συνθήκες διατροφής, οι συνθήκες ενσταυλισμού, κ.ά. Καθένας από τους παράγοντες αυτούς επιφέρει ένα θετικό ή αρνητικό αποτέλεσμα και όλοι μαζί αθροιστικά συντελούν στη διαμόρφωση του τελικού αποτελέσματος. Τέτοια χαρακτηριστικά (μεταβλητές), εμφανίζονται σε πολλά φαινόμενα και πειράματα. Το *Κεντρικό Οριακό Θεώρημα* λει ότι αυτά ακριβώς τα χαρακτηριστικά περιγράφονται ικανοποιητικά από την *κανονική κατανομή*. Επιπλέον, το *Κεντρικό Οριακό Θεώρημα* **συνδέει** την *κανονική κατανομή* με **οποιαδήποτε άλλη κατανομή** (αφού δεν προϋποθέτει να ακολουθούν οι παρατηρήσεις την *κανονική κατανομή*), γεγονός το οποίο, απαντάει, επίσης, στο ερώτημα, γιατί η *κανονική κατανομή* βρίσκει εφαρμογή σε μεγάλο πλήθος φαινομένων και πειραμάτων.

Πρέπει να τονίσουμε ότι για να αποδειχθεί ότι ένα συγκεκριμένο χαρακτηριστικό (μεταβλητή) προσεγγίζεται ικανοποιητικά από την *κανονική κατανομή*, πρέπει να γίνουν μετρήσεις που να επαληθεύουν ένα τέτοιο συμπέρασμα. Μια από τις πρώτες εφαρμογές της *κανονικής κατανομής*, έγινε το 1809 από το μεγάλο Γερμανό Μαθηματικό *Carl F. Gauss* ο οποίος διαπίστωσε ότι τα σφάλματα που γίνονται σε αστρονομικές παρατηρήσεις μπορούν να περιγραφούν ικανοποιητικά από την *κανονική κατανομή*. Στη συνέχεια, διαπιστώθηκε επίσης, ότι τα τυχαία σφάλματα (όχι τα συστηματικά) που εμφανίζονται σε διάφορες μετρήσεις ακολουθούν με ικανοποιητική προσέγγιση *κανονική κατανομή*. Για το λόγο αυτό, η *κανονική κατανομή* ονομάζεται και **κατανομή των σφαλμάτων (*law of errors*)**. Επίσης, είναι γνωστή ως **κατανομή του Gauss (*Gaussian distribution*)**, για τη μεγάλη συνεισφορά του *Gauss* στην ανάδειξη των ιδιοτήτων και της σημασίας της.



### 3.2 Η MDAV

Πολλοί από τους συνήθεις μετασχηματισμούς εκτελούνται σε σύνολο δεδομένων (data mining, ανακάλυψη της γνώσης, κτλ) μπορούν να προβληθούν σαν διαδικασίες ομαδοποίησης με διαφορετικά είδη περιορισμών [96, 69, 104, 108]. Σε αυτό το σημείο, παρουσιάζουμε την αντιμετώπιση του πρόβληματος της μικροσυσσωμάτωσης, ένα ιδιαίτερο είδος του προβλήματος ομαδοποίησης όπου υπάρχουν περιορισμοί σχετικά με το ελάχιστο μέγεθος των clusters ή των ομάδων, αλλά όχι με τον αριθμό τους, και η εντός-ομάδων ομοιογένεια πρέπει να μεγιστοποιηθεί. Η μικροσυσσωμάτωση είναι ένα πρόβλημα που εμφανίζεται στα statistical disclosure control (SDC), όπου χρησιμοποιείται για την ομαδοποίηση ενός συνόλου εγγραφών σε ομάδες τουλάχιστον  $k$  εγγραφών, με το  $k$  να είναι μία παράμετρος προσδιορίσιμη από τον χρήστη. Η συλλογή από τις ομάδες λέγεται  $k$ -τμήματα του συνόλου δεδομένων. Το {microaggregated} μικροσυνολικό σύνολο δεδομένων έχει κατασκευαστεί αντικαθιστώντας κάθε αυθεντική εγγραφή από το {centroid} κέντρο βάρους της ομάδας που ανήκει. Το {microaggregated} μικροσυνολικό σύνολο δεδομένων μπορεί να απελευθερωθεί χωρίς να θέσει σε κίνδυνο την προστασία προσωπικών δεδομένων των ατόμων τα οποία συγκροτούν το αυθεντικό σύνολο δεδομένων: οι εγγραφές μέσα σε μια ομάδα είναι δυσδιάκριτες στα {released} απελευθερωμένα σύνολα δεδομένων. Όσο υψηλότερη η εντός-ομάδων ομοιογένεια στο αυθεντικό σύνολο δεδομένων, τόσο μικρότερη η απώλεια των πληροφοριών που προκύπτει όταν αντικαθιστούμε εγγραφές σε μια ομάδα από το {centroid} κέντρο βάρους, ως εκ τούτου η εντός-ομάδων ομοιογένεια είναι αντιστρόφως ανάλογη με την χαμένη πληροφορία που προκαλείται από την μικροσυσσωμάτωση. Η αντιστροφή από το εντός-ομάδων άθροισμα των τετραγώνων  $SSE$  είναι η πιο συνηθισμένη μέτρηση για την ομοιογένεια ομαδοποίησης [98, 100, 101, 105, 109]. Το  $SSE$  μπορεί να υπολογιστεί ως

$$SSE = \sum_{i=1}^s \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)' (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)$$

Όπου

$s$ : ο αριθμός των παραγόμενων συνόλων

$n_i$ : ο αριθμός των εγγραφών σε  $i$ -th σετ

$x_{ij}$ : οι  $j$ -th εγγραφές

$\bar{x}_i$ : το {centroid} κέντρο βάρους του  $i$ -th σετ.

Σε όρους  $SSE$  το πρόβλημα της μικροσυσσωμάτωσης συνίσταται στην εύρεση των  $K$  - τμημάτων με το μικρότερο  $SSE$ .

Η Μικροσυσσωμάτωση έχει χρησιμοποιηθεί για πολλά χρόνια σε διάφορες χώρες: ξεκίνησε στην Eurostat [95] στις αρχές της δεκαετίας του ενενήντα, και έκτοτε χρησιμοποιείται στη Γερμανία [107] και πολλές άλλες χώρες [99]. Αυτά τα χρόνια εμπειρίας και πρακτικής που διατίθενται για μικροσυσσωμάτωση μας κληροδότησαν μια ποικιλία προσεγγίσεων, την οποία εμείς έπειτα εν συντομία θα συνοψίσουμε.[70]

- Βέλτιστες μεθόδους:

- Η Μονοπαραγοντική περίπτωση: Στο [102] παρουσιάστηκε ένας πολυώνυμος αλγόριθμος για βέλτιστη μονομεταβλητή μικροσυσσωμάτωση . Αλλά η βέλτιστη πολυμεταβλητή μικροσυσσωμάτωση απεδείχθη ότι είναι NP-hard στο [106].

- Η πολυμεταβλητή περίπτωση: Η Βέλτιστη πολυμεταβλητή μικροσυσσωμάτωση φάνηκε να είναι NP-hard στο [106]. Έτσι, οι μόνες πρακτικές μέθοδοι πολυμεταβλητής μικροσυσσωμάτωσης είναι οι ευρετικές.

- ευριστικές μεθόδους:

- Ευριστικές σταθερού μεγέθους: Το πιο γνωστό παράδειγμα αυτής της κατηγορίας είναι η μέγιστη Απόσταση προς Μέση Vector (MDAV) [64, 69, 103]. Η MDAV παράγει ομάδες σταθερής πληθικότητας  $k$  και, όταν ο αριθμός των εγγραφών δεν διαιρείται με  $k$ , μία ομάδα με πληθικότητα μεταξύ  $k$  και  $2k - 1$ . Η MDAV έχει αποδειχθεί ότι έχει την καλύτερη επίδοση από άποψη χρόνου και όσον αφορά την ομοιογένεια των προκύπτων ομάδων.

- Ευριστικές μεταβλητού μεγέθους: Η απόδοση ποικίλει μεταξύ των  $k$  και  $2k - 1$  με μεγέθη ομάδας  $k$ -τμημάτων . Μια τέτοια ευελιξία μπορεί να αξιοποιηθεί για την επίτευξη υψηλότερης ομοιογένεια μέσα σε μία ομάδα. Οι Μέθοδοι μεταβλητού μεγέθους περιλαμβάνουν τη γενετική εμπνευσμένη προσέγγιση για μικρά σύνολα δεδομένων [108], καθώς και [96, 104, 97].

Για να συνοψίσουμε, η πρόκληση στη μικροσυσσωμάτωση είναι να σχεδιάσει καλές ευριστικές για τη πολυμεταβλητή περίπτωση, όπου το ιδανικότερο αναφέρεται στο



συνδυασμό υψηλής ομοιογένειας σε μια ομάδα και στην υπολογιστική αποδοτικότητα. Αυτό το χαρτί είναι για μια νέα μέθοδο κατά μήκος αυτής της γραμμής.

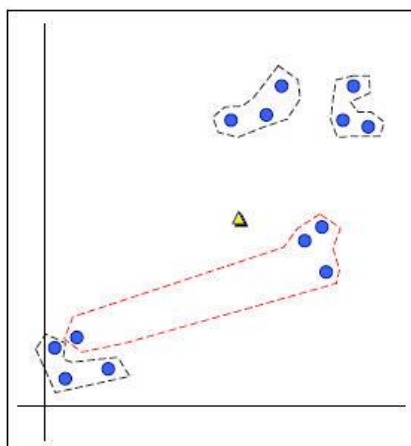
Η MDAV είναι, μία από τις καλύτερες μεθόδους ευρετικής πολυμεταβλητής μικροσυσσωμάτωση. Ωστόσο, παραμένει μια σταθερού μεγέθους ευρετική μέθοδος, υπάρχουν καταστάσεις στις οποίες δίδει  $k$ -τμήμα μακριά από τη βέλτιστη τιμή. Αυτό απεικονίζεται από το επόμενο παράδειγμα με παιχνίδια.

**Παράδειγμα 1.** Έστω ότι έχω ένα  $D$  σύνολο δεδομένων που αποτελείται από 13 βιβλία που έχουν 2 ιδιότητες.

$D = \{(2.4, 3), (1.68, 4.9), (3.18, 5.54), (5.32, 3.6), (18.68, 11.49), (20.14, 9.56), (19.85, 12.33),$

$(13.67, 18.9), (17.11, 21), (16.07, 19.23), (21.28, 18.9), (22, 21), (23, 18.5)\}$

Θεωρώντας  $k = 3$ , το Σχήμα 1 απεικονίζει την έξοδο του αλγορίθμου MDAV όταν εφαρμόζεται πάνω στο παράδειγμα μας. Το σχήμα απεικονίζει τις εγγραφές στο  $D$  (κύκλοι) και το κέντρο βάρους (τρίγωνο). Η ομάδα που σημειώνεται με κόκκινο χρώμα είναι πολύ διάσπαρτη και τελικά τα γενικά 3-τμήματα(ομάδες) δεν είναι αντιπροσωπευτικά. Αυτό το παράδειγμα δείχνει ότι το σταθερό μέγεθος που ορίζει η MDAV μπορεί να αποτύχει να προσαρμόσει κατάλληλα τα προκύπτον  $k$ -τμήματα στα συγκεκριμένα δεδομένα που έχουν οριστεί.[70]



**Σχήμα 3.1** Η έξοδος του αλγορίθμου MDAV για το παράδειγμα των παιχνιδιών

### 3.3 H SVD

#### 3.3.1 Εφαρμογή της SVD στο συνεργατικό φιλτράρισμα

Η αδυναμία των μεθόδων εύρεσης για τον πλησιέστερο γείτονα, σε μεγάλες και αραιές βάσεις δεδομένων οδήγησε στο να διερευνηθούν εναλλακτικοί αλγόριθμοι συστήματος συστάσεων. Μια προσέγγιση επιχείρησε να γεφυρώσει την αραιότητα με την ενσωμάτωση ημι-ευφών πρακτόρων στο σύστημα φιλτραρίσματος (Sarwar et al., 1998, Good et al., 1999). Αυτοί οι παράγοντες αξιολογούν και κατατάσσουν το κάθε προϊόν, χρησιμοποιώντας συντακτικά χαρακτηριστικά, παρέχοντας ένα πυκνό σύνολο από βαθμολογίες που, αυτές βοηθούν στην κάλυψη και τη βελτίωση της ποιότητας. Η λύση αυτή, ωστόσο, δεν μπόρεσε να αντιμετωπίσει το θεμελιώδες πρόβλημα των σχέσεων μεταξύ των χρηστών με λίγα κοινά σημεία, αλλά με αραιή βαθμολογία από τους πελάτες. Η ερευνητική κοινότητα είχε εκτεταμένη εμπειρία μάθησης σε αραιές βάσεις δεδομένων. Μετά την εξέταση σε αρκετές τεχνικές KDD, αποφάσισαν να δοκιμάσουν την εφαρμογή της Λανθάνουσας Σημασιολογικής Ευρετηρίασης (LSI) προς μείωση των διαστάσεων του πελάτη -προϊόν μιας αξιολόγησης σε μια μήτρα.[2]

