



**ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ**  
**ΤΜΗΜΑ ΨΗΦΙΑΚΩΝ ΣΥΣΤΗΜΑΤΩΝ**

**Δημιουργία μεθόδου λήψης απόφασης με χρήση  
αντίστροφων ερωτημάτων k κορυφαίων σημείων  
σε συστήματα πραγματικού χρόνου**

Διπλωματική διατριβή για το  
Π.Μ.Σ. «Διδακτική της Τεχνολογίας και Ψηφιακών Συστημάτων»  
του  
Νικητόπουλου Παναγιώτη

Επιβλέπων: Δουλκερίδης Χρήστος

Πειραιάς, Οκτώβριος 2014

Πανεπιστήμιο Πειραιώς

# Ευχαριστίες

Η παρούσα διπλωματική εργασία ολοκληρώνει το μεταπτυχιακό κύκλο σπουδών μου στο τμήμα Ψηφιακών Συστημάτων του Πανεπιστημίου Πειραιώς.

Οφείλω να ευχαριστήσω την οικογένειά μου, που με στήριξε κατά την διάρκεια των σπουδών μου και ήταν δίπλα μου σε κάθε ανάγκη.

Φυσικά ένα μεγάλο ευχαριστώ οφείλω και στο διδακτικό προσωπικό του τμήματος Ψηφιακών Συστημάτων του Πανεπιστημίου Πειραιώς και ιδιαίτερα στους καθηγητές μου Χρήστο Δουλκερίδη και Μαρίνο Θεμιστοκλέους, για το δύσκολο έργο που επιτελούν στον τομέα της εκπαίδευσης, και για τις πολύτιμες γνώσεις και συμβουλές, που θα με βοηθήσουν στη μετέπειτα σταδιοδρομία μου.

Πανεπιστήμιο Πειραιώς

Πανεπιστήμιο Πειραιώς

ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

## Περίληψη

Σχολή Τεχνολογιών Πληροφορικής και Επικοινωνιών  
Τμήμα Ψηφιακών Συστημάτων

Π.Μ.Σ. «Διδακτική της Τεχνολογίας και Ψηφιακών Συστημάτων»

Νικητόπουλος Παναγιώτης

Τα ερωτήματα  $k$  κορυφαίων σημείων, αποτελούν ένα πολύ χρήσιμο υποστηρικτικό εργαλείο στη διαδικασία λήψης αποφάσεων. Στην επιστήμη των βάσεων δεδομένων θεωρούνται από τα πιο σημαντικά ερωτήματα κατάταξης, αφού αναζητούν τα  $k$  πλήθους καλύτερα αντικείμενα βασιζόμενα στις προτιμήσεις του εκάστοτε ενδιαφερόμενου. Τα αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων αναζητούν τους ενδιαφερόμενους εκείνους που κρίνουν, βάσει των προτιμήσεών τους, ένα δοθέν αντικείμενο ως ένα από τα  $k$  πλήθους καλύτερα.

Στην παρούσα ερευνητική εργασία, γίνεται για πρώτη φορά μια προσπάθεια προσέγγισης του προβλήματος των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων σε συστήματα ροών δεδομένων. Προτείνεται ένα μοντέλο επίλυσης των ερωτημάτων τέτοιου τύπου, που περιλαμβάνει μια καινοτόμα αρχιτεκτονική συστήματος και τρεις πιθανούς αλγόριθμους που μπορούν να χρησιμοποιηθούν για την επίλυση του συγκεκριμένου προβλήματος. Οι προτεινόμενοι αλγόριθμοι συγκρίνονται μεταξύ τους μέσω εκτεταμένων εργαστηριακών πειραμάτων, όπου λαμβάνονται υπόψη (i) ο ρυθμός επεξεργασίας των εισερχόμενων δεδομένων, (ii) η υπολογιστική πολυπλοκότητα των αλγορίθμων και (iii) οι απαιτήσεις τους σε κύρια μνήμη. Από τα αποτελέσματα που συλλέχτηκαν, αναγνωρίζεται σημαντική βελτίωση στις επιδόσεις του συστήματος, υιοθετώντας τις προτεινόμενες βελτιώσεις.

**Θεματική περιοχή:** Αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων σε ροές δεδομένων

**Λέξεις κλειδιά:** ερώτημα  $k$  κορυφαίων σημείων, υπολογιστικό νέφος, πραγματικού χρόνου, μεγάλα δεδομένα, λήψη απόφασης

# Περιεχόμενα

	Σελ.
Ευχαριστίες	i
Περίληψη	iii
Περιεχόμενα	iv
Ευρετήριο σχημάτων	vi
Ευρετήριο ορισμών	vii
<b>1 Εισαγωγή</b>	<b>1</b>
1.1 Ορισμός του προβλήματος	2
1.1.1 Πολυκριτηριακή ανάλυση προβλημάτων	2
1.1.2 Νεφοϋπολογιστική	4
1.1.3 Ερωτήματα k κορυφαίων σημείων	5
1.2 Αντικείμενο και ερευνητικοί στόχοι	6
1.3 Διάρθρωση εγγράφου	6
<b>2 Θεωρητικό υπόβαθρο</b>	<b>9</b>
2.1 Κατανεμημένα συστήματα	10
2.1.1 Συστήματα υπολογιστικών συστοιχιών	11
2.1.2 Συστήματα υπερυπολογιστών	11
2.1.3 Συστήματα υπολογιστικών πλεγμάτων	11
2.1.4 Συστήματα υπολογιστικών νεφών	12
2.2 Συστήματα διαχείρισης μεγάλων δεδομένων	17
2.2.1 Διαχείριση συναλλαγών δεδομένων	17
2.2.2 Διαχείριση ανάλυσης δεδομένων	18
2.2.3 Στόχοι ανάλυσης μεγάλων δεδομένων	19
2.3 Συστήματα πραγματικού χρόνου	21
2.3.1 Χαρακτηριστικά συστημάτων πραγματικού χρόνου	22
2.3.2 Ταξινόμηση συστημάτων πραγματικού χρόνου	24
2.4 Πολυκριτηριακή ανάλυση αποφάσεων	25
2.4.1 Πολυκριτηριακή θεωρία αξίας ή χρησιμότητας	27
2.4.2 Προσέγγιση σχέσεων υπεροχής	28
2.4.3 Πολυκριτηριακός μαθηματικός προγραμματισμός	30
2.5 Ερωτήματα κορυφογραμμής	30
2.5.1 Ορισμοί	31

2.5.2	Χαρακτηριστικά κορυφογραμμής και κυριαρχίας Pareto	35
2.6	Ερωτήματα $k$ κορυφαίων σημείων	36
2.6.1	Κατηγοριοποίηση ερωτημάτων $k$ κορυφαίων σημείων	37
2.6.2	Αντίστροφα ερωτήματα $k$ κορυφαίων σημείων	38
<b>3</b>	<b>Ανασκόπηση βιβλιογραφίας</b>	<b>41</b>
3.1	Ερωτήματα κορυφογραμμών σε ροές δεδομένων	42
3.2	Ερωτήματα $k$ κορυφαίων σημείων σε ροές δεδομένων	44
3.2.1	Κατανομημένη διαχείριση ερωτημάτων	44
3.2.2	Αλγόριθμοι TMA και SMA	45
3.2.3	Διαχείριση ερωτημάτων σε αβέβαιες ροές	46
3.3	Αντίστροφα ερωτήματα $k$ κορυφαίων σημείων	48
3.3.1	Αλγόριθμος RTA	48
3.3.2	Αλγόριθμος branch and bound	50
3.3.3	Αλγόριθμοι DRT και DRT*	51
3.4	Ανοικτές ερευνητικές περιοχές	52
3.5	Συμπεράσματα	53
<b>4</b>	<b>Μοντέλο επίλυσης αντίστροφων ερωτημάτων <math>k</math> κορυφαίων σημείων</b>	<b>55</b>
4.1	Ορισμοί και σύμβολα	56
4.2	Ιδιότητες	58
4.3	Αρχιτεκτονική συστήματος	60
4.4	Αλγόριθμοι επεξεργασίας	62
4.4.1	Αλγόριθμος 1	62
4.4.2	Αλγόριθμος 2	65
4.4.3	Αλγόριθμος 3	70
<b>5</b>	<b>Πειραματική αξιολόγηση</b>	<b>75</b>
5.1	Παράμετροι αξιολόγησης	76
5.2	Οργάνωση πειραμάτων	76
5.3	Αποτελέσματα	80
5.3.1	Μεταβάλλοντας το πλήθος των διαστάσεων	80
5.3.2	Μεταβάλλοντας το μέγεθος του $Q$	81
5.3.3	Μεταβάλλοντας το μέγεθος του συρόμενου χρονικού παραθύρου	82
5.3.4	Μεταβάλλοντας το μέγεθος του $W$	84
5.4	Συμπεράσματα αξιολόγησης	85
<b>6</b>	<b>Συμπεράσματα και μελλοντική έρευνα</b>	<b>87</b>
6.1	Σύνοψη	87
6.2	Συμπεράσματα	89
6.3	Μελλοντική έρευνα	90
<b>A</b>	<b>Δείγμα πειραματικών μετρήσεων</b>	<b>93</b>
	<b>Βιβλιογραφία</b>	<b>95</b>

# Ευρετήριο σχημάτων

	Σελ.
1.1 Πρόβλεψη κίνησης δεδομένων διαδικτύου . . . . .	1
2.1 Κατανεμημένα συστήματα και τα υποσύνολά τους . . . . .	10
2.2 Διαφορές μεταξύ υπολογιστικού πλέγματος και υπολογιστικού νέφους . . . . .	12
2.3 Μοντέλα εφαρμογής της νεφοϋπολογιστικής . . . . .	15
2.4 Τεχνολογικές δυνατότητες της νεφοϋπολογιστικής . . . . .	16
2.5 Διαδικασία πολυκριτηριακής ανάλυσης . . . . .	26
2.6 Αναπαράσταση ενός συνόλου δεδομένων στο καρτεσιανό σύστημα . . . . .	33
2.7 Κορυφογραμμή ενός συνόλου δεδομένων στο καρτεσιανό σύστημα . . . . .	34
2.8 Παράδειγμα αντίστροφου ερωτήματος $k$ κορυφαίων σημείων . . . . .	40
4.1 Αναπαράσταση προτεινόμενου μοντέλου . . . . .	60
4.2 Περιοχές αμοιβαίας κυριαρχίας και αντι-κυριαρχίας . . . . .	65
4.3 Μέγιστη και ελάχιστη τιμή περιοχής πλέγματος . . . . .	66
5.1 Κατανομές δεδομένων . . . . .	78
5.2 Απόδοση του συστήματος μεταβάλλοντας τον αριθμό των νημάτων . . . . .	79
5.3 Μέσος ρυθμός επεξεργασίας (Διαστάσεις) . . . . .	80
5.4 Πολυπλοκότητα υπολογιστικής επεξεργασίας (Διαστάσεις) . . . . .	80
5.5 Μέγιστος απαιτούμενος χώρος κύριας μνήμης (Διαστάσεις) . . . . .	81
5.6 Μέσος ρυθμός επεξεργασίας ( $Q$ ) . . . . .	81
5.7 Πολυπλοκότητα υπολογιστικής επεξεργασίας ( $Q$ ) . . . . .	82
5.8 Μέγιστος απαιτούμενος χώρος κύριας μνήμης ( $Q$ ) . . . . .	82
5.9 Μέσος ρυθμός επεξεργασίας (Συρόμενο παράθυρο) . . . . .	82
5.10 Πολυπλοκότητα υπολογιστικής επεξεργασίας (Συρόμενο παράθυρο) . . . . .	83
5.11 Μέγιστος απαιτούμενος χώρος κύριας μνήμης (Συρόμενο παράθυρο) . . . . .	83
5.12 Μέσος ρυθμός επεξεργασίας ( $W$ ) . . . . .	84
5.13 Πολυπλοκότητα υπολογιστικής επεξεργασίας ( $W$ ) . . . . .	84
5.14 Μέγιστος απαιτούμενος χώρος κύριας μνήμης ( $W$ ) . . . . .	85



# Ευρετήριο ορισμών

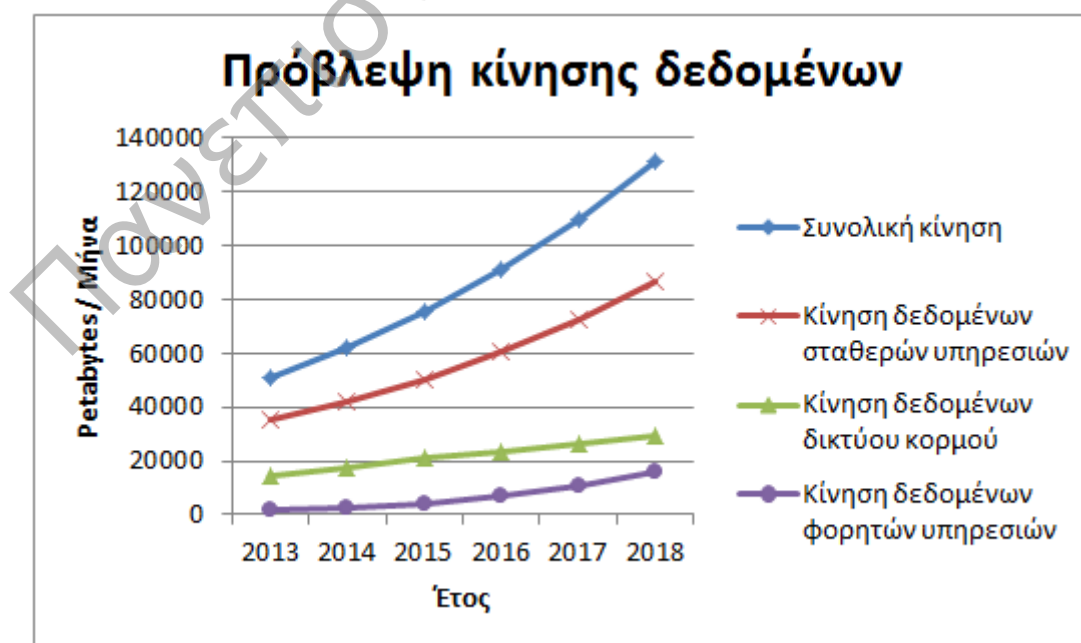
	Σελ.
Ορισμός 1.1 (Ανάλυση Μεγάλων Δεδομένων) . . . . .	2
Ορισμός 1.2 (Πολυκριτηριακή Ανάλυση Αποφάσεων) . . . . .	3
Ορισμός 1.3 (Μοντέλο Νεφοϋπολογιστικής) . . . . .	4
Ορισμός 2.1 (Κατανεμημένο Σύστημα) . . . . .	10
Ορισμός 2.2 (Σύστημα Υπολογιστικής Συστοιχίας) . . . . .	11
Ορισμός 2.3 (Σύστημα Υπερυπολογιστή) . . . . .	11
Ορισμός 2.4 (Σύστημα Υπολογιστικού Πλέγματος) . . . . .	11
Ορισμός 2.5 (Σύστημα Υπολογιστικού Νέφους) . . . . .	12
Ορισμός 2.6 (Σύστημα Πραγματικού Χρόνου) . . . . .	21
Ορισμός 2.7 (Κυριαρχία Pareto) . . . . .	31
Ορισμός 2.8 (Ερώτημα Κορυφογραμμής) . . . . .	32
Ορισμός 2.9 (Ερώτημα k Κορυφαίων Σημείων) . . . . .	36
Ορισμός 2.10 (Αντίστροφα Ερωτήματα k Κορυφαίων Σημείων) . . . . .	38
Ορισμός 4.1 (Περιοχή Αμοιβαίας Κυριαρχίας) . . . . .	65
Ορισμός 4.2 (Περιοχή Αμοιβαίας Αντι-Κυριαρχίας) . . . . .	65

Πανεπιστήμιο Πειραιώς

# Κεφάλαιο 1

## Εισαγωγή

Η εξάπλωση του διαδικτύου και της χρήσης διαφόρων φορητών δικτυακών συσκευών, οδήγησαν στην ανάγκη δημιουργίας νέας γενιάς εφαρμογών, οι οποίες θα είναι ικανές να παρέχουν αναβαθμισμένες υπηρεσίες στους χρήστες τους. Οι εφαρμογές αυτές παράγουν και επεξεργάζονται δεδομένα, το πλήθος των οποίων χρόνο με το χρόνο αυξάνεται σημαντικά. Σύμφωνα με έναν από τους μεγαλύτερους κατασκευαστές δικτυακών IP δρομολογητών, σε παγκόσμια βάση, προβλέπεται αύξηση της κίνησης δεδομένων του διαδικτύου κατά 20% μέσο όρο ετησίως ως το 2018 [1]. Ένα πολύ σημαντικό ποσοστό της αυξανόμενης κίνησης, αναμένεται για τα δεδομένα που κινούνται από και προς φορητές συσκευές (Σχήμα 1.1).



Σχήμα 1.1: Πρόβλεψη κίνησης δεδομένων διαδικτύου ως το 2018 [1]

Η ύπαρξη διαφόρων συσκευών που συνδέονται μεταξύ τους μέσω του διαδικτύου, δημιουργεί ένα νέο δίκτυο, που ονομάζεται Διαδίκτυο των Πραγμάτων - Internet of Things (IoT) [2]. Στην εποχή του Διαδικτύου των Πραγμάτων, τα δεδομένα που διακινούνται στο διαδίκτυο έχουν πολύ μεγάλο όγκο και μπορεί να διαφέρουν πολύ ως προς την ομοιογένειά τους, το περιεχόμενό τους, τη φύση τους και την προέλευσή τους. Από τα παραπάνω προκύπτει ότι η αποδοτική αποθήκευση, ανάκτηση και ανάλυσή των δεδομένων, αποτελεί πρόκληση για κάθε σύγχρονο πληροφοριακό σύστημα.

## 1.1 Ορισμός του προβλήματος

Ο ερευνητικός χώρος των βάσεων δεδομένων δεν μπορεί να μείνει ανεπηρέαστος από αυτές τις εξελίξεις, στο χώρο των σύγχρονων πληροφοριακών συστημάτων. Ο όρος Ανάλυση Μεγάλων Δεδομένων (Big Data Analytics) χρησιμοποιείται προκειμένου να περιγράψει αυτήν την τάση προς τη χρήση νέων εργαλείων και τεχνικών, για την αποδοτική διαχείριση πολύ μεγάλου όγκου δεδομένων. Πιο συγκεκριμένα:

### **Ορισμός 1.1 (Ανάλυση Μεγάλων Δεδομένων) :**

*Η ανάλυση μεγάλων δεδομένων είναι η διαδικασία εξέτασης μεγάλου όγκου διαφόρων τύπων δεδομένων, με σκοπό την αποκάλυψη κρυφών μοτίβων, άγνωστων συσχετίσεων και άλλων χρήσιμων πληροφοριών [3].*

Ένα πολύ σύνηθες χαρακτηριστικό των Μεγάλων Δεδομένων, αποτελεί η ύπαρξη πολλών διαστάσεων δεδομένων, οι οποίες μπορεί να μην είναι εξ αρχής γνωστές ή σαφώς καθορισμένες.

### 1.1.1 Πολυκριτηριακή ανάλυση προβλημάτων

Από τα πιο γνωστά προβλήματα της ύπαρξης πολλών διαστάσεων στα δεδομένα, είναι ότι αυτά καταλήγουν να απομακρύνονται πολύ το ένα από το άλλο. Το πρόβλημα αυτό είναι γνωστό και ως η κατάρα της διάστασης (curse of dimensionality) [4]. Αυτό έχει ως συνέπεια σε πολύ υψηλές διαστάσεις, η έννοια της ομοιότητας να καταρρέει, καθώς όλα τα σημεία απέχουν σημαντικά μεταξύ τους. Τα προβλήματα αυτά, σε συνδυασμό με την υπολογιστική πολυπλοκότητα που έχουν οι πράξεις στα πολυδιάστατα δεδομένα, κάνουν εμφανή την ανάγκη για ύπαρξη αποδοτικών μηχανισμών, υποστηριζόμενους από επαρκή επεξεργαστική ισχύ. Οι παραδοσιακές τεχνικές ανάλυσης δεδομένων, δεν μπορούν να έχουν εφαρμογή στο σύγχρονο χώρο της πολυκριτηριακής ανάλυσης αποφάσεων στα Μεγάλα Δεδομένα.

**Ορισμός 1.2 (Πολυκριτηριακή Ανάλυση Αποφάσεων) :**

Η Πολυκριτηριακή Ανάλυση Αποφάσεων (*Multicriteria Decision Making, MCDM*) αποτελεί κλάδο της Επιχειρησιακής Έρευνας (*Operational Research*), ο οποίος ασχολείται με την υποστήριξη της διαδικασίας λήψης αποφάσεων, λαμβάνοντας υπόψη περισσότερα του ενός κριτήρια [5].

Η διαδικασία λήψης απόφασης αποβλέπει στην επιλογή ενός υποσυνόλου βέλτιστων λύσεων, μέσα από ένα σύνολο εναλλακτικών επιλογών, οι οποίες από μαθηματική οπτική, μπορούν να θεωρηθούν εξίσου αποδεκτές. Οι εν λόγω λύσεις αποκαλούνται μη κυριαρχούμενες ή μη κατώτερες ή βέλτιστες Pareto. Η τελική λήψη της απόφασης, γίνεται συνήθως από κάποιον αποφασίζοντα (ανθρώπινος παράγοντας), ο οποίος συγκρίνει και αξιολογεί το υποσύνολο των βέλτιστων λύσεων, ώστε να εξαχθεί τελικά η καταλληλότερη λύση για το πρόβλημα. Ο τρόπος προσέγγισης του προβλήματος, μέσω της εισαγωγής περισσότερων του ενός κριτηρίων στη διαδικασία λήψης αποφάσεων, δίνει την δυνατότητα σε μια πιο ρεαλιστική επεικόνιση του πραγματικού προβλήματος, και άρα της καλύτερης προσέγγισης της βέλτιστης λύσης. Για το λόγο αυτό η Πολυκριτηριακή Ανάλυση Αποφάσεων, αποτελεί έναν από τους πιο δημοφιλείς και ταχύτερα αναπτυσσόμενους κλάδους της επιχειρησιακής έρευνας.

Η σημαντικότερη διαφορά της πολυκριτηριακής ανάλυσης από αντίστοιχες εναλλακτικές προσεγγίσεις, αποτελεί το γεγονός ότι η πολυκριτηριακή ανάλυση, δεν είναι απλή σύνθεση όλων των παραμέτρων ενός προβλήματος. Το βασικό χαρακτηριστικό γνώρισμα, που την διαφοροποιεί από τις υπόλοιπες, είναι η πραγματοποίηση της αναγκαίας σύνθεσης υπό το πρίσμα της πολιτικής λήψης των αποφάσεων και του συστήματος προτιμήσεων και αξιών, το οποίο συνειδητά ή ασυνείδητα χρησιμοποιεί ο αποφασίζων. Εφαρμογές της πολυκριτηριακής ανάλυσης υπάρχουν σχεδόν παντού, πχ: εφαρμογές που σχετίζονται με το νερό και, εν γένει, το περιβάλλον, από την αναπάρσταση των φυσικών διεργασιών μέχρι τον σχεδιασμό των υδραυλικών έργων και τη διαχείριση των υδροσυστημάτων.

Κατά κανόνα, τα πολυκριτηριακά προβλήματα αντιμετωπίζονται με τη μέθοδο του καθολικού μέτρου επίδοσης (*universal performance measure*), δηλαδή ενός βαθμωτού δείκτη που αποτιμά μονοσήμαντα την επίδοση του προβλήματος ως προς κάθε επικτική του λύση. Ο συγκεκριμένος δείκτης περιλαμβάνει είτε ένα και μόνο από τα κριτήρια, το οποίο θεωρείται μείζον, ή ένα συνδυασμό κριτηρίων, αυθαίρετα εντεταγμένων σε μια ενιαία μαθηματική έκφραση. Μια τέτοια προσέγγιση μπορεί, στη γενική της μορφή, να αποκρύψει σημαντικές πτυχές του προβλήματος και να οδηγήσει σε μια υποκειμενική ή μεροληπτική αντιμετώπισή του, με δεδομένο ότι, με κατάλληλη στάθμιση των διαφόρων κριτηρίων, μπορεί κάποιος να κατευθύνει εκ των προτέρων

τη διαδικασία βελτιστοποίησης προς μια υποκειμενική λύση. Έτσι τέτοιου είδους λύσεις προτείνονται μόνο σε περιπτώσεις που επιθυμείται η υποκειμενικότητα της εκάστοτε προτεινόμενης λύσης, με χρήση προτιμήσεων των χρηστών για τη δημιουργία της μαθηματικής συνάρτησης. Από τα μέσα της δεκαετίας του 1990, έχει ξεκινήσει η ανάπτυξη σχημάτων ταυτόχρονης αναζήτησης βέλτιστων Pareto λύσεων οι οποίες έκτοτε έχουν υποστεί σημαντικές βελτιώσεις. Ωστόσο, η έρευνα είναι ακόμα ανοιχτή.

### 1.1.2 Νεφοϋπολογιστική

Όπως αναφέρθηκε, η Πολυκριτηριακή Ανάλυση Αποφάσεων, λόγω της υπολογιστικής της πολυπλοκότητας, είναι ιδιαίτερα απαιτητική σε υπολογιστικούς πόρους, με αποτέλεσμα να απαιτείται υπολογιστική υποδομή που να μπορεί να υποστηρίξει αποτελεσματικά τη λειτουργία της.

#### **Ορισμός 1.3 (Μοντέλο Νεφοϋπολογιστικής) :**

*Η νεφοϋπολογιστική είναι ένα πολύ επιτυχημένο μοντέλο που επιτρέπει την εύκολη, συνεχή, και κατ' απαίτηση δικτυακή πρόσβαση σε ένα κοινόχρηστο παραμετροποιήσιμο σύνολο υπολογιστικών πόρων (πχ. δίκτυα, εξυπηρετητές, αποθηκευτικά μέσα, εφαρμογές, υπηρεσίες), το οποίο μπορεί πολύ γρήγορα να τροφοδοτηθεί και να διατεθεί με ελάχιστο κόστος διαχείρισης ή αλληλεπίδρασης με τον πάροχο της υπηρεσίας [6].*

Τα κύρια πλεονεκτήματα της χρήσης νεφοϋπολογιστικής αρχιτεκτονικής, είναι η ελαστικότητα, το μοντέλο πληρωμής με βάση τη χρήση, το χαμηλό κόστος εκκίνησης, η ευκολία πρόσβασης στην αγορά και η αποποίηση ευθυνών [7]. Η νεφοϋπολογιστική προσφέρεται παραδοσιακά στα εξής τρία μοντέλα: (i) Υποδομή υπό μορφή υπηρεσίας (IaaS), (ii) πλατφόρμα υπό μορφή υπηρεσίας (PaaS) και (iii) λογισμικό υπό μορφή υπηρεσίας (SaaS). Επιπλέον μπορεί να επεκταθεί και στην παροχή άλλων μοντέλων, όπως βάσεις δεδομένων υπό μορφή υπηρεσίας (DBaaS). Το τελευταίο μπορεί να εκμεταλλευτεί πολλά από τα χαρακτηριστικά και τις δυνατότητες που προσφέρονται σε ένα περιβάλλον νεφοϋπολογιστικής, όπως η επεκτασιμότητα και η κατανεμημένη διαχείριση των δεδομένων.

Ο ερευνητικός χώρος των βάσεων δεδομένων μελετά εδώ και τρεις δεκαετίες [7] τρόπους αποδοτικής και αποτελεσματικής σχεδίασης επεκτάσιμων, παράλληλων βάσεων δεδομένων. Τα αποτελέσματα της πολύχρονης μελέτης στον τομέα αυτό, περιλαμβάνουν πολλές προτάσεις και μοντέλα λύσεων: Το MapReduce [8] (με την ανοιχτού κώδικα υλοποίησή του, το Apache Hadoop) και το Apache Storm, αποτελούν τα τελευταία χρόνια τις πιο δημοφιλείς λύσεις σε συστήματα παράλληλων και κατανεμημένων βάσεων δεδομένων και ροών δεδομένων [9].

Οι απαιτήσεις στα σύγχρονα πληροφοριακά συστήματα, προϋποθέτουν τη συνεχή παρακολούθηση (monitoring) της ροής των δεδομένων, ώστε να είναι δυνατή η εξαγωγή συμπερασμάτων οποιαδήποτε στιγμή απαιτηθεί. Η διαδικασία αυτή είναι γνωστή υπό τον όρο ανάλυση δεδομένων σε πραγματικό χρόνο (real time analytics) [10]. Το πιο δημοφιλές εργαλείο διαχείρισης ροών δεδομένων σε πραγματικό χρόνο, αποτελεί αυτή τη στιγμή το Apache Storm [11]. Το εργαλείο αυτό μπορεί να χρησιμοποιηθεί για την υλοποίηση, τον έλεγχο και την αξιολόγηση της προτεινόμενης προσέγγισης των ερωτημάτων κορυφογραμμής [12].

Η υποδομή που απαιτείται για τη γρήγορη εκτέλεση ερωτημάτων (queries) στα σύγχρονα (state of the art) δημοφιλή μοντέλα καταναμημένων συστημάτων μεγάλων βάσεων δεδομένων (Hadoop, Storm), περιλαμβάνει την ύπαρξη ενός ή περισσότερων δικτυωμένων υπολογιστικών συμπλεγμάτων (clusters). Κάθε cluster μπορεί να αποτελείται εξ ολοκλήρου από υπολογιστικές μονάδες φθηνού υλικού (commodity hardware). Το μοντέλο της νεφοϋπολογιστικής, μπορεί να αποδώσει πολύ καλά στην αποθήκευση πολύ μεγάλου όγκου δεδομένων (Big Data Storage) και στην εφαρμογή μεθόδων πολυκριτηριακής ανάλυσης στα δεδομένα αυτά [12].

### 1.1.3 Ερωτήματα k κορυφαίων σημείων

Την τελευταία δεκαετία η υποστήριξη ερωτημάτων κατάταξης πλειάδων (ή σημείων) σε συστήματα βάσεων δεδομένων έχει τραβήξει το έντονο ενδιαφέρον της επιστημονικής κοινότητας [13–21]. Τα ερωτήματα k κορυφαίων σημείων επιτρέπουν την ανάκτηση k πλήθους πλειάδων, οι οποίες ταιριάζουν βέλτιστα στις προτιμήσεις ενός χρήστη, αποφεύγοντας με αυτόν τον τρόπο, την παραγωγή πολύ μεγάλου συνόλου αποτελεσμάτων. Έτσι τα ερωτήματα τέτοιου είδους, μπορούν να χρησιμοποιηθούν ως ένα εργαλείο υποστήριξης της διαδικασίας λήψης αποφάσεων [22].

Από τα παραπάνω προκύπτει ότι είναι πολύ σημαντικό για έναν κατασκευαστή, τα προϊόντα ή οι υπηρεσίες του να αποτελούν τμήμα των αποτελεσμάτων των k κορυφαίων ερωτημάτων [23]. Ωστόσο η παροχή προϊόντος ή υπηρεσίας που καλύπτει πλήρως τις προτιμήσεις και τα κριτήρια όλων των πιθανών αποφασιζόντων, είναι πολύ δύσκολη ίσως και αδύνατη. Έτσι δημιουργήθηκε η τεχνική αντίστροφων ερωτημάτων k κορυφαίων σημείων, η οποία επιτρέπει στους κατασκευαστές, να εντοπίσουν τις προτιμήσεις εκείνες, για τις οποίες το προϊόν ή η υπηρεσία τους αποτελεί τμήμα ενός ερωτήματος k κορυφαίων σημείων, συμβάλλοντας με αυτόν τον τρόπο στην αποτελεσματικότερη προώθησή τους [23].

## 1.2 Αντικείμενο και ερευνητικοί στόχοι

Αντικείμενο της παρούσας ερευνητικής εργασίας είναι (i) η μελέτη, αξιολόγηση και σύγκριση των υπάρχουσών προσεγγίσεων για τα ερωτήματα k κορυφαίων σημείων, και (ii) η προσπάθεια σχεδιασμού και ανάπτυξης μιας καινοτόμας πρότασης αποδοτικής εκτέλεσης των αντίστροφων ερωτημάτων k κορυφαίων σημείων, που να προσαρμόζεται και να ανταποκρίνεται στις ανάγκες των σύγχρονων πληροφοριακών συστημάτων, χρησιμοποιώντας state of the art εργαλεία διαχείρισης δεδομένων.

Μέσα από την επίτευξη των παραπάνω, μπορούν παράλληλα να εκπληρωθούν και οι εξής ερευνητικοί στόχοι:

**Ερευνητικός στόχος 1:** Μελέτη της βιβλιογραφίας, αξιολόγηση των υπάρχουσών λύσεων και σύγκριση των αναγκών που καλύπτει η κάθε μία.

**Ερευνητικός στόχος 2:** Προσδιορισμός και μελέτη ανοιχτών ερευνητικών θεμάτων.

**Ερευνητικός στόχος 3:** Σχεδιασμός και υλοποίηση καινοτόμας λύσης στο θέμα της αποδοτικής εκτέλεσης αντίστροφων ερωτημάτων k κορυφαίων σημείων.

**Ερευνητικός στόχος 4:** Αξιολόγηση της προτεινόμενης λύσης με πειραματική πραγμάτωση της.

**Ερευνητικός στόχος 5:** Βελτίωσή της σχεδιαζόμενης πρότασης με βάση τα δεδομένα που εξάγονται από τα πειραματικά αποτελέσματα.

**Ερευνητικός στόχος 6:** Προσδιορισμός πιθανών μελλοντικών επεκτάσεων και άλλων ανοιχτών ερευνητικών προβλημάτων.

## 1.3 Διάρθρωση εγγράφου

Το παρόν τεύχος αποτελείται από 7 κεφάλαια που καλύπτουν πλήρως την ανάπτυξη της ερευνητικής εργασίας.

**Στο Κεφάλαιο 1, Εισαγωγή:** ορίστηκε η θεματική περιοχή του ερευνητικού προβλήματος και τέθηκαν οι ερευνητικοί στόχοι για την επιτυχή εκπόνηση της παρούσας ερευνητικής εργασίας.

**Στο Κεφάλαιο 2, Θεωρητικό υπόβαθρο:** παρουσιάζεται το θεωρητικό υπόβαθρο των ερευνητικών περιοχών της παρούσας εργασίας. Συγκεκριμένα αναλύεται η



έννοια των κατανεμημένων συστημάτων και η εξέλιξη αυτών σε υποδομές υπολογιστικών νεφών, παρουσιάζονται οι προκλήσεις από την διαχείριση του μεγάλου όγκου των σύγχρονων βάσεων δεδομένων, περιγράφονται τα βασικά χαρακτηριστικά των συστημάτων πραγματικού χρόνου, γίνεται εισαγωγή στο χώρο των πολυκριτηριακών προβλημάτων, παρουσιάζονται οι έννοιες της κυριαρχίας Pareto και των ερωτημάτων κορυφογραμμής ως συναφή θέματα έρευνας και ορίζονται τα ερωτήματα (αντίστροφα και μη)  $k$  κορυφαίων σημείων αποδεικνύοντας τη χρησιμότητά τους στο πλαίσιο των πολυκριτηριακών προβλημάτων.

**Στο Κεφάλαιο 3, Ανασκόπηση βιβλιογραφίας:** παρουσιάζονται υπάρχουσες τεχνικές επίλυσης των ερωτημάτων (αντίστροφων και μη)  $k$  κορυφαίων σημείων και γίνεται εκτενής αναφορά στις τεχνικές εκείνες, που εκμεταλλεύονται την κατανεμημένη τοπολογία ώστε να αυξήσουν την αποδοτικότητά τους. Συγκρίνονται οι υπάρχουσες προσεγγίσεις και παρουσιάζονται ανοικτά ερευνητικά θέματα που υπάρχουν στο χώρο των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων.

**Στο Κεφάλαιο 4, Μοντέλο επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων:** παρουσιάζεται και αναλύεται το προτεινόμενο μοντέλο επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Προτείνονται επίσης τεχνικές και εργαλεία που επιτρέπουν την εφαρμογή του μοντέλου αυτού, σε συνθήκες συστημάτων πραγματικού χρόνου. Παρουσιάζονται και αναλύονται τρεις αλγόριθμοι επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων, οι οποίοι μπορούν να εκτελεστούν ακολουθώντας το προτεινόμενο μοντέλο.

**Στο Κεφάλαιο 5, Πειραματική αξιολόγηση:** παρουσιάζεται η διαδικασία αξιολόγησης των αλγορίθμων και του προτεινόμενου μοντέλου. Επιλέγονται τα κριτήρια αξιολόγησης και παρατίθεται η οργάνωση των πειραμάτων. Τα αποτελέσματα της διαδικασίας αξιολόγησης, παρατίθενται και αναλύονται.

**Στο Κεφάλαιο 6, Συμπεράσματα και μελλοντική έρευνα:** παρατίθενται τα συμπεράσματα από το σύνολο της διαδικασίας εκπόνησης της παρούσας ερευνητικής εργασίας. Αναφέρονται ανοικτά θέματα προς μελλοντική έρευνα και προτείνονται πιθανές τροποποιήσεις για την προτεινόμενη προσέγγιση.

Πανεπιστήμιο Πειραιώς

## Κεφάλαιο 2

# Θεωρητικό υπόβαθρο

Για την ορθότερη ερευνητική προσέγγιση της προτεινόμενης λύσης στην εκτέλεση των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων, απαραίτητη είναι η ανασκόπηση της υπάρχουσας βιβλιογραφίας που αφορά τα επιστημονικά πεδία που καλύπτει η συγκεκριμένη προσέγγιση. Η μελέτη αυτή συνεισφέρει σημαντικά στην εκπλήρωση των ερευνητικών στόχων που τέθηκαν στο Κεφάλαιο 1. Έτσι, σε αυτό το κεφάλαιο, γίνεται εισαγωγή στις βασικές έννοιες για την πραγματικού χρόνου εκτέλεση αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων, σε κατανομημένα συστήματα διαχείρισης Μεγάλων Δεδομένων.

Πιο συγκεκριμένα τα πεδία ερευνητικού ενδιαφέροντος της παρούσας εργασίας, είναι: (i) τα κατανομημένα συστήματα και τα συστήματα υπολογιστικών νεφών, (ii) η διαχείριση και ανάλυση Μεγάλων Δεδομένων, (iii) τα συστήματα πραγματικού χρόνου, (iv) τα πολυκριτηριακά προβλήματα, (v) τα ερωτήματα κορυφογραμμής και (vi) τα ερωτήματα  $k$  κορυφαίων σημείων. Για αυτές τις έννοιες, αναλύεται η σημασία, δικαιολογείται η χρησιμότητά τους και παρουσιάζεται το θεωρητικό υπόβαθρο στο οποίο στηρίζονται.

Στόχοι του κεφαλαίου αυτού είναι:

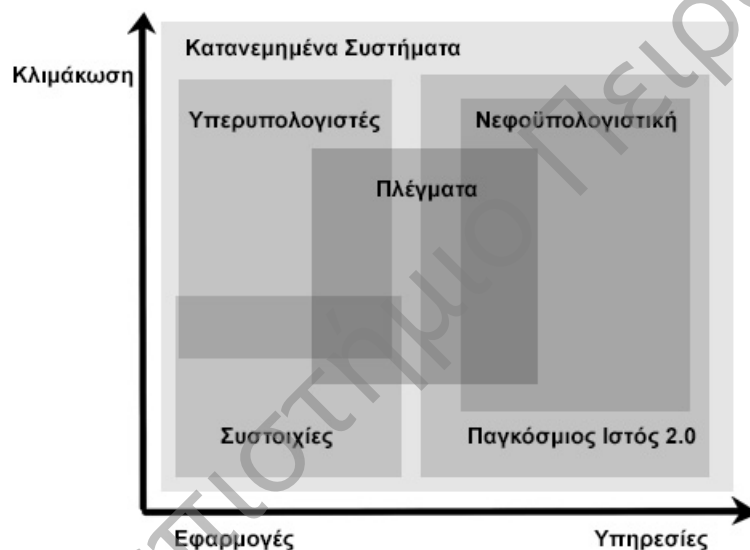
- Η εισαγωγή στις βασικές έννοιες των ερευνητικών περιοχών της παρούσας εργασίας.
- Η ιστορική αναδρομή και παρουσίαση της προέλευσης των εννοιών αυτών.
- Η ανάλυση του θεωρητικού υποβάθρου των επιστημονικών πεδίων που πραγματεύονται στη συνέχεια.

## 2.1 Κατανεμημένα συστήματα

### Ορισμός 2.1 (Κατανεμημένο Σύστημα) :

Τα κατανεμημένα υπολογιστικά συστήματα αποτελούν ένα σύνολο υπολογιστικών συστημάτων, τα οποία συνδέονται μεταξύ τους μέσω ενός ή περισσότερων υπολογιστικών δικτύων, με τρόπο τέτοιο, ώστε να εξυπηρετούν ένα συγκεκριμένο στόχο [24].

Η επικοινωνία των υπολογιστών αυτών, μπορεί να είναι ομοιογενής ή ετερογενής ενώ η κατανομή των υπολογιστικών συστημάτων μπορεί να είναι τοπική ή παγκόσμια. Τα κατανεμημένα συστήματα διακρίνονται σε υπερυπολογιστές (supercomputers), πλέγματα (grids), συστοιχίες (clusters), παγκόσμιο ιστό 2.0 (Web 2.0) και υπολογιστικά νέφη (clouds). Για καλύτερη κατανόηση, παρουσιάζεται σχηματικά η διασύνδεση των εννοιών αυτών στο Σχήμα 2.1.



Σχήμα 2.1: Κατανεμημένα συστήματα και τα υποσύνολά τους [24]

Ένα χαρακτηριστικό παράδειγμα κατανεμημένης υπολογιστικής, αποτελεί το έργο αναζήτησης εξωγήινης νοϋμοσύνης (Search for Extra-Terrestrial Intelligence, SETI) το οποίο βασίζεται σε ετερογενείς, παγκόσμια συνδεδεμένους υπολογιστές. Κάθε συμμετέχων έχει την δυνατότητα, μεταφορτώνοντας ένα μικρού μεγέθους λογισμικό, να συνδεθεί στους κεντρικούς εξυπηρετητές του έργου. Οι εξυπηρετητές αυτοί, παρέχουν στους συμμετέχοντες, τις καταγραφές των αισθητήρων τους. Τα δεδομένα των αισθητήρων έχουν πολύ μεγάλο όγκο συμβάλλοντας με αυτό τον τρόπο στην αύξηση της πολυπλοκότητας επεξεργασίας τους. Η επεξεργασία των δεδομένων αυτών, αποτελεί πρόκληση για κάθε αυτόνομο υπολογιστή, ενώ με τη συμμετοχή εκατομμυρίων υπολογιστικών συστημάτων, η διαδικασία γίνεται πιο εύκολη.

### 2.1.1 Συστήματα υπολογιστικών συστοιχιών

**Ορισμός 2.2 (Σύστημα Υπολογιστικής Συστοιχίας) :**

Οι υπολογιστικές συστοιχίες είναι κατανεμημένα συστήματα, το χαρακτηριστικό γνώρισμα των οποίων, αποτελεί το γεγονός ότι οι υπολογιστικές τους μονάδες είναι κατανεμημένες τοπικά και αποτελούνται από το ίδιο υλικό και λειτουργικό σύστημα. Έτσι μια διασυνδεδεμένη υπολογιστική συστοιχία μπορεί πιθανά να χρησιμοποιηθεί ή να νοηθεί σαν ένας υπερυπολογιστής [25].

Παράδειγμα μιας υπολογιστικής συστοιχίας, αποτελεί η διασύνδεση πολλών κονσολών Playstation 3, από την Αμερικανική Αεροπορία (US Air Force). Η συστοιχία αυτή χρησιμοποιείται για να υποστηρίξει υπολογιστικές πράξεις που απαιτούν υψηλής αποδοτικότητας υπολογιστικές μονάδες.

### 2.1.2 Συστήματα υπερυπολογιστών

**Ορισμός 2.3 (Σύστημα Υπερυπολογιστή) :**

Οι υπερυπολογιστές είναι κατανεμημένα συστήματα τοπικής κατανομής, που οι υπολογιστικές τους μονάδες αποτελούνται από παρόμοιο υλικό και λειτουργικό σύστημα και παρέχονται συγχωνευμένες σε ένα φυσικό πλαίσιο [25].

Η IBM κατασκευάζει υπερυπολογιστές αποτελούμενους από πολλούς επεξεργαστές, οι οποίοι παρέχονται υπό μορφή μιας υπολογιστικής μονάδας με πολύ υψηλές δυνατότητες επιδόσεων. Το μειονέκτημα χρήσης υπερυπολογιστών αποτελεί το γεγονός ότι συνήθως είναι πολύ ακριβοί και απαιτούν μεγάλα ποσά ενέργειας για τη λειτουργία τους.

### 2.1.3 Συστήματα υπολογιστικών πλεγμάτων

**Ορισμός 2.4 (Σύστημα Υπολογιστικού Πλέγματος) :**

Τα υπολογιστικά πλέγματα είναι κατανεμημένα συστήματα παγκόσμιας κατανομής, που πιθανά αποτελούνται από ανόμοια υποσυστήματα υλικού και λειτουργικού συστήματος [25].

Οι υπολογιστικές μονάδες των πλεγμάτων που διασυνδέονται μέσω του διαδικτύου, μπορούν να βρίσκονται οπουδήποτε ενώ συνήθως το κόστος διασύνδεσής τους είναι σημαντικά μικρό ή μηδαμινό. Είναι προφανές λοιπόν ότι τα πλέγματα δεν είναι τόσο κοστοβόρα όπως οι υπερυπολογιστές που περιγράφηκαν νωρίτερα. Το έργο SETI που περιγράφηκε νωρίτερα χρησιμοποιεί χαλαρά συνδεδεμένα πλέγματα. Κάθε τμήμα του

συγκεκριμένου συστήματος μπορεί να συνδέεται ή να αποσυνδέεται δυναμικά και να επεξεργάζεται δεδομένα συνεργατικά προς ένα κοινό σκοπό.

### 2.1.4 Συστήματα υπολογιστικών νεφών

#### Ορισμός 2.5 (Σύστημα Υπολογιστικού Νέφους) :

Τα υπολογιστικά νέφη είναι συστήματα παγκόσμιας κατανομής, των οποίων οι υπολογιστικοί πόροι παρέχονται υπό τη μορφή του μοντέλου της νεφοϋπολογιστικής (Ορισμός 1.3) [26].

Σε μια πρώτη ματιά, τα συστήματα υπολογιστικών νεφών μπορεί να φαίνονται παρόμοια με τα συστήματα υπολογιστικών πλεγμάτων, στην πραγματικότητα όμως, διαφέρουν αρκετά. Η σημαντικότερη διαφορά μεταξύ των δύο αυτών συστημάτων, αποτελεί το γεγονός ότι στα συστήματα υπολογιστικών νεφών, οι εξυπηρετητές αποτελούνται από εικονικούς πόρους, ενώ στα συστήματα υπολογιστικών πλεγμάτων, από φυσικούς. Μια πιο εκτεταμένη σύγκριση μεταξύ των δύο τεχνολογιών παρουσιάζεται στο Σχήμα 2.2

	Συστήματα υπολογιστικών πλεγμάτων	Συστήματα υπολογιστικών νεφών
<b>Χρησιμοποιούμενα μέσα</b>	Κατανομή μιας απλής διαδικασίας σε πολλαπλούς εξυπηρετητές	Εκτέλεση πολλών διεργασιών ταυτόχρονα σε έναν εξυπηρετητή, με χρήση εικονικών πόρων
<b>Τυπικό μοντέλο χρήσης</b>	Εκτέλεση διεργασιών σε ένα περιορισμένο χρονικό διάστημα	Τακτική χρήση για υπηρεσίες μακροχρόνιας υποστήριξης
<b>Επίπεδο χρήσης εικονικών μέσων</b>	Περιορισμένη χρήση εικονικών μέσων	Εκτεταμένη χρήση εικονικών μέσων

Σχήμα 2.2: Διαφορές μεταξύ υπολογιστικού πλέγματος και υπολογιστικού νέφους

Στη συνέχεια παρουσιάζονται τα πέντε βασικά χαρακτηριστικά των συστημάτων υπολογιστικών νεφών, σύμφωνα με τον ορισμό του μοντέλου της νεφοϋπολογιστικής κατά NIST [6]:

**Αυτοεξυπηρέτηση κατ' απαίτηση (On-demand self-service):** Οι υπολογιστικοί πόροι προσφέρονται στους καταναλωτές υπό μορφή υπηρεσίας, έπειτα από δική τους απαίτηση, χωρίς να είναι απαραίτητη η ενδιάμεση αλληλεπίδραση με τον πάροχο της υπηρεσίας.

**Ευρεία πρόσβαση στο δίκτυο (Broad network access):** Οι υπολογιστικοί πόροι είναι διαθέσιμοι μέσω του διαδικτύου και η πρόσβαση γίνεται μέσω τυποποιημένων μηχανισμών από συσκευές-πελάτες (π.χ. κινητά τηλέφωνα, ταμπλέτες, φορητοί υπολογιστές και σταθμούς εργασίας).

**Κοινή διάθεση πόρων (Resource pooling):** Η χρήση του μοντέλου πολλαπλών μισθωτών (multi-tenant), επιτρέπει στους παρόχους την εξυπηρέτηση πολλαπλών καταναλωτών, με την δυναμική ανάθεση υπολογιστικών πόρων σε αυτούς, σύμφωνα με τη ζήτηση που υπάρχει κάθε στιγμή.

**Ταχεία ελαστικότητα (Rapid elasticity):** Οι υπολογιστικοί πόροι μπορούν να δεσμεύονται και να αποδεσμεύονται ελαστικά και σε ορισμένες περιπτώσεις, αυτόματα, ώστε το υπολογιστικό νέφος να κλιμακώνεται ανάλογα με τη ζήτηση. Έτσι δίνεται η αίσθηση στον καταναλωτή ότι οι υπολογιστικοί πόροι είναι απεριόριστοι και μπορούν να διατεθούν σε οποιαδήποτε ποσότητα κι ανά πάσα στιγμή.

**Μετρήσιμα επίπεδα παροχής υπηρεσιών (Measured service):** Τα συστήματα υπολογιστικού νέφους ελέγχονται και βελτιστοποιούνται αυτόματα, αξιοποιώντας μια δυνατότητα μέτρησης σε κάποιο επίπεδο αφαίρεσης που είναι κατάλληλο για το είδος της υπηρεσίας (π.χ. αποθήκευση, επεξεργασία, χρήση δικτύου, ενεργοί λογαριασμοί χρηστών). Η χρήση των πόρων μπορεί να παρακολουθείται, να ελέγχεται, και να παρουσιάζεται υπό μορφή αναφορών, παρέχοντας διαφάνεια και ευκολία τόσο στον πάροχο όσο και στον καταναλωτή της υπηρεσίας.

Η λίστα των χαρακτηριστικών των υπολογιστικών συστημάτων κατά NIST, παρόλο που δεν είναι πλήρης, αποτελεί μια πολύ καλή αναφορά στα βασικά στοιχεία της νεφοϋπολογιστικής. Στη συνέχεια, για λόγους πληρότητας, παρουσιάζονται πρόσθετα χαρακτηριστικά, που δεν αποτελούν μέρος της λίστας του NIST:

**Κλιμάκωση (Scalability):** Στο υπολογιστικό νέφος είναι δυνατή η χειροκίνητη ή αυτόματη προσθήκη και αφαίρεση κόμβων για εκτέλεση εργασιών επεξεργασίας ή αποθήκευσης δεδομένων, ανάλογα με την αυξομείωση των απαιτήσεων, χωρίς να αλλοιώνεται η υπηρεσία που παρέχεται στους καταναλωτές.

**Εικονικοποίηση (Virtualization):** Οι λεπτομέρειες υλοποίησης και κατάστασης των υπολογιστικών πόρων αποκρύπτονται από τον καταναλωτή. Οι εικονικές μηχανές (Virtual Machines) δίνουν την εντύπωση μιας ολοκληρωμένης φυσικής υπολογιστικής μονάδας, ενώ στην πραγματικότητα είναι ένα σύνολο αρχείων και εφαρμογών που υπάρχουν και επεξεργάζονται σε μια άλλη (ή άλλες) φυσικές υπολογιστικές μονάδες. Τα τεχνικά χαρακτηριστικά των φυσικών υπολογιστικών μονάδων, μπορεί να διαφέρουν από αυτά των εικονικών, τόσο όσον αφορά το λογισμικό, όσο και το υλικό.

**Αξιοπιστία (Reliability):** Το μοντέλο του υπολογιστικού νέφους μπορεί να εγγυηθεί την ορθή λειτουργία των εφαρμογών των καταναλωτών, την μόνιμη αποθήκευση των δεδομένων τους και την αξιόπιστη δικτυακή μεταφορά τους. Για την εξασφάλιση αυτών, τα δεδομένα αντιγράφονται σε περισσότερες από μία υπολογιστικές μονάδες στο νέφος. Ο συνηθισμένος αριθμός αντιγράφων είναι 3 (replication level three).

**Συντήρηση (Maintenance):** Ανάλογα με το μοντέλο παροχής υπηρεσιών, η αποσφαλμάτωση των βασικών εφαρμογών και η αναβάθμιση των εκδόσεών τους, μπορεί να μην είναι πλέον ευθύνη του χρήστη, αλλά του παρόχου υπολογιστικού νέφους. Με τον τρόπο αυτό αποτελεί πλέον δικαιοδοσία του παρόχου να προσφέρονται συμβατές και αναβαθμισμένες εκδόσεις των βασικών εφαρμογών που είναι απαραίτητες για την ομαλή λειτουργία των εφαρμογών του καταναλωτή.

**Πολυελαστικότητα (Multi-tenancy):** Οι καταναλωτές των υπηρεσιών ενός υπολογιστικού νέφους δε χρειάζεται να έχουν αντίγραφα των εφαρμογών τους. Αρκεί ένα μοναδικό στιγμιότυπο (instance) το οποίο μπορεί να προσαρμοστεί στις ανάγκες του κάθε καταναλωτή. Αυτό έχει ως αποτέλεσμα την εξοικονόμηση πόρων στο νέφος και την ευκολότερη συντήρηση των εφαρμογών τους.

Ως μοντέλα πρακτικής εφαρμογής της νεφοϋπολογιστικής, όπως απεικονίζονται στο Σχήμα 2.3, μπορούν να θεωρηθούν τα ακόλουθα [6, 27, 28]:

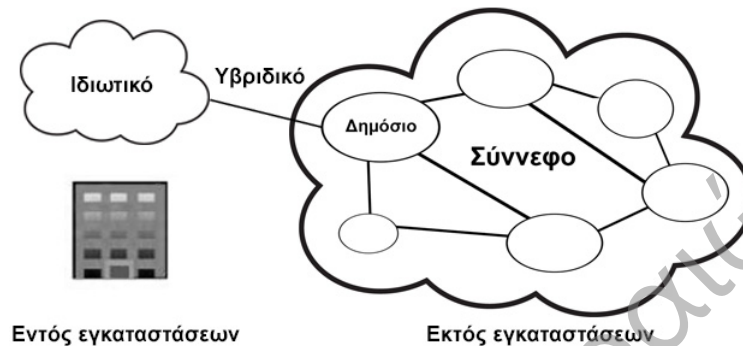
**Ιδιωτικό νέφος (Private cloud):** Η υποδομή υπολογιστικού νέφους λειτουργεί αποκλειστικά και μόνο για έναν οργανισμό. Η διαχειριστικές διαδικασίες μπορεί να εκτελούνται από τον ίδιο ή από τρίτους, ενώ η υποδομή μπορεί να βρίσκεται εντός ή εκτός των εγκαταστάσεών του.

**Κοινοτικό νέφος (Community cloud):** Η υποδομή υπολογιστικού νέφους μοιράζεται μεταξύ πολλών οργανισμών, υποστηρίζοντας μια συγκεκριμένη κοινότητα που έχει κοινές απαιτήσεις (π.χ. απαιτήσεις ασφαλείας, πολιτική και θέματα συμμόρφωσης). Η διαχειριστικές διαδικασίες μπορεί να εκτελούνται από τον ίδιο τον οργανισμό ή από τρίτους και η υποδομή μπορεί να βρίσκεται εντός ή εκτός των εγκαταστάσεων των οργανισμών που συμμετέχουν στην κοινότητα.

**Δημόσιο νέφος (Public cloud):** Η υποδομή υπολογιστικού νέφους διατίθεται στο ευρύ κοινό ή σε μια μεγάλη ομάδα εταιρειών και ανήκει σε έναν οργανισμό που πουλά υπηρεσίες νεφοϋπολογιστικής. Οι διαχειριστικές διαδικασίες εκτελούνται συνήθως από τον πάροχο του υπολογιστικού νέφους.



**Υβριδικό νέφος (Hybrid cloud):** Η υποδομή υπολογιστικού νέφους είναι μια σύνθεση από δύο ή περισσότερα μοντέλα εφαρμογής υπολογιστικών νεφών (ιδιωτικό, κοινοτικό ή δημόσιο), τα οποία παραμένουν μοναδικές οντότητες, αλλά συνδέονται μεταξύ τους με τυποποιημένη ή ιδιοκτήτη τεχνολογία που επιτρέπει τη φορητότητα δεδομένων και εφαρμογών (π.χ. εξισορρόπηση φόρτου εργασίας μεταξύ των υπολογιστικών νεφών).



Σχήμα 2.3: Μοντέλα εφαρμογής της νεφοϋπολογιστικής

Η νεφοϋπολογιστική περιλαμβάνει διάφορους τύπους υπηρεσιών. Υπάρχουν τρεις βασικές κατηγορίες των τεχνολογικών δυνατοτήτων που προσφέρονται ως υπηρεσία μέσω του μοντέλου της νεφοϋπολογιστικής, οι οποίες παρουσιάζονται σχηματικά στο Σχήμα 2.4.

**Υποδομή υπό μορφή υπηρεσίας (Infrastructure as a Service, IaaS):** Στατιστικά αποδεδειγμένα στοιχεία δείχνουν ότι το 80% της υπολογιστικής ισχύος δε χρησιμοποιείται αποδοτικά, ενώ το ίδιο συμβαίνει και με το 65% του αποθηκευτικού χώρου των εξυπηρετητών. Έτσι υπάρχει μεγάλη δυνατότητα κοινής χρήσης πόρων, με σκοπό την αποτελεσματικότερη, από πλευράς κόστους, χρήση τους [26].

Αντί της επένδυσης σε ιδιόκτητο εταιρικό εξυπηρετητή ή σε εύρος δικτύου, οι εταιρίες έχουν την δυνατότητα αγοράς των πόρων υπό μορφή χρέωσης ενοικίου και να της χρησιμοποιήσουν αυτών κατ' απαίτηση. Οι πάροχοι λαμβάνουν την ευθύνη παροχής και συντήρησης των εξυπηρετητών, των αποθηκευτικών μέσων και των ρυθμίσεων δικτύου ενώ ο καταναλωτής χρησιμοποιεί εικονικά τους πόρους αυτούς [25, 26].

**Πλατφόρμα υπό μορφή υπηρεσίας (Platform as a Service, PaaS):** Η πλατφόρμα υπό μορφή υπηρεσίας παρέχει τη δυνατότητα υποστήριξης όλου του κύκλου ζωής ανάπτυξης ενός λογισμικού: τον σχεδιασμό, την υλοποίηση, τη διόρθωση, τη δοκιμή, την εγκατάσταση, τη λειτουργία και την υποστήριξή του [25, 26].

Με τη χρήση του μοντέλου αυτού, το σύνολο ενός περιβάλλοντος λογισμικού μπορεί να βρίσκεται στον πάροχο υπηρεσιών, χωρίς να υπάρχει ανησυχία για

τις τεχνολογίες που χρησιμοποιούνται. Η μόνη ευθύνη που αναλαμβάνει ο καταναλωτής είναι ο σχεδιασμός και η δημιουργία των βάσεων δεδομένων και των εξατομικευμένων εφαρμογών του [26].

**Λογισμικό υπό μορφή υπηρεσίας (Software as a Service, SaaS):** Οι εταιρίες μπορούν να χρησιμοποιήσουν λογισμικό το οποίο είναι διαθέσιμο στο διαδίκτυο υπό μορφή χρέωσης ενοικίου, αντί να αγοράσουν το πλήρες πακέτο λογισμικού, επιτρέποντας στον καταναλωτή την κατάργηση του κινδύνου μη απόσβεσης της επένδυσής του. Σε αυτό το μοντέλο ο καταναλωτής δεν έχει την ευθύνη συντήρησης ή αναβάθμισης του παρεχόμενου λογισμικού. Την ευθύνη αυτή την αναλαμβάνει ο πάροχος της υπηρεσίας [26].

Τύπος νεφοϋπολογιστικής υπηρεσίας	Επιπτώσεις καταναλωτών
<b>Λογισμικό υπό μορφή υπηρεσίας</b>	<ul style="list-style-type: none"> <li>• Υπαρξη επιπέδου αφαίρεσης για τη λογική των εφαρμογών, της πλατφόρμας και της υποδομής.</li> <li>• Σημαντική μείωση δυσκολίας και κόστους για την εκτέλεση, την διαχείριση και την δημοσιοποίηση εφαρμογών.</li> <li>• Οι εφαρμογές μπορεί να ρυθμιστούν αλλά δεν μπορούν να καλύψουν πλήρως τις ανάγκες του καταναλωτή.</li> </ul>
<b>Πλατφόρμα υπό μορφή υπηρεσίας</b>	<ul style="list-style-type: none"> <li>• Υπαρξη επιπέδου αφαίρεσης για τη λογική της πλατφόρμας και της υποδομής.</li> <li>• Εξειδικευμένες εφαρμογές μπορούν να δημιουργηθούν πολύ γρήγορα και φθηνά.</li> <li>• Οι εξειδικευμένες εφαρμογές πρέπει να συντηρούνται και να υποστηρίζονται με την ευθύνη του καταναλωτή.</li> </ul>
<b>Υποδομή υπό μορφή υπηρεσίας</b>	<ul style="list-style-type: none"> <li>• Υπαρξη επιπέδου αφαίρεσης για την υπολογιστική υποδομή.</li> <li>• Υπαρξη δυνατότητας δυναμικής κλιμάκωσης, σύμφωνα με τις ανάγκες του καταναλωτή.</li> <li>• Ανάγκη ελέγχου και διαχείρισης των χρησιμοποιούμενων πόρων.</li> <li>• Τα ανώτερα επίπεδα της στοίβας πρέπει να διαχειρίζονται και να υποστηρίζονται υπό την ευθύνη του καταναλωτή.</li> </ul>

Σχήμα 2.4: Τεχνολογικές δυνατότητες της νεφοϋπολογιστικής

## 2.2 Συστήματα διαχείρισης μεγάλων δεδομένων

Το μοντέλο της νεφοϋπολογιστικής, όπως περιγράφηκε, αποτελεί όραμα μιας γενικής μετατόπισης της φυσικής θέσης των υπολογιστικών πόρων, εκτός των εγκαταστάσεων των επιχειρήσεων και, μέσω του δικτύου, προς τα διάφορα κέντρα υποδομών επόμενης γενιάς, τα οποία φιλοξενούνται σε μεγάλες εταιρίες, όπως Amazon, Google, Yahoo, Microsoft, Sun κλπ. Το υπολογιστικό νέφος θεωρείται η επανάσταση της πληροφορικής, επιτρέποντας στις επιχειρήσεις την αποφυγή λήψης υψηλού ρίσκου, από πραγματοποίηση μεγάλων κεφαλαιακών επενδύσεων σε υποδομές, και δίνοντάς τους τη δυνατότητα, να συνδέονται σε εξαιρετικά ισχυρούς υπολογιστικούς πόρους, μέσω της δικτυακής υποδομής τους.

Τα συστήματα διαχείρισης δεδομένων αποτελούν μια πολύ δημοφιλή κατηγορία εφαρμογών, οι οποίες είναι δυνητικά υποψήφιες για πιθανή μετάβασή τους σε υποδομές υπολογιστικού νέφους. Αυτό μπορεί να εξηγηθεί λαμβάνοντας υπόψη ένα ιδιόκτητο επιχειρησιακό σύστημα βάσης δεδομένων, το οποίο συνήθως προσεγγίζει ένα απαγορευτικά μεγάλο αρχικό κόστος, για την αγορά του απαιτούμενου υλικού και λογισμικού. Για πολλές εταιρίες το μοντέλο της νεφοϋπολογιστικής, υπό την έννοια της ενοικίασης πόρων ανάλογα με τη χρήση, αποτελεί ένα ιδιαίτερα ελκυστικό μοντέλο. Εντούτοις δεν ενδείκνυται για εφαρμογή σε όλες τις περιπτώσεις χρήσης των συστημάτων διαχείρισης δεδομένων. Στη συνέχεια παρουσιάζονται οι λόγοι για τους οποίους συνίσταται, ή μη, η μετάβαση των συστημάτων αυτών στο υπολογιστικό νέφος [29].