Η LSI είναι μια τεχνική μείωσης διάστασης που έχει χρησιμοποιηθεί ευρέως στην ανάκτηση πληροφοριών (IR) προς λύσει σε προβλήματα της συνωνυμίας και της πολυσημίας (Deerwester et al., 1990). Λαμβάνοντας υπόψη μια συχνότητα όρος-έγγραφο μήτρας, η LSI χρησιμοποιείται για την κατασκευή δύο πινάκων με μειωμένες διαστάσεις. Στην ουσία, αυτές μήτρες αντιπροσωπεύουν λανθάνουσες ιδιότητες των όρων, όπως αντανακλάται από την εμφάνισή τους σε έγγραφα, καθώς και των έγγραφων, όπως φαίνεται από τους όρους που εμφανίζονται στο εσωτερικό τους. Προσπαθεί να συνδέσει τις σχέσεις μεταξύ των πελατών με βάση τις αξιολογήσεις τους στα προϊόντα. Με τη μείωση της διάστασης του χώρου του προϊόντος, μπορούμε να αυξήσουμε την πυκνότητα και με αυτόν τον τρόπο να βρείτε περισσότερες ειδικότητες. Η Ανακάλυψη της λανθάνουσας σχέσης σε μια βάση δεδομένων μπορεί δυναμικά να λύσει το πρόβλημα της συνωνυμίας στα συστήματα συστάσεων. Στο Berry et al. (1995) επισημαίνουν ότι η μείωση των διαστάσεων του ορθογωνίου που προκύπτουν από τη SVD είναι λιγότερο θορυβώδεις από τα αρχικά δεδομένα και σύλληψη

των λανθάνοντα ενώσεων μεταξύ των όρων και των έγγραφων. Το (Billsus et al., 1998) πήρε τα πλεονέκτημα αυτής της σημασιολογικής για να μειωθεί η διάσταση του χώρου. Η μειωμένη λειτουργία διαστήματος χρησιμοποιήθηκε για να εκπαιδεύσει ένα νευρωνικό δίκτυο ώστε να δημιουργήσει προβλέψεις. Η κατασκευή του SVD-based αλγορίθμου συστάσεων έγινε για το σκοπό της δημιουργίας προβλέψεων και top-N συστάσεων.[2]

### 3.3.2 Singular Value Decomposition (SVD)

Η SVD είναι μία πολύ γνωστή τεχνική παραγοντοποίησης μιας μήτρας όπου οι παράγοντες είναι ένας  $m \times n$  πίνακας  $R$  που αποτελείται από τρεις πίνακες όπως το ακόλουθο:

$$R = U \cdot S \cdot V'$$

Όπου,  $U$  και  $V$  είναι δύο ορθογώνιες μήτρες μεγέθους  $m \times r$  και  $n \times r$  αντίστοιχα,  $r$  είναι οι γραμμές του πίνακα  $R$ . Το  $S$  είναι ένας διαγώνιος πίνακας του μεγέθους  $r \times r$  που έχει όλες τις ιδιόμορφες τιμές του πίνακα  $R$  ως διαγώνια στοιχεία του. Όλες οι εγγραφές του πίνακα  $S$  είναι θετικοί και αποθηκεύονται σε φθίνουσα σειρά μεγέθους. Οι μήτρες που λαμβάνεται με την εκτέλεση SVD είναι ιδιαίτερες χρήσιμες στην εφαρμογή μας, λόγω της ιδιότητας της SVD να παρέχει την καλύτερη προσέγγιση με τη χαμηλότερη κατάταξη του πρωτότυπου πίνακα  $R$ . Είναι δυνατό να μειωθεί η  $r \times r$  μήτρα του  $S$  και να έχει μόνο  $k$  μεγαλύτερες διαγώνιες τιμές για να ληφθεί μία μήτρα  $S_k$ , με  $k < r$ . Αν οι μήτρες  $U$  και  $V$  μειώνονται ανάλογα, τότε η ανακατασκευασμένη μήτρα  $R_k = U_k \cdot S_k \cdot V'_k$  είναι το πιο κοντινό rank- $k$  μήτρα για τον πίνακα  $R$ . Με άλλα λόγια, το  $R_k$  ελαχιστοποιεί την Frobenius νόρμα  $\|R - R_k\|$  πάνω από όλες τις μήτρες rank- $k$ . [2]

Χρησιμοποιούμε την SVD σε συστήματα συστάσεων για την εκτέλεση δύο διαφορετικών ασκήσεων: Πρώτον, χρησιμοποιούμε την SVD για να συλλάβει τις λανθάνουσες σχέσεις μεταξύ των πελατών και των προϊόντων που μας επιτρέπουν να υπολογίσουμε την προβλεπόμενη πιθανότητα ενός ορισμένου προϊόντος από τον πελάτη. Δεύτερον, χρησιμοποιούμε την SVD για να παραχθεί μια χαμηλού διαστάσεων αναπαράσταση του αρχικού χώρου του πελάτη-προϊόν και στη συνέχεια να υπολογίσει τη

γειτονιά στο μειωμένο χώρο. Στη συνέχεια χρησιμοποιείται για να δημιουργηθεί μια λίστα των top-N προϊόν για συστάσεις στους πελάτες. [2]

Πανεπιστήμιο Πειραιώς

## 4. PPCF με k-Ανωνυμία μέσω μικροσυσσωμάτωσης

### 4.0 Εισαγωγή

Η μέθοδος του συνεργατικού φιλτραρίσματος -*Collaborative Filtering* -(CF) είναι ένα εισηγητικό σύστημα το οποίο έχει ολοένα και μεγαλύτερη σημασία για τη βιομηχανία.

Η τρέχουσα εργασία εστιάζει στην Διατήρηση Προστασίας Προσωπικών Δεδομένων Συνεργατικού Φιλτραρίσματος -*Privacy Preserving Collaborative Filtering*- (PPCF), σκοπός της οποίας είναι να επιλύσει τα ζητήματα που αντιμετωπίζονται με την συστηματική συλλογή των προσωπικών πληροφοριών. Σε αυτή την εργασία, παρουσιάζεται μια μέθοδος PPCF που βασίζεται στη μικροσυσσωμάτωση, νοθεύει δεδομένα για την παροχή K-ανωνυμίας, ενώ ταυτόχρονα κάνει ακριβείς συστάσεις. Τα πειραματικά αποτελέσματα αποδεικνύουν ότι η προτεινόμενη μέθοδος διαταράσσει τα δεδομένα πιο αποτελεσματικά από ό,τι η γνωστή και ευρέως χρησιμοποιούμενη μέθοδος που βασίζεται στην παραμόρφωση Gaussian.

Οι περισσότεροι άνθρωποι έχουν την τάση να λαμβάνουν σοβαρά υπόψη άλλες συστάσεις ανθρώπων, όταν θέλουν να αγοράσουν ένα προϊόν. Η διαδικασία αξιολόγησης πριν από την αγορά ενός προϊόντος είναι δύσκολο όταν οι απόψεις που παίρνουμε είναι διαφορετικές ή / και από πολλαπλές πηγές. Αυτό έχει ιδιαίτερη σημασία όταν η πηγή των συστάσεων είναι το Διαδίκτυο. Από τη μία το Διαδίκτυο προσφέρει έναν πλούτο πληροφοριών σχετικά με μια τεράστια ποικιλία προϊόντων και υπηρεσιών που μπορούν να είναι χρήσιμες σε πιθανούς αγοραστές. Από την άλλη πλευρά, αυτός ο πλούτος των πληροφοριών μπορεί να γίνει ένα πρόβλημα και όχι η λύση, διότι μπορεί να παρεμποδίσει τη λήψη αποφάσεων.

Τα εισηγητικά συστήματα [1], προέρχονται από την ανακάλυψη της γνώσης σε βάσεις δεδομένων (KDD) [2], [3]. Τα συστήματα KDD είναι διαδικασίες που χρησιμοποιούνται από τις εταιρείες για να αναλύουν και να ανακαλύψουν κατανοητά πρότυπα μέσα από μεγάλες συλλογές δεδομένων που μπορεί να βοηθήσουν, για παράδειγμα, να εξοικονομήσουν χρήματα ή να πουλήσουν περισσότερα προϊόντα στους πελάτες. Το συνεργατικό φιλτράρισμα (CF) [4], [5] είναι ένα σύστημα συστάσεων το οποίο περιλαμβάνει μια μεγάλη οικογένεια των μεθόδων σύστασης. Ο στόχος της CF είναι να κάνει υποδείξεις για τα στοιχεία (π.χ. βιβλία, μουσική ή ταινίες), με βάση τις προτιμήσεις των χρηστών που έχουν

ήδη αποκτήσει ή / και βαθμολόγησαν αυτά τα στοιχεία. Το CF εμφανίστηκε με σκοπό την αυτόματη σύσταση σε ένα ψηφιακό περιβάλλον.

Προκειμένου να κάνουν προβλέψεις, οι μέθοδοι CF χρησιμοποιούν μεγάλες βάσεις δεδομένων που αποθηκεύουν πληροφορίες σχετικά με τις σχέσεις μεταξύ των σερτ των χρηστών και των αντικειμένων. Αυτά τα δεδομένα όπως έχει αναφερθεί και σε προηγούμενο κεφάλαιο, μοντελοποιούνται ως μήτρες που αποτελούνται από  $n$  χρήστες και  $m$  στοιχεία, καθώς και κάθε κελί  $(i, j)$  αποθηκεύει την αξιολόγηση του χρήστη  $i$  στη θέση  $j$ . Ως εκ τούτου, μια τιμή μπορεί να είναι σε ένα εύρος τιμών (π.χ. μεταξύ 0 και 10) ή απλά με το δυαδικό ψηφίο (θετικό / αρνητικό, ή αγόρασαν / δεν αγόρασαν) όπως και σε βάσεις δεδομένων με το καλάθι αγοράς.

Οι μέθοδοι CF βασίζονται στην υπόθεση ότι παρόμοιοι χρήστες, με την έννοια ότι έχουν παρόμοια ενδιαφέροντα ή συμπεριφορές, θα τους ενδιαφέρουν τα ίδια προϊόντα. Ως εκ τούτου, τα αντικείμενα που αγοράζονται από ένα χρήστη  $U_a$  μπορεί να συστηθούν σε άλλο χρήστη  $U_b$ , αν  $U_a$  και  $U_b$  έχουν ανάλογο ενδιαφέρον ή παρόμοια συμπεριφορά. Άλλες προσεγγίσεις, ανατρέπουν αυτή τη διαδικασία, κάνοντας συστάσεις με βάση την ομοιότητα μεταξύ των στοιχείων. Σε αυτό το πλαίσιο, μια ομάδα παρόμοιων χρηστών ή είδη, αποτελούν μια γειτονιά.

Κατά τη διάρκεια των τελευταίων ετών, η χρήση των συστημάτων εμπιστοσύνης έχει γνωρίσει μια σημαντική αύξηση στο διαδίκτυο. Μια δήλωση εμπιστοσύνης ορίζεται ως η ρητή γνώμη που εκφράζεται από ένα χρήστη σε έναν άλλο χρήστη, όσον αφορά την αντίληψη για την ποιότητα ορισμένων χαρακτηριστικών του χρήστη [24]. Η έννοια της εμπιστοσύνης χρησιμοποιείται ευρέως, για παράδειγμα, σε μηχανές αναζήτησης όπως το Google, το οποίο χρησιμοποιεί την παγκόσμια μέτρηση εμπιστοσύνης [25], καθώς και στο ηλεκτρονικό εμπόριο (π.χ. μέσω ebay) στο οποίο οι χρήστες εκφράζουν το επίπεδο ικανοποίησής τους από την αγορά ενός προϊόντος.

Με την άθροιση των δηλώσεων εμπιστοσύνης που εκφράζεται από κάθε χρήστη, είναι δυνατόν να παραχθεί μια κοινότητα ή μια γειτονιά [27], όπως φαίνεται, για παράδειγμα, στα κοινωνικά δίκτυα. Με αυτό τον τρόπο, οι πληροφορίες που παρέχονται από τα δίκτυα της εμπιστοσύνης μπορεί να συνδυαστούν με πίνακες για τη δημιουργία συστημάτων συστάσεων

εμπιστοσύνης με επίγνωση [24], η οποία μπορεί να χειριστεί τα προβλήματα όπως η σπανιότητα των δεδομένων και τεχνητές ταυτότητες, πιο αποτελεσματικά.

Οι μέθοδοι CF μπορούν να ταξινομηθούν σε τρεις κύριες κατηγορίες ανάλογα με τα δεδομένα που χρησιμοποιούν για τον υπολογισμό συστάσεων: Α) μέθοδοι που βασίζονται σε μνήμη, οι οποίοι χρησιμοποιούν την πλήρη μήτρα με όλες τις αξιολογήσεις Β) μοντέλα με βάση τις μεθόδους, οι οποίες χρησιμοποιούν στατιστικά μοντέλα και τις λειτουργίες του πίνακα δεδομένων, αλλά όχι τα πλήρη στοιχεία της μήτρας και Γ) υβριδικές μεθόδους, οι οποίες συνδυάζουν την προηγούμενη μέθοδο με βάση το περιεχόμενο των μεθόδων σύστασης [6], οι οποίες και αναλύονται σε προηγούμενο κεφάλαιο.