### 2.2.1 Διαχείριση συναλλαγών δεδομένων

Ο όρος συναλλακτική διαχείριση δεδομένων χρησιμοποιείται προκειμένου να περιγράψει τον τρόπο που πραγματοποιείται ο χειρισμός των δεδομένων στα συστήματα επεξεργασίας συναλλαγών (Online transaction processing, OLTP). Οι εφαρμογές που στηρίζονται στα συστήματα αυτά, συνήθως βασίζονται στο ACID μοντέλο [30] (Ατομικότητα, Συνεκτικότητα, Ακεραιότητα, Αντοχή - Atomocity, Consistensy, Integrity, Durability) ενώ τείνουν να πραγματοποιούν εγγραφές στη βάση δεδομένων σε εντατικούς ρυθμούς. Οι εφαρμογές συναλλακτικής διαχείρισης δεδομένων, αποτελούν τις λιγότερο πιθανές υποψήφιες για επιτυχή μετάβασή τους στο υπολογιστικό νέφος, για τους εξής λόγους:

**Μη χρήση αρχιτεκτονικής μηδενικού επιμερισμού:** Η αγορά των συστημάτων επεξεργασίας συναλλαγών κυριαρχείται από τις Microsoft, Oracle, Sybase, IBM και

DB2. Τα προϊόντα των εταιριών αυτών αποτυγχάνουν ή αδυνατούν να υποστηρίξουν την αρχιτεκτονική μηδενικού επιμερισμού (shared-nothing architecture). Η δημιουργία ενός συστήματος επεξεργασίας συναλλαγών με τέτοια αρχιτεκτονική, αποτελεί πρόκληση, καθώς ο τρόπος οργάνωσης και τοποθέτησης των δεδομένων, σε αυτά τα συστήματα, δυσκολεύει πολύ τη διαδικασία εκτέλεσης ταυτόχρονων ερωτημάτων. Το κύριο πλεονέκτημα της αρχιτεκτονικής μηδενικού επιμερισμού, αποτελεί η ευκολία στην επεκτασιμότητά της. Το πλεονέκτημα αυτό, έχει μικρότερη σημασία για τα δεδομένα συναλλαγών, για τα οποία η συντριπτική πλειοψηφία των εγκαταστάσεων έχει μέγεθος μικρότερο του 1 TB.

**Μη υποστήριξη του ACID μοντέλου:** Σύμφωνα με το θεώρημα CAP [31], κάθε σύστημα διαμοιρασμού δεδομένων, μπορεί να υποστηρίξει μόνο δύο από τις εξής ιδιότητες: (i) Συνεκτικότητα, (ii) Διαθεσιμότητα και (iii) Ανοχή σε Καταμήσεις. Επειδή τα δεδομένα στο υπολογιστικό νέφος είναι συνήθως κατανομημένα σε πολλές γεωγραφικές περιοχές, οι κατασκευαστές συστημάτων βάσεων δεδομένων υποχρεώνονται να θυσιάσουν μια ιδιότητα εκ των Συνεκτικότητα και Διαθεσιμότητα. Η Συνεκτικότητα είναι αυτή που επιλέγεται συνήθως να θυσιαστεί, με σκοπό την αύξηση της Διαθεσιμότητας του συστήματος.

**Αποθήκευση δεδομένων σε μη αξιόπιστο παροχέα υπηρεσιών:** Οι συναλλακτικές βάσεις δεδομένων περιέχουν συνήθως το πλήρες σύνολο των δεδομένων που είναι απαραίτητο για την ορθή εκτέλεση των κρίσιμων επιχειρησιακών διαδικασιών. Τα δεδομένα αυτά μπορεί ακόμη να είναι και απόρρητα, όπως πληροφορίες πελατών, ή αριθμοί πιστωτικών καρτών. Οποιαδήποτε αύξηση του κινδύνου ασφαλείας των δεδομένων για τα συστήματα αυτά, θεωρείται γενικά μη επιτρεπτή.

## 2.2.2 Διαχείριση ανάλυσης δεδομένων

Ο όρος αναλυτική διαχείριση δεδομένων αναφέρεται σε εφαρμογές διαχείρισης δεδομένων που χρησιμοποιούνται στον επιχειρηματικό σχεδιασμό, στη λύση προβλημάτων και στην υποστήριξη λήψης αποφάσεων. Στην αναλυτική διαδικασία χρησιμοποιούνται συνήθως δεδομένα ιστορικού ενδιαφέροντος και δεδομένα από συναλλακτικές βάσεις δεδομένων. Για το λόγο αυτό, τα δεδομένα που αποθηκεύονται σε αναλυτικές βάσεις δεδομένων έχουν γενικά πολύ μεγαλύτερο μέγεθος από τα δεδομένα των συναλλακτικών. Επιπλέον οι εφαρμογές ανάλυσης δεδομένων επικεντρώνονται κυρίως στην ανάγωση των δεδομένων και πολύ λιγότερο ή και καθόλου στην εγγραφή. Οι εφαρμογές αναλυτικής διαχείρισης δεδομένων, αποτελούν ένα πολύ καλό παράδειγμα επιτυχούς μετάβασής τους στο υπολογιστικό νέφος, για τους εξής λόγους:

**Εύκολη χρήση αρχιτεκτονικής μηδενικού επιμερισμού:** Το διαρκώς αυξανόμενο πλήθος δεδομένων στο οποίο εκτελούνται τα ερωτήματα ανάλυσης, αποτελεί την κύρια αιτία επιλογής της αρχιτεκτονικής μηδενικού επιμερισμού, αφού με τη χρήση της γίνεται δυνατή η καλύτερη κλιμάκωση. Επίσης ο φόρτος εργασίας της ανάλυσης δεδομένων τείνει να αποτελείται από ολοένα και πιο πολλές εργασίες μεγάλων σαρώσεων, πολυδιάστατων ομαδοποιήσεων και αστεροειδών ενώσεων, οι οποίες λειτουργούν καλύτερα σε ένα περιβάλλον παράλληλων κόμβων με χρήση δικτύου μηδενικού επιμερισμού. Επιπροσθέτως η δυνατότητα μη συχνής πραγματοποίησης εγγραφών στη βάση δεδομένων, περιορίζει την ανάγκη του περίπλοκου κλειδώματος των καταμερισμένων δεδομένων.

**Το ACID μοντέλο συνήθως δεν είναι απαραίτητο:** Η δυνατότητα μη συχνών εγγραφών σε αναλυτικές βάσεις δεδομένων, μαζί με το γεγονός ότι είναι συνήθως αποδεκτό να πραγματοποιείται η ανάλυση σε ένα πρόσφατο στιγμιότυπο της κατάστασης των δεδομένων, επιτρέπουν την εύκολη υιοθέτηση των ιδιοτήτων Ατομικότητα, Συνεκτικότητα και Ακεραιότητα του ACID μοντέλου. Έτσι η μείωση της Συνεκτικότητας που απαιτείται συνήθως να εφαρμοστεί στα καταναμημένα δεδομένα, δεν αποτελεί πρόβλημα για τις αναλυτικές βάσεις δεδομένων.

**Τα ευαίσθητα δεδομένα μπορούν να μείνουν εκτός ανάλυσης:** Στις περισσότερες περιπτώσεις, είναι δυνατό να αναγνωριστούν τα δεδομένα που θα μπορούσαν να προκαλέσουν τη μεγαλύτερη ζημιά, αν αποκτούσε πρόσβαση σε αυτά κάποιος τρίτος. Τα δεδομένα αυτά μπορούν να αφεθούν εκτός της διαδικασίας ανάλυσης, είτε να συμπεριληφθούν μόνο μετά την εφαρμογή μιας συνάρτησης ανωνυμίας, ή να αποθηκευτούν κρυπτογραφημένα. Επίσης τα δεδομένα που περιέχονται συνήθως στις αναλυτικές βάσεις δεδομένων, περιέχουν λιγότερη λεπτομέρεια από αυτά των συναλλακτικών βάσεων δεδομένων και άρα είναι λιγότερο ευαίσθητα από τα δεδομένα χαμηλού επιπέδου.

### 2.2.3 Στόχοι ανάλυσης μεγάλων δεδομένων

Η μεταφορά των δεδομένων σε μια νέα, πιο ισχυρή υπολογιστική υποδομή, δημιουργεί νέου είδους απαιτήσεις. Το υπολογιστικό νέφος, μπορεί να αποδειχθεί πολύ χρήσιμο εργαλείο στα συστήματα που ικανοποιούν τους ακόλουθους στόχους [32–34]:

**Διαθεσιμότητα:** Τα δεδομένα πρέπει να είναι πάντα διαθέσιμα ακόμα και σε περιπτώσεις δικτυακών σφαλμάτων. Για το σκοπό αυτό αναδείχτηκε η έννοια του μοντέλου της επικοινωνίας ως υπηρεσία (Communication as a Service, CaaS),

ώστε να υποστηριχθούν τέτοιου είδους απαιτήσεις. Το μοντέλο αυτό υποστηρίζει επίσης τη βελτιωμένη δικτυακή ασφάλεια, τη δυναμική πρόβλεψη της χρήσης εικονικών δικτυακών πόρων, την κρυπτογράφηση των επικοινωνιών και την παρακολούθηση του δικτύου.

**Επεκτασιμότητα:** Πρέπει να υποστηρίζονται πολύ μεγάλες βάσεις δεδομένων με πολύ υψηλό ρυθμό εξυπηρέτησης ερωτημάτων σε πολύ μικρό χρόνο. Πρέπει να υποστηρίζεται η προσθήκη νέων πελατών ή η αύξηση μεγέθους των υπαρχόντων, χωρίς να απαιτείται ιδιαίτερη προσπάθεια, πέραν της προσθήκης επιπλέον υπολογιστικών υποδομών. Συγκεκριμένα, απαιτείται να είναι δυνατή η αυτόματη ανακατανομή των δεδομένων, ώστε να εκμεταλλεύεται η προσθήκη νέου υλικού.

**Ελαστικότητα:** Οι απαιτήσεις των συστημάτων βάσεων δεδομένων σε υπολογιστικούς πόρους μπορούν να μεταβάλλονται ανάλογα με τη χρήση. Ζητούμενο είναι να ικανοποιείται η αλλαγή και προς τις δύο κατευθύνσεις (κλιμάκωση προς τα πάνω ή προς τα κάτω). Το σύστημα πρέπει να είναι ικανό να επανέρχεται πολύ γρήγορα σε κατάσταση ετοιμότητας, έπειτα από πραγματοποίηση τέτοιου είδους αλλαγών.

**Υψηλές επιδόσεις:** Σε πλατφόρμες δημόσιου υπολογιστικού νέφους, η τιμολόγηση των παρεχόμενων υπηρεσιών είναι δομημένη με τρόπο τέτοιο, ώστε κάθε πελάτης να χρεώνεται με βάση τη χρήση. Το κόστος για τον πελάτη αυξάνεται γραμμικά ανάλογα με το πλήθος των υπολογιστικών πόρων που χρησιμοποιεί. Έτσι οι υπολογιστικές επιδόσεις του συστήματος, μπορούν να έχουν άμεση επίδραση στην τιμολόγηση του πελάτη, και αποτελούν σημαντικό στοιχείο για την εξοικονόμηση χρημάτων.

**Υποστήριξη πολυ-πελατειακού μοντέλου:** Τα συστήματα βάσεων δεδομένων στο υπολογιστικό νέφος, πρέπει να είναι ικανά να υποστηρίξουν πολλές διαφορετικού είδους εφαρμογές, με τη χρήση παρόμοιου υλικού και λογισμικού. Οι επιδόσεις της κάθε εφαρμογής πρέπει να μην επηρεάζονται από τον τρόπο που χρησιμοποιούν άλλες εφαρμογές τους υπολογιστικούς τους πόρους.

**Δίκαιη εξισορρόπηση φορτίου:** Ο φόρτος εργασίας πρέπει να μεταφέρεται αυτόματα μεταξύ των διαφόρων κόμβων, ούτως ώστε οι πόροι να χρησιμοποιούνται αποτελεσματικά και να αποφεύγονται φαινόμενα υπερφόρτωσης.



**Ανοχή στα σφάλματα:** Κατά την εκτέλεση ερωτημάτων προς αναλυτικές βάσεις δεδομένων, απαραίτητο είναι το σύστημα να μπορεί να συνεχίσει την ορθή εκτέλεση του ερωτήματος, ακόμα και αν κάποιος κόμβος του συστήματος παρουσιάσει πρόβλημα δυσλειτουργικότητας. Αυτό συνήθως πραγματοποιείται με την ύπαρξη πολλών κόμβων που περιέχουν αντίγραφα των δεδομένων.

**Δυνατότητα εκτέλεσης σε ετερογενή περιβάλλοντα:** Στις πλατφόρμες υπολογιστικού νέφους υπάρχει η τάση της διαρκούς αύξησης των κόμβων που συμμετέχουν στην εκτέλεση ερωτημάτων. Είναι όμως σχεδόν αδύνατο να επιτευχθούν παρόμοιες επιδόσεις σε εκατοντάδες χιλιάδες υπολογιστικούς κόμβους. Οι αστοχίες υπολογιστικών εξαρτημάτων που δεν προκαλούν ολική αστοχία του κόμβου, αλλά οδηγούν σε μείωση των επιδόσεών του, είναι πολύ συχνές σε υψηλά επίπεδα κλιμάκωσης. Ένα σύστημα διαχείρισης βάσεων δεδομένων στο υπολογιστικό νέφος, πρέπει να λαμβάνει υπόψη του τέτοιου είδους θέματα και να υιοθετεί τις απαραίτητες προφυλάξεις που θα αποτρέπουν τη μείωση των επιδόσεών του.

**Ευέλικτη διεπαφή ερωτημάτων:** Πρέπει να υποστηρίζεται η χρήση SQL γλώσσας ερωτημάτων και η μη χρήση αυτής (noSQL). Ακόμη πρέπει να επιτρέπεται η δημιουργία προσαρμοσμένων συναρτήσεων και η εύκολη εκτέλεση αυτών σε περιβάλλοντα παράλληλης επεξεργασίας με χρήση ερωτημάτων.

## 2.3 Συστήματα πραγματικού χρόνου

### **Ορισμός 2.6 (Σύστημα Πραγματικού Χρόνου) :**

*Ως σύστημα πραγματικού χρόνου, ορίζεται ένα σύστημα επεξεργασίας πληροφορίας, του οποίου η αποτελεσματική λειτουργία δεν καθορίζεται μόνο από την ορθότητα των λογικών αποτελεσμάτων του, αλλά και από το φυσικό χρόνο που απαιτείται ώστε να παραχθούν τα αποτελέσματα αυτά [35].*

Ο υπολογισμός και οι επικοινωνίες πραγματικού χρόνου αποτελούν τεχνολογίες που καθιστούν εφικτή την ύπαρξη πολλών τρεχουσών και μελλοντικών περιοχών εφαρμογών. Τα πληροφοριακά συστήματα και τα ενσωματωμένα συστήματα (embedded systems), για τα οποία η εγγύηση της ασφάλειας αποτελεί καθοριστικό παράγοντα, αποτελούν στην πραγματικότητα συστήματα πραγματικού χρόνου. Οι τεχνολογίες πραγματικού χρόνου γίνονται όλο και περισσότερο διεισδυτικές (pervasive) και αναμένεται να καλύψουν ένα ευρύ φάσμα εφαρμογών, με διαφορετικούς χρονικούς περιορισμούς [35].

Τα συστήματα πραγματικού χρόνου πρέπει να είναι ικανά να ανταποκριθούν σε ένα ευρύ φάσμα γεγονότων του πραγματικού κόσμου. Έτσι ο σχεδιασμός τους γίνεται συνήθως ικανοποιώντας ένα υψηλό επίπεδο επεκτασιμότητας, προκειμένου να μπορούν να εξελιχθούν σε ένα περιβάλλον διαρκώς μεταβαλλόμενων απαιτήσεων. Από τα παραπάνω προκύπτει ότι τα συστήματα πραγματικού χρόνου έχουν ιδιαίτερες απαιτήσεις σχεδιασμού και υλοποίησης. Για το λόγο αυτό έχουν αναπτυχθεί εξειδικευμένες γλώσσες προγραμματισμού, οι οποίες μπορούν να συμβάλλουν στην ομαλή και αξιόπιστη λειτουργία του συστήματος. Ο σχεδιασμός των γλωσσών προγραμματισμού συστημάτων πραγματικού χρόνου έχει γίνει με τρόπο τέτοιο, ώστε να μην απαιτείται η παρέμβαση του λειτουργικού συστήματος τη στιγμή λειτουργίας του συστήματος. Οι συγκεκριμένες γλώσσες γνωρίζουν μεγάλη δημοφιλία ειδικά στα ενσωματωμένα υπολογιστικά συστήματα, στα οποία μέσω τεχνικών προγραμματισμού χαμηλού επιπέδου, δίνεται η δυνατότητα άμεσου χειρισμού και ελέγχου του συστήματος.

### 2.3.1 Χαρακτηριστικά συστημάτων πραγματικού χρόνου

Τα συστήματα πραγματικού χρόνου έχουν ορισμένα χαρακτηριστικά που τα διακρίνουν από τα μη πραγματικού χρόνου. Στη συνέχεια παρουσιάζεται η πλήρης λίστα των χαρακτηριστικών τους, η οποία όμως δεν είναι αντιπροσωπευτική για κάθε ένα από αυτά. Η έννοια του συστήματος πραγματικού χρόνου είναι πολύ γενική, και η εύρεση κοινών χαρακτηριστικών, αποτελεί πρόκληση [36].

**Χρονικός περιορισμός:** Κάθε εργασία πραγματικού χρόνου σχετίζεται με κάποιου είδους χρονικό περιορισμό. Η πιο συνήθης μορφή χρονικού περιορισμού είναι η προθεσμία (deadline), η οποία θέτει περιορισμό στο φυσικό χρόνο που πρέπει να εκτελεστεί μια διεργασία. Άλλοι τύποι χρονικών περιορισμών αποτελούν η καθυστέρηση και η χρονική διάρκεια.

**Νέο κριτήριο ορθότητας:** Η έννοια της ορθότητας στα συστήματα πραγματικού χρόνου είναι διαφορετική από εκείνη που χρησιμοποιείται στα παραδοσιακά συστήματα. Στα συστήματα πραγματικού χρόνου, η ορθότητα συνεπάγεται όχι μόνο τη λογική ορθότητα των αποτελεσμάτων, αλλά και την τήρηση των χρονικών περιορισμών κατά την εξαγωγή των αποτελεσμάτων. Ένα λογικά ορθό αποτέλεσμα που παράγεται αποτυγχάνοντας στους χρονικούς περιορισμούς που έχουν τεθεί, μπορεί να θεωρηθεί ως εσφαλμένο αποτέλεσμα.

**Κρισιμότητα ασφάλειας:** Για τα παραδοσιακά συστήματα μη πραγματικού χρόνου η ασφάλεια και η αξιοπιστία είναι ανεξάρτητα ζητήματα. Ωστόσο, σε πολλά συστήματα πραγματικού χρόνου, τα δύο αυτά θέματα είναι άρρηκτα συνδεδεμένα μεταξύ τους, καθιστώντας την ασφάλεια κρίσιμο παράγοντα για την αξιοπιστία



τους. Ασφαλές σύστημα είναι αυτό που δεν προκαλεί φθορά, ακόμη και σε περίπτωση αποτυχίας. Αξιόπιστο σύστημα, από την άλλη πλευρά, είναι αυτό που μπορεί να λειτουργεί για μεγάλη χρονική διάρκεια, χωρίς να παρουσιάζει ενδεχόμενες αποτυχίες.

**Παράλληλισμός:** Ένα σύστημα πραγματικού χρόνου επεξεργάζεται και αναλύει δεδομένα που πολλές φορές προέρχονται από γεγονότα του πραγματικού κόσμου. Επειδή τα γεγονότα αυτά μπορεί να λαμβάνουν χώρα ταυτόχρονα, θα πρέπει ένα σύστημα πραγματικού χρόνου να είναι ικανό να τα επεξεργάζεται παράλληλα.

**Κατανεμημένη δομή:** Τα δομικά στοιχεία πολλών συστημάτων πραγματικού χρόνου μπορεί να είναι κατανεμημένα σε διάφορες γεωγραφικές περιοχές. Έτσι, τα γεγονότα που προκύπτουν στις τοποθεσίες αυτές, συχνά αντιμετωπίζονται σε τοπικό επίπεδο, ώστε να μπορούν να εξαχθούν αποτελέσματα χωρίς να επιβαρύνεται το υποκείμενο δίκτυο επικοινωνίας.

**Ανατροφοδοτούμενη δομή:** Πολλά κατανεμημένα συστήματα πραγματικού χρόνου, μπορεί να έχουν μια δομή ανατροφοδότησης. Σε αυτά τα συστήματα, οι αισθητήρες ανιχνεύουν τα δεδομένα του περιβάλλοντος περιοδικά και τα επεξεργάζονται για να καθορίσουν τις απαραίτητες διορθωτικές ενέργειές τους.

**Κρισιμότητα διαδικασίας:** Η κρισιμότητα μιας διαδικασίας (task criticality) είναι ο δείκτης μέτρησης του κόστους αποτυχίας της διαδικασίας αυτής. Η κρισιμότητα μιας διαδικασίας προσδιορίζεται εξετάζοντας την κρισιμότητα των αποτελεσμάτων που παράγονται από τη διαδικασία αυτή. Ένα σύστημα πραγματικού χρόνου μπορεί να έχει διάφορα καθήκοντα διαφορετικής κρισιμότητας. Είναι, ως εκ τούτου, φυσική η ανάγκη, κατά το σχεδιασμό ενός αποσφαλματωμένου συστήματος, να λαμβάνεται υπόψη η κρισιμότητα των διαφόρων καθυκόντων του. Όσο μεγαλύτερη είναι η κρισιμότητα μιας εργασίας, τόσο πιο αξιόπιστη θα πρέπει να γίνει. Επιπλέον, σε περίπτωση σφάλματος μιας πολύ κρίσιμης εργασίας, είναι σημαντική η άμεση ανίχνευση του σφάλματος και αποκατάστασή του.

**Προσαρμοσμένο υλικό:** Ένα σύστημα πραγματικού χρόνου συχνά λειτουργεί χρησιμοποιώντας προσαρμοσμένο υλικό που έχει σχεδιαστεί και αναπτυχθεί ειδικά για το σκοπό αυτό.

**Μηχανισμός αντίδρασης:** Τα συστήματα πραγματικού χρόνου χρησιμοποιούν συχνά μηχανισμούς αντίδρασης (reactive mechanisms). Τα συστήματα αυτά διατηρούν μια συνεχή αλληλεπίδραση μεταξύ του υπολογιστικού συστήματος και του περιβάλλοντός του.

**Σταθερότητα:** Ακόμα και υπό συνθήκες υπερφόρτωσης, τα συστήματα πραγματικού χρόνου τηρούν τους χρονικούς περιορισμούς τους για τις πιο κρίσιμες εργασίες,

εις βάρος των μη κρίσιμων. Αυτό έρχεται σε αντίθεση με την απαίτηση αμεροληψίας των παραδοσιακών συστημάτων, ακόμη και υπό συνθήκες υπερφόρτωσης.

**Αυτόματος χειρισμός σφαλμάτων:** Πολλά συστήματα πραγματικού χρόνου λειτουργούν όλο το εικοσιτετράωρο, χωρίς να απαιτείται ανθρώπινη παρέμβαση. Έτσι η αυτόματη λήψη διορθωτικών μέτρων, σε περίπτωση σφάλματος, είναι επιβεβλημένη. Ακόμη και σε περίπτωση μη δυνατής λήψης άμεσων διορθωτικών μέτρων, επιθυμητό είναι το σφάλμα να μην οδηγεί σε καταστροφικές συνέπειες για το σύστημα. Τα σφάλματα ανιχνεύονται αυτόματα, και το σύστημα συνεχίζει να λειτουργεί σε υποβαθμισμένη κατάσταση χωρίς να τερματίζεται με βίαιο τρόπο.

### 2.3.2 Ταξινόμηση συστημάτων πραγματικού χρόνου

Για τα συστήματα πραγματικού χρόνου, η απώλεια ή μη τήρηση μιας προθεσμίας εκτέλεσης των λειτουργιών τους, μπορεί να επιφέρει αρνητικές επιπτώσεις. Έτσι πολλά συστήματα περιλαμβάνουν μια συνάρτηση μέτρησης επιπτώσεων, η οποία συσχετίζεται με τις μη τηρούμενες προθεσμίες. Ανάλογα με τις απαιτήσεις χρόνου εκτέλεσης που υπάρχουν, τα συστήματα πραγματικού χρόνου διακρίνονται σε υποκατηγορίες. Ένα σύστημα μπορεί να αποτελείται από υποσυστήματα τα οποία εμπίπτουν σε περισσότερες της μιας, υποκατηγορίες. Οι υποκατηγορίες παρουσιάζονται σχηματικά και εξηγούνται στη συνέχεια [37]:

**Αυστηρά πραγματικού χρόνου (Hard real-time):** Αποτελούν συστήματα στα οποία είναι απολύτως επιτακτικό, οι απαντήσεις να δίνονται μέσα στην απαιτούμενη προθεσμία π.χ. συστήματα ελέγχου πτήσεων.

**Χαλαρά πραγματικού χρόνου (Soft real-time):** Αποτελούν συστήματα στα οποία οι προθεσμίες είναι σημαντικές, θα συνεχίσουν όμως να λειτουργούν ορθά αν περιστασιακά οι προθεσμίες δεν τηρούνται π.χ. συστήματα πολυμέσων, συστήματα δρομολόγησης/μεταγωγής σε δίκτυα, συστήματα ελέγχου, παιχνίδια.

**Καθολικά πραγματικού χρόνου (Real real-time):** Αποτελούν συστήματα που είναι αυστηρά πραγματικού χρόνου, και στα οποία οι χρόνοι απόκρισης είναι πολύ μικροί π.χ. συστήματα καθοδήγησης πυραύλων.

**Σταθερά πραγματικού χρόνου (Firm real-time):** Αποτελούν συστήματα που είναι χαλαρά πραγματικού χρόνου, στα οποία δεν υπάρχει όμως όφελος από καθυστερημένη παράδοση.

## 2.4 Πολυκριτηριακή ανάλυση αποφάσεων

Η πολυκριτηριακή ανάλυση μπορεί να θεωρηθεί ως μία συστηματική και μαθηματικά τυποποιημένη προσπάθεια επίλυσης προβλημάτων που προκύπτουν από αντικρουόμενους στόχους. Ο Ορισμός 1.2 αποτελεί τον πλήρη επιστημονικά τεκμηριωμένο ορισμό της διαδικασίας.

Η επιστημονική περιοχή της πολυκριτηριακής ανάλυσης στηρίζεται σε ένα θεωρητικό υπόβαθρο, το οποίο αναπτύσει τη βασική λογική για την προσέγγιση προβλημάτων πολλαπλών κριτηρίων. Η διαδικασία της προσέγγισης περιλαμβάνει τον προσδιορισμό των κύριων δομικών στοιχείων του προβλήματος και την ανάλυση των βασικών τους ιδιοτήτων. Με βάση την διαδικασία αυτή, έχει αναπτυχθεί ένα πλήθος τεχνικών, κατάλληλων για την αντιμετώπιση ενός μεγάλου εύρους πολυκριτηριακών προβλημάτων. Αν και η ταξινόμηση των τεχνικών αυτών σε κατηγορίες δεν είναι αυστηρή, διακρίνονται τρεις βασικές ομάδες μεθόδων [5, 38, 39]:

**Αναλυτική-συνθετική προσέγγιση:** Εφαρμόζεται σε πολυκριτηριακά προβλήματα στα οποία εξετάζεται ένα πεπερασμένο σύνολο διακριτών επιλογών.

**Θεωρία των σχέσεων υπεροχής:** Επίσης εφαρμόζεται καλύτερα σε πολυκριτηριακά προβλήματα πεπερασμένων διακριτών επιλογών.

**Πολυκριτηριακός μαθηματικός προγραμματισμός:** Εφαρμόζεται σε προβλήματα πολυκριτηριακής φύσης, όπου εξετάζεται ένα συνεχές σύνολο άπειρου πλήθους επιλογών για τις οποίες, κατά αναλογία με τα προβλήματα γραμμικού μονοκριτηριακού προγραμματισμού, τα κριτήρια απόφασης μπορούν να έχουν οποιαδήποτε τιμή εντός ενός καθορισμένου επιτρεπτού πεδίου.

**Πολυκριτηριακή θεωρία χρησιμότητας:** Εφαρμόζεται σε πολυκριτηριακά προβλήματα με συνεχές ή διακριτό σύνολο επιλογών και στηρίζεται στη λογική της αναγωγής του πολυκριτηριακού σε μονοκριτηριακό πρόβλημα, μέσω του προσδιορισμού μιας καθολικής συνάρτησης χρησιμότητας που συνθέτει τις επιμέρους (ανά κριτήριο) προτιμήσεις του αποφασίζοντα σε ένα ενιαίο μέτρο, με βάση το οποίο προχωράει στη λήψη της τελικής απόφασης.

Κάθε πρόβλημα πολυκριτηριακής ανάλυσης προσδιορίζεται από ορισμένα δομικά χαρακτηριστικά, που απορρέουν είτε από την ίδια τη φύση του προβλήματος είτε από τις απόψεις και τις προτιμήσεις του αποφασίζοντα. Η ταυτοποίηση του αντικειμένου της πολυκριτηριακής ανάλυσης ως προς τα χαρακτηριστικά αυτά, αποτελεί ένα πρώτο στάδιο της αναλυτικής διαδικασίας, που διευκολύνει την κατανόηση του προβλήματος και επιτρέπει την επιλογή της κατάλληλης μεθόδου επίλυσης.

Η διαδικασία της πολυκριτηριακής ανάλυσης (Σχήμα 2.5), ως μαθηματικό μοντέλο επίλυσης πολυκριτηριακών προβλημάτων, μπορεί να βοηθήσει τον αποφασίζοντα στην αναζήτηση της βέλτιστης λύσης και στην καλύτερη κατανόηση των συνεπειών των αποφάσεών του. Αποτελείται από δύο βασικά στάδια: (i) το στάδιο δόμησης του προβλήματος, όπου γίνεται ο καθορισμός του προβλήματος και επιλέγονται τα πιθανά σενάρια λύσης του, και (ii) το στάδιο ανάλυσης των αποτελεσμάτων, στο οποίο ελέγχονται οι πιθανές λύσεις [40]. Το πρώτο στάδιο περιλαμβάνει την επιλογή των κριτηρίων απόφασης, τη μέτρηση των επιδόσεών τους και την ταξινόμηση αυτών, ώστε να εκτιμηθεί η βαρύτητά τους. Επίσης καθορίζονται οι πιθανές περιοριστικές παράμετροι, ανάλογα με το αντικείμενο του εξεταζόμενου προβλήματος. Δημιουργείται ένα μοντέλο βαθμωτής αξιολόγησης των σεναρίων, το οποίο χρησιμοποιείται για την ταξινόμησή τους και την επιλογή της τελικής λύσης. Στο δεύτερο στάδιο αναλύεται η ευαισθησία της λύσης και προσδιορίζονται οι συγκρούσεις των κριτηρίων απόφασης.



Σχήμα 2.5: Διαδικασία πολυκριτηριακής ανάλυσης

Το μαθηματικό μοντέλο επίλυσης των πολυκριτηριακών προβλημάτων περιλαμβάνει ορισμένα στοιχεία [41] που πρέπει να αναφερθούν για την αποσαφήνιση των βασικών αρχών του. Ένα πρόβλημα αποτελείται από:

**Τη μήτρα αξιολόγησης** που περιλαμβάνει ένα σύνολο διακριτών επιλογών, ένα σύνολο κριτηρίων αξιολόγησης και την επίδοση της κάθε επιλογής στο αντίστοιχο κριτήριο.

**Το σύστημα προτιμήσεων του αποφασίζοντα** που εμπεριέχει τη σχετική βαρύτητα των κριτηρίων, την κατεύθυνση προτίμησης των επιδόσεων (ελάχιστο ή μέγιστο) και τα όρια ανοχής.

Ωστόσο ζητούμενο του κάθε προβλήματος παραμένει ο προσδιορισμός της σχετικά βέλτιστης λύσης, μέσα από την ταξινόμηση όλων των πιθανών λύσεων σε ομάδες και την ιεράρχησή τους με βάση τα κριτήρια και τη βαρύτητά τους. Ένα πρόβλημα μπορεί

να επιλυθεί είτε με την εφαρμογή μεθόδων σύνθεσης των επιδόσεων, όπου γίνεται αναγωγή σε μονοκριτηριακό πρόβλημα, είτε με την εφαρμογή μεθόδων ιεράρχησης των επιλογών όπου γίνεται δυαδική σύγκριση σε κάθε κριτήριο και διατυπώνονται σχέσεις επικράτησης. Στη συνέχεια παρουσιάζονται οι δύο αυτές μέθοδοι επίλυσης των πολυκριτηριακών προβλημάτων.

### 2.4.1 Πολυκριτηριακή θεωρία αξίας ή χρησιμότητας

Σε αυτό το σύστημα της πολυκριτηριακής ανάλυσης, η συγκριτική αξιολόγηση των εναλλακτικών σεναρίων ακολουθεί τα εξής στάδια:

**1ο Στάδιο:** Αρχικά, γίνεται η επιλογή των κριτηρίων, τα οποία θα πρέπει να καλύπτουν όλες τις πλευρές του εξεταζόμενου προβλήματος και να μπορούν να βαθμολογηθούν σε κατάλληλη κλίμακα. Μετά, ακολουθεί η ταξινόμηση των κριτηρίων σε ομάδες. Καθεμιά από αυτές τις ομάδες χαρακτηρίζεται από ένα συντελεστή βαρύτητας, που δηλώνει το “βάρος” της στο κάθε σενάριο και προσδιορίζεται μετά από συζητήσεις με όλους τους εμπλεκόμενους φορείς, λαμβάνοντας υπόψη και δεδομένα ανάλογων περιπτώσεων. Το άθροισμα των συντελεστών αυτών θα πρέπει να είναι ίσο με 100%. Κατόπιν, βάσει των παραπάνω προκύπτει η αντίστοιχη αθροιστική συνάρτηση, η οποία θα έχει τη μορφή:

$$f(O) = \sum A_i * O_i$$

όπου:

- $O_i$  είναι οι επιμέρους ομάδες κριτηρίων
- $A_i$  είναι ο συντελεστής βαρύτητας κάθε μίας από τις ομάδες κριτηρίων  $O_i$  και
- το άθροισμα των συντελεστών βαρύτητας πρέπει να ισούται με 1 (100%),  $\sum A_i = 1$ .

**2ο Στάδιο:** Οι ομάδες κριτηρίων αναλύονται στα επιμέρους κριτήρια αξιολόγησης, για τα οποία επίσης καθορίζεται η σχετική σπουδαιότητά τους μέσα στην ομάδα κριτηρίων με τη βοήθεια κατάλληλων συντελεστών βαρύτητας. Το άθροισμα των συντελεστών βαρύτητας των επιμέρους κριτηρίων μέσα σε κάθε ομάδα είναι επίσης 100%.

**3ο Στάδιο:** Πραγματοποιείται ανάλυση όλων των εναλλακτικών χαρακτηριστικών κάθε επιμέρους κριτηρίου τα οποία στη συνέχεια ποσοτικοποιούνται βάσει κλίμακας 1-10, όπου οι μικρότερες τιμές αφορούν στις δυσμενέστερες αποδόσεις των χαρακτηριστικών του κριτηρίου και οι μεγαλύτερες τιμές στις ευνοϊκότερες (καλύπτοντας με τον τρόπο αυτό όλες τις πιθανές περιπτώσεις).

**4ο Στάδιο:** Αρχικά γίνεται αποτύπωση των χαρακτηριστικών κάθε επιμέρους κριτηρίου για κάθε εναλλακτικό σενάριο και αφού γίνει σύγκριση τους με την κλίμακα που αναπτύσσεται στο 3ο στάδιο, λαμβάνει μία συγκεκριμένη τιμή απόδοσης σε κλίμακα από 1 έως 10. Στη συνέχεια, οι τιμές που προκύπτουν, πολλαπλασιάζονται με το σχετικό συντελεστή βαρύτητας που έχει καθένα από τα κριτήρια σε κάθε ομάδα. Ακολούθως, προστίθενται τα αντίστοιχα γινόμενα για την κάθε ομάδα και με τον τρόπο αυτό ποσοτικοποιείται κάθε ομάδα κριτηρίων. Μετά, ο βαθμός κάθε ομάδας πολλαπλασιάζεται με τον αντίστοιχο συντελεστή βαρύτητάς της, κι έτσι προκύπτει μέσω της αθροιστικής συνάρτησης ένα μέτρο της συνολικής αποτελεσματικότητας κάθε επιλογής. Με βάση τη βαθμολογία αυτή γίνεται κατάταξη των εναλλακτικών σεναρίων, με ευνοϊκότερο, αυτό που έχει την υψηλότερη επίδοση.

#### 2.4.2 Προσέγγιση σχέσεων υπεροχής

Η προσέγγιση των σχέσεων υπεροχής βασίζεται στην ανά ζεύγη σύγκριση των επιλογών σε κάθε μεμονωμένο κριτήριο με βάση τις επιδόσεις τους και τις ενδοκριτηριακές προτιμήσεις του αποφασίζοντα, όπως αυτές εκφράζονται με τα κατώφλια αδιαφορίας ή/και προτίμησης. Χαρακτηριστικό των μεθόδων υπεροχής είναι ότι η σύγκριση γίνεται στην αρχική κλίμακα μέτρησης των επιδόσεων (ποσοτική ή ποιοτική) χωρίς αναγωγή στο διάστημα  $[0, 1]$ . Ο δείκτης που προκύπτει από την ανά κριτήριο σύγκριση συντίθεται στη συνέχεια σε ένα συνολικό δυαδικό δείκτη λαμβάνοντας υπόψη τους συντελεστές βαρύτητας των κριτηρίων.

Οι δυαδικοί δείκτες χαρακτηρίζουν ζεύγη επιλογών  $(a, b)$  και προσδιορίζουν στο διάστημα  $[0, 1]$  το βαθμό στον οποίο ισχύει η υπόθεση: «η λύση  $a$  είναι τουλάχιστον τόσο καλή όσο και η λύση  $b$ ». Ανάλογα με την μέθοδο και τον ακριβή τρόπο υπολογισμού τους, οι δείκτες αυτοί ονομάζονται δείκτες προτίμησης ή δείκτες συμφωνίας (ως προς την υπόθεση). Μια λύση  $a$  που εμφανίζει υψηλές τιμές δεικτών προτίμησης σε σχέση με τις υπόλοιπες εναλλακτικές λύσεις χαρακτηρίζεται από μία σχετική υπεροχή, ενώ

αντίθετα άλλες λύσεις που δεν επιβεβαιώνουν την υπόθεση σε σημαντικό βαθμό, κρίνονται ως υποδεέστερες. Επομένως, το τελικό στάδιο στις μεθόδους υπεροχής είναι η επεξεργασία των δυαδικών δεικτών έτσι ώστε να προκύψουν σχέσεις υπεροχής και η τελική κατάταξη των εναλλακτικών λύσεων.

**1ο Στάδιο:** Αρχικά, γίνεται η επιλογή των κριτηρίων, τα οποία θα πρέπει να καλύπτουν όλες τις πλευρές του εξεταζόμενου προβλήματος και να μπορούν να βαθμολογηθούν σε κατάλληλη κλίμακα.

**2ο Στάδιο:** Για όλα τα κριτήρια αξιολόγησης καθορίζεται η σπουδαιότητά τους με τη βοήθεια κατάλληλων συντελεστών βαρύτητας. Το άθροισμα των συντελεστών βαρύτητας των κριτηρίων είναι 100%.

**3ο Στάδιο:** Πραγματοποιείται ανάλυση όλων των εναλλακτικών χαρακτηριστικών κάθε επιμέρους κριτηρίου τα οποία στη συνέχεια ποσοτικοποιούνται βάσει κλίμακας 1-10, όπου οι μικρότερες τιμές αφορούν στις δυσμενέστερες αποδόσεις των χαρακτηριστικών του κριτηρίου και οι μεγαλύτερες τιμές στις ευνοϊκότερες (καλύπτοντας με τον τρόπο αυτό όλες τις πιθανές περιπτώσεις).

**4ο Στάδιο:** Αρχικά γίνεται αποτύπωση των χαρακτηριστικών κάθε επιμέρους κριτηρίου για κάθε εναλλακτικό σενάριο και αφού γίνει σύγκριση τους με την κλίμακα που αναπτύσσεται στο 3ο στάδιο, λαμβάνει μία συγκεκριμένη τιμή απόδοσης σε κλίμακα από 1 έως 10.

**5ο Στάδιο:** Εφαρμογή του μοντέλου πολυκριτηριακής ανάλυσης.

Αυτό που διαφοροποιεί τις μεθόδους υπεροχής από τις μεθόδους πολυκριτηριακής ανάλυσης αθροιστικής συνάρτησης, είναι ότι το μέτρο χαρακτηρισμού και αξιολόγησης των λύσεων δεν είναι μία συνολική σταθμισμένη «επίδοση», αλλά ένας δείκτης σύνθεσης των προτιμήσεων του αποφασίζοντα. Αυτό σημαίνει ότι και οι συντελεστές βαρύτητας στις μεθόδους υπεροχής παίζουν ένα διαφορετικό ρόλο. Ειδικότερα, δεν έχουν το χαρακτήρα των συντελεστών αντιστάθμισης μεταξύ των επιδόσεων στα επιμέρους κριτήρια, γι' αυτό και δεν χρησιμοποιείται η μέθοδος αντιστάθμισης για την εξαγωγή τους. Αντίθετα, υποδηλώνουν το βαθμό συμβολής κάθε κριτηρίου στη διαμόρφωση του συνολικού δείκτη προτίμησης ή συμφωνίας.

### 2.4.3 Πολυκριτηριακός μαθηματικός προγραμματισμός

Όταν το σύνολο των δυνατών επιλογών του αποφασίζοντα δεν δίδεται ρητά, αλλά έμμεσα μέσω των τιμών των μεταβλητών απόφασης ενός προβλήματος μαθηματικού προγραμματισμού, τότε το πρόβλημα ανήκει στον Πολυκριτηριακό Μαθηματικό Προγραμματισμό (Multiple Objective Mathematical Programming). Οι μαθηματικές σχέσεις μεταξύ των μεταβλητών απόφασης, που πρέπει να ικανοποιούνται, αποτελούν τους περιορισμούς του προβλήματος, ενώ οι συναρτήσεις εκείνες των μεταβλητών απόφασης που πρέπει να αριστοποιηθούν ονομάζονται αντικειμενικές συναρτήσεις. Με τον όρο λύση του προβλήματος εννοείται κάθε συνδυασμός τιμών που μπορούν να λάβουν οι μεταβλητές απόφασης.

Ο Πολυκριτηριακός Μαθηματικός Προγραμματισμός επιλύει το πρόβλημα της διανυσματικής βελτιστοποίησης, που αποτελεί μία επέκταση της βαθμωτής βελτιστοποίησης (scalar optimization), με την οποία ασχολείται ο συμβατικός Μαθηματικός Προγραμματισμός χρησιμοποιώντας μια αντικειμενική συνάρτηση. Όταν επιπλέον οι περιορισμοί και οι αντικειμενικές συναρτήσεις είναι γραμμικές συναρτήσεις των μεταβλητών απόφασης τότε το πρόβλημα ανήκει στον Πολυκριτηριακό Γραμμικό Προγραμματισμό (Multiple Objective Linear Programming) που αποτελεί με τη σειρά του επέκταση του συμβατικού Γραμμικού Προγραμματισμού.

Η εισαγωγή πολλών αντικειμενικών συναρτήσεων σ' ένα πρόβλημα Γραμμικού Προγραμματισμού, δημιουργεί το πρόβλημα της γραμμικής διανυσματικής μεγιστοποίησης (linear vector maximum problem). Με τον όρο μεγιστοποίηση εννοείται γενικότερα η αριστοποίηση, αφού και η ελαχιστοποίηση μπορεί εύκολα να μετατραπεί σε μεγιστοποίηση με αλλαγή προσήμου. Η ανάλυση των προβλημάτων αυτών αποτελεί το αντικείμενο του Πολυκριτηριακού Γραμμικού Προγραμματισμού.

## 2.5 Ερωτήματα κορυφογραμμής

Τα πολυκριτηριακά προβλήματα αποτελούν γενικά μια κατηγορία προβλημάτων, στην οποία όσο αυξάνει το πλήθος των κριτηρίων απόφασης, τόσο δυσχεραίνει η εκτέλεση της διαδικασίας της πολυκριτηριακής ανάλυσης. Για τη λύση αυτού του προβλήματος, έχουν αναπτυχθεί διάφορες τεχνικές, οι οποίες επιτρέπουν την εκτέλεση τμημάτων της διαδικασίας σε υπολογιστικές υποδομές. Οι περισσότερες τεχνικές τέτοιου είδους, αυτοματοποιούν το τμήμα εκείνο που αφορά στη στάθμιση και εξαγωγή αποτελεσμάτων, όπως περιγράφηκε στο Σχήμα 2.5. Η πιο συχνά χρησιμοποιούμενη τεχνική εξαγωγής αποτελεσμάτων, αποτελεί η μέθοδος του καθολικού μέτρου επίδοσης, που



περιγράφηκε στο Κεφάλαιο 1. Ωστόσο στη συνέχεια, για λόγους πληρότητας, παρουσιάζεται μια παρεμφερής τεχνική που επιτρέπει την παροχή αποτελεσμάτων, χωρίς να είναι απαραίτητο να καθοριστεί από τον αποφασίζοντα το καθολικό μέτρο επίδοσης.

### 2.5.1 Ορισμοί

Τα ερωτήματα κορυφογραμμής αποτελούν εργαλείο της διαδικασίας λήψης απόφασης, τα οποία επιτρέπουν την αντικειμενική θεώρηση των κριτηρίων απόφασης. Η τεχνική αυτή, αποδίδει πολύ καλά όταν ο αποφασίζων αδυνατεί να ορίσει τις βαρύτητες των κριτηρίων, υστερεί ωστόσο να προσαρμοστεί αποτελεσματικά στις προσωπικές του προτιμήσεις. Πριν αναλυθεί η έννοια των ερωτημάτων κορυφογραμμής, είναι απαραίτητο να παρουσιαστεί η έννοια της κυριαρχίας Pareto.

#### Ορισμός 2.7 (Κυριαρχία Pareto) :

Έστω  $\mathcal{A} = \{A_1, \dots, A_d\}$  ένα πεπερασμένο σύνολο χαρακτηριστικών (ένα σχεσιακό σχήμα). Το νούμερο  $d$  είναι η διάσταση της κυριαρχίας Pareto. Κάθε χαρακτηριστικό  $A_i \in \mathcal{A}$  σχετίζεται με ένα άπειρο πεδίο τιμών  $\mathcal{D}_{A_i}$  (π.χ. πραγματικούς αριθμούς). Το διάστημα των πλειάδων ορίζεται ως  $\mathcal{U} = \prod_{A_i \in \mathcal{A}} \mathcal{D}_{A_i}$ . Η σχέση κυριαρχίας αποτελεί υποσύνολο του  $\mathcal{U} \times \mathcal{U}$ . Δοθείσης μιας πλειάδας  $t \in \mathcal{U}$ , η τιμή του χαρακτηριστικού της  $A_i$  συμβολίζεται ως  $t[A_i]$ . Ο συμβολισμός αυτός μπορεί εύκολα να γενικευθεί για σύνολα χαρακτηριστικών. Σε κάθε χαρακτηριστικό  $A_i \in \mathcal{A}$  αντιστοιχίζεται μια μέθοδος ταξινόμησης των τιμών του (αύξουσα ή φθίνουσα), ανάλογα με την κατεύθυνση προτίμησής του, η οποία συμβολίζεται ως  $>_{A_i}$ . Η κυριαρχία Pareto  $\succ^{pto}$  ορίζεται ως εξής [42]:

$$t \succ^{pto} s \equiv t \neq s \wedge \bigwedge_{A_i \in \mathcal{A}} t[A_i] \geq_{A_i} s[A_i]$$

Ο παραπάνω μαθηματικός τύπος μπορεί εύκολα να επεξηγηθεί: Μια πλειάδα  $t$  κυριαρχεί έναντι μιας άλλης πλειάδας  $s$  όταν και οι δύο ανήκουν στο διάστημα  $\mathcal{U}$ , και για κάθε χαρακτηριστικό  $A_i$  η τιμή της πλειάδας  $t$  ( $t[A_i]$ ) είναι καλύτερη, σύμφωνα με τη μέθοδο ταξινόμησης του χαρακτηριστικού ( $>_{A_i}$ ) από την αντίστοιχη τιμή του  $s$  ( $s[A_i]$ ). Στην πράξη, πολλά χαρακτηριστικά από πραγματικές εφαρμογές, είναι άνευ σημασίας προς την κυριαρχία Pareto, όπως για παράδειγμα τα πεδία ταυτότητας (id). Τα πεδία αυτά μπορούν να μη λαμβάνονται υπόψη όταν υπολογίζεται η κυριαρχία Pareto.

Από την οπτική της άλγεβρας, η κυριαρχία Pareto μπορεί να οριστεί χρησιμοποιώντας τον δυαδικό τελεστή  $\otimes$  (που επίσης ονομάζεται τελεστής Pareto) σε συνδυασμό με την

ταξινόμηση των χαρακτηριστικών  $\succ_{A_i}, i = 1, \dots, d$ :

$$\succ_{\mathcal{A}}^{pto} = \succ_{A_1} \otimes \succ_{A_2} \otimes \dots \otimes \succ_{A_d}$$

όπου οι εξισώσεις  $\succ_{A_i}^{pto}$  και  $\succ_{XY}^{pto} = \succ_X^{pto} \otimes \succ_Y^{pto}$  μπορούν να ορισθούν ως εξής:

$$\begin{aligned} t[XY] \succ_{XY}^{pto} s[XY] \equiv \\ t[X] \succ_X^{pto} s[X] \wedge t[Y] \succ_Y^{pto} s[Y] \\ \vee t[X] \preceq_X^{pto} s[X] \wedge t[Y] \succ_Y^{pto} s[Y] \end{aligned}$$

όταν  $XY \subseteq \mathcal{A}$  και  $X \cap Y = \emptyset$ . Ο τελεστής Pareto είναι προσεταιριστικός και αντιμεταθετικός [42].

Η έννοια των ερωτημάτων κορυφογραμμής, αναφέρθηκε για πρώτη φορά το 2001 από τους Borzsony κ.ά. οπότε και δόθηκε ο πρώτος ορισμός τους:

**Ορισμός 2.8 (Ερώτημα Κορυφογραμμής) :**

Δοθέντος ενός συνόλου πλειάδων  $S$  στο χώρο δεδομένων  $\mathcal{U}$ , το ερώτημα κορυφογραμμής είναι η πράξη εκείνη που υπολογίζει το σύνολο  $SKY(S) \subseteq S$  των πλειάδων, οι οποίες δεν κυριαρχούνται Pareto από καμία πλειάδα του  $S$ . Τα σημεία (ή πλειάδες) του  $SKY(S)$  ονομάζονται σημεία κορυφογραμμής [43].

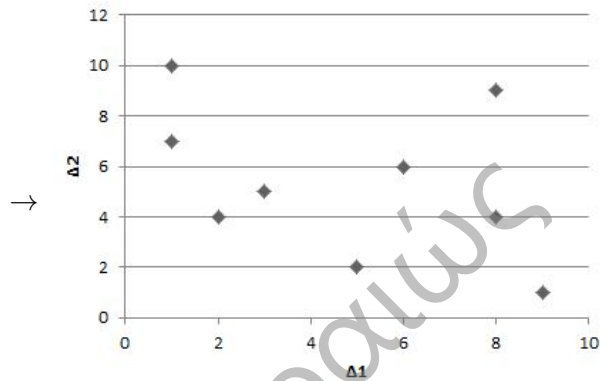
Είναι σημαντικό να επισημανθεί η διάκριση μεταξύ των τεσσάρων: (i) κυριαρχία Pareto, (ii) ερωτήματα κορυφογραμμής που επιστρέφουν όλα τα μη κυριαρχούμενα Pareto δεδομένα, (iii) κορυφογραμμές που είναι τα αποτελέσματα των ερωτημάτων κορυφογραμμής και (iv) αλγόριθμοι κορυφογραμμής που υπολογίζουν κορυφογραμμές. Η τεχνική των ερωτημάτων κορυφογραμμής, προϋπήρχε ήδη από το 1975 με την ονομασία «πρόβλημα μεγίστων ανυσμάτων» (maximum vector problem) ή βέλτιστη Pareto (Pareto optimum) [44, 45], όπως αυτή προτάθηκε από τους Preparata κ.ά., και αντιμετώπιζε πολυκριτηριακά προβλήματα, στα πλαίσια της υπολογιστικής γεωμετρίας. Τα ερωτήματα κορυφογραμμής (skyline queries) είναι ένας από τους πολλούς τύπους ερωτημάτων βάσεων δεδομένων, που μπορούν να πραγματοποιηθούν πάνω σε πολυδιάστατα δεδομένα. Τα ερωτήματα αυτά έχουν προσελκύσει το ενδιαφέρον μεγάλου μέρους της επιστημονικής κοινότητας, από διάφορους τομείς (Βάσεις Δεδομένων, Υπολογιστική Γεωμετρία, Συστήματα Λήψης Αποφάσεων κ.ά.), λόγω της ιδιαίτερης χρησιμότητας και χρηστικότητας που παρουσιάζουν.

Στην επιστήμη των βάσεων δεδομένων, τα σύνολα δεδομένων (datasets) μελετώνται καλύτερα όταν αποτυπωθούν στο καρτεσιανό σύστημα αξόνων. Αυτό συμβαίνει γιατί οι διαδικασίες ανάλυσης των δεδομένων, γίνονται πιο εύκολα κατανοητές όταν αποτυπωθούν και αναλυθούν σε ένα σύστημα δύο ή τριών διαστάσεων, και στη συνέχεια

γενικευθούν για μεγαλύτερο πλήθος διαστάσεων. Για την αναπαράσταση των δεδομένων στο καρτεσιανό σύστημα, αρκεί κάθε άξονάς του να αποτελεί ένα χαρακτηριστικό (ή διάσταση) των δεδομένων αυτών. Οι πλειάδες των δεδομένων αναπαριστώνται ως σημεία του καρτεσιανού συστήματος, τα οποία έχουν τιμές στους άξονες, ίδιες με αυτές των χαρακτηριστικών τους. Η τεχνική αυτή παρουσιάζεται στο Σχήμα 2.6.

$\Delta 1$	$\Delta 2$
5	2
8	4
3	5
2	4
1	7
8	9
9	1
6	6
1	10

(α) Πίνακας δεδομένων



(β) Δισδιάστατη αναπαράσταση δεδομένων

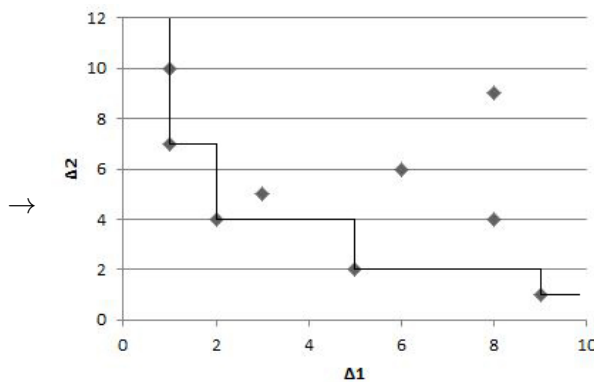
Σχήμα 2.6: Αναπαράσταση ενός συνόλου δεδομένων στο καρτεσιανό σύστημα

Όπως αναφέρθηκε, η κορυφογραμμή (skyline) ενός συνόλου δεδομένων αποτελείται από τα σημεία εκείνα, τα οποία δεν κυριαρχούνται Pareto από κανένα άλλο. Ένα παράδειγμα που χρησιμοποιείται πολύ συχνά στη βιβλιογραφία προκειμένου να περιγράψει τη χρηστικότητα των ερωτημάτων κορυφογραμμής, είναι το ακόλουθο: Έστω ότι διατίθεται ένα σύνολο δεδομένων, που αποτελείται από σημεία στο δισδιάστατο χώρο. Το κάθε σημείο αντιπροσωπεύει ένα ξενοδοχείο για το οποίο υπάρχει (ως πληροφορία) η τιμή ενός δωματίου του και η απόστασή του από τη θάλασσα. Έχοντας τις πληροφορίες αυτές, είναι δυνατή η απεικόνιση των ξενοδοχείων ως σημεία στο καρτεσιανό σύστημα αξόνων. Το σύνηθες πρόβλημα στα δεδομένα αυτά, αποτελεί η αναζήτηση των ξενοδοχείων που ικανοποιούν ταυτόχρονα το κριτήριο της χαμηλής τιμής και το κριτήριο της μικρής απόστασης από τη θάλασσα. Η κορυφογραμμή ενός τέτοιου συνόλου δεδομένων, αποτελείται από τα ξενοδοχεία εκείνα που δεν έχουν χειρότερη τιμή και χειρότερη απόσταση από τη θάλασσα από οποιοδήποτε άλλο ξενοδοχείο.

Το Σχήμα 2.7 παρουσιάζει ένα παράδειγμα κορυφογραμμής, που αντιστοιχεί στα σημεία του σχήματος 2.6.

Με την αναπαράσταση των σημείων στο καρτεσιανό σύστημα, γίνεται αμέσως εμφανής η ευκολία με την οποία μπορούν να εξαχθούν εμπειρικά τα σημεία κορυφογραμμής. Στη δισδιάστατη αναπαράσταση των σημείων του συνόλου δεδομένων του σχήματος 2.7, η συνεχόμενη μαύρη γραμμή, που ενώνει τα σημεία κορυφογραμμής, αποτελεί ουσιαστικά το όριο κυριαρχίας: όλα τα σημεία που βρίσκονται στο χώρο

$\Delta 1$	$\Delta 2$
5	2
8	4
3	5
2	4
1	7
8	9
9	1
6	6
1	10



(α) Πίνακας με επισημασμένο skyline

(β) Δισδιάσταση αναπαράσταση του skyline

Σχήμα 2.7: Κορυφογραμμή ενός συνόλου δεδομένων στο καρτεσιανό σύστημα

πάνω και δεξιά από τη γραμμή αυτή, κυριαρχούνται Pareto από ένα ή περισσότερα σημεία κορυφογραμμής.

Η πρώτη εφαρμογή των ερωτημάτων κορυφογραμμής, στα πλαίσια των υπολογιστικών βάσεων δεδομένων, είναι αυτή των Borzsonyi κ.ά. [43], στην οποία ορίζεται ο τελεστής κορυφογραμμής (skyline operator). Επίσης εκεί παρουσιάζονται δύο αλγόριθμοι υπολογισμού της κορυφογραμμής των δεδομένων που είναι αποθηκευμένα σε μία βάση δεδομένων, προτείνοντας επιπλέον, μια επέκταση της γλώσσας SQL με τον τελεστή SKYLINE OF, ώστε οι υπολογισμοί να πραγματοποιούνται κατά τη διάρκεια της ανάκτησης των δεδομένων και όχι αργότερα. Η σύνταξη των ερωτημάτων κορυφογραμμής, όπως αυτή προτάθηκε από τους Borzsonyi κ.ά. [43], παρουσιάζεται στον Αλγόριθμο 2.1:

```
SELECT ... FROM ...
WHERE ...
GROUP BY ... HAVING ...
SKYLINE OF [DISTINCT]  $d_1$  [MIN|MAX|DIFF], ...,  $d_m$  [MIN|MAX|DIFF]
ORDER BY ...
```

Αλγόριθμος 2.1: SQL σύνταξη ερωτήματος κορυφογραμμής κατά Borzsony κ.ά.

Τα γνωρίσματα  $d_1, \dots, d_m$ , αντιστοιχούν στις διαστάσεις ενδιαφέροντος, για τον υπολογισμό του ερωτήματος κορυφογραμμής (skyline), όπως η τιμή ενός δωματίου, η απόστασή του από τη θάλασσα κ.λ.π. Οι τελεστές MIN, MAX και DIFF ορίζουν την κατεύθυνση προτίμησης των κριτηρίων (ελάχιστο, μέγιστο, διαφορετικό). Στο παράδειγμα με τα ξενοδοχεία, το κριτήριο της τιμής ενός δωματίου ακολουθεί την κατεύθυνση του ελαχίστου (MIN επισημείωση), ενώ ο βαθμός αξιολόγησής του, αποτελεί κριτήριο μεγιστοποίησης (MAX επισημείωση). Το προαιρετικό όρισμα DISTINCT καθορίζει την διαχείριση των διπλότυπων. Η σημασιολογία του SKYLINE OF είναι ξεκάθαρη. Η SKYLINE OF πρόταση εκτελείται μετά το SELECT ... FROM ... WHERE ... GROUP BY

... HAVING ... μέρος του ερωτήματος, αλλά πριν από την ORDER BY πρόταση ενώ είναι πιθανό να ακολουθούν και άλλες προτάσεις μετά από αυτήν. Το SKYLINE OF επιλέγει τις ενδιαφέρουσες πλειάδες, δηλαδή εκείνες που δεν κυριαρχούνται Pareto από καμία άλλη.

### 2.5.2 Χαρακτηριστικά κορυφογραμμής και κυριαρχίας Pareto

Η κυριαρχία Pareto και, κατ' επέκταση, οι κορυφογραμμές ενός συνόλου δεδομένων, έχουν ορισμένες πολύ ενδιαφέρουσες ιδιότητες [42, 43]. Κατά την παρουσίαση των παρακάτω ιδιοτήτων, θεωρείται η ύπαρξη ενός χώρου δεδομένων  $D$  που ορίζεται από ένα σύνολο  $d$  διαστάσεων  $\{d_1, d_2, \dots, d_d\}$ , ένα σύνολο δεδομένων  $S$  στο χώρο  $D$  πλήθους  $|S|$  και η κορυφογραμμή  $SKY(S) \subseteq S$ :

- Έστω μια μονότονη συνάρτηση βαθμολόγησης  $S \Rightarrow R$ . Αν υπάρχει πλειάδα (ή σημείο)  $p \in S$  που μεγιστοποιεί το αποτέλεσμα της συνάρτησης αυτής, τότε η  $p$  αποτελεί οπωσδήποτε τμήμα της κορυφογραμμής, δηλαδή  $p \in SKY(S)$ .
- Για κάθε πλειάδα (ή σημείο)  $p \in SKY(S)$  υπάρχει τουλάχιστον μια μονότονη συνάρτηση βαθμολόγησης, τέτοια ώστε το αποτέλεσμά της να μεγιστοποιείται για την πλειάδα  $p$ .
- Όταν το  $S$  αποτελείται από μια διάσταση δεδομένων, ο υπολογισμός της μονοδιάστατης κορυφογραμμής  $SKY(S)$ , δεν παρουσιάζει κανένα ενδιαφέρον, καθώς είναι ισοδύναμος με τον υπολογισμό του MIN, MAX ή DISTINCT της συγκεκριμένης διάστασης.
- Για οποιαδήποτε πλειάδα  $p \in S$  ισχύει ότι  $p \not\prec^{pto} p$ . Η σχέση κυριαρχίας είναι μη αυτοπαθής.
- Έστω πλειάδες  $p, q \in S$  για τις οποίες ισχύει  $p \succ^{pto} q$ . Τότε πρέπει να ισχύει ότι  $q \not\prec^{pto} p$ . Η σχέση κυριαρχίας είναι αντισυμμετρική.
- Έστω πλειάδες  $p, q, v \in S$  για τις οποίες ισχύει  $p \succ^{pto} q$  και  $q \succ^{pto} v$ . Τότε πρέπει να ισχύει ότι  $p \succ^{pto} v$ . Η σχέση κυριαρχίας είναι μεταβατική.
- Έστω πλειάδες  $p, q, v \in S$  για τις οποίες ισχύει  $p \not\prec^{pto} q$  και  $q \not\prec^{pto} v$ . Τότε δεν είναι απαραίτητο να ισχύει ότι  $p \not\prec^{pto} v$ . Η σχέση κυριαρχίας δεν είναι αρνητικά μεταβατική. Αυτό έχει ως αποτέλεσμα η κυριαρχία Pareto  $\succ^{pto}$  να μην μπορεί να αναπαρασταθεί χρησιμοποιώντας οποιαδήποτε συνάρτηση βαθμολόγησης [42].

## 2.6 Ερωτήματα $k$ κορυφαίων σημείων

Τα ερωτήματα  $k$  κορυφαίων σημείων αποτελούν μια τεχνική υποστήριξης της διαδικασίας λήψης απόφασης. Όπως και τα ερωτήματα κορυφογραμμής, η τεχνική αυτή, εμπίπτει στο τμήμα εκείνο που αφορά τη στάθμιση και εξαγωγή αποτελεσμάτων, όπως περιγράφηκε στο Σχήμα 2.5. Η ουσιαστική διαφορά της με την τεχνική των ερωτημάτων κορυφογραμμής, αποτελεί το γεγονός ότι η τεχνική  $k$  κορυφαίων σημείων, απαιτεί τον ορισμό της βαρύτητας των κριτηρίων απόφασης. Η διαδικασία του ορισμού της βαρύτητας του κάθε κριτηρίου απόφασης, μπορεί να πραγματοποιηθεί είτε από τον ίδιο τον αποφασίζοντα, είτε από κάποιον ενδογενή ή εξωγενή παράγοντα. Τα ερωτήματα  $k$  κορυφαίων σημείων, αναφέρονται συχνά στη βιβλιογραφία και ως γραμμικά ερωτήματα (linear queries) [46].

### **Ορισμός 2.9 (Ερώτημα $k$ Κορυφαίων Σημείων) :**

*Τα ερωτήματα  $k$  κορυφαίων σημείων αναζητούν τις κορυφαίες  $k$  πλήθους πλειάδες, οι οποίες μεγιστοποιούν ή ελαχιστοποιούν το γραμμικό σταθμισμένο άθροισμα των τιμών των διαστάσεων των σημείων, ενός συνόλου δεδομένων. [46].*

Το γραμμικό σταθμισμένο άθροισμα των τιμών των διαστάσεων των πλειάδων του συνόλου δεδομένων, μπορεί να θεωρηθεί και ως μια συνάρτηση βαθμολόγησης των πλειάδων, ώστε αυτές να είναι δυνατόν να ταξινομηθούν μονοσήμαντα, ως προς την τιμή της συνάρτησης βαθμολόγησης. Η βαθμονόμηση και παράλληλα η ταξινόμηση των πλειάδων, μπορεί να βοηθήσει στην επιλογή των  $k$  καλύτερων πλειάδων, που μεγιστοποιούν ή ελαχιστοποιούν, ανάλογα την περίπτωση, την τιμή της συνάρτησης βαθμολόγησης. Η τεχνική του αθροίσματος των σταθμισμένων τιμών των διαστάσεων μιας πλειάδας, είναι σύμφωνη με την έννοια της προτίμησης, όπως αυτή αναφέρεται εκτενώς στη βιβλιογραφία [47, 48], και ακολουθεί τους κανόνες της πολυκριτηριακής θεωρίας αξίας ή χρησιμότητας που περιγράφηκε νωρίτερα. Η στάθμιση ενός κριτηρίου απόφασης, καθορίζει ουσιαστικά το βαθμό της σπουδαιότητάς του.