Στη CF που βασίζεται στη μνήμη, οι συστάσεις της διατυπώνονται σε δύο στάδια που είναι: (i) Αναζήτηση γείτονα και (ii) πρόβλεψη σύστασης. Δεδομένου ενός χρήστη  $U_a$ , οι λειτουργίες συσχέτιση και απόσταση χρησιμοποιούνται για τον υπολογισμό γείτονα. Η πιο κοινή συσχέτιση και η απόσταση λειτουργίας που χρησιμοποιούνται είναι η συσχέτιση Pearson [5], η ομοιότητα συνημιτόνου [8] και η Ευκλείδεια απόσταση. Η ομοιότητα μεταξύ των χρηστών μπορεί επίσης να υπολογιστεί με ένα πολύ πιο αποτελεσματικό τρόπο, σύμφωνα με τη συμπεριφορά τους όταν ψηφίζουν. Τα παραδείγματα αυτών φαίνονται στο [12], όπου οι τάσεις των χρηστών υπολογίζονται, ή στο [13], όπου υπολογίζονται σύμφωνα με την προστασία προσωπικών δεδομένων. Μόλις καθοριστεί η γειτονιά του χρήστη  $U_a$ , οι συστάσεις μπορούν να υπολογιστούν με τη χρήση, για παράδειγμα, των μεθόδων που περιγράφονται στο [5], [14]. Αυτές οι μέθοδοι μπορούν να χρησιμοποιηθούν για την πρόβλεψη ψηφοφορίας ή να προτείνουν top-n αντικείμενα για τον χρήστη  $U_a$ .

Με βάση το μοντέλο CF και τις μεθόδους [2], [6], [12], [21], θα οικοδομηθεί ένα μοντέλο από το πλήρες πλέγμα επί του οποίου θα γίνουν προβλέψεις. Η εμφάνιση αυτών των μεθόδων δικαιολογείται από τους περιορισμούς ότι η μνήμη βασίζεται σε μεθόδους CF όσον αφορά την προσαρμοστικότητα, τις υπολογιστικές δαπάνες και την σπανιότητα.

Οι σύνθετες μέθοδοι CF συνδυάζουν τη μνήμη και τις modelbased μεθόδους, διατηρώντας τα πλεονεκτήματα των αλγορίθμων που εμπλέκονται, ενώ ταυτόχρονα ελαχιστοποιούν τα μειονεκτήματά και τις ανεπάρκειες τους. Παραδείγματα αυτών των μεθόδων είναι η Διάγνωση της προσωπικότητας [22] και το μοντέλο της πιθανοτικής μνήμης [23]. Οι υβριδικές μέθοδοι μπορούν επίσης να ληφθούν με το συνδυασμό των μεθόδων που



βασίζονται σε μνήμη και του μοντέλου που βασίζεται στο περιεχόμενο των συστημάτων συστάσεων. Ορισμένα γνωστά παραδείγματα [6] είναι: Filterbots, Content-booster, Fab and Ripper

#### **4.1 Η Κ – ανωνυμία με μικροσυσσωμάτωση**

Ανεξάρτητα από τη μέθοδο CF, υπάρχουν διάφοροι περιορισμοί σε αυτό το είδος των συστημάτων συστάσεων. Οι πιο σημαντικοί περιορισμοί [6], [24], [65], είναι η αραιότητα, η επεκτασιμότητα, cold start, η συνωνυμία, η δωροδοκία, οι επιθέσεις αντιγραφής του προφίλ και η έλλειψη της ιδιωτικής ζωής.

Μεταξύ όλων των προαναφερθέντων ανοικτών προβλημάτων, αυτή η εργασία θα επικεντρωθεί στην προστασία της ιδιωτικής ζωής των χρηστών που συμμετέχουν στις διαδικασίες της CF. Η κύρια έννοια αυτού του κεφαλαίου είναι η παρουσίαση μιας νέας μεθόδου με μικροσυσσωμάτωση με βάση την Διατήρηση Προστασίας Προσωπικών Δεδομένων Συνεργατικού Φιλτραρίσματος που εγγυάται Κ-ανωνυμία στους χρήστες. Η μέθοδος αυτή έχει αποδειχθεί ότι είναι πιο αποτελεσματική όσον αφορά την προστασία της ιδιωτικής ζωής και την απώλεια των πληροφοριών από μια χρησιμοποιούμενη ευρέως μέθοδο όπως η Gaussian με προσθήκη θορύβου. Επιπλέον, για τη βελτίωση της Gaussian με την προσθήκη του θορύβου, η παραπάνω μέθοδος παραμένει απλή, διευκολύνοντας την ευρεία υιοθέτησή της

#### **4.2 Το Υπόβαθρο**

##### **4.2.1 Η Προστασία των Προσωπικών Δεδομένων και η Διατήρηση του Συνεργατικού Φιλτραρίσματος**

Η ευρεία χρήση της μεθόδου CF στο Διαδίκτυο παρέχει μεγάλες ευκαιρίες για οφέλη τόσο στις επιχειρήσεις όσο και τους χρήστες. Ωστόσο, ένα σημαντικό μειονέκτημα που τίθεται είναι η έλλειψη της ιδιωτικής ζωής του χρήστη. Η σημασία της προστασίας της ιδιωτικής ζωής σε συστήματα CF τονίζεται από το αυξανόμενο ρυθμό με τον οποίο οι πληροφορίες του κάθε χρήστη συλλέγονται και αποθηκεύονται. Η απρόσεκτη διαχείριση των προσωπικών πληροφοριών, μπορεί να είναι παράνομη σε πολλές χώρες, αλλά επιπλέον έχει σοβαρές συνέπειες για τους χρήστες των οποίων γνωστοποιούνται πληροφορίες, καθώς και



για τις επιχειρήσεις. Ένα από τα κύρια προβλήματα στην CF είναι ότι οι πελάτες πιστεύουν ότι οι προτιμήσεις / προφίλ τους μπορεί να εκτεθούν, έτσι αποφασίζουν είτε να μην δώσουν την εκτίμησή τους για ένα συγκεκριμένο στοιχείο ή δεν το κάνουν σωστά ή το κάνουν ανακριβώς [28]. Ως εκ τούτου, η αίσθηση της έλλειψης της ιδιωτικής ζωής, οδηγεί σε μείωση τόσο στον αριθμό των αξιολογήσεων καθώς και στην ποιότητα τους. Η Προστασία Προσωπικών Δεδομένων Διατήρησης συνεργατικού φιλτραρίσματος (PPCF) αποσκοπεί στην επίλυση των θεμάτων προστασίας της ιδιωτικής ζωής που θέτει η συστηματική συλλογή των ιδιωτικών πληροφοριών, που απαιτούνται για την ορθή χρήση των μεθόδων CF.

Στις δυναμικές αγορές, οι εταιρείες μπορούν να ενδιαφέρονται να συνεργαστούν για να επιτύχουν καλύτερες συστάσεις για τους πελάτες τους. Λόγω της ιδιωτικής ζωής και τις ανησυχίες των επιχειρήσεων, τα δεδομένα δεν θα πρέπει να γνωστοποιούνται μεταξύ των εταιρειών. Σε αυτό το πλαίσιο, τα δεδομένα μπορούν να είναι καταναμημένα σε διάφορα κομμάτια με διάφορους τρόπους:

**Κάθετος διαχωρισμός (VP)** : στην οποία οι εταιρείες κατέχουν ασυνεχές σύνολα στοιχείων, αλλά με τους ίδιους χρήστες.

**Οριζόντιο διαχωρισμό (HP)** : στις οποίες κατέχουν διάφορα μέρη ξένα μεταξύ τους σύνολα των χρηστών για τις απόψεις των ίδιων ειδών.

**Αυθαίρετος διαχωρισμός (AP)** : στην οποία δεν υπάρχει κανόνας για το πώς διανέμονται τα δεδομένα [30]. Αν το σύνολο είναι καθορισμένο από  $m \times n$  αντικείμενα ανά χρήστη, ένα μέρος A κατέχει ένα υποσύνολο των χρηστών  $m_a \leq m$ , ενώ ένα άλλο μέρος B κατέχει το υπόλοιπο  $m_b = m - m_a$  και το ίδιο ισχύει και για τα υπόλοιπα στοιχεία.

Υπάρχουν διάφοροι τρόποι στη βιβλιογραφία για την προστασία της ιδιωτικής ζωής στις βάσεις δεδομένων. Για τα ανώνυμα προφίλ σε μεγάλες βάσεις δεδομένων μπορούμε χρησιμοποιήσουμε μεθόδους για την παροχή της K-ανωνυμίας [31], [66]. Εξαιτίας αυτών των ενδιαφερόντων ιδιοτήτων, το δημόσιο κλειδί πρέπει να κρυπτογραφείται [32], [33], με ασφαλείς πολυτμηματοποιημένους υπολογισμούς [34] και πρωτόκολλα κρυπτογράφησης που χρησιμοποιούνται συχνά στο διαδίκτυο και στα συστήματα ψηφοφορίας [67]. Άλλες μέθοδοι προσθέτουν θόρυβο στα δεδομένα για να νοθεύσει τις τιμές τους, με έναν τρόπο που επηρεάζει όσο το δυνατόν λιγότερο τη στατιστική ιδιότητα της μήτρας, όπως η μέση βαθμολογία των χρηστών και τα στοιχεία.

Σύμφωνα με το πώς οι πληροφορίες αποθηκεύονται και πώς οι συστάσεις υπολογίζονται, ταξινομούμε τις PPCF προσεγγίσεις σε κεντρικές ή αποκεντρωμένες. Θεωρούμε ότι μια PPCF μέθοδος είναι κεντρική αν χρησιμοποιεί ένα τρίτο μέρος για την εκτέλεση ενδιάμεσων υπολογισμών μεταξύ των χρηστών ή οντοτήτων. Μια μέθοδος θεωρείται επίσης ως κεντρική εάν αποθηκεύονται βαθμολογίες σε έναν server όπου οι συστάσεις και οι προβλέψεις υπολογίζονται. Περιπτώσεις στις οποίες μοιράστηκαν τα δεδομένα δεν θεωρούνται ως κεντρική διότι τα δεδομένα διανέμονται μεταξύ των διαφόρων μερών. Συνήθως, η κεντρική μέθοδος PPCF προσφέρει υψηλότερη απόδοση από τις αποκεντρωμένες ομολόγους της σε σχέση με την ομοιότητα και την πρόβλεψη υπολογισμών γιατί αποφεύγει πολλά γενικά έξοδα επικοινωνίας. Ωστόσο, στις κεντρικές μεθόδους, τα δεδομένα διαχειρίζονται από ένα μόνο τμήμα που έχει τον πλήρη έλεγχο πάνω τους, με διατήρηση των θεμάτων ιδιωτικής ζωής αν και τα δεδομένα έχουν ένα χαμηλό επίπεδο προστασίας. Οι περισσότερες κεντρικές μεθόδους PPCF προσθέτουν θόρυβο με διάφορους τρόπους για να διαταράξουν τα δεδομένα. Στο [35], [36], οι συγγραφείς προτείνουν τη διατάραξη των δεδομένων ακολουθώντας ένα στοιχείο αμετάβλητης Gaussian / ομοιόμορφης κατανομής του θορύβου. Στο itembased σύστημα η PPCF που προτείνει το Zhang et al. [37], οι διαταραχές προστίθενται ανάλογα με τη σημασία των αξιολογήσεων κατά τη διαδικασία της σύστασης. Ένας άλλος τρόπος για να διαταράξει τα δεδομένα παρουσιάζεται στο [38], όπου οι συγγραφείς προτείνουν μεταθέσεις αντικείμενων και γεωμετρικούς μετασχηματισμούς για να συσκοτίσουν τα δεδομένα.

Από την άλλη πλευρά, έχουμε τις αποκεντρωμένες μεθόδους PPCF. Αυτές οι μέθοδοι χρησιμοποιούν τα μέλη των κατανεμημένων δικτύων, θεωρούνται στις περισσότερες περιπτώσεις ως χρήστες, και διενεργεί ενδιάμεσους υπολογισμούς και προβλέψεις. Στα αποκεντρωμένα προγράμματα γενικά συνεπάγεται λιγότερη αποκάλυψη πληροφοριών από ότι η κεντρική μέθοδος στους ομολόγους της, αλλά συνεπάγεται τη χρήση των πρωτοκόλλων και πιο περίπλοκους υπολογισμούς. Συνήθως, στις αποκεντρωμένες PPCF μεθόδους, οι χρήστες αποθηκεύουν τις δικές τους αξιολογήσεις. Αυτό οδηγεί σε μια σειρά ελλείψεων, όπως η απαίτηση της ενεργούς συμμετοχής των χρηστών, η οποία είναι απαραίτητη για να μοιράζονται τα στοιχεία τους και να εκτελούν ενδιάμεσους υπολογισμούς. Εάν οι χρήστες δεν είναι ενεργοί, το ποσό των στοιχείων επί των οποίων εφαρμόζεται CF μειώνεται δραστικά αλλάζοντας την ακρίβεια των συστάσεων. Διάφορες προσεγγίσεις σε συστήματα βάσεων δεδομένων του καλαθιού της αγοράς [40], [41] έχουν προταθεί στην βιβλιογραφία. Αυτό το

είδος των βάσεων δεδομένων είναι κατάλληλο για να κάνει συστάσεις top-N με υψηλή ακρίβεια και το χαμηλό υπολογιστικό κόστος οφείλεται σε δυαδικές αξιολογήσεις του περιεχομένου του. Στο πλαίσιο αριθμητικής βαθμολογίας, έχουμε αρκετές προσεγγίσεις με στεγανά συστήματα δεδομένων, όπως αυτά που προτείνονται στο [45] [47] και [51]. Τέλος, τα συστήματα στα οποία οι χρήστες αποθηκεύουν τις δικές τους αξιολογήσεις μπορούν να βρεθούν στο [52], [54], [56], [42] και [59].