Ο βαθμός σπουδαιότητας των εφαρμοζόμενων κριτηρίων απόφασης για την αξιολόγηση των διαφόρων εναλλακτικών σεναρίων, καθορίζεται από το συντελεστή βαρύτητας που αποδίδεται στα κριτήρια αυτά. Ανάλογα με την περίπτωση, χρησιμοποιούνται είτε άμεσοι συντελεστές βαρύτητας, είτε έμμεσοι. Οι άμεσοι συντελεστές βαρύτητας χρησιμοποιούνται στην περίπτωση που ο αριθμός των κριτηρίων είναι μικρός και η επιλογή των συντελεστών βαρύτητας είναι δυνατή και εύκολη. Οι έμμεσοι συντελεστές βαρύτητας προσδιορίζονται με την ταξινόμηση των κριτηρίων κατά σειρά σπουδαιότητας, την απόδοση ενός συνολικού συντελεστή βαρύτητας, ή ενός μέγιστου συντελεστή βαρύτητας και στη συνέχεια τον ακριβή προσδιορισμό αυτών, σε

σχέση με το μέτρο που επιλέχθηκε. Επιπλέον, είναι δυνατή η χρήση κριτηρίων, στα οποία δεν έχει αποδοθεί συντελεστής βαρύτητας.

Οι συντελεστές βαρύτητας αντικατοπτρίζουν το σύστημα αξιών και προτιμήσεων του αποφασίζοντα. Δηλαδή, ο προσδιορισμός της σπουδαιότητας του κάθε κριτηρίου βασίζεται στην ιδιαίτερη σημασία που δίνουν οι ενδιαφερόμενοι φορείς για κάθε κριτήριο. Συνεπώς, ανάλογα με το είδος του προβλήματος, είναι δυνατό να παρουσιάζουν μεγαλύτερη σημασία για τους ενδιαφερόμενους φορείς τα περιβαλλοντικά κριτήρια σε σχέση με τα οικονομικά ή και το αντίστροφο. Έτσι, για τον προσδιορισμό των συντελεστών βαρύτητας απαιτείται η προσεκτική ιεραρχική ταξινόμηση των διαφόρων κριτηρίων από τον αποφασίζοντα.

### 2.6.1 Κατηγοριοποίηση ερωτημάτων k κορυφαίων σημείων

Το ερευνητικό ενδιαφέρον για τα ερωτήματα k κορυφαίων σημείων, είναι έντονο εδώ και τουλάχιστον μια δεκαετία. Αυτό έχει ως αποτέλεσμα, την ύπαρξη μεγάλου πλήθους τεχνικών προσέγγισης και λύσης ερωτημάτων τέτοιου είδους. Για το λόγο αυτό, στη συνέχεια παρουσιάζονται διάφορες πιθανές κατηγοριοποιήσεις που μπορούν να αναγνωριστούν στη βιβλιογραφία [49] για τους αλγορίθμους επίλυσης ερωτημάτων k κορυφαίων σημείων:

**Μοντέλο ερωτήματος:** Οι τεχνικές επεξεργασίας ερωτημάτων k κορυφαίων σημείων, μπορούν να ταξινομηθούν με βάση το μοντέλο ερωτήματος που υιοθετούν. Κάποιες τεχνικές πραγματεύονται ένα μοντέλο ερωτήματος επιλογής (selection query model), στο οποίο η βαθμολόγηση κάθε πλειάδας συνδέεται άμεσα με τις πλειάδες. Άλλες τεχνικές υποθέτουν τη χρήση ερωτημάτων συνένωσης (join query model), όπου η βαθμολόγηση των πλειάδων πραγματοποιείται σε συνενωμένα αποτελέσματα πλειάδων. Μια τρίτη κατηγορία, ασχολείται με το μοντέλο συγκεντρωτικών ερωτημάτων (aggregate query model), όπου οι εργασίες βαθμολόγησης και ταξινόμησης, πραγματοποιούνται σε ομάδες πλειάδων.

**Μέθοδος προσπέλασης δεδομένων:** Οι τεχνικές επεξεργασίας ερωτημάτων k κορυφαίων σημείων, μπορούν να ταξινομηθούν με βάση τις μεθόδους προσπέλασης που υποθέτουν ότι διατίθενται στις υποφαινόμενες πηγές δεδομένων. Για παράδειγμα ορισμένες τεχνικές υποθέτουν την ύπαρξη δυνατότητας για τυχαία προσπέλαση των δεδομένων, ενώ άλλες περιορίζονται μόνο στην ταξινομημένη σειριακή προσπέλαση.



**Επίπεδο εφαρμογής:** Οι τεχνικές επεξεργασίας ερωτημάτων  $k$  κορυφαίων σημείων, μπορούν να ταξινομηθούν με βάση το επίπεδο ενσωμάτωσης στα συστήματα βάσεων δεδομένων. Για παράδειγμα κάποιες τεχνικές έχουν υλοποιηθεί στο επίπεδο εφαρμογής, πάνω από το επίπεδο του συστήματος της βάσης δεδομένων, ενώ άλλες έχουν υλοποιηθεί ως τελεστές ερωτημάτων βάσεων δεδομένων.

**Αβεβαιότητα δεδομένων και ερωτημάτων:** Οι τεχνικές επεξεργασίας ερωτημάτων  $k$  κορυφαίων σημείων, μπορούν να ταξινομηθούν με βάση το επίπεδο της αβεβαιότητας που εμπλέκεται στα μοντέλα δεδομένων ή/και ερωτημάτων τους. Κάποιες τεχνικές παράγουν ακριβή αποτελέσματα, ενώ κάποιες άλλες επιτρέπουν την παραγωγή αποτελεσμάτων κατά προσέγγιση, ή πραγματεύονται μη βέβαια δεδομένα.

**Συνάρτηση κατάταξης:** Οι τεχνικές επεξεργασίας ερωτημάτων  $k$  κορυφαίων σημείων, μπορούν να ταξινομηθούν με βάση τον ορισμό που δίνουν στη συνάρτηση ταξινόμησης. Οι περισσότερες τεχνικές υιοθετούν τη χρήση μονότονων συναρτήσεων, ενώ λίγες είναι αυτές που επιτρέπουν τη χρήση γενικών συναρτήσεων.

### 2.6.2 Αντίστροφα ερωτήματα $k$ κορυφαίων σημείων

Η ερευνητική δραστηριότητα γύρω από τα ερωτήματα  $k$  κορυφαίων σημείων, παραδοσιακά αναλύει τη διαδικασία του σεναρίου χρήσης του ερωτήματος, από τη σκοπιά του καταναλωτή προϊόντων ή υπηρεσιών. Ωστόσο ενδιαφερόμενος για τα ερωτήματα τέτοιου τύπου δεν είναι μόνο ο καταναλωτής, αλλά και ο πάροχος. Η ερευνητική δραστηριότητα ανάλυσης των ερωτημάτων  $k$  κορυφαίων σημείων από τη σκοπιά του παρόχου, αναλύθηκε για πρώτη φορά το 2010 από τη Βλάχου κ.ά. [23]. Η πρόταση της συγκεκριμένης εργασίας περιελάμβανε τον ορισμό των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων, ως εξής:

**Ορισμός 2.10 (Αντίστροφα Ερωτήματα  $k$  Κορυφαίων Σημείων) :**

*Ένα αντίστροφο ερώτημα  $k$  κορυφαίων σημείων, ορίζεται από μια δοθείσα πλειάδα  $p$ , και αναζητά τις βαρύτητες των κριτηρίων απόφασης  $w$ , για τις οποίες η πλειάδα  $p$ , αποτελεί τμήμα του ερωτήματος  $k$  κορυφαίων σημείων [23].*

Ουσιαστικά τα αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων απαντούν στο ερώτημα: «Δοθέντος ενός πιθανού προϊόντος, ποιες είναι οι προτιμήσεις χρηστών για τις οποίες το προϊόν αυτό αποτελεί τμήμα του αποτελέσματος ενός ερωτήματος  $k$  κορυφαίων σημείων;». Τέτοιου είδους ερωτήματα μπορούν να αποκτήσουν πολύ μεγάλο ενδιαφέρον για κατασκευαστές προϊόντων ή παρόχους υπηρεσιών, οι οποίοι θέλουν



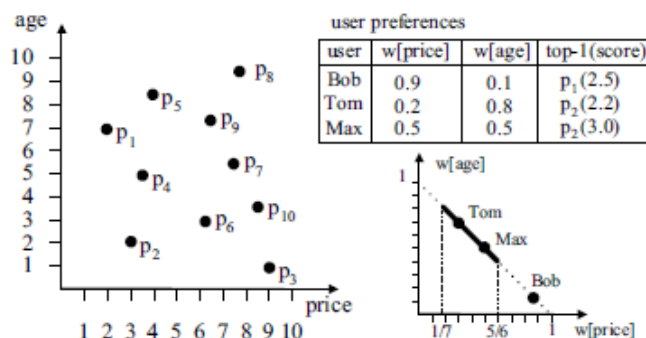
να μελετήσουν την επίδραση που μπορεί να έχει το προϊόν τους στην αγορά. Επίσης μπορεί να βοηθήσει τους κατασκευαστές να καθορίσουν τις προδιαγραφές του υπό ανάπτυξη προϊόντος τους, έτσι ώστε το προϊόν αυτό να έχει αυξημένες πιθανότητες επιτυχίας στην αγορά.

Τα αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων διαφέρουν από τα αντίστροφα ερωτήματα κοντινότερου γείτονα (reverse nearest neighbor query, RNN) [50]. Ένα RNN ερώτημα αναζητά τις πλειάδες εκείνες οι οποίες έχουν την πλειάδα αναφοράς ως την πλησιέστερη πλειάδα. Αντίθετα τα αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων αναζητούν τις συναρτήσεις απόστασης (σε όρους βαρύτητας) για τις οποίες η πλειάδα αναφοράς πληροί τις προϋποθέσεις ως ο  $k$  πλησιέστερος γείτονας του σημείου εκκίνησης του χώρου δεδομένων. Για το λόγο αυτό, οι υπάρχουσες προσεγγίσεις του προβλήματος των RNN ερωτημάτων, δεν μπορούν να έχουν εφαρμογή στα αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων.

Τα αντίστροφα ερωτήματα κορυφογραμμής (reverse skyline queries) [51] στοχεύουν στην αναγνώριση πελατών που ενδιαφέρονται για ένα προϊόν βασιζόμενα στη σχέση κυριαρχίας. Παρόλα αυτά οι προτιμήσεις των χρηστών εκφράζονται ως σημεία με τις ίδιες ιδιότητες των προϊόντων. Στην περίπτωση των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων όμως, οι προτιμήσεις των χρηστών μοντελοποιούνται με έναν πιο γενικό τρόπο (μόνο όσον αφορά τις βαρύτητες των κριτηρίων) και δε χρειάζεται να αντιστοιχιστούν μοναδικά σε ένα σημείο στο χώρο δεδομένων.

Στην ερευνητική εργασία των Βλάχου κ.ά [23] παρουσιάζονται δύο εκδοχές αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων: Η monochromatic και η bichromatic. Στην πρώτη δεν υπάρχει καμιά πληροφορία για τις προτιμήσεις των πελατών, και ο κατασκευαστής στοχεύει στην αναγνώριση της επίπτωσης που θα είχε στην αγορά ένα πιθανό προϊόν. Στη δεύτερη δίνεται ένα σύνολο δεδομένων που περιέχει τις προτιμήσεις των χρηστών και το αντίστροφο ερώτημα  $k$  κορυφαίων σημείων αναζητά τις προτιμήσεις εκείνες, για τις οποίες ένα πιθανό προϊόν κατατάσσεται εντός του αποτελέσματος του ερωτήματος  $k$  κορυφαίων σημείων.

Στο Σχήμα 2.8 το αντίστροφο ερώτημα 1 κορυφαίου σημείου του  $p_1$  (το  $p_1$  αποτελεί σημείο αναφοράς) έχει ως αποτέλεσμα το σύνολο βαρύτητας  $(0.9, 0.1)$  που έχει οριστεί από τον Bob. Για το σημείο  $p_2$  όμως δύο σύνολα βαρύτητας αποτελούν τμήμα του αποτελέσματος του αντίστροφου ερωτήματος 1 κορυφαίου σημείου: του Tom και του Max. Στην πραγματικότητα όλα τα σύνολα βαρύτητας που διαθέτουν τιμή  $w[price]$  στο τμήμα  $[\frac{1}{7}, \frac{5}{6}]$  ανήκουν στο αποτέλεσμα του αντίστροφου ερωτήματος 1 κορυφαίου σημείου του  $p_2$ . Αυτό το τμήμα της γραμμής  $w[price] + w[age] = 1$  αντιστοιχεί στο αποτέλεσμα της monochromatic εκδοχής του αντίστροφου ερωτήματος



Σχήμα 2.8: Παράδειγμα αντίστροφου ερωτήματος k κορυφαίων σημείων [23]

1 κορυφαίου σημείου για το  $p_2$ , ενώ το σύνολο  $\{(0.5, 0.5), (0.2, 0.8)\}$  είναι το αποτέλεσμα της bichromatic εκδοχής του αντίστροφου ερωτήματος 1 κορυφαίου σημείου για το  $p_2$  για το δοθέν σύνολο δεδομένων των προτιμήσεων των χρηστών.

Διαισθητικά ο χώρος λύσεων των αντίστροφων ερωτημάτων k κορυφαίων σημείων, είναι ο χώρος που ορίζεται από τις βαρύτητες  $w[price]$  και  $w[age]$ . Η μονοχρωματική εκδοχή επιστρέφει τμήματα του χώρου λύσεων και είναι χρήσιμη στην επιχειρηματική ανάλυση και πιο συγκεκριμένα στη διαδικασία εκτίμησης της επίδρασης ενός προϊόντος, όταν δεν έχει δοθεί σύνολο προτιμήσεων των χρηστών, αλλά η κατανομή αυτών είναι γνωστή. Στο παράδειγμα που δόθηκε, με την προϋπόθεση ότι οι προτιμήσεις των χρηστών ακολουθούν ομοιόμορφη κατανομή, η επίδραση ενός πιθανού προϊόντος  $p_2$ , στην αγορά μπορεί να εκτιμηθεί ως:  $(\frac{5}{6} - \frac{1}{7}) \times 100\% = 69\%$ . Από την άλλη, η bichromatic εκδοχή των αντίστροφων ερωτημάτων k κορυφαίων σημείων, έχει ακόμα ευρύτερη εφαρμογή, αφού αναγνωρίζει τους χρήστες που ενδιαφέρονται για ένα συγκεκριμένο προϊόν, δοθέντος ενός συνόλου προτιμήσεων των χρηστών. Για παράδειγμα, η καλύτερη στρατηγική για προώθηση του προϊόντος με βάση τα προφίλ των πελατών θα ήταν η διαφήμιση του προϊόντος  $p_1$  στον Bob και του προϊόντος  $p_2$  στον Tom και στον Max. Πρέπει ωστόσο να επισημανθεί ότι ένα κενό αποτέλεσμα ερωτήματος για ένα προϊόν (πχ. για το προϊόν  $p_3$ ), υποδεικνύει ότι το προϊόν αυτό δεν είναι ενδιαφέρον για τους πελάτες με βάση τις προτιμήσεις τους. Η bichromatic εκδοχή μπορεί να χρησιμοποιηθεί σε πρακτικές εφαρμογές, και είναι ευκολότερο να ενσωματωθεί σε συστήματα βάσεων δεδομένων, ενώ η μονοχρωματική εκδοχή, παρέχει κυρίως μια γεωμετρική ερμηνεία και βοηθά στην διαισθητική κατανόηση του προβλήματος.

## Κεφάλαιο 3

# Ανασκόπηση βιβλιογραφίας

Στο προηγούμενο κεφάλαιο έγινε εισαγωγή στις βασικές έννοιες που πραγματεύεται η παρούσα ερευνητική εργασία. Αναλύθηκαν εκτενώς όλα τα βασικά σημεία από κάθε ερευνητική περιοχή, ώστε να επεξηγηθούν όλες οι τεχνολογίες και οι μέθοδοι που χρησιμοποιούνται στη συνέχεια. Ωστόσο απαραίτητο είναι να παρουσιαστεί και μια σειρά από ερευνητικές εργασίες, οι οποίες προτείνουν τεχνικές υποστήριξης της διαδικασίας λήψης αποφάσεων.

Έτσι, στο κεφάλαιο αυτό παρουσιάζεται ένα σύνολο τέτοιων μεθόδων, οι οποίες υποβοηθούνται από τη χρήση υπολογιστικών συστημάτων. Αναλύονται διεξοδικά όλες οι υπάρχουσες μέθοδοι επίλυσης αντίστροφων ερωτημάτων κ κορυφαίων σημείων, ενώ παράλληλα παρουσιάζονται τα πλεονεκτήματα και τα μειονεκτήματά τους. Μέσα από τη διαδικασία αυτή, αναδεικνύονται διάφορες ανοικτές ερευνητικές περιοχές, που μπορούν μελλοντικά να μελετηθούν. Με αυτόν τον τρόπο πραγματοποιείται επιλογή μιας ανοικτής ερευνητικής περιοχής, ώστε να αποτελέσει αντικείμενο της παρούσας ερευνητικής εργασίας.

Στόχοι του κεφαλαίου αυτού είναι:

- Η παρουσίαση, κριτική ανάλυση και σύγκριση των συναφών ερευνητικών εργασιών.
- Η αναγνώριση των ανοικτών ερευνητικών περιοχών.
- Η επιλογή μιας ανοικτής ερευνητικής περιοχής, ως αντικείμενο της παρούσας εργασίας.

### 3.1 Ερωτήματα κορυφογραμμών σε ροές δεδομένων

Πριν παρουσιαστούν οι ερευνητικές εργασίες που αφορούν τα ερωτήματα  $k$  κορυφαίων σημείων, απαραίτητο είναι να μελετηθούν πρώτα οι εργασίες που αφορούν σε ένα παρεμφερές ερευνητικό θέμα: τα ερωτήματα κορυφογραμμών σε ροές δεδομένων. Με τον τρόπο αυτό, γίνεται δυνατή η σύγκριση των προσεγγίσεων σε σχέση με τις αντίστοιχες των ερωτημάτων  $k$  κορυφαίων σημείων. Στην εργασία αυτή μελετάται μόνο η περίπτωση του μοντέλου αποκλειστικής προσάρτησης (append only model) [52] για τις ροές δεδομένων. Σε αυτό το πλαίσιο, οι πλειάδες εισέρχονται στο σύστημα, και θεωρούνται ενεργές μόνο για τη χρονική διάρκεια που αυτές ανήκουν σε ένα ορισμένο συρόμενο χρονικό παράθυρο (sliding window). Υπάρχουν δύο ειδών συρόμενα χρονικά παράθυρα: τα βασιζόμενα στο πλήθος των πλειάδων, όπου περιέχουν τις πιο προσφάτες  $N$  πλήθους πλειάδες, και τα βασιζόμενα στο χρόνο, τα οποία περιέχουν τις πλειάδες εκείνες που εισήχθησαν το σύστημα σε ένα ορισμένο χρονικό διάστημα, καλύπτοντας τις πιο πρόσφατες χρονικές στιγμές.

Το 2005 πραγματοποιήθηκε η πρώτη ερευνητική προσέγγιση στο συγκεκριμένο χώρο [53] από τους Lin κ.ά. [54]. Μελετήθηκε η δυνατότητα υπολογισμού κορυφογραμμών  $n$  μεγέθους σε συρόμενα χρονικά παράθυρα μεγέθους  $N$ , όπου  $n < N$ . Στην εργασία αυτή χρησιμοποιείται μια δομή που ονομάζεται γράφος κυριαρχίας (dominance graph), για την αναπαράσταση των κρίσιμων κυριαρχιών. Η κυριαρχία ενός σημείου  $e$  προς ένα άλλο σημείο  $\acute{e}$  είναι κρίσιμη μόνο εφόσον το  $e$  είναι το νεότερο σημείο (αλλά αρχαιότερο του  $\acute{e}$ ) το οποίο κυριαρχεί επί του  $\acute{e}$ . Ο γράφος κυριαρχίας είναι κατευθυνόμενος έτσι ώστε να ορίζει τις κρίσιμες κυριαρχίες και αναπαρίσταται από κόμβους που περιγράφουν την στιγμή άφιξης του εκάστοτε σημείου. Παράλληλα διατηρείται μια δομή R-Tree στην κύρια μνήμη, ώστε να επιτυγχάνεται γρήγορος υπολογισμός της κυριαρχίας ενός νεοεισερχόμενου σημείου επί των σημείων που υπάρχουν ήδη στο συρόμενο χρονικό παράθυρο.

Η δεύτερη εργασία στο χώρο των ερωτημάτων κορυφογραμμής σε ροές δεδομένων πραγματοποιήθηκε από τους Tao κ.ά. [55]. Σε αυτήν αρχικά ορίζεται ένα σύνολο κρίσιμων ιδιοτήτων οι οποίες βοηθούν στο σχεδιασμό του συστήματος παρακολούθησης ερωτημάτων κορυφογραμμής. Συντηρούνται δύο ξεχωριστά σύνολα δεδομένων: το  $DB_{REST}$  όπου αποθηκεύονται οι πλειάδες που δεν αποτελούν τμήμα της κορυφογραμμής και  $DB_{SKY}$ , όπου βρίσκεται το σύνολο των πλειάδων κορυφογραμμής. Το προτεινόμενο σύστημα αποτελείται από δύο μονάδες (modules): την προεπεξεργαστική μονάδα (pre-processing module, PM) και τη μονάδα συντήρησης (maintenance module, MM). Παράλληλα προτείνονται δύο πιθανές προσεγγίσεις για τη λύση του

προβλήματος. Η χαλαρή μέθοδος (lazy method) εκτελεί το PM κάθε φορά που καταφθάνει νέο σημείο και το MM κάθε φορά που λήγει κάποιο σημείο από το  $DB_{SKY}$ . Το PM ελέγχει αν τα νεοεισερχόμενα σημεία ανήκουν στην  $DB_{SKY}$ . Το MM αφαιρεί από την  $DB_{SKY}$  τα σημεία που έχουν λήξει και μεταφέρει σημεία από την  $DB_{REST}$  στην  $DB_{SKY}$ . Ωστόσο δεν έχει μηχανισμό αφαίρεσης των περιττών σημείων από την  $DB_{REST}$ . Η πρόθυμη μέθοδος (eager method) εκτελεί το PM κάθε φορά που καταφθάνει νέο σημείο και το MM κάθε φορά που λαμβάνει χώρα ένα γεγονός (event). Τα γεγονότα αποθηκεύονται σε μια λίστα της μορφής  $\langle r, r.t_{ev}, tag \rangle$ , όπου  $r$  είναι ένας pointer προς το σημείο το οποίο αφορά το συγκεκριμένο entry της λίστας,  $r.t_{ev}$  είναι το χρονικό σημείο εκτέλεσης του event, και το tag, είναι ο τύπος του event: EX σε περίπτωση που το σημείο είναι ήδη skyline, SK σε αντίθετη περίπτωση.

Βασιζόμενοι στο μοντέλο κατανεμημένων ροών δεδομένων, οι Sun κ.ά. [56] πρότειναν έναν αλγόριθμο (BOCS) για ερωτήματα κορυφογραμμών σε ροές δεδομένων. Σε ένα τέτοιο σύστημα, θεωρείται η ύπαρξη ενός συνόλου από εξυπηρετητές, καθένας από τους οποίους παράγει ένα υποσύνολο των δεδομένων της ροής. Επιπρόσθετα υπάρχει ένας κεντρικός εξυπηρετητής που επικοινωνεί με τους απομακρυσμένους εξυπηρετητές και είναι υπεύθυνος για τον υπολογισμό του τελικού αποτελέσματος ενός ερωτήματος. Αυτή η τοπολογία είναι παρόμοια με αυτήν που ακολουθούν τα συστήματα πολύ μεγάλης κατανομής, και πιο συγκεκριμένα στην πλήρως συνδεδεμένη τοπολογία δικτύου. Η κύρια διαφορά, είναι ότι στην περίπτωση των ροών δεδομένων, στόχος αποτελεί η διαρκής παρακολούθηση της κορυφογραμμής με την πάροδο του χρόνου, αντί του στόχου του υπολογισμού της κορυφογραμμής σε μια δοθείσα χρονική στιγμή. Στον BOCS κάθε απομακρυσμένος εξυπηρετητής υπολογίζει την τοπική του κορυφογραμμή, χρησιμοποιώντας έναν κεντρικοποιημένο αλγόριθμο. Έτσι σε κάθε χρονική στιγμή, μόνο τα σημεία που προστέθηκαν στην εκάστοτε τοπική κορυφογραμμή αποστέλλονται στον κεντρικό εξυπηρετητή. Με βάση την ιδιότητα της προσθετικότητας του τελεστή κορυφογραμμής, τα σημεία που δε συμπεριλαμβάνονται στην τοπική κορυφογραμμή, δεν αποστέλλονται στον κεντρικό εξυπηρετητή, μειώνοντας με τον τρόπο αυτό το συνολικό κόστος επικοινωνίας. Έτσι ο κεντρικός εξυπηρετητής εφαρμόζει εκ νέου έναν κεντρικοποιημένο αλγόριθμο υπολογισμού κορυφογραμμής στο σύνολο των σημείων που συνέλεξε από τους απομακρυσμένους και υπολογίζει αποτελεσματικά την τελική κορυφογραμμή για κάθε χρονική στιγμή.

Άλλες εργασίες στον ερευνητικό αυτό χώρο, αποτελούν η εργασία των Sarka κ.ά. [57] οι οποίοι ασχολήθηκαν με τον υπολογισμό κορυφογραμμών σε ροές κατηγορικών δεδομένων (κατηγορικά δεδομένα είναι αυτά που δεν μπορούν να λάβουν ποσοτική τιμή, παρά μόνο ποιοτική), η εργασία των Wenlin κ.ά [58] όπου προτείνεται καινοτόμος μηχανισμός επίλυσης ερωτημάτων τέτοιου τύπου σε χώρους δεδομένων δύο διαστάσεων, και η εργασία των Fang κ.ά. [59] στην οποία προτείνεται αλγόριθμος

υπολογισμού κορυφογραμμών σε ροές δεδομένων που περιέχουν κατηγορικά και μη δεδομένα.

## 3.2 Ερωτήματα k κορυφαίων σημείων σε ροές δεδομένων

### 3.2.1 Κατανεμημένη διαχείριση ερωτημάτων

Το 2003 η δουλειά των Babcock κ.ά [14] προσέγγισε για πρώτη φορά την ανάλυση ερωτημάτων k κορυφαίων σημείων σε ροές δεδομένων και μάλιστα σε κατανεμημένο περιβάλλον. Ωστόσο το μοντέλο που υιοθέτησαν στην εργασία τους, δεν ακολουθεί το μοντέλο των συρόμενων χρονικών παραθύρων (sliding window). Αντίθετα υιοθέτησαν το μοντέλο των «κατανεμημένων ταμειακών μηχανών» [60]. Πιο συγκεκριμένα στο μοντέλο αυτό, διατίθενται ένα σύνολο από εξυπηρετητές παρακολούθησης και ένας κεντρικός εξυπηρετητής συντονισμού.

Στο σύστημα υπάρχουν ένα πεπερασμένο σύνολο από πλειάδες, οι οποίες ορίζονται από ένα σύνολο αριθμητικών τιμών  $V = \{V_1, V_2, \dots, V_n\}$ . Κάθε εξυπηρετητής παρακολούθησης μπορεί να λάβει μέσω της ροής δεδομένων αλλαγές στις τιμές οποιασδήποτε υπάρχουσας πλειάδας (θετική ή αρνητική μεταβολή), καθώς και πιθανές προσθήκες νέων πλειάδων στο σύστημα. Το προτεινόμενο μοντέλο παρέχει τη δυνατότητα υπολογισμού των k κορυφαίων πλειάδων με 100% ακρίβεια ή με ένα ποσοστό ακρίβειας οριζόμενο δυναμικά από το χρήστη. Ο κεντρικός εξυπηρετητής είναι υπεύθυνος για τη συνεχή παρακολούθηση των αποτελεσμάτων ενός ερωτήματος k κορυφαίων σημείων σε μια οριοθετημένη από το χρήστη περιοχή ανοχής σφάλματος.

Η συνολική προσέγγιση που ακολουθείται στην παραπάνω εργασία έχει ως στόχο τον υπολογισμό και τη συντήρηση στον κεντρικό εξυπηρετητή, ένα αρχικά 100% ορθό σύνολο  $T$  με τα k κορυφαία σημεία ενός ερωτήματος. Αυτό επιτυγχάνεται θέτοντας τον κεντρικό εξυπηρετητή υπεύθυνο για τον ορισμό αριθμητικών περιορισμών για τα επιμέρους δεδομένα σε κάθε εξυπηρετητή παρακολούθησης. Καθώς αλλαγές στις τιμές των δεδομένων λαμβάνουν χώρα, οι κόμβοι παρακολούθησης παρακολουθούν τις αλλαγές στα επιμέρους δεδομένα τους, διασφαλίζοντας ότι κάθε αριθμητικός περιορισμός εξακολουθεί να ικανοποιείται. Για όσο καιρό οι αριθμητικοί περιορισμοί παραμένουν απαραβίαστοι για όλους τους κόμβους παρακολούθησης, δεν απαιτείται κανένας είδους επικοινωνία με τον κεντρικό εξυπηρετητή, για την επιβεβαίωση της ορθότητας του συνόλου  $T$ . Από την άλλη, εάν ένας ή περισσότεροι περιορισμοί παραβιαστούν, λαμβάνει χώρα μια κατανεμημένη διαδικασία αποκαλούμενη ως «Επίλυση», μεταξύ του κεντρικού εξυπηρετητή και των κόμβων παρακολούθησης, με σκοπό να

καθοριστεί αν το σύνολο  $T$  είναι ακόμα ορθό, ή να πραγματοποιήσει τις απαραίτητες αλλαγές. Στη συνέχεια, αν το σύνολο  $T$  έχει μεταβληθεί, ο κεντρικός εξυπηρετητής θέτει νέους αριθμητικούς περιορισμούς στους κόμβους παρακολούθησης, με σκοπό τη συνεχιζόμενη ορθή παρακολούθηση του συνόλου  $T$ , και δεν πραγματοποιείται οποιαδήποτε άλλη ενέργεια μέχρι να παραβιαστεί κάποιος από τους νέους αριθμητικούς περιορισμούς.

Η μέθοδος που περιγράφηκε είναι πολύ αποδοτική σε περιβάλλοντα υψηλής κατανομής. Αυτό συμβαίνει γιατί οι επικοινωνίες μεταξύ των διαφόρων κόμβων, έχουν μειωθεί στο ελάχιστο: Μια αλληλουχία από αμφίδρομες επικοινωνίες μπορεί να εκκινήσει, μόνο έπειτα από την παραβίαση ενός ή περισσότερων αριθμητικών περιορισμών. Η συγκεκριμένη εργασία αποτελεί ένα πολύ καλό παράδειγμα συνεχούς παρακολούθησης ερωτημάτων  $k$  κορυφαίων σημείων σε κατανομημένες ροές δεδομένων.

### 3.2.2 Αλγόριθμοι TMA και SMA

Το 2006 ο Μουρατίδης κ.ά. [13] πρότειναν έναν αποδοτικό αλγόριθμο για τον υπολογισμό ερωτημάτων  $k$  κορυφαίων σημείων σε ροές δεδομένων. Ο προτεινόμενος αλγόριθμος, χρησιμοποιεί συρόμενα χρονικά παράθυρα βασιζόμενα στο πλήθος των πλειάδων (count based sliding window), ενώ μπορεί να υποστηρίξει το ίδιο αποδοτικά και τους δύο τύπους συρόμενων χρονικών παραθύρων. Υποστηρίζεται ακόμη η δυναμική εισαγωγή στο σύστημα νέων ερωτημάτων  $q$ , οι οποίες ορίζονται από μια μονότονη συνάρτηση  $f$  και ένα  $k$  πλήθος αναμενόμενων αποτελεσμάτων.

Αρχικά διατυπώνεται ένα σύνολο ιδιοτήτων, οι οποίες στη συνέχεια συμβάλλουν στο σχεδιασμό και τη βελτίωση του προτεινόμενου συστήματος. Να σημειωθεί ότι για τη συγκεκριμένη εργασία οι υψηλότερες βαθμολογίες είναι προτιμότερες.