#### **4.2.2 Έλεγχος αποκάλυψης στατιστικών στοιχείων και μικροσυσσωμάτωση**

Ο στατιστικός έλεγχος αποκάλυψης [68], γνωστός και ως ανωνυμοποίηση των δεδομένων, επιδιώκει να μετατρέψει σειρές μικροδεδομένων (δηλαδή σύνολα δεδομένων που αποτελούνται από αρχεία που αντιστοιχούν σε μεμονωμένους ερωτηθέντες) πριν τη δημοσίευση, με τέτοιο τρόπο που δεν είναι δυνατόν να αναγνωρίσει εκ νέου ο εναγόμενος που αντιστοιχεί σε κάθε συγκεκριμένη εγγραφή σε ανώνυμες δημοσιεύσεις μικροστοιχείων του -identity αποκάλυψη- ούτε είναι δυνατόν να ανακαλύψουν την αξία του εμπιστευτικού χαρακτηριστικού (π.χ. μισθός) για ένα συγκεκριμένο εναγόμενο -attribute αποκάλυψη-.

Πριν από οποιαδήποτε διαδικασία της ανωνυμίας, τα άμεσα αναγνωριστικά (όνομα, ταυτότητα, κ.λπ.) πρέπει φυσικά να κατασταλούν από το σύνολο δεδομένων. Ωστόσο, μερικά από τα χαρακτηριστικά που παραμένουν στο ανώνυμο σύνολο δεδομένων μπορεί να είναι ημι-αναγνωριστικά, δηλαδή, χαρακτηριστικά τα οποία μπορούν να διευκολύνουν τον έμμεσο εκ νέου προσδιορισμό των ερωτώμενων από εξωτερικές πηγές δεδομένων (διαθέσιμο ως εισβαλλόμενο φόντο γνώσης) που συνδυάζουν αυτά τα χαρακτηριστικά με άμεσα αναγνωριστικά.

Η Μικροσυσσωμάτωση [64] είναι μια οικογένεια αλγορίθμων για ανωνυμοποίηση συνόλων δεδομένων και λειτουργεί σε δύο στάδια:

1) Το σύνολο των εγγραφών σε ένα σύνολο δεδομένων είναι συγκεντρωμένα έτσι ώστε:

i) για κάθε συστάδα περιέχει τουλάχιστον  $k$  εγγραφές .

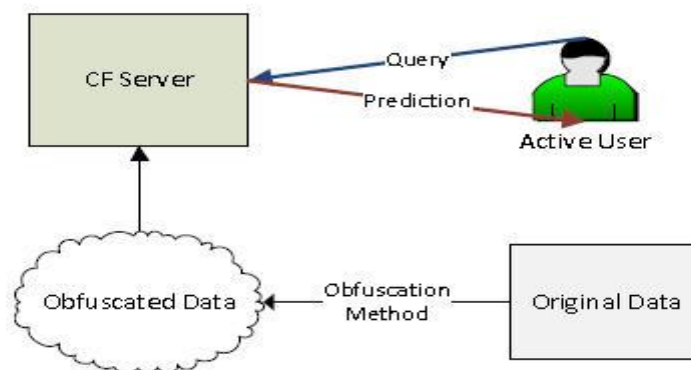
ii) τα αρχεία μέσα σε μια συστάδα είναι όσο το δυνατόν παρόμοια.

2) Τα αρχεία της κάθε ομάδας αντικαθίστουνται από έναν εκπρόσωπο του συμπλέγματος, συνήθως είναι η εγγραφή με το κέντρο βάρους (δηλαδή ο μέσος όρος του συμπλέγματος).

Όταν η μικροσυσσωμάτωση εφαρμόζεται στην πρόβλεψη των εγγράφων σχετικά με τα χαρακτηριστικά του σχεδόν- αναγνωριστικό τους, το προκύπτον σύνολο δεδομένων είναι k- ανώνυμο, δηλαδή, σε έναν εισβολέα κάθε εγγραφή στο σύνολο δεδομένων δεν μπορεί να διακριθεί σε μια ομάδα εγγραφών k το σχεδόν – αναγνωριστικό τους. Στο [69], περιγράφεται μια απλή μικροσυσσωμάτωση που ονομάζεται MDAV, στην οποία όλα τα συμπλέγματα έχουν ακριβώς τα αρχεία k, εκτός από το τελευταίο, το οποίο θα μπορούσε να έχει μεταξύ k και  $2k - 1$  εγγραφές. Πιο εξελιγμένες προσεγγίσεις μπορεί να βρεθεί στο [70] και [71].

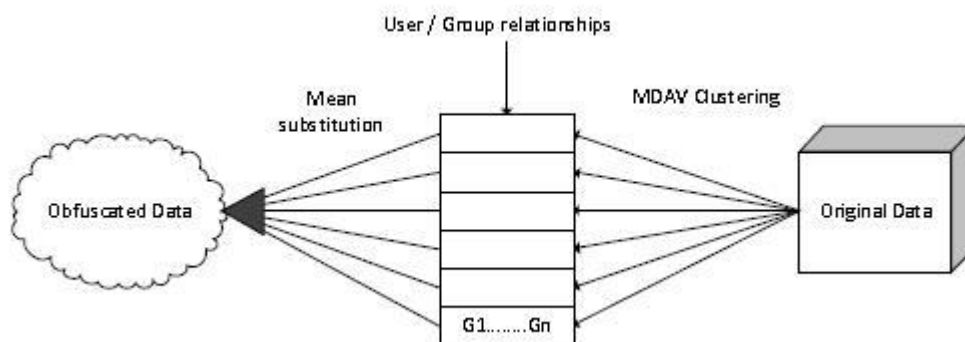
### 4.3 Η Προτεινόμενη Μέθοδος

Στην ενότητα αυτή, περιγράφουμε μία νέα προτεινόμενη μέθοδο. Το σύστημα αυτό θα μπορούσε να χαρακτηριστεί ως μια κεντρική μέθοδος PPCF. Η αρχιτεκτονική του φαίνεται στο Σχήμα 2.[73]



Σχήμα 4.2 Σχήμα Κεντρικής PPCF. Ο χρήστης κάνει μια αίτηση σε ένα στοιχείο στο διακομιστή, ο οποίος απαντά με μια εξατομικευμένη πρόβλεψη.

Αυτή η νέα προσέγγισή βασίζεται στην προαναφερθείσα μικροσυσσωμάτωση MDAV και απεικονίζεται στο Σχήμα 3. με ελαφρώς τροποποιημένη MDAV με τον τρόπο που να έχουν απομείνει αρχεία διαχείρισης. Ο αρχικός αλγόριθμος MDAV ορίζει ότι αν στο τέλος της διαδικασίας ομαδοποίησης υπάρχουν εγγραφές  $p$  μεταξύ  $k$   $p < 2k$  που δεν ανήκουν σε καμία ομάδα, θα πρέπει να αποτελούν μια τελική ομάδα  $C_f$  από μόνοι τους. Στην προσέγγισή αυτή, προκειμένου να μειωθεί η απώλεια πληροφοριών, πρέπει πρώτα να υπολογιστεί η μέση τιμή των  $C_f$ , που συμβολίζεται με  $M_f$  και υπολογίζουμε την απόσταση μεταξύ κάθε εγγραφής  $C_f$  με  $M_f$ . Μετά από αυτό, συγκρίνουμε την απόσταση ανάμεσα σε κάθε μέλος της  $C_f$  με όλες τις ήδη σχηματιζόμενες ομάδες. Αν περισσότερο από το ήμισυ των εγγραφών στην  $C_f$  είναι πιο κοντά στην  $M_f$  από ό,τι σε οποιαδήποτε άλλη ομάδα, τότε έχουμε σχηματίσει ένα τελικό σύμπλεγμα με τα στοιχεία  $C_f$ . αλλιώς, κάθε εγγραφή έχει ομαδοποιηθεί με την πλησιέστερη ομάδα μεταξύ εκείνων που έχουν ήδη σχηματιστεί.



**Σχήμα 4.3** Σχήμα ομαδοποίησης MDAV. Τα βήματα ροής από δεξιά (αρχικό σύνολο δεδομένων) προς τα αριστερά (ασαφές σύνολο δεδομένων).

Το σύστημα αυτό, που απεικονίζεται στο Σχήμα 4[73], λειτουργεί ως εξής:

1) Βεβαιωθείτε ότι το σύνολο δεδομένων δεν περιέχει τιμές που λείπουν για οποιαδήποτε ιδιότητα σε οποιοδήποτε εγγραφή. Αυτό είναι αναγκαίο προκειμένου να υπολογίσουμε την Ευκλείδεια απόσταση μεταξύ των εγγραφών. Η μεθόδους καταλογισμού [55], [62] ή μη εξατομικευμένων τιμών μπορούν να χρησιμοποιηθούν για να γεμίσουν τα κενά πεδία του συνόλου δεδομένων της μήτρας.

2) Όταν η μήτρα είναι εντελώς γεμάτη, υπολογίζουμε το Z-score της κάθε στήλης (στοιχείο) του συνόλου δεδομένων, χρησιμοποιώντας την ακόλουθη έκφραση

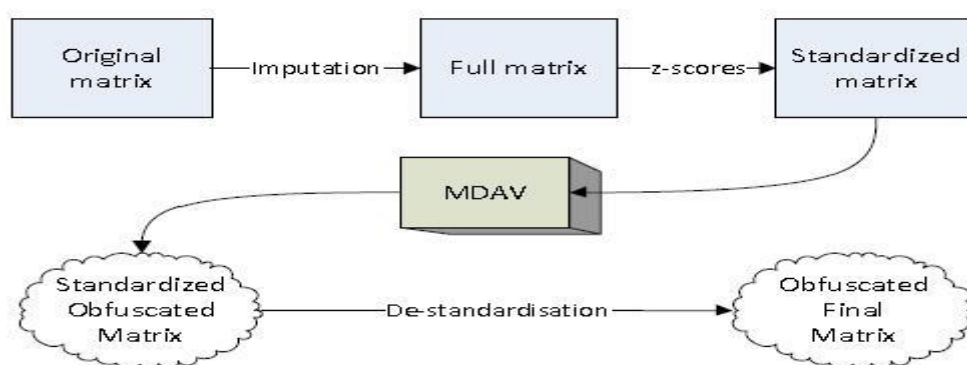
$$z\text{-score} = \frac{x_i - \mu}{\sigma}$$

όπου  $x_i$  είναι η  $i$ -οστή τιμή του στοιχείου  $x$  και  $\mu$  και  $\sigma$  είναι ο μέσος όρος και η τυπική απόκλιση του στοιχείου  $x$ , αντίστοιχα. Με τον τρόπο αυτό, η μέση και η τυπική απόκλιση του μετασχηματισμένου στοιχείου είναι 0 και 1, αντίστοιχα.

3) Όταν έχει τυποποιηθεί η μήτρα, είμαστε σε θέση να εφαρμόσουμε την ομαδοποίηση MDAV. Οι χρήστες θα ομαδοποιηθούν σε  $n$  συστάδες, με κάθε συστάδα  $C_i$  που αποτελείται από τις  $k$  όμοιους χρήστες, σύμφωνα με Ευκλείδεια απόσταση. Η τιμή του  $k$  υποδηλώνει το πλήθος της ομάδας. Με την επιλογή των πιο όμοιων χρηστών, μεγιστοποιούμε την ομοιογένεια της ομάδας και ως εκ τούτου μειώνουμε την απώλεια πληροφοριών. Μόλις οι σχέσεις προς ομάδας καθοριστούν, οι μέσες τιμές του κάθε  $C_i$ , συμβολίζεται ως  $M_i$ , και υπολογίζονται. Στη συνέχεια, κάθε τιμή του  $C_i$  αντικαθίσταται από το αντίστοιχο  $M_i$ .

4) Η διαδικασία ομαδοποίησης MDAV θα οδηγήσει σε ένα νέο σύνολο δεδομένων στο οποίο τα μέλη προς στο ίδιο σύμπλεγμα  $C_i$  θα έχουν τα ίδια χαρακτηριστικά. Ως εκ τούτου, μετά την εφαρμογή MDAV, αυτό το σύνολο δεδομένων θα ικανοποιήσει  $K$ -ανωνυμία.

5) Τελικά, σε περίπτωση που θέλουμε να γίνουν προβλέψεις, τα αποτελέσματα είναι μη-τυποποιημένα ώστε να ληφθεί το τελικό ασαφές σύνολο δεδομένων.



**Σχήμα 4.4** Προτεινόμενη βήμα προς βήμα μέθοδος



### 4.3.1 Τα Πειραματικά Αποτελέσματα

Στην ενότητα αυτή, δείχνουμε τα πειραματικά αποτελέσματα της μεθόδου αυτής και τα συγκρίνουμε με αυτά που λαμβάνονται από την ευρέως χρησιμοποιούμενη μέθοδος προσθήκης θορύβου Gaussian (GNA), η οποία χρησιμοποιεί μια κατανομή Gauss με μηδενική μέση τιμή και την τυπική απόκλιση  $\sigma$  (δηλαδή  $N(0, \sigma)$ ) για να διαταράξει το σύνολο δεδομένων.

Πρώτον, δείχνουμε τα αποτελέσματα που σχετίζονται με την προστασία της ιδιωτικής ζωής και το βοηθητικό πρόγραμμα που παρέχεται από τις αναλυτικές μεθόδους στο τμήμα IV-A. Εμείς τότε θα αξιολογήσει την ποιότητα των προβλέψεων του τμήματος IV-B.