**Περιοχή επιρροής:** Για ένα ερωτήμα  $q$ , σε μια συγκεκριμένη χρονική στιγμή  $t$ , το σημείο  $p_k$  διαθέτει την  $k$ -οστή χαμηλότερη βαθμολογία, σύμφωνα με τη συνάρτηση του  $q$ , έστω  $f(x_1, x_2) = x_1 + 2 \times x_2$ . Η γραμμή που ορίζεται από το  $score(p_k) = f(x_1, x_2) = x_1 + 2 \times x_2$  χωρίζει το χώρο δεδομένων σε δύο κομμάτια. Το κομμάτι που περιέχει το όριο του χώρου των δεδομένων, αποτελεί την περιοχή επιρροής. Αν προστεθεί κάποια πλειάδα στην περιοχή αυτή, θα προκαλέσει την μείωση του μεγέθους του χώρου αυτού, ενώ αν αφαιρεθεί κάποια πλειάδα, θα προκαλέσει την αύξησή του.

**Άνω όριο βαθμολόγησης:** Όλες οι πλειάδες που ανήκουν σε ένα ορθογώνιο  $R$ , μπορούν να λάβουν βαθμολόγηση το πολύ  $maxscore(R)$ , η οποία αντιστοιχεί στη βαθμολογία του πάνω δεξιά άκρου του ορθογωνίου.



**k skyband ερωτήματα:** Το αποτέλεσμα ενός οποιουδήποτε ερωτήματος  $k$  κορυφαίων σημείων, περιέχεται οπωσδήποτε και σε ένα ερώτημα  $k$  skyband.

Όλα τα σημεία  $p$  που εισέρχονται στο σύστημα, αποθηκεύονται στην κύρια μνήμη σε μια δομή FIFO (First In, First Out) ενώ παράλληλα συντηρείται μια απλή δομή πλέγματος, ως ευρετήριο αυτών των σημείων. Για κάθε κελί του πλέγματος, συντηρείται η λίστα με τα ερωτήματα εκείνα που το συγκεκριμένο κελί (ή χώρος δεδομένων) αποτελεί περιοχή επιρροής.

Προτείνονται δύο πιθανοί αλγόριθμοι: ο TMA και ο SMA. Στον TMA αλγόριθμο, υπολογίζεται αρχικά το σύνολο των πλειάδων που αποτελούν τμήμα του αποτελέσματος, για τις πλειάδες εκκίνησης. Για κάθε νεοεισερχόμενη πλειάδα  $p$ , συγκρίνεται η βαθμολογία της σε σχέση με τη βαθμολογία του  $k$ -οστού στοιχείου για ένα ερώτημα  $q$ . Αν η πρώτη βαθμολογία είναι καλύτερη της δεύτερης, τότε εκτελείται ένας αλγόριθμος εκ νέου υπολογισμού του αποτελέσματος του ερωτήματος  $q$ . Παρόμοια διαδικασία εκτελείται σε περίπτωση λήξης κάποιου σημείου. Αν το λήξαν σημείο ανήκει σε ένα κελί  $c$ , και η περιοχή αυτή αποτελεί περιοχή επιρροής για κάποιο σημείο  $q$ , τότε εκτελείται ο αλγόριθμος εκ νέου υπολογισμού του αποτελέσματος του ερωτήματος  $q$ .

Ο αλγόριθμος SMA προσπαθεί να βελτιώσει την απόδοση του TMA, με την αξιοποίηση της ιδιότητας των  $k$  skyband ερωτημάτων. Για κάθε ερώτημα  $q$  συντηρείται στην κύρια μνήμη το αποτέλεσμα του skyband ερωτήματος μεγέθους  $k$ . Κάθε νεοεισερχόμενο σημείο  $p$ , αντιστοιχίζεται σε ένα κελί  $c$  του πλέγματος ευρετηρίου. Αν το  $c$  αποτελεί περιοχή επιρροής για κάποιο ερώτημα  $q$ , τότε συγκρίνεται η βαθμολογία της  $p$  σε σχέση με τη βαθμολογία του  $k$ -οστού στοιχείου για ένα ερώτημα  $q$ . Αν η πρώτη βαθμολογία είναι μεγαλύτερη της δεύτερης, ενημερώνεται το skyband του  $q$ . Αντίστοιχα για κάθε σημείο  $p$  που λήγει, ελέγχεται αν το κελί στο οποίο ανήκει είναι χώρος επιρροής για κάποιο  $q$ , και αν είναι, αφαιρείται το σημείο  $p$  από το skyband του. Τελικά για κάθε skyband που επηρεάστηκε υπολογίζονται τα  $k$  πλήθους σημεία με τη μεγαλύτερη βαθμολογία. Αυτά αποτελούν το τρέχον αποτέλεσμα του αλγορίθμου. Εφόσον το skyband ενός σημείου  $q$ , έχει λιγότερα από  $k$  πλήθους αποτελέσματα, τότε εκτελείται ο αλγόριθμος εκ νέου υπολογισμού του αποτελέσματος του ερωτήματος  $q$ .

### 3.2.3 Διαχείριση ερωτημάτων σε αβέβαιες ροές

Στην εργασία των Jin κ.ά. το 2008 [61], παρουσιάστηκε για πρώτη φορά μέθοδος για τον υπολογισμό των ερωτημάτων  $k$  κορυφαίων σημείων σε ροές αβέβαιων (πιθανοτικών) δεδομένων. Αβέβαια δεδομένα καλούνται αυτά, τα οποία η πληροφορία τους



περιέχουν κάποια πιθανότητα σφάλματος. Για παράδειγμα μια εφαρμογή παρακολούθησης της κίνησης με χρήση ραντάρ για τον εντοπισμό της ταχύτητας των κινούμενων αυτοκινήτων, μπορεί να περιέχει σφάλματα στις μετρήσεις που προκαλούνται από κοντινές γραμμές υψηλής τάσης, από παρεμβολές άλλων αυτοκινήτων, ανθρώπινα σφάλματα κλπ. Σε τέτοιου είδους σενάρια, τα δεδομένα θεωρούνται ορθά με μια συγκεκριμένη πιθανότητα σφάλματος.

Στόχος της παραπάνω εργασίας αποτέλεσε η δημιουργία μεθόδου, η οποία θα ήταν αποδοτική τόσο όσον αφορά το χρόνο εκτέλεσης όσο και το μέγεθος του διαθέσιμου χώρου δεδομένων που απαιτείται. Οι δύο αυτοί στόχοι έχουν πολύ μεγάλη σημασία σε συστήματα ροών δεδομένων, διότι ο όγκος των εισερχόμενων ή/και των εξερχόμενων δεδομένων μπορεί να είναι πολύ μεγάλος. Οι Soliman κ.ά. [22] πρότειναν δύο ορισμούς για τα ερωτήματα τέτοιου τύπου, τον  $U\text{-Top}k$  και τον  $U\text{-}k\text{Ranks}$ . Οι Hua κ.ά. [62] πρότειναν έναν ακόμα ορισμό με όνομα  $PT\text{-}k$ . Η ανάλυση των δοθέντων ορισμών, είναι εκτός του σκοπού της παρούσης εργασίας.

Οι Jin κ.ά. πρότειναν τη δημιουργία ενός ενοποιημένου πλαισίου για τη διαρκή επεξεργασία ερωτημάτων  $k$  κορυφαίων σημείων σε ροές αβεβίων δεδομένων, καλύπτοντας όλους τους προαναφερθέντες ορισμούς. Το προτεινόμενο πλαίσιο αποτελείται από πολλά αποδοτικά στοιχεία, τόσο όσον αφορά το χρησιμοποιούμενο χώρο, όσο και τον απαιτούμενο χρόνο. Όπως αποδεικνύεται από την εργασία τους, παρόλο που είναι αρκετά εύκολος ο χειρισμός των εισερχόμενων πλειάδων, οι εξερχόμενες πλειάδες απαιτούν πολύ δύσκολους χειρισμούς. Από τους στόχους που τέθηκαν για τη δημιουργία του μηχανισμού, απαιτείται ο προσεκτικός σχεδιασμός των στοιχείων του, ώστε να αποθηκεύεται μόνο το ελάχιστο πλήθος της πληροφορίας που χρειάζεται για την ορθή παρακολούθηση της εκτέλεσης του ερωτήματος.

Περιγράφεται ο τρόπος με τον οποίο μπορούν να προσαρμοστούν οι ήδη υπάρχουσες τεχνικές για τον υπολογισμό του ερωτήματος σε ροές δεδομένων, οι οποίες εισάγουν δεδομένα στο σύστημα. Με τον τρόπο αυτό, παρέχεται λύση στο πρόβλημα των ροών δεδομένων που δεν έχουν όρια (unbounded data streams - landmark windows), δηλαδή των ροών που τροφοδοτούν το σύστημα με πλειάδες που δε λήγουν ποτέ. Ωστόσο η υποστήριξη ροών δεδομένων με όρια, και ο χειρισμός των διαγραφών των ληγμένων πλειάδων, αποτελεί μια πολύ πιο δύσκολη διαδικασία από την εισαγωγή πλειάδων. Στην πραγματικότητα αν οι διαγραφές είναι απαραίτητες, δεν υπάρχει καλύτερος τρόπος από την αποθήκευση όλων των δεδομένων που παραμένουν ενεργά στο συρόμενο χρονικό παράθυρο, αφού κάθε πλειάδα θα μπορούσε να έχει την ευκαιρία να αποτελέσει τμήμα του αποτελέσματος. Έτσι, με σκοπό τη μείωση του χώρου που απαιτείται για την αποθήκευση των πλειάδων αυτών, αξιοποιείται η πολύ χρήσιμη ιδιότητα των συρόμενων χρονικών παραθύρων, η οποία αναφέρεται στο γεγονός ότι

οι πλειάδες που εισέρχονται πρώτες, εξέρχονται πρώτες. Προτείνονται επίσης ένα σύνολο απο τεχνικές, οι οποίες βελτιώνουν σημαντικά την πολυπλοκότητα του χώρου και του χρόνου. Οι τεχνικές αυτές περιλαμβάνουν τη συμπίεση των δεδομένων, το buffering και τα εκθετικά ιστογράμματα. Αποδεικνύεται αναλυτικά ότι παρόλο που η χρήση των τεχνικών αυτών μειώνει τις απαιτήσεις σε χώρο δεδομένων σε σχέση με το μέγεθος του συρόμενου χρονικού παραθύρου, επιτυγχάνεται ο μη επηρεασμός του χρόνου εκτέλεσης του προτεινόμενου αλγορίθμου. Παρατηρούνται ακόμη, με τα εργαστηριακά τους πειράματα, σημαντικές βελτιώσεις κατά μία τάξη μεγέθους, στην απόδοση του προτεινόμενου αλγορίθμου, σε σχέση με άλλες λύσεις που χρησιμοποιούν το μοντέλο του συρόμενου χρονικού παραθύρου. Αναλύεται επίσης διεξοδικά ο τρόπος με τον οποίο μπορούν να προσαρμοστούν όλοι οι ορισμοί των ερωτημάτων που περιγράφηκαν πιο πάνω, στο προτεινόμενο πλαίσιο της εργασίας.

### 3.3 Αντίστροφα ερωτήματα k κορυφαίων σημείων

#### 3.3.1 Αλγόριθμος RTA

Το 2011 οι Βλάχου κ.ά. [63] μελέτησαν εκτενώς τα αντίστροφα ερωτήματα k κορυφαίων σημείων, και συγκεκριμένα τις μονοχρωματικές και δίχρωματικές εκδοχές του. Παρουσιάστηκαν εκ νέου οι ορισμοί για καθένα από τις δύο εκδοχές, και προτάθηκαν μέθοδοι αντιμετώπισης του εκάστοτε προβλήματος. Στην εν λόγω εργασία, θεωρείται πως οι χαμηλότερες τιμές βαθμολόγησης είναι προτιμότερες.

Για τη μονοχρωματική εκδοχή, αναγνωρίζονται αρχικά ορισμένες ιδιότητες που ισχύουν για τα ερωτήματα αυτά, και στη συνέχεια προτείνεται ένας αλγόριθμος για τη λύση του προβλήματος σε χώρο δύο διαστάσεων. Μια από αυτές τις ιδιότητες, αφορά το κριτήριο της κυριαρχίας Pareto. Πιο συγκεκριμένα, όταν ένα σημείο του δισδιάστατου χώρου  $p$  κυριαρχεί έναντι του σημείου αναφοράς  $q$  ( $p \succ_{pto} q$ ), τότε το σημείο  $p$  θα έχει πάντα καλύτερη βαθμολόγηση από το  $q$  για οποιαδήποτε μονότονη συνάρτηση βαθμολόγησης. Και αντιστρόφως, όταν το σημείο αναφοράς  $q$  κυριαρχεί έναντι του σημείου  $p$ , τότε το σημείο  $q$  θα έχει πάντα καλύτερη βαθμολόγηση από το  $p$ , για οποιαδήποτε μονότονη συνάρτηση βαθμολόγησης. Έτσι ο προτεινόμενος αλγόριθμος για τα μονοχρωματικά ερωτήματα, δεν εξετάζει τα σημεία  $p$  για τα οποία ισχύει μία από τις παραπάνω συνθήκες, καθώς αυτά δεν επηρεάζουν το τελικό αποτέλεσμα. Ο προτεινόμενος αλγόριθμος για κάθε εξεταζόμενο σημείο  $p$ , ορίζει το μια διακριτή βαρύτητα  $w_i$  η οποία αντιστοιχεί στην κλίση του ευθύγραμμου τμήματος που είναι κάθετο στο τμήμα  $p-q$ . Το τελικό αποτέλεσμα υπολογίζεται με βάση όλους τους πιθανούς συνδυασμούς που μπορούν να πραγματοποιηθούν χρησιμοποιώντας τα  $w_i$  που εξήχθησαν,

με τρόπο τέτοιο ώστε να δημιουργούνται διαστήματα βαρυτήτων (π.χ. το διάστημα  $[w_0, w_1]$ ), για τα οποία επιβεβαιώνεται ότι το σημείο  $q$  αποτελεί τμήμα του top- $k$  αποτελέσματος. Το πρόβλημα των μονοχρωματικών ερωτημάτων, σε πλήθος διαστάσεων μεγαλύτερο του δύο, αναμένεται να μελετηθεί μελλοντικά.

Για τη bichromatic εκδοχή του αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων, θα μπορούσε να χρησιμοποιηθεί ο αφελής (naïve) τρόπος, της εξέτασης όλων των υπαρκτών βαρυτήτων  $w_i$ , ώστε να υπολογιστεί για κάθε μία από αυτές, το top- $k$  αποτέλεσμα. Με τον τρόπο αυτό, μπορεί να απαντηθεί το ερώτημα, αν το σημείο  $q$ , αποτελεί τμήμα του top- $k$  αποτελέσματος για καθεμιά από τις βαρύτητες  $w_i$ . Ωστόσο για τη βελτίωση της αποδοτικότητας ενός τέτοιου συστήματος, προτάθηκε η χρήση του αλγορίθμου Reverse top- $k$  Threshold Algorithm (RTA). Σκοπός του RTA είναι η μείωση των ερωτημάτων top- $k$  που απαιτείται να εκτελεστούν, με βάση την παρατήρηση ότι παρόμοιες βαρύτητες, μπορούν να επιστρέψουν παρόμοια αποτελέσματα. Πιο συγκεκριμένα, ο αλγόριθμος RTA, εκτελεί για το πρώτο εξεταζόμενο  $w_i$ , τον υπολογισμό του top- $k$  ερωτήματος. Το αποτέλεσμα αυτό, αποθηκεύεται σε ένα buffer, και συγκρίνεται η βαθμολόγηση του σημείου  $q$  σε σχέση με τη μεγαλύτερη βαθμολογία που συναντάται στα σημεία του buffer. Αν το  $q$  διαθέτει καλύτερη βαθμολογία, τότε το  $w_i$  προστίθεται στο τελικό αποτέλεσμα. Στη συνέχεια υπολογίζεται η μεγαλύτερη βαθμολογία των σημείων που βρίσκονται στο buffer (κατώφλι), με βάση την επόμενη βαρύτητα  $w_{i+1}$  και συγκρίνεται με τη βαθμολογία του σημείου  $q$  χρησιμοποιώντας την ίδια βαρύτητα. Εφόσον το  $q$  έχει χειρότερη βαθμολογία από το κατώφλι, τότε η βαρύτητα  $w_i$  μπορεί να απορριφθεί.

Ο αλγόριθμος RTA μπορεί να βελτιωθεί περαιτέρω, αν προηγουμένως πραγματοποιηθεί ταξινόμηση των βαρυτήτων, με τρόπο τέτοιο ώστε να μεγιστοποιείται η πιθανότητα της απόρριψης βαρυτήτων από το κατώφλι. Για το λόγο αυτό, προτείνεται ένας αλγόριθμος που μεγιστοποιεί την ομοιότητα όλων των διαδοχικών ζευγαριών  $w_i$ . Αποδεικνύεται ότι το συγκεκριμένο πρόβλημα βελτιστοποίησης είναι NP-hard, και προτείνεται ένας άπληστος (greedy) αλγόριθμος του πλησιέστερου γείτονα, για την υλοποίησή του. Προτείνεται ακόμη η χρήση μιας ευρετηριακής δομής δεδομένων (RTOP-Grid) βασιζόμενη στη χωρική τμηματοποίηση του χώρου δεδομένων. Αναλύονται εκτενώς η δομή και οι ιδιότητες του RTOP-Grid, ενώ παρουσιάζεται ο τρόπος με τον οποίο μπορεί να χρησιμοποιηθεί, για τη βελτίωση της απόδοσης του RTA αλγορίθμου.

### 3.3.2 Αλγόριθμος branch and bound

Το 2013 οι Βλάχου κ.ά. [64] πρότειναν ένα νέο τρόπο ανάλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Σε αυτήν την ερευνητική εργασία αναγνωρίστηκαν για πρώτη φορά ορισμένες χρήσιμες ιδιότητες, που επιτρέπουν τη σημαντική βελτίωση της συνολικής απόδοσης ενός συστήματος: Ορίστηκαν οι έννοιες του ελαχίστου ( $\ell_V(p)$ ) και μεγίστου ( $u_V(p)$ ) ορίου βαθμολόγησης, για σημεία και χώρους δεδομένων. Ορίστηκε επίσης η ιδιότητα της υπερίσχυσης, έτσι ώστε: Δοθέντος ενός συνόλου  $V \subseteq W$  με βαρύτητες, ένα σημείο αναφοράς ερωτήματος  $q$  και ένα χώρο δεδομένων  $m \subseteq S$ , αν ισχύει  $u_V(q) < \ell_V(m)$  τότε ορίζεται ότι το σημείο  $q$  υπερισχύει έναντι του  $m$  ή αλλιώς  $q \prec_V m$ . Παρομοίως αν  $u_V(m) < \ell_V(q)$  τότε ορίζεται ότι η περιοχή  $m$  υπερισχύει έναντι του  $p$  ή αλλιώς  $m \prec_V q$ . Αν δεν ισχύει τίποτα από τα δύο, τότε τα  $q$  και  $m$  είναι μη συγκρίσιμα ή αλλιώς  $m \asymp_V q$ .

Η ιδιότητα της υπερίσχυσης μπορεί πολύ εύκολα να χρησιμοποιηθεί ώστε να δημιουργηθεί η ιδιότητα της απόρριψης: δοθέντος ενός συνόλου από βαρύτητες  $V \subseteq W$ , οι οποίες αναπαριστώνται από ένα χώρο δεδομένων  $m_V$ , και ενός αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων  $RTOP_k(q)$ , εάν  $k$  πλήθους αντικείμενα δεδομένων (χώροι ή πλειάδες) υπερισχύουν έναντι του  $q$  με βάση το  $V$ , τότε το  $m_V$  μπορεί να απορριφθεί. Η ιδιότητα αυτή είναι πολύ χρήσιμη καθώς επιτρέπει την απόρριψη ενός συνόλου από βαρύτητες ( $m_V$ ) χωρίς να είναι απαραίτητη η εξέταση όλων των επιμέρους βαρυτήτων.

Στην παραπάνω εργασία προτείνεται ένας αλγόριθμος branch and bound ( $BBR$ ), ο οποίος αρχικά δημιουργεί ένα R-tree ευρετήριο για τις βαρύτητες  $W$ . Ο αλγόριθμος  $INTOP_k$  χρησιμοποιεί τις ιδιότητες που περιγράφηκαν πιο πάνω, ώστε να απορρίπτει πολύ γρήγορα ομάδες βαρυτήτων, οι οποίες δεν μπορούν να αποτελούν τμήμα του  $RTOP_k$  αποτελέσματος. Ο αλγόριθμος  $INTOP_k$  χρησιμοποιείται από τον  $BBR$ , προκειμένου να παραχθεί το τελικό αποτέλεσμα. Η απόδοση του  $BBR$  εξαρτάται άμεσα από το πλήθος των κλήσεων προς τον  $INTOP_k$ . Έτσι προτάθηκε ο αλγόριθμος  $BBR^*$  που αποθηκεύει προηγούμενα αποτελέσματα του  $INTOP_k$ , ώστε να μην απαιτείται η κλήση του για κάθε κόμβο του R-tree. Προτάθηκε επίσης η χρήση μιας εναλλακτικής δομής δεδομένων για τις βαρύτητες: το συγκεντρωτικό R-tree [65]. Με τις δύο προτεινόμενες επεκτάσεις, είναι δυνατόν να βελτιωθεί σημαντικά ο χρόνος εκτέλεσης του αλγορίθμου  $BBR$ , κατατάσσοντάς τον στους γρηγορότερους αλγόριθμους που υπάρχουν αυτή τη στιγμή για την επίλυση των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων [64].

### 3.3.3 Αλγόριθμοι DRT και DRT\*

Το 2011 οι Βλάχου κ.ά. [66] εξέτασαν για πρώτη φορά το πρόβλημα των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων σε συστήματα που βασίζονται στη φυσική θέση. Έτσι εισάγεται η έννοια των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων βασιζόμενα στην απόσταση, ως εξής: Έστω μια πλειάδα αναφοράς  $q$ , ένας θετικός ακέραιος αριθμός  $k$ , ένα σύνολο  $S$  από πλειάδες και ένα σύνολο  $M$  από φορητές συσκευές  $m_i \in M$ , η κάθε μία από τις οποίες συνδέεται με ένα σύνολο βαρυτήτων  $w_i \in W$ . Το αποτέλεσμα ενός αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων βασιζόμενο στην απόσταση ( $RTOP_k(q)$ ), περιέχει μια φορητή συσκευή  $m_i$  αν και μόνο αν  $\exists p \in TOP_k(w_i)$  τέτοιο ώστε  $f(m_i, q) \leq f(m_i, p)$ .

Η βασική διαφορά μεταξύ των ερωτημάτων αυτών με τα παραδοσιακά αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων, αποτελεί το γεγονός ότι για δύο πανομοιότυπα σύνολα βαρυτήτων  $w_1 \equiv w_2$  που ανήκουν στις φορητές συσκευές  $m_1$  και  $m_2$  αντίστοιχα, είναι πιθανό να ισχύει  $m_1 \in RTOP_k(q)$  ενώ  $m_2 \notin RTOP_k(q)$ . Ο λόγος που συμβαίνει είναι επειδή οι συσκευές μπορεί να βρίσκονται σε διαφορετική φυσική θέση, και έτσι το σύνολο δεδομένων  $S$  να είναι διαφορετικό για κάθε μια συσκευή. Η δυναμική φύση του συγκεκριμένου προβλήματος, καθιστά τις παραδοσιακές τεχνικές μη αποδοτικές στη διαχείριση τέτοιου τύπου ερωτημάτων. Στην εν λόγω εργασία προτείνονται αποδοτικές μέθοδοι για την ορθή παρακολούθηση τέτοιων ερωτημάτων για όσο οι συσκευές μετακινούνται, με την προϋπόθεση ότι διατίθεται ένα αποτέλεσμα ενός αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων του παρελθόντος.

Ο αλγόριθμος DRT αποτελεί μια πρώτη προσέγγιση στην προσπάθεια επίλυσης του συγκεκριμένου προβλήματος: Έστω μια φορητή συσκευή στη φυσική θέση  $m$  με αποτέλεσμα αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων  $RTOP_k(q)$ , η οποία μετακινείται στη φυσική θέση  $m'$  με αντίστοιχο αποτέλεσμα  $RTOP'_k(q)$ . Έστω επίσης ότι για τη συσκευή αυτή, η βαρύτητα της απόστασης ορίζεται από το  $w[n]$ . Τότε η μέγιστη δυνατή μεταβολή που μπορεί να πραγματοποιηθεί για τη συγκεκριμένη συσκευή, στη βαθμολογία του  $q$ , είναι η εξής:  $\Delta f = w[n] \times d(m, m')$ , όπου  $d(m, m')$  είναι η απόσταση του σημείου  $m$  από το σημείο  $m'$ . Ορίζεται επίσης ως  $f_k$  η βαθμολογία του  $k$ -οστού σημείου στο top- $k$  αποτέλεσμα του σημείου  $m$ . Από τα παραπάνω προκύπτει ότι αν ένα σύνολο βαρυτήτων  $w$  ανήκει στο  $RTOP_k(q)$ , τότε το  $w$  θα αποτελεί επίσης τμήμα του  $RTOP'_k(q)$  αν και μόνο αν ισχύει ότι  $f(m', q) \leq f_k - 2 \times w[n] \times d(m, m')$ . Από την ιδιότητα αυτή, εξάγεται ο αλγόριθμος DRT, ο οποίος αποφεύγει τον υπολογισμό top- $k$  ερωτημάτων για τις περιπτώσεις που αυτό δεν είναι απαραίτητο.

Παρόλο που ο αλγόριθμος DRT βελτιώνει δραστικά την απόδοση εκτέλεσης ερωτημάτων τέτοιου τύπου σε σχέση με την αφελή μέθοδο, στην εργασία τους οι Βλάχου

κ.ά. πρότειναν μια επέκταση του υπάρχοντος μοντέλου, ώστε να βελτιωθούν περαιτέρω οι επιδόσεις του συστήματος. Δοθέντος ενός σημείου αναφοράς  $q$  και μιας φορητής συσκευής  $m$ , η κεντρική ιδέα περιλαμβάνει τον ορισμό μιας «αφαλούς περιοχής» (safe area) γύρω από το  $q$  (για κάθε φορητή συσκευή) με την ακόλουθη ενδιαφέρουσα ιδιότητα: για όσο η συσκευή  $m$  βρίσκεται εντός της ασφαλούς περιοχής, το σημείο  $q$  αποτελεί τμήμα του top-k αποτελέσματος για το  $m$ , ή αλλιώς  $w \in RTOP_k(q)$ . Με άλλα λόγια, μπορεί να αποφευχθεί οποιοσδήποτε top-k υπολογισμός για το  $m$  όσον αφορά το  $q$ , για όσο η συσκευή  $m$  δεν ξεπερνά τα όρια της ασφαλούς περιοχής. Η ασφαλής περιοχή  $S(q, R, m)$  αποτελεί ουσιαστικά την περιοχή που ορίζεται από έναν κύκλο με κέντρο το σημείο  $q$  και ακτίνα  $R$ . Πρέπει να σημειωθεί ότι η ακτίνα  $R$  της ασφαλούς περιοχής  $S(q, R, m)$ , εξαρτάται μόνο από το  $w$  και είναι ανεξάρτητη από την πραγματική μετατόπιση του  $m$ . Έτσι μπορεί να υπολογιστεί μια φορά για κάθε συσκευή  $m$  και να παραμένει ορθή για όσο το σύνολο βαρυτήτων  $w$  για τη συσκευή αυτή, δε μεταβάλλεται. Ο αλγόριθμος που υπολογίζει την ακτίνα για κάθε συσκευή  $m$ , ονομάστηκε DRT\*, και στηρίζεται στη λογική της περιοχής προτεραιότητας. Η περιοχή προτεραιότητας  $PP(q, r, m)$  ενός σημείου  $q$  για μια συσκευή  $m$  και ένα σημείο ενδιαφέροντος  $p$  (που δεν είναι ίδιο με το  $q$ ), αποτελεί ουσιαστικά την περιοχή του χώρου δεδομένων, η οποία όταν περιέχει τη συσκευή  $m$ , το σημείο  $q$  λαμβάνει καλύτερη βαθμολογία από το  $p$  για το ερώτημα του  $m$ . Ορισμός της περιοχής αυτής βοηθάει τον αλγόριθμο DRT\* στον ακριβή ορισμό της ασφαλούς περιοχής μιας συσκευής  $m$ .

### 3.4 Ανοικτές ερευνητικές περιοχές

Μέσα από την ανασκόπηση της υπάρχουσας βιβλιογραφίας και τη μελέτη άλλων παρεμφερών εργασιών στο χώρο της υποστήριξης της διαδικασίας λήψης αποφάσεων, κατέστη δυνατή η αναγνώριση ενός συνόλου από ερευνητικές περιοχές, οι οποίες δεν έχουν μελετηθεί ακόμα επαρκώς σε ερευνητικό επίπεδο.

Η αναγνώριση των ανοικτών ερευνητικών περιοχών, που πραγματοποιείται στην παρούσα εργασία, περιστρέφεται κυρίως γύρω από την περιοχή των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Τα ερωτήματα τέτοιου τύπου είναι σχετικά καινούργια στη βιβλιογραφία (η πρώτη εμφάνισή τους έγινε το 2010 [23]), και έκτοτε έχουν τραβήξει αρκετά το ενδιαφέρον της επιστημονικής κοινότητας. Έτσι προέκυψε ένα σύνολο από ερευνητικές εργασίες που μελετάνε τα ερωτήματα αυτά με διάφορους τρόπους και σε διάφορες παραλλαγές τους. Ωστόσο, αν πραγματοποιηθεί σύγκριση των υπάρχουσών εργασιών με αυτές άλλων ερευνητικών πεδίων, γίνεται σαφές πως τα ερωτήματα τέτοιου τύπου επιδέχονται περαιτέρω ερευνητικής μελέτης. Στη συνέχεια παρουσιάζονται οι ανοικτές ερευνητικές περιοχές που αναγνωρίστηκαν έπειτα

από την ανασκόπηση που πραγματοποιήθηκε στο πλαίσιο της παρούσας εργασίας. Για καλύτερη κατανόηση αυτών, οι περιοχές χωρίζονται σε τέσσερις πιθανές κατηγορίες:

**Μοντέλο ερωτήματος:** Το σύνολο των υπαρχουσών τεχνικών θεωρεί την υιοθέτηση ενός ερωτήματος επιλογής, στο οποίο η βαθμολόγηση της εκάστοτε πλειάδας, συνδέεται άμεσα με τις ίδιες τις πλειάδες. Μελλοντικά θα μπορούσαν να εξεταστούν ερωτήματα συνένωσης (join queries) ή και συγκεντρωτικά ερωτήματα (aggregate queries) για τη βαθμολόγηση των πλειάδων. Εδώ πρέπει επίσης να αναφερθεί η απουσία αποτελεσματικής μεθόδου για την επίλυση του μονοχρωματικού ερωτήματος σε χώρο δεδομένων άνω των δύο διαστάσεων.

**Μέθοδος προσπέλασης δεδομένων:** Οι υπάρχουσες τεχνικές υποθέτουν την ύπαρξη των δεδομένων σε ένα σύστημα διαχείρισης, το οποίο επιτρέπει την τυχαία προσπέλαση των δεδομένων. Ωστόσο απουσιάζουν μελέτες που θα μπορούσαν να γίνουν σε συστήματα σειριακής προσπέλασης των δεδομένων (π.χ. μέσω ροών δεδομένων) ή για δεδομένα υψηλής κατανομής.

**Τύποι δεδομένων:** Οι υπάρχουσες προσεγγίσεις επικεντρώνονται κυρίως σε σενάρια όπου τα δεδομένα είναι ποσοτικά μετρήσιμα και για τα οποία υπάρχει βεβαιότητα της τιμής τους. Παρεμφερείς εργασίες από άλλα επιστημονικά πεδία, δείχνουν ότι τα δεδομένα μπορεί να μην είναι ποσοτικά μετρήσιμα (κατηγορικά δεδομένα), να υπάρχει ένα ποσοστό αβεβαιότητας της τιμής τους (probabilistic data) ή ακόμα και να είναι ελλιπή (incomplete data). Για κάθε μία από τις παραπάνω περιπτώσεις, υπάρχει χώρος για μελλοντικές έρευνες.

**Συνάρτηση κατάταξης:** Όλες οι τεχνικές που εξετάστηκαν, χρησιμοποιούν το μοντέλο της μονότονης συνάρτησης αθροίσματος των βαθμωτών τιμών, για τη βαθμολόγηση κάθε πλειάδας. Ωστόσο το μοντέλο αυτό δεν είναι το μοναδικό, και θα μπορούσαν μελλοντικά να μελετηθούν επιπλέον τύποι πιθανών συναρτήσεων βαθμολόγησης.

### 3.5 Συμπεράσματα

Η μελέτη της υπάρχουσας βιβλιογραφίας, επικεντρώθηκε κυρίως στα πολυκριτηριακά ερωτήματα ανάλυσης δεδομένων, για την υποστήριξη συστημάτων λήψης αποφάσεων πραγματικού χρόνου. Τέτοιου είδους συστήματα έχουν μελετηθεί εκτενώς σε ερωτήματα κορυφογραμμής, ή και σε ερωτήματα  $k$  κορυφαίων σημείων. Οι δύο αυτοί τύποι ερωτημάτων παρουσιάζουν αρκετές ομοιότητες όσον αφορά τη διαχείρισή

τους και έτσι λύσεις και τεχνικές που προτείνονται για κάποιο τύπο από αυτούς, μπορούν συνήθως να προσαρμοστούν ώστε να υποστηρίξουν και τον άλλο τύπο.

Πιο συγκεκριμένα μελετήθηκαν τρόποι υποστήριξης των ερωτημάτων κορυφογραμμής σε διάφορες παραλλαγές που παρουσιάζουν αυτά, σε συστήματα πραγματικού χρόνου: συστήματα αποκλειστικής προσάρτησης με χρήση χρονικού συρόμενου παραθύρου, ιδιότητες των συστημάτων αυτών, συστήματα που λαμβάνουν υπόψη υποσύνολα του χρονικού συρόμενου παραθύρου, καθώς και συστήματα κατανεμημένων ροών δεδομένων. Από τη μελέτη και την ανάλυση αυτών, γίνεται σαφές πως ο χώρος για περαιτέρω ερευνητικές προσεγγίσεις στα ερωτήματα κορυφογραμμής είναι περιορισμένος.

Εξετάστηκαν επίσης τα ερωτήματα  $k$  κορυφαίων σημείων σε συστήματα ροών δεδομένων πραγματικού χρόνου, όπου αναλύθηκαν οι ιδιότητες των συστημάτων τέτοιου τύπου, παρουσιάστηκαν οι πιο αποδοτικές μέθοδοι επίλυσης του προβλήματος, μελετήθηκαν οι προτεινόμενοι τρόποι διαχείρισης ερωτημάτων σε κατανεμημένα συστήματα και παρουσιάστηκαν τρόποι διαχείρισης του προβλήματος των ροών αβέβαιων δεδομένων (uncertain data). Ο χώρος των ερωτημάτων  $k$  κορυφαίων σημείων έχει κεντρίσει το έντονο ενδιαφέρον της επιστημονικής κοινότητας τα τελευταία χρόνια, και για το λόγο αυτό υπάρχει πληθώρα από προτεινόμενες ερευνητικές μεθόδους.

Τέλος μελετήθηκαν υπάρχουσες προσεγγίσεις στο χώρο των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων, οι οποίες πραγματεύονται το πρόβλημα σε ένα πιο σταθερό περιβάλλον, όπου τα δεδομένα είναι όλα διαθέσιμα, τη στιγμή που χρειάζονται από το σύστημα. Παρουσιάστηκε ένα σύνολο από ανοικτές ερευνητικές περιοχές για τα ερωτήματα αυτά, οι οποίες χρήζουν μελέτης. Στην παρούσα εργασία ωστόσο επιλέγεται να μελετηθεί η περίπτωση των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων σε ροές δεδομένων, με χρήση συρόμενων χρονικών παραθύρων. Πιο συγκεκριμένα επιλέγεται η μελέτη των ερωτημάτων επιλογής - όπου κάθε πλειάδα συνδέεται άμεσα με τη βαθμολογία της - σε συστήματα ροών δεδομένων σειριακής προσπάθειας, όπου τα δεδομένα είναι ποσοτικά μετρήσιμα και δεν είναι αβέβαια, ενώ για τη συνάρτηση βαθμολόγησης επιλέγεται το κλασικό μοντέλο της συνάρτησης αθροίσματος των βαθμωτών τιμών της εκάστοτε πλειάδας. Το υπόλοιπο της εργασίας θα επικεντρωθεί στην παρουσίαση πιθανών προσεγγίσεων στο πρόβλημα αυτό.



## Κεφάλαιο 4

# Μοντέλο επίλυσης αντίστροφων ερωτημάτων $k$ κορυφαίων σημείων

Στο προηγούμενο κεφάλαιο αναλύθηκαν οι υπάρχουσες τεχνικές επίλυσης των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Έτσι αναδείχθηκε η απουσία μιας μεθόδου επίλυσης τέτοιων ερωτημάτων σε ροές δεδομένων. Τα ερωτήματα τέτοιου τύπου απαιτούν διαφορετικού είδους προσέγγιση, η οποία θα ανταποκρίνεται στους περιορισμούς και θα εκμεταλλεύεται τα πλεονεκτήματα των ροών δεδομένων.

Στο κεφάλαιο αυτό προτείνεται μια καινοτόμα μέθοδος επίλυσης των ερωτημάτων τέτοιου τύπου. Η μέθοδος αυτή, περιλαμβάνει τη δημιουργία ενός συστήματος, το οποίο είναι ικανό να χειρίζεται αποδοτικά μεγάλου όγκου πληροφορίες. Παράλληλα είναι απαραίτητο να αναγνωριστούν ορισμένες ιδιότητες οι οποίες ισχύουν στη λειτουργία τέτοιου τύπου συστημάτων, και μπορούν να βοηθήσουν στο σχεδιασμό της προτεινόμενης μεθόδου.

Στόχοι του κεφαλαίου αυτού είναι:

- Η διατύπωση των παραδοχών οι οποίες λαμβάνονται υπόψη στη συνέχεια.
- Η αναγνώριση και διατύπωση ενός συνόλου ιδιοτήτων που ισχύουν στο περιβάλλον των παραδοχών που διατυπώνονται.
- Η παρουσίαση του προτεινόμενου μοντέλου επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων σε ροές δεδομένων.
- Η μελέτη πιθανών μεθόδων βελτίωσης του προτεινόμενου μοντέλου με χρήση των αναγνωρισμένων ιδιοτήτων.

## 4.1 Ορισμοί και σύμβολα

Σύμβολο	Περιγραφή
$D$	Χώρος δεδομένων
$d_d$	Διάσταση $d$ του χώρου δεδομένων $D$
$Q$	Σύνολο πλειάδων αναφοράς
$n$	Μέγεθος συνόλου πλειάδων αναφοράς $Q$
$q_i$	Πλειάδα αναφοράς $i$ του συνόλου $Q$
$q_i[j]$	Αξία της πλειάδας $q_i$ στη διάσταση $j$
$S$	Σύνολο δεδομένων ροής
$p$	Πλειάδα του συνόλου $S$
$p[j]$	Αξία της πλειάδας $p$ στη διάσταση $j$
$p.t_{arr}$	Λογικός χρόνος άφιξης πλειάδας $p$
$p.t_{exp}$	Λογικός χρόνος λήξης πλειάδας $p$
$p_{in}$	Εισερχόμενη πλειάδα $p$
$p_{exp}$	Εξερχόμενη πλειάδα $p$
$T$	Λογικό χρονικό διάστημα ζωής πλειάδας $p$
$W$	Σύνολο προτιμήσεων
$m$	Μέγεθος συνόλου προτιμήσεων $W$
$W_a$	Προτίμηση $a$ του συνόλου $W$
$W_a[j]$	Αξία της προτίμησης $W_a$ στη διάσταση $j$
$f_{W_a}(p)$	Συνάρτηση βαθμολόγησης της πλειάδας $p$ με βάση την προτίμηση $W_a$
$k$	Πλήθος των αποτελεσμάτων του top-k ερωτήματος
$RTOP_k(q_i), W'(q_i)$	Αποτέλεσμα αντίστροφου top-k ερωτήματος για το σημείο $q_i$
$W_{inf}(p, q_i)$	Σύνολο προτιμήσεων που επηρεάζουν το $W'(q_i)$ για τη διάρκεια ισχύος του $p$
$BUFFER$	Περιοχή στην κύρια μνήμη, μεγέθους $T$ , που αποθηκεύονται τα ενεργά $p$
$SCORE(W, Q)$	Περιοχή στην κύρια μνήμη που αποθηκεύονται οι βαθμολογίες των $q_i$ για κάθε $W_a$
$COUNTER(W, Q)$	Περιοχή στην κύρια μνήμη που συντηρείται το πλήθος των $p$ που έχουν καλύτερη βαθμολογία από κάποιο $q_i$ για την προτίμηση $W_a$
$V$	Υποσύνολο των προτιμήσεων $W$
$ev$	Περιοχή - κόμβος του RTree που αντιστοιχεί στο υποσύνολο των προτιμήσεων $V$
$M$	Υποσύνολο των πλειάδων αναφοράς $Q$
$m_M$	MBR του υποσυνόλου των πλειάδων αναφοράς $M$
$\ell_V(m)$	Η ελάχιστη βαθμολογία που λαμβάνει το $m$ για προτιμήσεις του $V$
$u_V(m)$	Η μέγιστη βαθμολογία που λαμβάνει το $m$ για προτιμήσεις του $V$

Πίνακας 4.1: Επισκόπηση συμβόλων

Για τη δημιουργία του μοντέλου επίλυσης των αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων σε ροή δεδομένων, θεωρείται η ύπαρξη του ακόλουθου σεναρίου: Έστω χώρος δεδομένων  $D$ , ο οποίος ορίζεται από ένα σύνολο  $d$  διαστάσεων  $d_1, d_2, \dots, d_d$  και ένα πεπερασμένο σύνολο δεδομένων  $Q = \{q_1, q_2, \dots, q_n\}$  στο  $D$  με μέγεθος  $n = |Q|$ . Ένα σημείο  $q_i \in Q \forall 1 \leq i \leq n$  μπορεί να αναπαρασταθεί ως  $q_i = \{q_i[1], q_i[2], \dots, q_i[d]\}$ , όπου  $q_i[j] \forall 1 \leq j \leq d$  είναι η τιμή του σημείου  $q_i$  στη διάσταση  $d_j$ . Κάθε διάσταση του χώρου  $D$ , θεωρείται ότι αναπαριστά ένα αριθμητικό χαρακτηριστικό, και έτσι οι τιμές  $q_i[j]$  σε οποιαδήποτε διάσταση  $d_j$  θεωρούνται ως αριθμητικές και μη αρνητικές. Στο σύστημα εισέρχονται μέσω μιας ροής δεδομένων σημεία  $p = \{p[1], p[2], \dots, p[d]\}$  που ανήκουν στο χώρο δεδομένων  $D$ . Σε αυτήν την εργασία θεωρείται ότι η ροή δεδομένων ακολουθεί τους κανόνες του μοντέλου της αποκλειστικής προσάρτησης (append only model) [52]. Σε κάθε εισερχόμενο σημείο  $p$ , ανατίθενται τιμές  $p.t_{arr}$  και  $p.t_{exp}$  οι οποίες αναπαριστούν το λογικό χρόνο άφιξης και το λογικό χρόνο λήξης του σημείου  $p$  αντίστοιχα. Ορίζεται σταθερό λογικό χρονικό διάστημα  $T = p.t_{exp} - p.t_{arr}$  το οποίο αναπαριστά το διάστημα στο οποίο κάθε εισερχόμενο σημείο  $p$  παραμένει ενεργό στο σύστημα, διατηρώντας τη δυνατότητά του να επηρεάσει το τελικό αποτέλεσμα. Όλα τα ενεργά σημεία  $p$  θεωρείται ότι ανήκουν στο σύνολο δεδομένων  $S$  του χώρου δεδομένων  $D$ .

Θεωρείται ακόμη η ύπαρξη ενός πεπερασμένου συνόλου δεδομένων  $W = \{W_1, W_2, \dots, W_m\}$  μεγέθους  $m = |W|$ , του οποίου κάθε σημείο  $W_a \forall 1 \leq a \leq m$  ορίζεται από διαστάσεις πλήθους  $d$ , ώστε  $W_a = \{W_a[1], W_a[2], \dots, W_a[d]\}$ . Το σύνολο δεδομένων  $W$  αναπαριστά προτιμήσεις με χρήση του όρου της βαρύτητας, ως προς τις διαστάσεις του χώρου δεδομένων των συνόλων  $S, Q$ . Κάθε προτίμηση  $W_a$  αντιπροσωπεύει τις βαρύτητες των διαστάσεων  $d_1, d_2, \dots, d_d$  για τη συγκεκριμένη προτίμηση, με τον εξής τρόπο: Η τιμή  $W_a[1]$  αντιστοιχεί στην ποσοστιαία βαρύτητα της διάστασης  $d_1$ , η τιμή  $W_a[2]$  αντιστοιχεί στην ποσοστιαία βαρύτητα της διάστασης  $d_2$  κ.ο.κ. Έτσι το πλήθος των διαστάσεων κάθε σημείου  $W_a$  είναι  $d$ . Χωρίς να υπάρξει απώλεια της γενικότητας του μοντέλου, θεωρείται ότι οι βαρύτητες  $W_a[j]$  κανονικοποιούνται στο πεδίο τιμών  $[0, 1]$  με τέτοιο τρόπο, ώστε  $\sum W_a[j] = 1$ . Η τεχνική αυτή βρίσκεται σε συμφωνία με την έννοια της προτίμησης που παρουσιάζεται στη βιβλιογραφία και έχει υιοθετηθεί ευρέως σε παρόμοιες εργασίες [23, 67, 68]. Ορίζεται ακόμη γραμμική συνάρτηση βαθμολόγησης ενός σημείου  $p : f_{W_a}(p) = \sum W_a[j] \times p[j]$ . Οι μικρότερες τιμές βαθμολόγησης ενός σημείου  $p$ , θεωρούνται προτιμότερες έναντι των μεγαλύτερων, για τον υπολογισμό των ερωτημάτων  $k$  κορυφαίων σημείων.

Δοθέντος ενός θετικού ακεραίου αριθμού  $k$ , για τον υπολογισμό ενός αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων, για ένα σημείο  $q_i$  της παραπάνω υπόθεσης, αρκεί να υπολογιστεί το αποτέλεσμα του ερωτήματος  $k$  κορυφαίων σημείων  $TOP_k$  για

κάθε προτίμηση  $W_a$  λαμβάνοντας υπόψη κάθε φορά μόνο τα σημεία  $p$  που παραμένουν ενεργά στο σύστημα (δεν έχουν λήξει). Το αποτέλεσμα του αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων για το σημείο  $q_i$  αποτελεί το σύνολο των προτιμήσεων  $RTOP_k(q_i) = W'(q_i) \subseteq W$ , για τις οποίες υπάρχει  $p \in TOP_k(W_a)$  τέτοιο ώστε  $f_{W_a}(q_i) \leq f_{W_a}(p)$ .

Στον Πίνακα 4.1 παρουσιάζονται συνοπτικά τα σύμβολα και η σημασία τους.

## 4.2 Ιδιότητες

Για τη δημιουργία του μοντέλου απαραίτητη είναι η καταγραφή ενός συνόλου ιδιοτήτων, οι οποίες ισχύουν στα αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων, όταν αυτά εξετάζονται σε περιβάλλοντα ρών δεδομένων:

**Ιδιότητα 1:** Αν  $|S| + 1 \leq k$ , τότε  $W'(q_i) = W \forall q_i \in Q$ .

**Απόδειξη:** Ισχύει διότι  $|TOP_k| < k$  για κάθε  $W_a \in W$ , με αποτέλεσμα το  $q_i$  να έχει πάντα τη δυνατότητα να αποτελεί τμήμα του  $TOP_k$  για όλα τα  $W_a$ .

**Ιδιότητα 2:** Αν  $f_{W_a}(p) \geq f_{W_a}(q_i) \forall W_a \in W$ , τότε το σημείο  $p$  δεν επηρεάζει το  $W'(q_i)$ .

**Απόδειξη:** Από τον ορισμό του bichromatic αντίστροφου ερωτήματος  $k$  κορυφαίων σημείων, προκύπτει ότι για να επηρεάσει ένα σημείο  $p$  το αποτέλεσμα  $W'(q_i)$  θα πρέπει  $f_{W_a}(p) \leq f_{W_a}(q_i)$ .

**Ιδιότητα 3:** Αν για κάθε  $q_i \in Q$  το σημείο  $p$  δεν επηρεάζει το  $W'(q_i)$  τότε το σημείο  $p$  μπορεί να απορριφθεί.

**Απόδειξη:** Προφανής.

**Ιδιότητα 4:** Αν  $q_i \succ^{pto} p$  τότε  $f_{W_a}(q_i) \leq f_{W_a}(p) \forall W_a \in W$ .

**Απόδειξη:** Με τη χρήση του ορισμού της γραμμικής συνάρτησης  $f_{W_a}(p) = \sum W_a[j] \times p[j]$ , προκύπτει ότι διατηρώντας σταθερές τις τιμές  $W_a[j]$ , η αύξηση της αξίας μιας ή περισσότερων τιμών  $p[j]$  προκαλεί αύξηση της τιμής της συνάρτησης  $f_{W_a}(p)$ . Και αντίστροφα, η μείωση της αξίας μιας ή περισσότερων τιμών  $p[j]$  προκαλεί μείωση της τιμής της συνάρτησης  $f_{W_a}(p)$ .

**Ιδιότητα 5:** Αν  $q_i \succ^{pto} p \forall q_i \in Q$  τότε το σημείο  $p$  μπορεί να απορριφθεί.

**Απόδειξη:** Προφανής από το συνδυασμό των ιδιοτήτων 2, 3 και 4.

**Ιδιότητα 6:** Αν  $\exists W_a \in W$  τέτοιο ώστε  $f_{W_a}(p) < f_{W_a}(q_i)$ , τότε το σημείο  $p$  μπορεί να επηρεάζει το  $W'(q_i)$  ως προς την προτίμηση  $W_a$ .

**Απόδειξη:** Αν το πλήθος  $c$  των σημείων  $p$  για τα οποία ισχύει  $f_{W_a}(p) < f_{W_a}(q_i)$  είναι μεγαλύτερο ή ίσο του  $k$ , τότε το σημείο  $q_i$  δεν αποτελεί τμήμα του  $TOP_k$  για την προτίμηση  $W_a$  και έτσι η προτίμηση  $W_a$  παύει να αποτελεί τμήμα του  $W'(q_i)$ . Η παύση ισχύος ενός σημείου  $p$  για το οποίο  $\exists W_a \in W$  τέτοιο ώστε  $f_{W_a}(p) < f_{W_a}(q_i)$ , μπορεί επίσης να προκαλέσει αλλαγή του  $W'(q_i)$  για τον ίδιο λόγο.

**Ιδιότητα 7:** Κατά την άφιξη ενός σημείου  $p$ , είναι δυνατός ο υπολογισμός ενός συνόλου προτιμήσεων  $W_{inf}(p, q_i) \in W$ , οι οποίες επηρεάζουν το  $W'(q_i)$  για τη διάρκεια ισχύος του  $p$ .

**Απόδειξη:** Αν για το σημείο  $p \exists W_a \in W$  τέτοιο ώστε  $f_{W_a}(p) < f_{W_a}(q_i)$  τότε, σύμφωνα με την ιδιότητα 6, το σημείο  $p$  είναι πιθανό να επηρεάζει το αποτέλεσμα  $W'(q_i)$  ως προς την προτίμηση  $W_a$  καθ'όλη τη διάρκεια ισχύος του. Έτσι το σύνολο  $W_{inf}$  αποτελείται από τις προτιμήσεις εκείνες, για τις οποίες ισχύει ότι  $f_{W_a}(p) < f_{W_a}(q_i)$  κατά την άφιξη του  $p$ .

**Ιδιότητα 8:** Αν  $u_V(q) < \ell_V(p)$  (το  $q$  υπερισχύει του  $p$  ή αλλιώς  $q \prec_V p$ ) τότε  $f_{W_a}(q) < f_{W_a}(p) \forall W_a \in V$ .

**Απόδειξη:** Εξ ορισμού ισχύει ότι  $f_{W_a}(q) \leq u_V(q) \forall W_a \in V$ . Ισχύει επίσης ότι  $\ell_V(p) \geq f_{W_a}(p) \forall W_a \in V$ . Από το συνδυασμό των δύο και με την υπόθεση ότι  $u_V(q) < \ell_V(p)$ , εξάγεται το συμπέρασμα ότι  $f_{W_a}(q) < f_{W_a}(p) \forall W_a \in V$ .

**Ιδιότητα 9:** Αν  $u_V(p) < \ell_V(q)$  (το  $p$  υπερισχύει του  $q$  ή αλλιώς  $p \prec_V q$ ) τότε  $f_{W_a}(p) < f_{W_a}(q) \forall W_a \in V$ .

**Απόδειξη:** Εξ ορισμού ισχύει ότι  $f_{W_a}(p) \leq u_V(p) \forall W_a \in V$ . Ισχύει επίσης ότι  $\ell_V(q) \geq f_{W_a}(q) \forall W_a \in V$ . Από το συνδυασμό των δύο και με την υπόθεση ότι  $u_V(p) < \ell_V(q)$ , εξάγεται το συμπέρασμα ότι  $f_{W_a}(p) < f_{W_a}(q) \forall W_a \in V$ .

Στη συνέχεια παρουσιάζεται η αρχιτεκτονική του προτεινόμενου συστήματος επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων σε ροές δεδομένων. Παρατίθενται επίσης τρεις σχετικοί αλγόριθμοι που μπορούν να λειτουργούν σε συστήματα της προτεινόμενης αρχιτεκτονικής, οι οποίοι παρέχουν τη δυνατότητα υπολογισμού αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων.

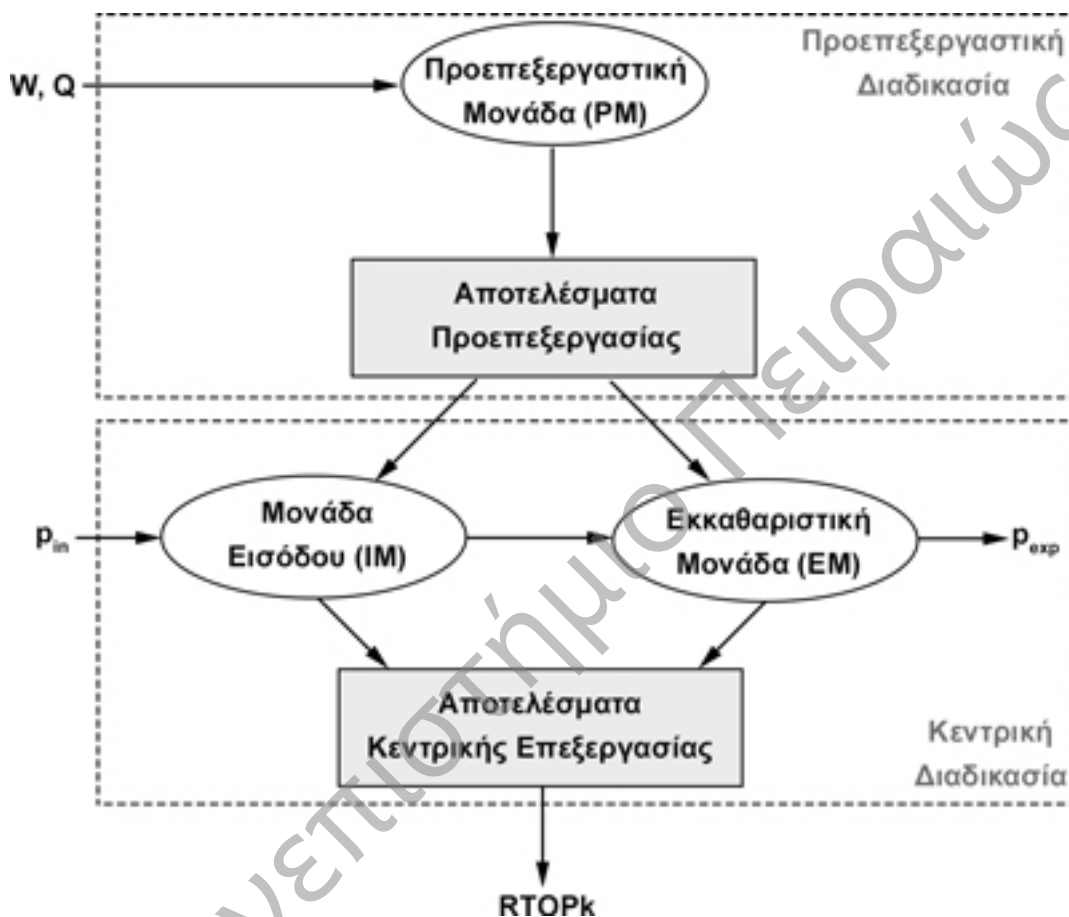
Για την καλύτερη κατανόηση της λειτουργίας του συστήματος, είναι απαραίτητο να σημειωθεί η εξής παρατήρηση: Το  $RTOP_k(q_i)$  για ένα δοθέν σημείο  $q_i \in Q$ , μπορεί να μεταβληθεί μόνο έπειτα από τις παρέλευση ενός εκ των ακόλουθων δύο γεγονότων:

1. Ένα νέο σημείο  $p_{in} \in S$  εισέρχεται στο σύστημα ( $p_{in}.t_{arr} = now$ ).

- Ένα νέο σημείο  $p_{exp} \in S$  εξέρχεται από το σύστημα ( $p_{exp}.t_{exp} = now$ ).

### 4.3 Αρχιτεκτονική συστήματος

Στο Σχήμα 4.1 παρουσιάζεται σχηματικά η αρχιτεκτονική του προτεινόμενου μοντέλου.



Σχήμα 4.1: Αναπαράσταση προτεινόμενου μοντέλου

Το μοντέλο του προτεινόμενου συστήματος, που περιγράφεται στη συνέχεια, μπορεί να υποστηρίξει και τα δύο είδη συρόμενων χρονικών παραθύρων: τα βασιζόμενα στο πλήθος των πλειάδων, όπου περιέχουν τις πιο πρόσφατες  $N$  πλήθους πλειάδες, και τα βασιζόμενα στο χρόνο, τα οποία περιέχουν τις πλειάδες εκείνες που εισήχθησαν το σύστημα σε ένα ορισμένο χρονικό διάστημα, καλύπτοντας τις πιο πρόσφατες χρονικές στιγμές. Ακόμη παρέχεται η δυνατότητα ταυτόχρονου υπολογισμού και παρακολούθησης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων για κάθε  $q_i \in Q$ . Δηλαδή για κάθε  $q_i$  υπολογίζεται και συντηρείται ξεχωριστό  $W'(q_i)$ . Τα σύνολα δεδομένων  $W$  και  $Q$  θεωρείται ότι είναι γνωστά εκ των προτέρων της εκκίνησης της

επεξεργασίας του συστήματος, ενώ το σύνολο δεδομένων  $S$  είναι άγνωστο, και εισέρχεται στο σύστημα κατά τη διάρκεια λειτουργίας του. Όλα τα σύνολα και οι δομές δεδομένων που περιγράφονται στη συνέχεια, αποθηκεύονται στην κύρια μνήμη του συστήματος ώστε να επιτυγχάνεται πολύ υψηλός ρυθμός επεξεργασίας της εισερχόμενης ροής δεδομένων.

Το σύστημα δέχεται ως είσοδο τα πεπερασμένα σύνολα δεδομένων  $W, Q$  καθώς και μια ροή δεδομένων  $S$ . Θεωρείται ότι οι πλειάδες της ροής δεδομένων  $S$ , αποθηκεύονται σε μια περιοχή της κύριας μνήμης,  $BUFFER$  μεγέθους  $T$ . Όταν το  $BUFFER$  ξεπεράσει το όριο  $T$ , αφαιρούνται οι αρχαιότερες πλειάδες της ροής  $S$ , έτσι ώστε το  $BUFFER$  να διατηρεί πάντα το πολύ μέγεθος  $T$ . Στο εξής κάθε νεοεισερχόμενη πλειάδα της ροής δεδομένων  $S$ , συμβολίζεται ως  $p_{in}$  ενώ κάθε πλειάδα που παύει να ισχύει στο  $BUFFER$ , συμβολίζεται ως  $p_{exp}$ . Η επεξεργασία των δεδομένων πραγματοποιείται από τρεις διακριτές επεξεργαστικές μονάδες (modules) ώστε να παραχθούν τα τελικά αποτελέσματα  $RTOP_k(q_i) \forall q_i \in Q$ . Η λειτουργία των μονάδων αυτών, διαφέρει από αλγόριθμο σε αλγόριθμο, ωστόσο στη συνέχεια παρουσιάζονται σε γενικές γραμμές τα καθήκοντά τους.

Η λειτουργία του συστήματος χωρίζεται σε δύο λογικές διαδικασίες: στην Προεπεξεργαστική Διαδικασία και στην Κεντρική Διαδικασία. Η προεπεξεργαστική διαδικασία περιέχει την Προεπεξεργαστική Μονάδα, η οποία αποτελεί τη μονάδα εκείνη που αναλαμβάνει την προεπεξεργασία των συνόλων  $W$  και  $Q$  ώστε να αυξηθεί η ταχύτητα λειτουργίας του συστήματος. Τα αποτελέσματα της προεπεξεργαστικής διαδικασίας, αποθηκεύονται στην κύρια μνήμη του συστήματος, ώστε να χρησιμοποιηθούν στη συνέχεια. Κατά τη διάρκεια εκτέλεσης της προεπεξεργαστικής διαδικασίας, η ροή δεδομένων  $S$  δεν έχει ξεκινήσει να εισέρχεται στο σύστημα, ενώ η λειτουργία της Προεπεξεργαστικής Διαδικασίας, έχει σαφή αρχή και τέλος, λόγω της πεπερασμένης φύσης των συνόλων  $W, Q$ .

Μετά το πέρας της εκτέλεσης της Προεπεξεργαστικής Μονάδας, στο σύστημα ξεκινά η εκτέλεση της Κεντρικής Διαδικασίας. Θεωρείται ότι με την εκκίνηση εκτέλεσης της Κεντρικής Διαδικασίας, η δομή  $BUFFER$  περιέχει ήδη  $T$  πλήθους πλειάδες του συνόλου  $S$ . Η Κεντρική Διαδικασία περιλαμβάνει την Μονάδα Εισόδου και την Εκκαθαριστική Μονάδα. Η Μονάδα Εισόδου αναλαμβάνει να επεξεργαστεί και να αποθηκεύσει στο  $BUFFER$  κάθε νεοεισερχόμενο σημείο  $p_{in}$ . Όταν το  $BUFFER$  γεμίσει με μεγαλύτερο του  $T$  πλήθους πλειάδες, η Εκκαθαριστική Μονάδα αναλαμβάνει να διαγράψει την αρχαιότερη πλειάδα  $p_{exp}$  από το  $BUFFER$ . Οι μονάδες αυτές υπολογίζουν και ενημερώνουν την κύρια μνήμη του συστήματος, με τα  $RTOP_k(q_i)$  σύνολα  $\forall q_i \in Q$ . Η λειτουργία της Κεντρικής Διαδικασίας, έχει σαφή αρχή αλλά μπορεί να μην

έχει τέλος. Η φύση των ερωτημάτων ρούν δεδομένων, απαιτεί την ύπαρξη δυνατότητας επ'αόριστον εκτέλεσης του συστήματος, ώστε να παρακολουθείται διαρκώς η εξέλιξη των αποτελεσμάτων του ερωτήματος.

Στη συνέχεια παρουσιάζονται οι τρεις προτεινόμενοι αλγόριθμοι επεξεργασίας αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Ο πρώτος αλγόριθμος αποτελεί τον αφελή αλγόριθμο και παρέχεται μόνο για δυνατότητα σύγκρισης με τους υπόλοιπους. Οι δύο επόμενοι αλγόριθμοι βελτιώνουν την αποδοτικότητα του πρώτου αλγορίθμου με χρήση των ιδιοτήτων που περιγράφηκαν νωρίτερα.

## 4.4 Αλγόριθμοι επεξεργασίας

### 4.4.1 Αλγόριθμος 1

Ο Αλγόριθμος 1 αποτελεί τον αφελή αλγόριθμο υπολογισμού αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Λόγω της απουσίας άλλης παρόμοιας μεθόδου στην ερευνητική βιβλιογραφία, ο αλγόριθμος αυτός θα αποτελέσει το μέτρο σύγκρισης και αξιολόγησης των αλγορίθμων που προτείνονται στη συνέχεια.

Ρόλος της Προεπεξεργαστικής Μονάδας είναι ο υπολογισμός της βαθμολόγησης όλων των πλειάδων  $q_i \in Q$ , σύμφωνα με τη συνάρτηση  $f_{W_a}(q_i) \forall W_a \in W$ . Ο διδιάστατος πίνακας  $SCORE(W, Q)$  στο εξής θα περιέχει τα αποτελέσματα όλων των υπολογισμών που πραγματοποιούνται στην Προεπεξεργαστική Μονάδα, με τρόπο τέτοιο ώστε η τιμή  $SCORE(W_a, q_i)$  να συμβολίζει την τιμή της συνάρτησης  $f_{W_a}(q_i)$ . Αυτή η προεπεξεργαστική διεργασία, μπορεί να αυξήσει σημαντικά τη συνολική απόδοση του συστήματος, καθώς κατά τη λειτουργία του, ελέγχονται πολύ συχνά οι τιμές  $f_{W_a}(q_i)$  για κάθε  $W_a \in W, q_i \in Q$ .