Πειράματα με GNA επαναλήφθηκαν 50 φορές με κάθε αξιολόγηση. Όπως φαίνεται ήδη στην παραπάνω πρότασή (βλέπε Σχήμα 4), οι τιμές των δεδομένων έχουν τυποποιηθεί πριν η Gauss να προσθέσει το θόρυβο. Προκειμένου να ελεγχθεί η ποιότητα της μεθόδου, χρησιμοποιείται το γνωστό σύνολο δεδομένων Movielens. Το Movielens αναπτύχθηκε με GroupLens [5], και είναι ένα από τα σύνολα αναφοράς CF. Εδώ, θα επικεντρωθεί στην Movielens 100k, το οποίο περιέχει 100.000 αξιολογήσεις από 943 χρήστες σε 1.682 ταινίες. Η 100k Movielens έχει τιμές εύρους που περιλαμβάνεται μεταξύ 1 και 5. Αυτή η βάση δεδομένων είναι πολύ αραιή, δεδομένου ότι περισσότερο από το 90% των πεδίων είναι άδεια. Έχουμε καθιερώσει τη μέση τιμή των τιμών εύρους (3) για να συμπληρώσουμε τα κενά πεδία της μήτρας. Μόλις γεμίσει πλήρως, η μήτρα περιέχει συνολικά 1.586.126 τιμές.

### 4.3.2 Η Προστασία της ιδιωτικής ζωής

Προκειμένου να μετρηθεί η ποιότητα της ιδιωτικής ζωής που παρέχεται από μια μέθοδο διαταραχής θεωρούμε δύο παράγοντες, και συγκεκριμένα η απώλεια πληροφοριών και ο κίνδυνος αποκάλυψης. Η απώλεια πληροφοριών συνδέεται γενικά με το άθροισμα των τετραγώνων των σφαλμάτων (SSE). Το SSE χρησιμοποιείται συνήθως ως μέτρο παραμόρφωσης που εισάγεται στα πρωτότυπα δεδομένα. Στην ειδική περίπτωση της μικροσυσσωμάτωσης, το SSE υπολογίζεται με διανυσματική σημειογραφία ως εξής:

$$SSE = \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 \quad (4.1)$$

όπου  $g$  είναι ο αριθμός των υποσυνόλων / συστάδων που δημιουργούνται από το αλγόριθμο,  $n_i$  ο αριθμός των στοιχείων σε κάθε ομάδα, η  $x_{ij}$  διάνυσμα του  $j$ -οστού χρήστη της  $i$ -οστής ομάδας και  $\bar{x}_i$  είναι ο μέσος όρος του διανύσματος της  $i$ -οστής ομάδας. Γενικότερα, δεδομένου ενός πρωτότυπο σύνολο δεδομένων  $O$  αντιπροσωπεύεται από μία μήτρα του  $n \times m$  στοιχεία  $o_{ij}$  και διαστρεβλωμένο / προστατευμένο σύνολο δεδομένων  $P$  που παριστάνεται από μία μήτρα  $n \times m$  στοιχεία  $p_{ij}$ , το SSE υπολογίζεται ως εξής:

$$SSE = \sum_{i=1}^n \sum_{j=1}^m (o_{ij} - p_{ij})^2 \quad (4.2)$$

Ο κίνδυνος αποκάλυψης (DR) μετρά την πιθανότητα που σωστά είναι σχετικά με μια εγγραφή καλυμμένη / προστατευμένη σε μια μήτρα δεδομένων με μια καταγραφή της αρχικής μήτρας. Είναι επίσης γνωστή ως πιθανότητα της εκ νέου αναγνώρισης, ή τον κίνδυνο εκ νέου ταυτοποίησης.

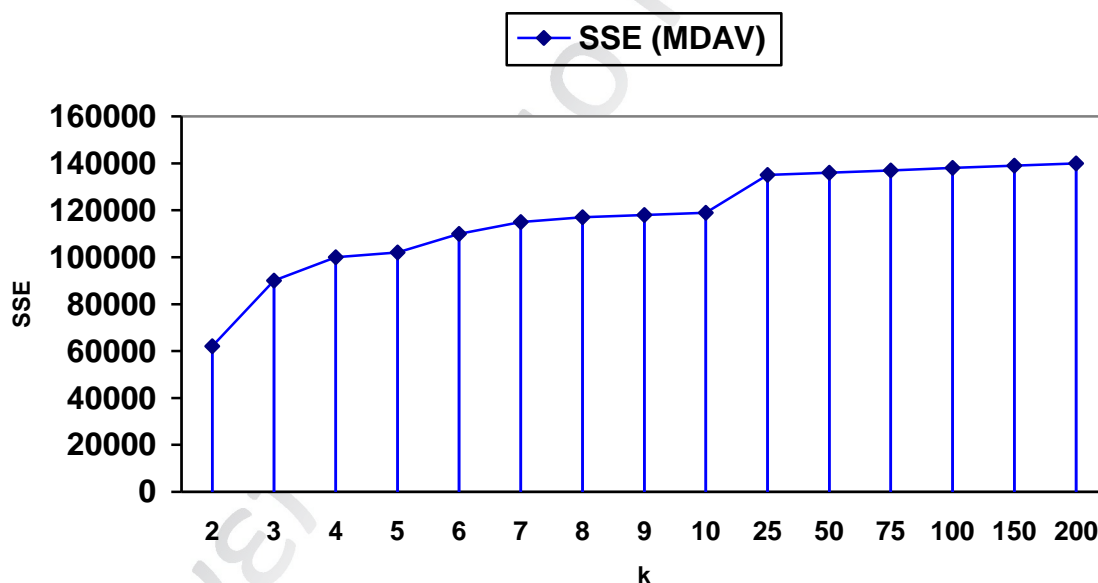
Για έναν εισβολέα, η διαδικασία της εκ νέου ταυτοποίησης συνίσταται τον υπολογισμό της απόστασης (π.χ. η Ευκλείδεια απόσταση) μεταξύ μιας δοσμένης προστατευόμενης εγγραφής  $p_i$  (που αντιστοιχεί στο  $i$  χρήστη), και σε μια εγγραφή στόχο  $o_j$ , που θα μπορούσαν να ληφθούν από τρίτες πηγές, όπως η απογραφή και άλλα παρόμοια. Στην περίπτωση μας, υποθέτουμε ότι το καλύτερο σενάριο είναι για έναν επιτιθέμενο (το σενάριο μαντείο), στο οποίο έχει το αρχικό σύνολο δεδομένων  $O$  και το διαστρεβλωμένο σύνολο δεδομένων  $P$  και ο ίδιος προσπαθεί να συνδέσει κάθε  $p_i$  εγγραφή του  $P$  με τα αρχεία στην  $o_j$  του  $O$ .

Για κάθε  $p_i$  εγγραφή στο  $P$  ο εισβολέας καθορίζει την πλησιέστερη εγγραφή  $o_j$  στο  $O$ . Αν καθοριστεί η πλησιέστερη εγγραφή  $o_j$  είναι στην πραγματικότητα η αυθεντική εγγραφή που ανήκει στην  $p_i$ , ο επιτιθέμενος τα έχει καταφέρει και λέμε ότι η  $p_i$  έχει ταχτοποιηθεί εκ νέου. Για τον υπολογισμό του κινδύνου αποκάλυψης, προσπαθούμε να αναγνωρίσουμε εκ νέου όλες τις εγγραφές και υπολογίζουμε το ποσοστό των σωστών των εκ νέου ταχτοποιημένων. Όσον αφορά την προστασία της ιδιωτικής ζωής και χρησιμότητα των δεδομένων, τόσο το SSE όσο και ο DR πρέπει να είναι σε χαμηλό επίπεδο.



Στους ακόλουθους πίνακες και σχήματα που δείχνουν τα αποτελέσματα του SSE και του DR για τις μεθόδους που αναλύθηκαν: η μέθοδος μικροσυσσωμάτωσης και η GNA. Ο πίνακας II δείχνει τα αποτελέσματα που ελήφθησαν με την προτεινόμενη μέθοδο μας για διαφορετικές τιμές του  $k$ , οι οποίες αντιπροσωπεύουν το πλήθος των στοιχείων του συνόλου των ομάδων, ενώ ο πίνακας III δείχνει τα αποτελέσματα που λαμβάνονται χρησιμοποιώντας την GNA με διαφορετικές τιμές του  $\sigma$ . Μπορεί να φανεί καθαρά ότι η σχέση μεταξύ SSE και DR είναι πολύ καλύτερη στη με βάση MDAV προσέγγιση.

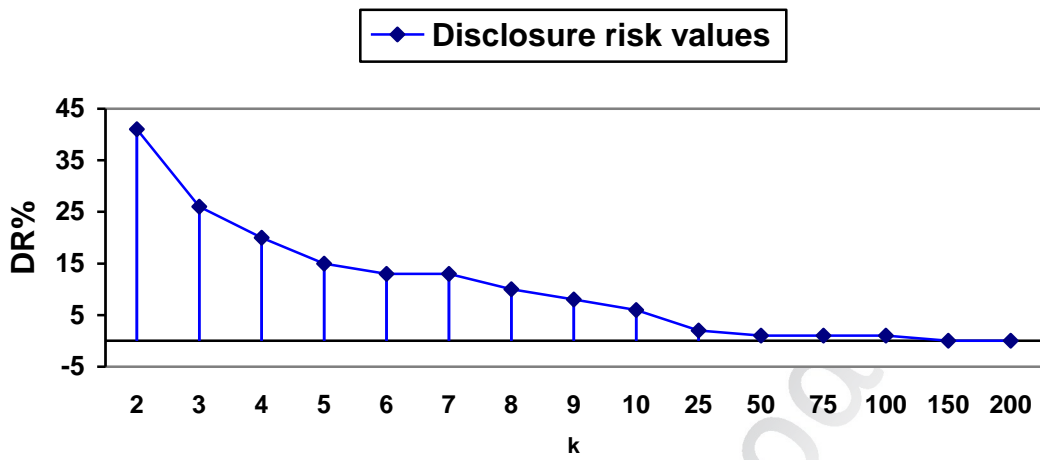
Το Σχήμα 5 και το Σχήμα 6 δείχνουν, αντίστοιχα, το SSE και το DR στη με βάση MDAV προσέγγιση του PPCF για διαφορετικές τιμές του  $k$ . Μπορεί να παρατηρηθεί ότι η συμπεριφορά τους είναι αρκετά ανταγωνιστική. Όταν το SSE αυξάνεται, μειώνεται αναλόγως και το DR.



Σχήμα 4.5 Οι SSE τιμές της μεθόδου μας για Monielens των 100k δεδομένων

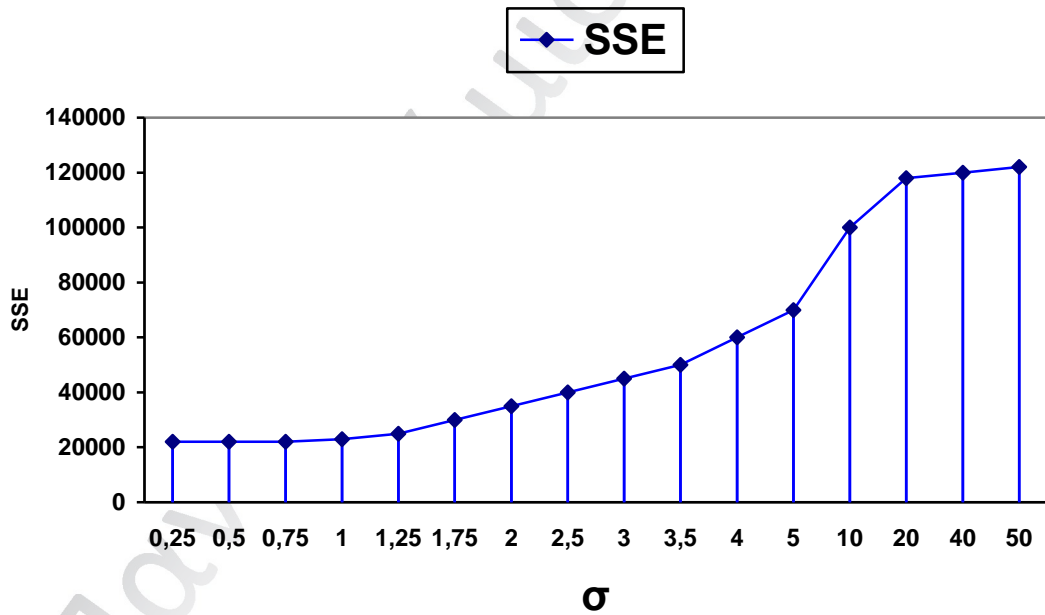
Στα Σχήμα 7 και Σχήμα 8 αντίστοιχως δείχνουν το SSE και DR για την προσέγγιση της GNA. Παρομοίως με την προσέγγιση που βασίζεται MDAV, όταν ο SSE μεγαλώνει το DR μειώνεται. Ωστόσο, η μέθοδος GNA πρέπει να προσθέσει πολλή περισσότερη παραμόρφωση στα δεδομένα (δηλαδή περισσότερο SSE) από την MDAV για την επίτευξη του ίδιου DR.

## Disclosure risk (MDAV)



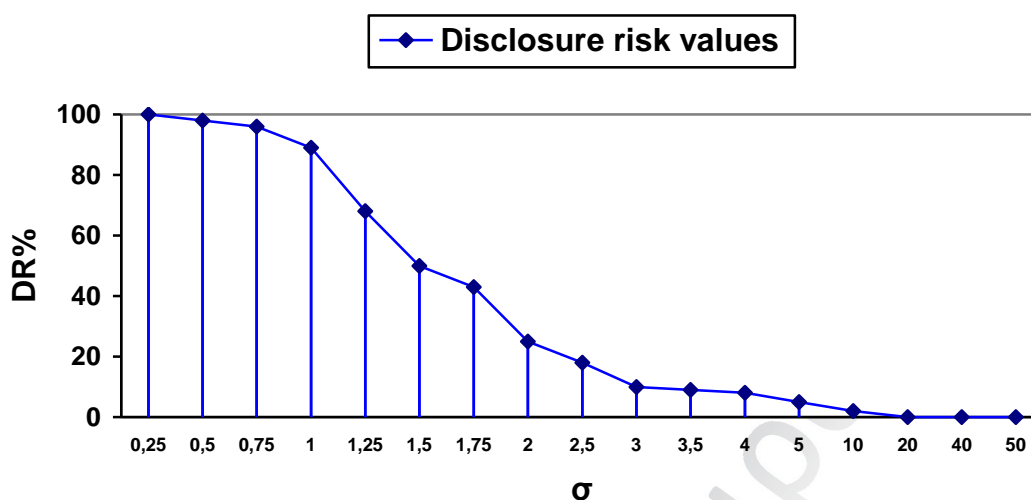
Σχήμα 4.6 Οι τιμές DR της μεθόδου μας, για Movielsens των 100k δεδομένων

## SSE (GNA)



Σχήμα 4.7 τιμές SSE της μεθόδου GNA για Movielsens των 100k δεδομένων

## Disclosure risk (GNA)



Σχήμα 4.8 Οι τιμές DR της μεθόδου GNA για Movielens των 100k δεδομένων

### 4.3.3 Η Προγνωστική ακρίβεια

Στην προηγούμενη ενότητα, έχουμε αναλύσει την ποιότητα της προστασίας της ιδιωτικής ζωής και της χρησιμότητας που παρέχεται από τη MDAV που βασίζεται και η GNA προσέγγιση. Ωστόσο, η ιδιωτικότητα και η χρησιμότητα είναι μόνο μια διάσταση του προβλήματος. Σημειώστε ότι τα προστατευμένα δεδομένα θα χρησιμοποιηθούν από εισηγητικά συστήματα για να γίνουν προβλέψεις σχετικά με τα στοιχεία που ένας χρήστης θα πρέπει να ενδιαφέρεται περισσότερο. Έτσι, είναι σημαντικό να μην μένουμε μόνο στην προστασία της ιδιωτικής ζωής, αλλά και να παρέχουμε ακριβείς προβλέψεις. Για να συγκρίνουμε την πρότασή μας με τη μέθοδο της GNA, έχουμε επιλέξει κωδικοποιημένα σύνολα δεδομένων με το ίδιο DR (π.χ. Στην περίπτωση αυτή, για λόγους απλότητας, έχουμε επιλέξει  $DR = 7,21\%$ , επειδή αντιστοιχεί στην τιμή που επιτυγχάνεται με  $k = 10$  για MDAV, και  $\sigma = 4$  χρησιμοποιώντας την GNA). Σημειώστε ότι κάθε άλλη τιμή DR θα μπορούσε να είχε επιλεγεί τόσο ώστε να είναι το ίδιο και για τις δύο μεθόδους.

Έχουμε ορίσει ένα σύνολο εκπαίδευσης με 80% των τιμών του στοιχείου και μια πρόβλεψη με το υπόλοιπο 20%. Οι προβλέψεις υπολογίζονται μόνο για τις αρχικές τιμές του κάθε χρήστη. Η πρόβλεψη των τιμών γίνεται σε δύο στάδια:

**Πίνακας 4.3 Αποτελέσματα του MDAV βασισμένα στο PPCF. Τα SSE αποτελέσματα εμφανίζονται στην κλίμακα  $10^3$**

ML 100k	MDAV														
k	2	3	4	5	6	7	8	9	10	25	50	75	100	150	200
SSE	64	87	99	105	110	114	117	119	120	180	184	186	186	188	189
DR%	40.82	26.51	19.93	15.9	12.19	12.19	9.65	7.95	7.21	2.33	0.63	0.21	0.21	0.1	0.1

**Πίνακας 4.4 Αποτελέσματα του GNA βασισμένα στο PPCF. Τα SSE αποτελέσματα εμφανίζονται στην κλίμακα  $10^3$**

ML 100k	GNA																
$\sigma$	0.25	0.5	0.75	1	1.25	1.5	1.75	2	2.5	3	3.5	4	5	10	20	40	50
SSE	246	248	257	275	302	336	376	418	509	592	663	727	830	1078	1221	1294	1889
DR%	100	100	98.51	89.28	68.5	50.58	44.53	27.99	18.76	10.49	8.58	7.21	4.24	1.4	0.42	0.31	0.1

- Εύρεση του πλησιέστερου γείτονα: Δεδομένου ενός  $u_i$  χρήστη για τον οποίο θέλουμε να προβλέψουμε κάποιες τιμές, θεωρούμε ένα σενετ εκπαίδευσης και βρίσκουμε τον πιο κοντινό γείτονά του, τον  $u_j$ .
- Εκχώρηση τιμής του / της: τις προβλεπόμενες τιμές για το χρήστη  $u_i$  είναι εκείνη που αντιστοιχεί σε  $u_j$  και αφορά την πρόβλεψη που έχει οριστεί.

Μόλις ολοκληρωθεί η πρόβλεψη για όλους τους χρήστες, υπολογίζουμε το σφάλμα μεταξύ των τιμών του αρχικού συνόλου δεδομένων (δηλαδή εκείνες τις τιμές δοκιμής που αποτελούν το 20%) και τις προβλεπόμενες τιμές. Για τον υπολογισμό του σφάλματος

εφαρμόζουμε το ευρέως χρησιμοποιούμενο μέσο για το απόλυτο σφάλμα (MAE), ορίζεται ως εξής:

$$MAE = \frac{\sum_{i=1}^n |p_i - r_i|}{n} \quad (4.3)$$

όπου  $n$  είναι ο αριθμός των προβλεπόμενων στοιχεία,  $P_i$  είναι η προβλεπόμενη τιμή που αφορά το στοιχείο  $i$  και  $r_i$  είναι η πραγματική αξία του  $i$ . Τα αποτελέσματα φαίνονται στον **Πίνακα 4.5**.

**Πίνακας 4.5** Το MAE περιέχει τιμές, συγκρίνοντας τις μήτρες πρόβλεψης με την αυθεντική βάση δεδομένων

Μέθοδος	MAE	%MAE
MDAV, $k=10$	0.89	22.25
GNA, $\sigma=4$	1.08	27

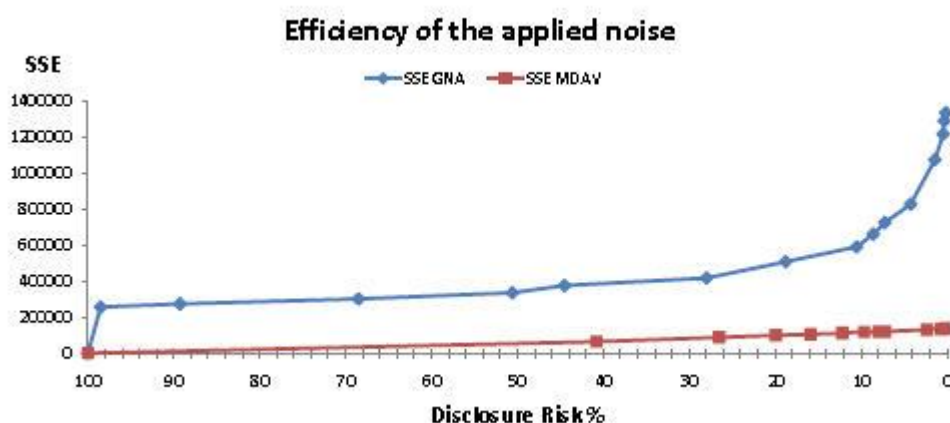
#### 4.4 Σύγκριση και συζήτηση

Στην προηγούμενη ενότητα παρουσιάστηκαν τα αποτελέσματα που προέκυψαν από την μικροσυσσωμάτωση με βάση την παραπάνω προσέγγισή και την κλασική GNA προσέγγιση. Σε αυτή την ενότητα θα συγκριθούν εν συντομία τα αποτελέσματα αυτά και θα δειχθεί ότι η MDAV προσέγγιση είναι ανώτερη, όσον αφορά τόσο την προστασία της ιδιωτικής ζωής όσο και την ακρίβεια πρόβλεψης.

Στο *Σχήμα 9*, μπορούμε να δούμε μια σύγκριση μεταξύ του SSE και τις τιμές DR και των δύο μεθόδων. Στον άξονα  $X$  εκπροσωπούμε το DR και στον άξονα  $y$  δείχνουμε το SSE. Αυτό το σχήμα μπορεί να χρησιμοποιηθεί για να διαβαστεί το ποσό του θορύβου, από την άποψη του SSE, που απαιτείται από κάθε μέθοδο για να επιτευχθεί ένα δεδομένο DR. Για παράδειγμα, μπορεί να παρατηρηθεί ότι, για μια σταθερή τιμή του  $DR = 30\%$ , η MDAV προσέγγιση εισάγει ένα σφάλμα της τάξης των 100K ενώ η GNA προσέγγιση απαιτεί 400K.

Η μικρότερη δυνατή τιμή του DR για την ανάλυση του συνόλου των δεδομένων είναι  $\frac{1}{943} \approx 0.1\%$ . Για να λάβει αυτή τη τιμή, η MDAV προσέγγιση θα πρέπει να σχηματιστούν ομάδες των  $k = 150$  στοιχεία, γεγονός που οδηγεί σε ένα SSE των 138.650. Αντιθέτως, η GNA αποκτά DR αξίας ίση με ένα SSE των 1.339.008, που είναι σχεδόν κατά μία τάξη μεγαλύτερου μεγέθους.

Αυτά τα αποτελέσματα δείχνουν σαφώς ότι η προτεινόμενη προσέγγιση διαταράσσει τα δεδομένα με ένα πολύ πιο αποτελεσματικό τρόπο. Επιπλέον, όπως έχει ήδη αναφερθεί, η μέθοδός μας εξασφαλίζει την προστασία της ιδιωτικής ζωής των χρηστών παρέχοντας Κ-ανωνυμία.



**Σχήμα 4.9** Σχέση μεταξύ της SSE και DR για τις μεθόδους που αναλύθηκαν και αφορούν την Movielen των 100k βάση δεδομένων

Όσον αφορά την ποιότητα της πρόβλεψης, ο Πίνακας IV δείχνει την ακρίβεια των προβλεπόμενων τιμών. Μπορεί να φανεί ότι όταν οι προβλέψεις γίνονται με βάση τα δεδομένα που προστατεύονται με MDAV, το MAE είναι 22,25% με τιμή DR 7,21%, που είναι ένα σημαντικό επίπεδο προστασίας της ιδιωτικής ζωής. Αντιθέτως, όταν οι τιμές προβλεφθούν με βάση τα δεδομένα που προστατεύονται από την GNA οδηγούμαστε σε σφάλμα του 27%, το οποίο είναι σχεδόν 5% υψηλότερο. Ως εκ τούτου, μπορούμε να συμπεράνουμε ότι τόσο η ποιότητα των προβλέψεων όσο και η ποιότητα της ιδιωτικής ζωής είναι καλύτερη στην προτεινόμενη μέθοδο που βασίζεται στην MDAV.

## **5. Συμπεράσματα**

### **5.0 Σύνοψη και Συμπεράσματα**

Το συνεργατικό φιλτράρισμα (CF) είναι ένα σύστημα συστάσεων που χρησιμοποιείται για να πραγματοποιηθούν αυτόματες συστάσεις προς τους χρήστες σε πολλαπλά πλαίσια. Παρά τα μεγάλα πλεονεκτήματα της χρήσης του CF, έχουμε υπογραμμίσει το σημαντικό αντίκτυπο που μπορεί να έχει για τους χρήστες στην προστασία της ιδιωτικής τους ζωής. Αν και ένα μεγάλο μέρος των μεθόδων CF προτείνεται, η μελέτη τους εξακολουθεί να είναι απαραίτητη και υπάρχουν πολλές προκλήσεις που πρέπει να ξεπεραστούν. Πιθανώς, η πιο σημαντική μεταξύ τους είναι η σωστή προστασία της ιδιωτικής ζωής των χρηστών.[73]

Σίγουρα, η προστασία της ιδιωτικής ζωής των χρηστών με την απόκρυψη όσον το δυνατόν περισσότερων προτιμήσεων τους, η ποιότητα των συστάσεων μειώνεται. Ως εκ τούτου, στο προηγούμενο κεφάλαιο, έχουμε παρουσιάσει τη μέθοδο PPCF με βάση την μικροσυσσωμάτωση. Τα αποτελέσματα που προέκυψαν κατά τη διάρκεια της αξιολόγησης της βάσης δεδομένων καταδεικνύουν ότι η παραπάνω μέθοδος διαταράσσει δεδομένα με ένα πολύ πιο αποδοτικό τρόπο από ό,τι άλλες γνωστές μεθόδους όπως η GNA. Επιπλέον, η πρότασή αυτή επιτυγχάνει K-ανωνυμία, η οποία αυξάνει την προστασία της ιδιωτικής ζωής των χρηστών με τέτοιο τρόπο ώστε η GNA δεν μπορεί να εγγυηθεί.[73]

### **5.1 Μελλοντικές επεκτάσεις**

Μελλοντικές εργασίες θα επικεντρωθούν σε δύο διαφορετικές κατευθύνσεις. Η πρώτη είναι να βελτιωθεί η αποτελεσματικότητα της μεθόδου αυτής, προκειμένου να μπορεί να εφαρμοστεί σε ένα αποκεντρωμένο σύστημα. Η δεύτερη κατεύθυνση είναι να αναλύσει την επίδραση των μεθόδων υπολογισμού που αφορούν την προστασία της ιδιωτικής ζωής και τις συστάσεις της ποιότητας, στη συνέχεια, τη μελέτη των δικτύων της εμπιστοσύνης και των αποτελεσματικών πολιτικών καταλογισμού.[73]

Πανεπιστήμιο Πειραιώς



## ***Βιβλιογραφία***

- [1] P. Resnick and H. Varian, “Recommender systems,” *Communications of the ACM*, vol. 40, no. 3, pp. 56–58, 1997. [Online]. Available: <http://dl.acm.org/citation.cfm?id=245121>
- [2] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, “Application of dimensionality reduction in recommender system-a case study,” *ACM Web KDD 2000 Web Mining for E Commerce Workshop*, vol. 1625, no. 1, pp. 264–8, 2000. [Online]. Available: <http://oai.dtic.mil/oai/oai?verb=getRecordn&metadataPrefix=htmln&identifier=ADA439541>
- [3] U. Fayyad and G. Piatetsky-Shapiro, “Advances in knowledge discovery and data mining,” AAI/MIT Press, 1996.
- [4] D. Goldberg, D. Nichols, B. M. Oki, and D. Terry, “Using collaborative filtering to weave an information tapestry,” *Communications of the ACM*, vol. 35, no. 12, pp. 61–70, 1992. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=138859.138867>
- [5] P. Resnick, N. Iacovou, and M. Suchak, “GroupLens: an open architecture for collaborative filtering of netnews,” *Proceedings of the ACM conference on Computer supported cooperative work CSCW*, vol. pp, no. 3, pp. 175–186, 1994. [Online]. Available: <http://dl.acm.org/citation.cfm?id=192905>
- [6] X. Su and T. M. Khoshgoftaar, “A Survey of Collaborative Filtering Techniques,” *Advances in Artificial Intelligence*, vol. 2009, no. Section 3, pp. 1–19, 2009. [Online]. Available: <http://www.hindawi.com/journals/aai/2009/421425/>
- [7] J. L. Herlocker, J. a. Konstan, L. G. Terveen, and J. T. Riedl, “Evaluating collaborative filtering recommender systems,” *ACM Transactions on Information Systems*, vol. 22, no. 1, pp. 5–53, Jan. 2004. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=963770.963772>
- [8] J. Breese, D. Heckerman, and C. Kadie, “Empirical Analysis of Predictive Algorithms for Collaborative Filtering,” *UAI 98*, pp. 43–52, 1998.
- [9] U. Shardanand, “Social information filtering for music recommendation,” Ph.D. dissertation, MIT, EECS Dept., 1994. [Online]. Available: [http://dl.acm.org/ftn\\_gateway.cfm?id=223931n&type=html](http://dl.acm.org/ftn_gateway.cfm?id=223931n&type=html)
- [10] U. Shardanand and P. Maes, “Social information filtering: algorithms for automating word of mouth,” *Proceedings of the SIGCHI conference on Human factors in computing systems CHI 95*, vol. 1, pp. 210–217, 1995. [Online]. Available: <http://dl.acm.org/citation.cfm?id=223931>

- [11] J. Herlocker and J. Konstan, "An algorithmic framework for performing collaborative filtering," Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval, vol. 54, no. 2, pp. 230–237, 1999. [Online]. Available: <http://portal.acm.org/citation.cfm?id=312682>
- [12] F. Cacheda, V. Carneiro, D. Fernandez, and V. Formoso, "Comparison of collaborative filtering algorithms," ACM Transactions on the Web, vol. 5, no. 1, pp. 1–33, Feb. 2011. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1921591.1921593>
- [13] N. Lathia, S. Hailes, and L. Capra, "Private Distributed Collaborative Filtering Using Estimated Concordance Measures," Proceedings of the 2007 ACM conference on Recommender systems RecSys 07, p. 1, 2007. [Online]. Available: <http://discovery.ucl.ac.uk/52802/>
- [14] B. Sarwar, G. Karypis, J. Konstan, and J. Reidl, "Item-based collaborative filtering recommendation algorithms," Proceedings of the 10th international conference on World Wide Web, vol. 1581133480, no. 15, pp. 285–295, 2001. [Online]. Available: <http://portal.acm.org/citation.cfm?id=371920.372071>
- [15] J. Wang, A. De Vries, and M. Reinders, "Unifying user-based and item-based collaborative filtering approaches by similarity fusion," Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval SIGIR 06, vol. 211, no. 1107, p. 501, 2006. [Online]. Available: <http://discovery.ucl.ac.uk/13465/>
- [16] S. A. Goldman and M. K. Warmuth, "Learning binary relations using weighted majority voting," Machine Learning, vol. 20, no. 3, pp. 245–271, 1995. [Online]. Available: <http://www.springerlink.com/index/10.1007/BF00994017>
- [17] Y. Koren and R. Bell, "Advances in collaborative filtering," Recommender Systems Handbook, pp. 43–52, 2011. [Online]. Available: <http://www.springerlink.com/index/X78806445324K172.pdf>
- [18] T. Hofmann, "Latent semantic models for collaborative filtering," ACM Transactions on Information Systems TOIS, vol. 22, no. 1, pp. 89–115, 2004. [Online]. Available: <http://portal.acm.org/citation.cfm?id=963774>
- [19] D. Lemire and A. Maclachlan, "Slope one predictors for online rating-based collaborative filtering," Society for Industrial Mathematics, vol. 05, no. 12, pp. 471–475, 2005.

- [20] S. Lloyd, "Least squares quantization in PCM," *IEEE Transactions on Information Theory*, vol. 28, no. 2, pp. 129–137, 1982. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1056489>
- [21] Z. Xia, Y. Dong, and G. Xing, "Support vector machines for collaborative filtering," in *Proceedings of the 44th annual Southeast regional conference*, ser. ACM-SE 44. New York, NY, USA: ACM, 2006, pp. 169–174. [Online]. Available: <http://doi.acm.org/10.1145/1185448.1185487>
- [22] D. Y. Pavlov and D. M. Pennock, "A Maximum Entropy Approach to Collaborative Filtering in Dynamic, Sparse, High-Dimensional Domains," *Advances in Neural Information Processing Systems* 15, vol. 15, no. 2/3, pp. 1441–1448, 2003.
- [23] A. Schwaighofer, V. Tresp, and H. Kriegel, "Probabilistic memory based collaborative filtering," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 1, pp. 56–69, Jan. 2004. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1264822>
- [24] P. Massa and P. Avesani, "Trust-aware collaborative filtering for recommender systems," *On the Move to Meaningful Internet Systems 2004: CoopIS, DOA, and ODBASE*, 492–508., vol. 3290, no. 8, pp. 492–508, 2004. [Online]. Available: <http://www.springerlink.com/index/8BAJ2BP1HATVFGKC.pdf>
- [25] —, "Trust metrics on controversial users: balancing between tyranny of the majority," *International Journal on Semantic Web and Information Systems*, pp. 1–21, 2007. [Online]. Available: <http://www.igi-global.com/article/trust-metrics-controversial-users/2830>
- [26] L. Page, S. Brin, R. Motwani, and T. Winograd, "The PageRank Citation Ranking: Bringing Order to the Web," *World Wide Web Internet And Web Information Systems*, vol. 54, no. 2, pp. 1–17, 1998. [Online]. Available: <http://ilpubs.stanford.edu:8090/422>
- [27] J. Golbeck, "FilmTrust: Movie Recommendations from Semantic Webbased Social Networks," in *ISWC2005 Posters Demonstrations*, 2006, pp. 1314–1315.
- [28] L. Cranor, J. Reagle, and M. Ackerman, "Beyond concern: Understanding net users' attitudes about online privacy," *The Internet Upheaval: Raising Questions, Seeking Answers in Communications Policy*, Tech. Rep., 2000.

[29] M. McPherson, L. Smith-Lovin, and J. M. Cook, "Birds of a Feather: Homophily in Social Networks," *Annual Review of Sociology*, vol. 27, no. 1, pp. 415–444, 2001. [Online]. Available: <http://arjournals.annualreviews.org/doi/abs/10.1146/annurev.soc.27.1.415>

[30] G. Jagannathan and R. Wright, "Privacy-preserving distributed kmeans clustering over arbitrarily partitioned data," *conference on Knowledge discovery in data*, 2005. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1081942>

[31] L. Sweeney, "k-anonymity: A model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-based Systems*, vol. 10, no. 5, pp. 557–570, 2002. [Online]. Available: <http://www.worldscientific.com/doi/abs/10.1142/S0218488502001648>

[32] P. Paillier, "Public-Key Cryptosystems Based on Composite Degree Residuosity Classes," *Advances in Cryptology EUROCRYPT 99*, vol. 1592, pp. 223–238, 1999. [Online]. Available: <http://www.springerlink.com/index/kwjvf0k8fqyy2h3d.pdf>

[33] T. ElGamal, "A public key cryptosystem and a signature scheme based on discrete logarithms," *IEEE Transactions on Information Theory*, vol. 31, no. 4, pp. 469–472, 1985. [Online].

Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1057074>

[34] A. C. Yao, "Protocols for secure computations," *23rd Annual Symposium on Foundations of Computer Science (sfcs 1982)*, pp. 160–164, Nov. 1982. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4568388>

[35] H. Polat, "Privacy-preserving collaborative filtering using randomized perturbation techniques," *Third IEEE International Conference on Data Mining*, pp. 625–628, 2003. [Online].

Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1250993>

[36] H. Polat and L. Hall, "SVD-based Collaborative Filtering with Privacy," *Proceedings of the 2005 ACM symposium on Applied computing SAC 05*, pp. 791–795, 2005.

[37] S. Zhang, J. Ford, and F. Makedon, "A privacy-preserving collaborative filtering scheme with two-way communication," *Proceedings of the 7th ACM conference on Electronic commerce - EC '06*, pp. 316–323, 2006. [Online].

Available: <http://portal.acm.org/citation.cfm?doid=1134707.1134742>

[38] R. Parameswaran and D. M. Blough, "Privacy Preserving Collaborative Filtering Using Data Obfuscation," 2007 IEEE International Conference on Granular Computing (GRC 2007), pp. 380–380, Nov. 2007. [Online].

Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4403128>

[39] R. Shokri, P. Pedarsani, G. Theodorakopoulos, and J.-P. Hubaux, "Preserving privacy in collaborative filtering through distributed aggregation of offline profiles," Proceedings of the third ACM conference on Recommender systems RecSys 09, p. 157, 2009. [Online].

Available: <http://portal.acm.org/citation.cfm?doid=1639714.1639741>

[40] H. Polat and W. Du, "Privacy-preserving top-N recommendation on distributed data," Journal of the American Society for Information Science, vol. 59, no. 7, pp. 1093–1108, 2008. [Online].

Available: <http://ejournals.ebsco.com/direct.asp?ArticleID=4530945EBCA169FD5C19>

[41] I. Yakut and H. Polat, "Estimating NBC-based recommendations on arbitrarily partitioned data with privacy," Knowledge-Based Systems, vol. 36, pp. 353–362, Dec. 2012. [Online].

Available: <http://linkinghub.elsevier.com/retrieve/pii/S0950705112002031>

[42] C. Kaleli and H. Polat, "P2P collaborative filtering with privacy," Turkish Journal of Electric Electrical Engineering and Compute Science, vol. 18, no. 1, pp. 101–116, 2010.

[43] M. Kantarcioglu and J. Vaidya, "Privacy preserving naive bayes classifier for horizontally partitioned data," In IEEE ICDM workshop on privacy preserving data mining, pp. 3–9, 2003. [Online]. Available: <http://www.cis.syr.edu/wedu/ppdm2003/papers/1.pdf>

[44] I. Yakut and H. Polat, "Privacy-preserving hybrid collaborative filtering on cross distributed data," Knowledge and Information Systems, vol. 30, no. 2, pp. 405–433, Apr. 2011. [Online]. Available: <http://www.springerlink.com/index/10.1007/s10115-011-0395-3>

[45] —, "Arbitrarily distributed data-based recommendations with privacy," Data & Knowledge Engineering, vol. 72, pp. 239–256, Feb. 2012. [Online]. Available: <http://linkinghub.elsevier.com/retrieve/pii/S0169023X11001467>

[46] H. Polat and W. Du, "Privacy-preserving collaborative filtering on vertically partitioned data," Knowledge Discovery in Databases: PKDD 2005, pp. 651–658, 2005.

[47] I. Yakut and H. Polat, "Privacy-Preserving Svd-Based Collaborative Filtering on Partitioned Data," International Journal of Information Technology & Decision Making, vol. 09, no. 03, pp. 473–

502, May 2010. [Online]. Available: <http://www.worldscientific.com/doi/abs/10.1142/S0219622010003919>

[48] N. Dokoohaki, C. Kaleli, H. Polat, and M. Matskin, "Achieving Optimal Privacy in Trust-Aware Social Recommender Systems," *Lecture Notes in Computer Science*, vol. 6430, no. Social Informatics, pp. 62–79, 2010. [Online]. Available: <http://www.springerlink.com/content/33817677756006q0>

[49] C. Kaleli and H. Polat, "Privacy-Preserving Trust-Based Recommendations on Vertically Distributed Data," pp. 376–379, 2011. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6061362>

[50] C.-L. Hsieh, J. Zhan, D. Zeng, and F. Wang, "Preserving Privacy in Joining Recommender Systems," pp. 561–566, 2008. [Online]. Available: <http://ieeexplore.ieee.org/xpl/freeabsnall.jsp?arnumber=4511628>

[51] J. Zhan, I. Wang, and C. Hsieh, "Towards efficient privacy-preserving collaborative recommender systems," *GrC 2008.*, pp. 778–783, 2008. [Online]. Available: <http://ieeexplore.ieee.org/xpls/absnall.jsp?arnumber=4664769>

[52] J. Canny, "Collaborative filtering with privacy," *Security and Privacy, 2002. Proceedings. 2002 IEEE*, pp. 45–57, 2002. [Online]. Available: <http://ieeexplore.ieee.org/xpls/absnall.jsp?arnumber=1004361>

[53] E. Polak., *Computational methods in optimization: a unified approach.*, New York, USA, 1971.

[54] J. Canny, "Collaborative Filtering with Privacy via Factor Analysis," *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*, vol. 14, pp. 238—245, 2002.

[55] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society Series BMethodological*, vol. 39, no. 1, pp. 1–38, 1977. [Online]. Available: <http://www.jstor.org/stable/10.2307/2984875>

[56] S. Berkovsky, F. Ricci, Y. Eytani, and T. Kuflik, "Enhancing Privacy and Preserving Accuracy of a Distributed Collaborative Filtering," *Proceedings of the 2007 ACM conference on Recommender systems RecSys 07*, pp. 9–16, 2007.



- [57] M. Tada, H. Kikuchi, and S. Puntheeranurak, "Privacy-Preserving Collaborative Filtering Protocol Based on Similarity between Items," pp. 573–578, 2010. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5474755>
- [58] H. Kikuchi, H. Kizawa, and M. Tada, "Privacy-Preserving Collaborative Filtering Schemes," pp. 911–916, 2009. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5066586>
- [59] A. Basu, J. Vaidya, H. Kikuchi, and T. Dimitrakos, "Privacy-preserving Collaborative Filtering for the Cloud," pp. 223–230, 2011. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6133147>
- [60] D. Bogdanov and R. Sassoon, "Privacy preserving collaborative filtering with Sharemind," Cybernetica research report, pp. T–4–2, 2008.
- [61] A. Shamir, "How to share a secret," Communications of the ACM, vol. 22, no. 11, pp. 612–613, 1979. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=359168.359176>
- [62] J. Montaquila and C. Ponikowski, "An evaluation of alternative imputation methods," In Proceedings of the Survey Research Methods, 1995.
- [63] J. Vaidya and C. Clifton, "Privacy-preserving outlier detection," pp. 233–240, 2004. [Online]. Available: <http://www.ncbi.nlm.nih.gov/pubmed/20595089>
- [64] J. Domingo-Ferrer and J. M. Mateo-Sanz, "Practical data-oriented microaggregation for statistical disclosure control," pp. 189–201, 2002. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=979982>
- [65] C. Kaleli and H. Polat, "Providing naive Bayesian classifier-based private recommendations on partitioned data," Knowledge Discovery in Databases: PKDD 2007, pp. 515–522, 2007. [Online]. Available: <http://www.springerlink.com/index/R37157975874308P.pdf>
- [66] J. Wang, Y. Luo, S. Jiang, and J. Le, "A Survey on Anonymity-Based Privacy Preserving," 2009 International Conference on EBusiness and Information System Security, pp. 1–4, 2009. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5137908>
- [67] B. Adida and D. Wikstrom, "How to shuffle in public," in Artificial Intelligence. Springer-Verlag, 2007, pp. 555–574. [Online].



Available:<http://www.springerlink.com/index/j6p730488x602r28.pdf>

- [68] A. Hundepool, J. Domingo-Ferrer, L. Franconi, S. Giessing, E. Schulte- Nordholt, K. Spicer, and P.-P. de Wolf, *Statistical Disclosure Control*. Wiley, 2012.
- [69] J. Domingo-Ferrer and V. Torra, "Ordinal, continuous and heterogeneous k-anonymity through microaggregation," *Data Mining and Knowledge Discovery*, vol. 11, no. 2, pp. 195–212, 2005.
- [70] A. Solanas and A. Martinez-Balleste, "V-MDAV: Variable group size multivariate microaggregation," in *COMPSTAT 2006*, 2006, pp. 917–925.
- [71] A. Solanas, A. Martinez-Balleste, and U. Gonzalez-Nicalas, "A variable- MDAV-based partitioning strategy to continuous multivariate microaggregation with genetic algorithms," in *International Joint Conference on Neural Networks(IJCNN)*, 2010, pp. 1–7.
- [72] Fran Casino, Constantinos Patsakis, Domenec Puig and Agusti Solanas "On Privacy Preserving Collaborative Filtering: Current Trends, Open Problems and New Issues "
- [73] Fran Casino, Josep Domingo-Ferrer, Constantinos Patsakis, Domenec Puig and Agusti Solanas "Privacy Preserving Collaborative Filtering with k-Anonymity through Microaggregation"
- [74] M. J. Pazzani, "A framework for collaborative, content-based and demographic filtering," *Artificial Intelligence Review*, vol. 13, no. 5-6, pp. 393–408, 1999.
- [75] T. Zhu, R. Greiner, and G. Haubl, "Learning a model of a web user's interests," in *Proceedings of the 9th International Conference on User Modeling (UM '03)*, vol. 2702, pp. 65–75, Johnstown, Pa, USA, June 2003.
- [76] M. Pazzani and D. Billsus, "Learning and revising user profiles: the identification of interesting web sites," *Machine Learning*, vol. 27, no. 3, pp. 313–331, 1997.
- [77] M. K. Condif, D. D. Lewis, D. Madigan, and C. Posse, "Bayesian mixed-effects models for recommender systems," in *Proceedings of ACM SIGIR Workshop of Recommender Systems: Algorithm and Evaluation*, 1999.
- [78] B. Krulwich, "Lifestyle finder: intelligent user profiling using large-scale demographic data," *Artificial Intelligence Magazine*, vol. 18, no. 2, pp. 37–45, 1997.

- [79] R. Burke, "Hybrid recommender systems: survey and experiments," *UserModelling and User-Adapted Interaction*, vol. 12, no. 4, pp. 331–370, 2002.
- [80] R. H. Guttman, *Merchant differentiation through integrative negotiation in agent-mediated electronic commerce*, M.S. thesis, School of Architecture and Planning, MIT, 1998.
- [81] M. Balabanovic and Y. Shoham, "Content-based, collaborative recommendation," *Communications of the ACM*, vol. 40, no. 3, pp. 66–72, 1997.
- [82] G. Adomavicius and A. Tuzhilin, "Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions," *IEEE Transactions on Knowledge and Data Engineering*, vol. 17, no. 6, pp. 734–749, 2005.
- [83] A. Ansari, S. Essegaiar, and R. Kohli, "Internet recommendation systems," *Journal of Marketing Research*, vol. 37, no. 3, pp. 363–375, 2000.
- [84] A. E. Gelfand and A. F. M. Smith, "Sampling-based approaches to calculating marginal densities," *Journal of the American Statistical Association*, vol. 85, pp. 398–409, 1990.
- [85] B. M. Sarwar, J. A. Konstan, A. Borchers, J. Herlocker, B. Miller, and J. Riedl, "Using filtering agents to improve prediction quality in the grouplens research collaborative filtering system," in *Proceedings of the ACM Conference on Computer Supported Cooperative Work (CSCW '98)*, pp. 345–354, Seattle, Wash, USA, 1998.
- [86] R. J. Mooney and L. Roy, "Content-based book recommendation using learning for text categorization," in *Proceedings of the Workshop on Recommender Systems: Algorithms and Evaluation (SIGIR '99)*, Berkeley, Calif, USA, 1999.
- [87] C. Basu, H. Hirsh, and W. Cohen, "Recommendation as classification: using social and content-based information in recommendation," in *Proceedings of the 15th National Conference on Artificial Intelligence (AAAI '98)*, pp. 714–720, Madison, Wis, USA, July 1998.

- [88] M. Claypool, A. Gokhale, T. Miranda, et al., “Combining content-based and collaborative filters in an online newspaper,” in *Proceedings of the SIGIR Workshop on Recommender Systems: Algorithms and Evaluation*, Berkeley, Calif, USA, 1999.
- [89] X. Su, R. Greiner, T. M. Khoshgoftaar, and X. Zhu, “Hybrid collaborative filtering algorithms using a mixture of experts,” in *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence (WI '07)*, pp. 645–649, Silicon Valley, Calif, USA, November 2007.
- [90] J. A. Delgado, *Agent-based information filtering and recommender systems on the internet*, Ph.D. thesis, Nagoya Institute of Technology, February 2000.
- [91] R. E. Schapire, “A brief introduction to boosting,” in *Proceedings of the 16th International Joint Conference on Artificial Intelligence (IJCAI '99)*, pp. 1401–1405, 1999.
- [92] B. Smyth and P. Cotter, “A personalized TV listings service for the digital TV age,” in *Proceedings of the 19th International Conference on Knowledge-Based Systems and Applied Artificial Intelligence (ES '00)*, vol. 13, pp. 53–59, Cambridge, UK, December 2000.
- [93] A. Popescul, L. H. Ungar, D. M. Pennock, and S. Lawrence, “Probabilistic models for unified collaborative and contentbased recommendation in sparse-data environments,” in *Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence (UAI '01)*, pp. 437–444, 2001.
- [94] D. M. Pennock, E. Horvitz, S. Lawrence, and C. L. Giles, “Collaborative filtering by personality diagnosis: a hybrid memory- and model-based approach,” in *Proceedings of the 16th Conference on Uncertainty in Artificial Intelligence (UAI '00)*, pp. 473–480, 2000.
- [95] D. Defays and P. Nanopoulos. Panels of enterprises and confidentiality: the small aggregates method. In Proc. of 92 Symposium on Design and Analysis of Longitudinal Surveys, pages 195–204, Ottawa, 1993. Statistics Canada.
- [96] J. Domingo-Ferrer, A. Martinez-Balleste, and J. M. Mateo-Sanz. Efficient multivariate data-oriented microaggregation. Manuscript, 2005.
- [97] J. Domingo-Ferrer, F. Sebe, and A. Solanas. A polynomial-time approximation to optimal multivariate microaggregation. Manuscript, 2005.

- [98] A. W. F. Edwards and L. L. Cavalli-Sforza. A method for cluster analysis. *Biometrics*, 21:362–375, 1965.
- [99] Economic Commission for Europe. Statistical data confidentiality in the transition countries: 2000/2001 winter survey. In *Joint ECE/Eurostat Work Session on Statistical Data Confidentiality*, 2001. Invited paper n.43.
- [100]. A. D. Gordon and J. T. Henderson. An algorithm for euclidean sum of squares classification. *Biometrics*, 33:355–362, 1977.
- [101] P. Hansen, B. Jaumard, and N. Mladenovic. Minimum sum of squares clustering in a low dimensional space. *Journal of Classification*, 15:37–55, 1998.
- [102] S. L. Hansen and S. Mukherjee. A polynomial algorithm for optimal univariate microaggregation. *IEEE Transactions on Knowledge and Data Engineering*, 15(4):1043–1044, July–August 2003.
- [103] A. Hundepool, A. Van de Wetering, R. Ramaswamy, L. Franconi, A. Capobianchi, P. DeWolf, J. Domingo-Ferrer, V. Torra, R. Brand, and S. Giessing.  $\mu$ -ARGUS version 4.0 Software and User’s Manual. Statistics Netherlands, Voorburg NL, may 2005. <http://neon.vb.cbs.nl/casc>.
- [104] M. Laszlo and S. Mukherjee. Minimum spanning tree partitioning algorithm for microaggregation. *IEEE Transactions on Knowledge and Data Engineering*, 17(7):902–911, 2005.
- [105] J. B. MacQueen. Some methods for classification and analysis of multivariate observations. In *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pages 281–297, 1967.
- [106] A. Oganian and J. Domingo-Ferrer. On the complexity of optimal microaggregation for statistical disclosure control. *Statistical Journal of the United Nations Economic Commission for Europe*, 18(4):345–354, 2001.
- [107] M. Rosemann. Erste ergebnisse von vergleichenden untersuchungen mit anonymisierten und nicht anonymisierten einzeldaten am beispiel der kostenstrukturerhebung und der

umsatzsteuerstatistik. In *G. Ronning and R. Gnoss (editors) Anonymisierung wirtschaftsstatistischer Einzeldaten, Wiesbaden: Statistisches Bundesamt*, pages 154–183, 2003.

[108] Agusti Solanas, Antoni Martinez-Balleste, Josep M. Mateo-Sanz, and Josep Domingo-Ferrer. Multivariate microaggregation based on a genetic algorithm. Manuscript, 2006.

[109] J. H. Ward. Hierarchical grouping to optimize an objective function. *Journal of the American Statistical Association*, 58:236–244, 1963.

[110] Χριστίνα Π. Χριστάκου Εισηγητικά Συστήματα Βασισμένα σε Μοντελοποίηση Προτιμήσεων Χρήστη και Μεθόδους Διήθησης της Πληροφορίας Εφαρμογή Επιλογής Κινηματογραφικών Ταινιών στο Διαδίκτυο, 2004.

Πανεπιστήμιο Πειραιώς