Η Μονάδα Εισόδου αναλαμβάνει να επεξεργαστεί όλες τις νεοεισερχόμενες πλειάδες  $p_{in}$ . Ρόλος της Μονάδας Εισόδου είναι να καταγράφει το πλήθος των πλειάδων  $p_{in}$ , οι οποίες έχουν τιμή βαθμολόγησης  $f_{W_a}(p_{in})$  μικρότερη της τιμής  $f_{W_a}(q_i)$ . Ο διδιάστατος πίνακας  $COUNTER(W, Q)$  στο εξής θα περιέχει το πλήθος των πλειάδων  $p \in S$ , οι οποίες έχουν μικρότερη τιμή βαθμολόγησης από κάποιο  $q_i \in Q$ . Ο πίνακας  $COUNTER(W, Q)$  κατά την αρχικοποίησή του, θεωρείται ότι θέτει όλες τις τιμές του ίσες με μηδέν. Οι τιμές αποθηκεύονται στον πίνακα με τρόπο τέτοιο, ώστε η τιμή  $COUNTER(W_a, q_i)$  να συμβολίζει το πλήθος των πλειάδων  $p \in S$ , οι οποίες έχουν τιμή βαθμολόγησης  $f_{W_a}(p)$  μικρότερη της τιμής  $f_{W_a}(q_i)$ . Κατ'επέκταση με χρήση των αποτελεσμάτων της προεπεξεργαστικής διαδικασίας, αρκεί να πραγματοποιηθεί σύγκριση της τιμής  $f_{W_a}(p)$  με την τιμή  $SCORE(W_a, q_i)$ , επιτρέποντας με τον τρόπο αυτό,



την αύξηση της αποδοτικότητας του συστήματος. Σε περίπτωση που η πρώτη τιμή είναι μικρότερη της δεύτερης, η Μονάδα Εισόδου αυξάνει την τιμή  $COUNTER(W_a, q_i)$  κατά 1.

Ρόλος της Εκκαθαριστικής Μονάδας είναι να συντηρεί ορθό το μετρητή του πλήθους των πλειάδων  $p \in S$ , οι οποίες έχουν τιμή βαθμολόγησης  $f_{W_a}(p)$  μικρότερη της τιμής  $f_{W_a}(q_i)$  ακόμα και μετά τη διαγραφή της πλειάδας  $p_{exp}$ . Για το σκοπό αυτό, η Εκκαθαριστική Μονάδα αναλαμβάνει να υπολογίσει όλες τις τιμές  $f_{W_a}(p_{exp})$  για κάθε  $W_a \in W$  και να τις συγκρίνει με τις αντίστοιχες τιμές  $SCORE(W_a, q_i)$ . Σε περίπτωση που η πρώτη τιμή είναι μικρότερη της δεύτερης, η Εκκαθαριστική Μονάδα μειώνει την τιμή  $COUNTER(W_a, q_i)$  κατά 1.

Με τη χρήση της παραπάνω μεθόδου, είναι δυνατό να παρακολουθείται διαρκώς η εξέλιξη πολλών αντίστροφων ερωτημάτων  $k$  κορυφαίων πλειάδων σε πραγματικό χρόνο. Το τελικό αποτέλεσμα  $RTOP_k(q_i)$  ενός ερωτήματος  $q_i$  μπορεί εύκολα να εξαχθεί, με τη χρήση του πίνακα  $COUNTER(W, Q)$ . Δοθέντος ενός ερωτήματος  $q_i$ , το σύστημα αρκεί να ελέγξει όλες τις τιμές του πίνακα  $COUNTER(W, Q)$  που αντιστοιχούν στο σημείο  $q_i$ , και να εξάγει τις προτιμήσεις  $W_a \in W$  για τις οποίες ισχύει ότι  $COUNTER(W_a, q_i) < k$ . Η ορθότητα της παραπάνω μεθόδου, μπορεί εύκολα να αποδειχθεί: Ο πίνακας  $COUNTER(W, Q)$  περιέχει κάθε στιγμή το πλήθος των πλειάδων εκείνων που έχουν τιμή βαθμολόγησης μικρότερη μικρότερη της τιμής βαθμολόγησης του  $q_i$ . Αν το πλήθος για μια ορισμένη προτίμηση  $W_a$  είναι μικρότερο της τιμής  $k$ , τότε η πλειάδα  $q_i$  δύναται να αποτελέσει τμήμα του  $TOP_k$  ερωτήματος, και άρα η προτίμηση αυτή αποτελεί τμήμα του  $RTOP_k(q_i)$  ερωτήματος.

Στη συνέχεια παρουσιάζονται συνοπτικά και με χρήση ψευδο-αλγορίθμων, η λειτουργία των τριών μονάδων επεξεργασίας:

#### Προεπεξεργαστική Μονάδα (pre-processing module, PM)

**Είσοδος:**  $W, Q$

**Έξοδος:**  $SCORE(W, Q)$

**Περιγραφή:** Για κάθε  $q_i \in Q, W_a \in W$  εφαρμόζεται η συνάρτηση  $f_{W_a}(q_i)$  και το αποτέλεσμα της αποθηκεύεται σε έναν δισδιάστατο πίνακα  $SCORE(W, Q)$ .

**Αλγόριθμος:** Αλγόριθμος 4.1

```
for (a = 0; a < m; a++) {  
  for (i = 0; i < n; i++) {  
    SCORE(W_a, q_i) = f_{W_a}(q_i);  
  }  
}
```

Αλγόριθμος 4.1: Αλγόριθμος προεπεξεργαστικής μονάδας

### Μονάδα εισόδου (input module, IM)

**Είσοδος:**  $p_{in}, SCORE(W, Q)$

**Έξοδος:**  $COUNTER(W, Q)$

**Περιγραφή:** Για κάθε νεοεισερχόμενη πλειάδα  $p_{in}$  εφαρμόζεται η συνάρτηση  $f_{W_a}(p_{in}) \forall W_a \in W$  και το αποτέλεσμα αυτής συγκρίνεται με τις τιμές  $SCORE(W_a, q_i)$  για κάθε  $q_i \in Q$ . Για τις περιπτώσεις που το αποτέλεσμα της συνάρτησης είναι μικρότερο από το αποτέλεσμα της τιμής  $SCORE(W_a, q_i)$ , αυξάνεται η τιμή του δείκτη  $COUNTER(W_a, q_i)$  κατά 1.

**Αλγόριθμος:** Αλγόριθμος 4.2

---

```

for (a = 0; a < m; a++) {
    score(pin) = fWa(pin);
    for (i = 0; i < n; i++) {
        if (score(pin) < SCORE(Wa, qi)) {
            COUNTER(Wa, qi)++;
        }
    }
}

```

---

Αλγόριθμος 4.2: Αλγόριθμος μονάδας εισόδου

### Εκκαθαριστική Μονάδα (exrunge module, EM)

**Είσοδος:**  $p_{exp}, SCORE(W, Q)$

**Έξοδος:**  $COUNTER(W, Q)$

**Περιγραφή:** Για κάθε εξερχόμενη πλειάδα  $p_{exp}$  εφαρμόζεται η συνάρτηση  $f_{W_a}(p_{exp}) \forall W_a \in W$  και το αποτέλεσμα αυτής συγκρίνεται με τις τιμές  $SCORE(W_a, q_i)$  για κάθε  $q_i \in Q$ . Για τις περιπτώσεις που το αποτέλεσμα της συνάρτησης είναι μικρότερο από το αποτέλεσμα της τιμής  $SCORE(W_a, q_i)$ , μειώνεται η τιμή του δείκτη  $COUNTER(W_a, q_i)$  κατά 1.

**Αλγόριθμος:** Αλγόριθμος 4.3

---

```

for (a = 0; a < m; a++) {
    scorep = fWa(pexp);
    for (i = 0; i < n; i++) {
        if (scorep < SCORE(Wa, qi)) {
            COUNTER(Wa, qi)--;
        }
    }
}

```

---

Αλγόριθμος 4.3: Αλγόριθμος εκκαθαριστικής μονάδας

#### 4.4.2 Αλγόριθμος 2

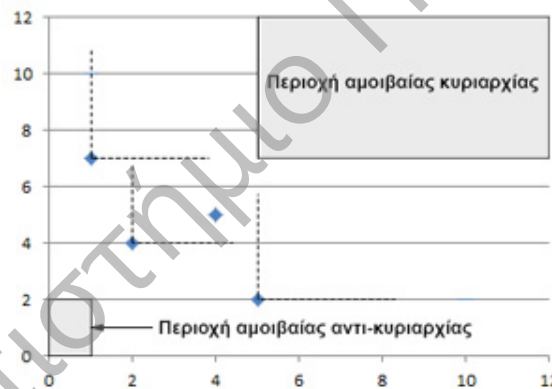
Με τη χρήση των Ιδιοτήτων 4 και 5, ορίζονται οι περιοχές αμοιβαίας κυριαρχίας και αντι-κυριαρχίας (mutual dominance region, MDR), (mutual anti-dominance area, MAR) ως εξής:

**Ορισμός 4.1 (Περιοχή Αμοιβαίας Κυριαρχίας) :**

Έστω σύνολο δεδομένων  $Q = \{q_1, \dots, q_n\}$ . Η περιοχή αμοιβαίας κυριαρχίας, αποτελεί την περιοχή διασταυρώσεως (ή περιοχή τομής) των περιοχών κυριαρχίας, όλων των σημείων του συνόλου  $Q$ . Ένα παράδειγμα μιας περιοχής αμοιβαίας κυριαρχίας, παρουσιάζεται στο Σχήμα 4.2

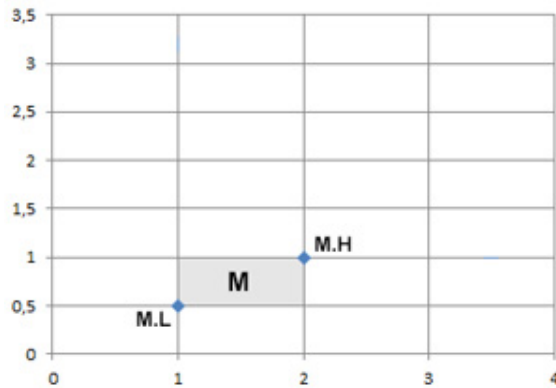
**Ορισμός 4.2 (Περιοχή Αμοιβαίας Αντι-Κυριαρχίας) :**

Έστω σύνολο δεδομένων  $Q = \{q_1, \dots, q_n\}$ . Η περιοχή αμοιβαίας κυριαρχίας, αποτελεί την περιοχή διασταυρώσεως (ή περιοχή τομής) των περιοχών αντι-κυριαρχίας, όλων των σημείων του συνόλου  $Q$ . Ένα παράδειγμα μιας περιοχής αμοιβαίας αντι-κυριαρχίας, παρουσιάζεται στο Σχήμα 4.2



Σχήμα 4.2: Περιοχές αμοιβαίας κυριαρχίας και αντι-κυριαρχίας

Έστω ακόμη ότι για το σύνολο των προτιμήσεων  $W$  υπάρχει ένα ευρετήριο τύπου  $RTree$   $RTreeW$ . Επιπλέον έστω ότι για το σύνολο δεδομένων  $Q$ , υπάρχει ένα ευρετήριο πλέγματος (grid index)  $GridQ$ . Ένα τμήμα  $M$  (κελί) του πλέγματος  $GridQ$ , ορίζεται από τα μέγιστα και ελάχιστα σημεία του  $M.h, M.l$  αντίστοιχα όπως φαίνεται στο Σχήμα 4.3. Με τον ίδιο τρόπο ένας κόμβος  $V$  (περιοχή) της δομής  $RTreeW$ , ορίζεται από τα μέγιστα και ελάχιστα σημεία του  $V.h, V.l$  αντίστοιχα. Η βαθμολογία  $f_{W_a}(q_i)$  ενός οποιουδήποτε  $q_i \in M$  χρησιμοποιώντας μια οποιαδήποτε προτίμηση  $W_a \in V$ , ορίζεται από ένα ελάχιστο όριο  $\ell_V(M) = f_{V.l}(M.l)$  και από ένα μέγιστο όριο  $u_V(M) = f_{V.h}(M.h)$ . Αντίστοιχα η ελάχιστη βαθμολογία που μπορεί να λάβει ένα νεοεισερχόμενο σημείο  $p_{in}$  για τις προτιμήσεις του  $V$  ορίζεται ως  $\ell_V(p_{in}) = f_{V.l}(p_{in})$  ενώ η μέγιστη τιμή ορίζεται ως  $u_V(p_{in}) = f_{V.u}(p_{in})$ .



Σχήμα 4.3: Μέγιστη και ελάχιστη τιμή περιοχής πλέγματος

Από τα παραπάνω και με τη χρήση των ιδιοτήτων 8 και 9, γίνεται σαφές πως αν  $u_V(M) < \ell_V(p_{in})$  τότε η περιοχή  $M$  υπερισχύει έναντι του σημείου  $p_{in}$  και το σημείο  $p_{in}$  δεν επηρεάζει το αποτέλεσμα  $RTOP_k(q_i) \forall q_i \in M$  όσον αφορά το σύνολο των προτιμήσεων  $W_a \in V$ . Έτσι ο αλγόριθμος μπορεί να αποφύγει τον υπολογισμό της βαθμολογίας για κάθε  $q_i \in M, W_a \in V$ . Αντιστοίχα αν ισχύει ότι  $u_V(p_{in}) < \ell_V(M)$  τότε το σημείο  $p_{in}$  υπερισχύει έναντι της περιοχής  $M$  και το σημείο αυτό, επηρεάζει το αποτέλεσμα  $RTOP_k(q_i) \forall q_i \in M$  για όλες τις προτιμήσεις του συνόλου  $W_a \in V$ . Έτσι ο αλγόριθμος μπορεί να συμπεριλάβει με ασφάλεια όλες τις προτιμήσεις  $W_a \in V$  για όλα τα σημεία  $q_i \in M$ .

Με βάση τους παραπάνω ορισμούς και με την εφαρμογή των κατάλληλων αλλαγών στη λειτουργία των Μονάδων Προεπεξεργασίας και Εισόδου, καθίσταται δυνατή η βελτίωση της αποδοτικότητας του βασικού μοντέλου. Στη συνέχεια αναλύεται η τροποποιημένη λειτουργία των μονάδων αυτών, επεξηγώντας τον τρόπο με τον οποίο οι αλλαγές συμβάλλουν στην αύξηση της αποδοτικότητας του συστήματος.

Η Μονάδα Προεπεξεργασίας αναλαμβάνει πλέον τον ορισμό των περιοχών αμοιβαίας κυριαρχίας και αντι-κυριαρχίας που ορίζονται από όλα τα σημεία  $q_i \in Q$ . Προς το σκοπό αυτό, ελέγχει όλα τα σημεία  $q_i$  αναζητώντας τη μέγιστη τιμή που εμφανίζεται σε κάθε διάσταση. Το σημείο που ορίζεται από τις μέγιστες τιμές που βρέθηκαν, αποτελεί το σημείο αρχής της Περιοχής Αμοιβαίας Κυριαρχίας και αποθηκεύεται στη μεταβλητή  $MDR$ . Αντίστοιχα αναζητά την ελάχιστη τιμή που εμφανίζεται σε κάθε διάσταση και το σημείο που ορίζεται από τις ελάχιστες τιμές που βρέθηκαν από το σημείο τέλους της Περιοχής Αμοιβαίας Αντι-Κυριαρχίας, το οποίο αποθηκεύεται στη μεταβλητή  $MAR$ . Επιπροσθέτως αναλαμβάνει να χτίσει ευρετήρια για το σύνολο των προτιμήσεων  $W$  με χρήση ενός  $RTree$ , και για το σύνολο των σημείων αναφοράς  $Q$  με χρήση ενός ευρετηρίου πλέγματος. Η δομή  $SCORE(W, Q)$  που περιγράφηκε στον Αλγόριθμο 1, δεν είναι πλέον απαραίτητη και η Μονάδα Επεξεργασίας δεν ασχολείται με τη δημιουργία της.

Η Μονάδα Εισόδου μπορεί πλέον να χρησιμοποιήσει τις τιμές των μεταβλητών  $MDR$  και  $MAR$ , ώστε να αποφευχθεί ο έλεγχος των πλειάδων  $p_{in} \in S$  που εμπίπτουν στους χώρους των περιοχής αμοιβαίας κυριαρχίας και αντι-κυριαρχίας. Πιο συγκεκριμένα για κάθε πλειάδα  $p_{in}$  που λαμβάνει η Μονάδα Εισόδου, εκτελεί έλεγχο κυριαρχίας pareto (pareto dominance) μεταξύ των πλειάδων  $p_{in}$  και  $MDR$ . Σε περίπτωση που η  $MDR$  κυριαρχεί pareto έναντι της  $p_{in}$ , τότε η πλειάδα  $p_{in}$  μπορεί να απορριφθεί με ασφάλεια από το  $BUFFER$ . Αντίστοιχα εκτελεί έλεγχο κυριαρχίας pareto (pareto dominance) μεταξύ των πλειάδων  $p_{in}$  και  $MAR$ . Σε περίπτωση που η  $p_{in}$  κυριαρχεί pareto έναντι της  $MAR$ , τότε όλες οι τιμές του πίνακα  $COUNTER(W, Q)$  πρέπει να αυξηθούν κατά 1.

Στη συνέχεια ο αλγόριθμος της Μονάδας Εισόδου ελέγχει ένα προς ένα τα κελιά της δομής  $GridQ$ , ξεκινώντας από το κελί στη θέση  $(0,0)$ . Συγκεκριμένα για κάθε κελί  $M \in GridQ$  πραγματοποιεί έλεγχο ύπαρξης κυριαρχίας pareto μεταξύ των  $M$  και  $p_{in}$ . Σε περίπτωση που το κελί  $M$  κυριαρχεί pareto έναντι της  $p_{in}$ , τότε η πλειάδα  $p_{in}$  δεν επηρεάζει το αποτέλεσμα  $RTOPk(q_i) \forall q_i \in M$ . Αντίστοιχα αν η πλειάδα  $p_{in}$  κυριαρχεί pareto έναντι του  $M$ , τότε όλες οι τιμές του πίνακα  $COUNTER(W, q_i)$  πρέπει να αυξηθούν κατά 1 για κάθε  $q_i \in M$ .

Στη συνέχεια σαρώνει τη δομή  $RTreeW$  από τη ρίζα έως τα φύλλα ελέγχοντας για κάθε κόμβο  $V$  του δέντρου, τις μέγιστες και τις ελάχιστες τιμές που λαμβάνει το σημείο  $p_{in}$ , όπως περιγράφηκε πιο πάνω. Για κάθε κόμβο  $V$  του  $RTreeW$  και κάθε κελί  $M$  του  $GridQ$ , πραγματοποιεί έλεγχο υπερίσχυσης του κελιού  $M$  σε σχέση με το σημείο  $p_{in}$ , με τον τρόπο που περιγράφηκε πιο πάνω. Έτσι, μπορεί να πραγματοποιήσει πολύ γρήγορα τον έλεγχο όλων των προτιμήσεων της δομής  $RTreeW$  για όλα τα σημεία αναφοράς.

Αντίστοιχα η Εκκαθαριστική Μονάδα εκτελεί τον ίδιο αλγόριθμο που περιγράφηκε στην Μονάδα Εισόδου, με τη μοναδική διαφορά, ότι της αύξησης των τιμών του πίνακα  $COUNTER(W, Q)$  όταν αυτό κρίνεται απαραίτητο, πραγματοποιεί μείωση αυτών κατά μια μονάδα.

Οι παραπάνω βελτιώσεις βασίζονται στις Ιδιότητες 4, 5, 8 και 9 που περιγράφηκαν νωρίτερα. Συγκεκριμένα λόγω της μεταβατικής ιδιότητας που έχει η πράξη της κυριαρχίας pareto, μπορούμε να ορίσουμε πως αφού  $q_i \succ^{pto} MDR \forall q_i \in Q$  και  $MDR \succ^{pto} p_{in}$ , τότε ισχύει και  $q_i \succ^{pto} p_{in} \forall q_i \in Q$ . Η Ιδιότητα 5, ορίζει ρητά στην περίπτωση αυτή, πως η πλειάδα  $p_{in}$ , μπορεί να απορριφθεί. Αντίστοιχα αποδεικνύεται και η ορθότητα της ιδιότητας της αντι-κυριαρχίας. Η ιδιότητα της υπερίσχυσης έχει ήδη αναλυθεί και αποδειχτεί ως προς την ορθότητά της.

Στη συνέχεια παρουσιάζονται συνοπτικά και με χρήση ψευδο-αλγορίθμων, η λειτουργία της Προεπεξεργαστικής Μονάδας, της Μονάδας Εισόδου και της Εκκαθαριστικής Μονάδας:

### Προεπεξεργαστική Μονάδα (pre-processing module)

**Είσοδος:**  $W, Q$

**Έξοδος:**  $MDR, MAR, RTreeW, GridQ$

**Περιγραφή:** Ελέγχονται οι τιμές των διαστάσεων για κάθε  $q_i \in Q$ . Αποθηκεύονται οι μέγιστες τιμές για κάθε διάσταση ξεχωριστά στη νέα πλειάδα  $MDR$  και οι ελάχιστες τιμές για κάθε διάσταση ξεχωριστά στην πλειάδα  $MAR$ . Δημιουργούνται τέλος τα ευρετήρια  $RTree$  για το σύνολο  $W$  και πλέγματος για το σύνολο  $Q$ .

**Αλγόριθμος:** Αλγόριθμος 4.4

---

```

MDR = [0,0,...,0];
MAR = [MAX,MAX,...,MAX];
for (i = 0; i < n; i++) {
  for (j = 0; j < d; j++) {
    if (qi[j] > MDR[j]) {
      MDR[j] = qi[j];
    }
    if (qi[j] < MAR[j]) {
      MAR[j] = qi[j];
    }
  }
}
RTreeW = makeRtreeIndex(W);
GridQ = makeGridIndex(Q);

```

---

Αλγόριθμος 4.4: Αλγόριθμος προεπεξεργαστικής μονάδας

### Μονάδα εισόδου (input module)

**Είσοδος:**  $p_{in}, MAR, MDR, GridQ, RTreeW$

**Έξοδος:**  $COUNTER(W, Q)$

**Περιγραφή:** Κάθε νεοεισερχόμενη πλειάδα  $p_{in}$ , ελέγχεται αν κυριαρχείται  $pareto$  από την πλειάδα  $MDR$ . Αν η συνθήκη είναι αληθής, η πλειάδα  $p_{in}$  απορρίπτεται. Αλλιώς ελέγχεται αν η πλειάδα  $p_{in}$  κυριαρχεί  $pareto$  το σημείο  $MAR$  και αν ισχύει, αυξάνονται όλοι οι δείκτες του πίνακα  $COUNTER$  κατά 1. Αν δεν ισχύει τίποτα από τα δύο, ο αλγόριθμος ελέγχει ένα προς ένα όλα τα κελιά  $M$  της δομής  $GridQ$ . Για κάθε κελί  $M$  ελέγχεται αν κυριαρχεί  $pareto$  την πλειάδα  $p_{in}$ . Αν η συνθήκη είναι αληθής, το κελί  $M$  απορρίπτεται για την πλειάδα  $p_{in}$ . Αλλιώς ελέγχεται αν η πλειάδα  $p_{in}$  κυριαρχεί  $pareto$  το κελί  $M$  και αν ισχύει, αυξάνονται όλοι οι δείκτες του πίνακα  $COUNTER(W, q_i) \forall q_i \in M$  κατά 1. Αν δεν ισχύει τίποτα από αυτά, ελέγχεται η σχέση υπερίσχυσης για κάθε κελί της δομής  $GridQ$

και κάθε κόμβο της δομής  $RTreeW$ . Για τις περιπτώσεις που διαπιστωθεί υπερίσχυση του σημείου  $p_{in}$  έναντι μιας ομάδας  $M$  σημείων αναφοράς με χρήση μιας ομάδας  $V$  προτιμήσεων, αυξάνονται κατά 1 οι τιμές του  $COUNTER$  που αφορούν τα  $W_a \in V, q_i \in M$ .

#### Αλγόριθμος: Αλγόριθμος 4.5

```
function increaseCounter(a, i) {
    COUNTER(W_a, q_i)++;
}

if (MDR  $\succ^{pto}$   $p_{in}$ ) {
    BUFFER.remove( $p_{in}$ );
    exit;
}

else if ( $p_{in}$   $\succ^{pto}$  MAR) {
    for (a = 0; a < m; a++) {
        for (i = 0; i < n; i++) {
            increaseCounter(a, i);
        }
    }
    exit;
}

foreach (Cell M in GridQ) {
    if (M  $\succ^{pto}$   $p_{in}$ ) {
        continue;
    }
    else if ( $p_{in}$   $\succ^{pto}$  M) {
        for (a = 0; a < m; a++) {
            foreach ( $q_i$  in M) {
                increaseCounter(a, i);
            }
        }
    }
    V = RTreeW.getRoot();
    heapW.enqueue(V);
    while (!heapW.isEmpty()) {
        V = heapW.dequeue();
        if ( $u_V(M) < \ell_V(p_{in})$ ) {
            continue; //  $p_{in}$  does not affect  $RTOP_k(M)$  by V preferences
        }
        else if ( $u_V(p_{in}) < \ell_V(M)$ ) {
            for each ( $W_a$  in V) {
                for each ( $q_i$  in M) {
                    increaseCounter(a, i); //  $p_{in}$  affects all  $RTOP_k(M)$ 
                }
            }
        }
        else if !(V is point) {
            heapW.enqueue(V.expand);
        }
        else { // V is a point
            score( $p_{in}$ ) =  $f_V(p_{in})$ 
            for each ( $q_i$  in M) {
```

```

    if (score( $p_{in}$ ) < SCORE( $V, q_i$ )) {
        increaseCounter( $a, i$ );
    }
}
}
}
}
}

```

Αλγόριθμος 4.5: Αλγόριθμος μονάδας εισόδου

### Εκκαθαριστική Μονάδα (expunge module)

**Είσοδος:**  $p_{exp}, MAR, MDR, GridQ, RTreeW$

**Έξοδος:**  $COUNTER(W, Q)$

**Περιγραφή:** Η λειτουργία της Εκκαθαριστικής Μονάδας είναι ακριβώς ίδια με τη λειτουργία της Μονάδας Εισόδου με τη διαφορά ότι αντί της αύξησης των τιμών του πίνακα  $COUNTER(W, Q)$ , πραγματοποιεί μείωση αυτών κατά 1. Για οικονομία χώρου, παρατίθενται μόνο η αλλαγή της συνάρτησης `increaseCounter` που χρησιμοποιείται και στη Μονάδα Εισόδου.

**Αλγόριθμος:** Αλγόριθμος 4.6

```

function increaseCounter( $a, i$ ) {
    COUNTER( $W_a, q_i$ )--;
}

```

Αλγόριθμος 4.6: Αλγόριθμος μονάδας εισόδου

#### 4.4.3 Αλγόριθμος 3

Επιπλέον των σχέσεων κυριαρχίας και υπερίσχυσης, προτείνεται μια ακόμη βελτίωση για την αύξηση της αποδοτικότητας του συστήματος. Με χρήση της Ιδιότητας 7 μπορούν να προταθούν σημαντικές βελτιώσεις στην απόδοση της Μονάδας Εκκαθάρισης, όπως περιγράφεται στη συνέχεια:

Κατά την επεξεργασία ενός νεοεισερχόμενου σημείου  $p_{in}$  στην Μονάδα Εισόδου, είναι δυνατόν να παρατηρηθούν οι προτιμήσεις εκείνες, οι οποίες επηρεάζονται από το σημείο  $p_{in}$ . Είναι αυτές που επηρεάζουν τον δείκτη  $COUNTER(W_a, q_i)$ . Η μεταβλητή  $W_{inf}(p, Q)$ , αποτελεί τον πίνακα που περιέχει τις λίστες επηρεαζόμενων προτιμήσεων, δηλαδή τις λίστες εκείνες στις οποίες αποθηκεύονται όλες οι προτιμήσεις  $W_a$  που επηρεάζονται από ένα σημείο  $p_{in}$  για το σύνολο των ερωτημάτων  $Q$ . Συγκεκριμένα η τιμή  $W_{inf}(p, q_i)$  περιέχει τη λίστα με τις προτιμήσεις  $W_a$  οι οποίες επηρεάζουν το τελικό αποτέλεσμα του ερωτήματος  $q_i$  για το χρόνο ζωής του σημείου  $p$ .



Με τη χρήση της λίστας επηρεαζόμενων προτιμήσεων, η Εκκαθαριστική Μονάδα είναι δυνατόν να ολοκληρώσει την ενημέρωση του πίνακα  $COUNTER(W, Q)$  σε σημαντικό μικρότερο χρονικό διάστημα, αφού αποκτά την δυνατότητα άμεσης πρόσβασης στις τιμές αυτές που πρέπει να μειώσει κατά 1. Πιο συγκεκριμένα, η Εκκαθαριστική Μονάδα λαμβάνει ως είσοδο την πλειάδα προς διαγραφή  $p_{exp}$ , και αναζητά στον πίνακα  $W_{inf}(p, Q)$  τις λίστες προτιμήσεων που αφορούν στην πλειάδα  $p_{exp}$ . Για κάθε υφιστάμενη τιμή  $W_a$  στις λίστες  $W_{inf}(p_{exp}, q_i)$ , η Εκκαθαριστική Μονάδα μειώνει κατά 1 την τιμή του πίνακα  $COUNTER(W, Q)$  στο κελί που αντιστοιχεί στις τιμές  $W_a, q_i$ .

Η ορθότητα της παραπάνω μεθόδου είναι προφανής, λαμβάνοντας υπόψη την Ιδιότητα 7: Κατά την άφιξη ενός σημείου  $p_{in}$  είναι δυνατόν να παρατηρηθούν οι προτιμήσεις που επηρεάζονται. Έτσι κατά τη λήξη του σημείου αυτού, αρκεί να ελεγχθούν οι προτιμήσεις που επηρεάστηκαν στην άφιξή του. Με τον τρόπο αυτό, δεν είναι απαραίτητο να ελεγχθούν εκ νέου όλες οι προτιμήσεις του συστήματος, βελτιώνοντας σημαντικά την απόδοσή του.

Στη συνέχεια παρουσιάζονται συνοπτικά και με χρήση ψευδο-αλγορίθμων, η τροποποιημένη λειτουργία των βελτιωμένων εκδόσεων της Μονάδας Εισόδου και της Εκκαθαριστικής Μονάδας:

#### Προεπεξεργαστική Μονάδα (pre-processing module)

**Είσοδος:**  $W, Q$

**Έξοδος:**  $MDR, MAR, RTreeW, GridQ$

**Περιγραφή:** Ελέγχονται οι τιμές των διαστάσεων για κάθε  $q_i \in Q$ . Αποθηκεύονται οι μέγιστες τιμές για κάθε διάσταση ξεχωριστά στη νέα πλειάδα  $MDR$  και οι ελάχιστες τιμές για κάθε διάσταση ξεχωριστά στην πλειάδα  $MAR$ . Δημιουργούνται τέλος τα ευρετήρια  $RTree$  για το σύνολο  $W$  και πλέγματος για το σύνολο  $Q$ .

**Αλγόριθμος:** Αλγόριθμος 4.7

```
MDR = [0, 0, ..., 0];
MAR = [MAX, MAX, ..., MAX];
for (i = 0; i < n; i++) {
    for (j = 0; j < d; j++) {
        if (qi[j] > MDR[j]) {
            MDR[j] = qi[j];
        }
        if (qi[j] < MAR[j]) {
            MAR[j] = qi[j];
        }
    }
}
RTreeW = makeRtreeIndex(W);
GridQ = makeGridIndex(Q);
```

Αλγόριθμος 4.7: Αλγόριθμος προεπεξεργαστικής μονάδας

**Μονάδα εισόδου (input module)**

**Είσοδος:**  $p_{in}, MAR, MDR, GridQ, RTreeW$

**Έξοδος:**  $COUNTER(W, Q), W_{inf}(p_{in}, Q)$

**Περιγραφή:** Κάθε νεοεισερχόμενη πλειάδα  $p_{in}$ , ελέγχεται αν κυριαρχείται *pareto* από την πλειάδα  $MDR$ . Αν η συνθήκη είναι αληθής, η πλειάδα  $p_{in}$  απορρίπτεται. Αλλιώς ελέγχεται αν η πλειάδα  $p_{in}$  κυριαρχεί *pareto* το σημείο  $MAR$  και αν ισχύει, αυξάνονται όλοι οι δείκτες του πίνακα  $COUNTER$  κατά 1. Αν δεν ισχύει τίποτα από τα δύο, ο αλγόριθμος ελέγχει ένα προς ένα όλα τα κελιά  $M$  της δομής  $GridQ$ . Για κάθε κελί  $M$  ελέγχεται αν κυριαρχεί *pareto* την πλειάδα  $p_{in}$ . Αν η συνθήκη είναι αληθής, το κελί  $M$  απορρίπτεται για την πλειάδα  $p_{in}$ . Αλλιώς ελέγχεται αν η πλειάδα  $p_{in}$  κυριαρχεί *pareto* το κελί  $M$  και αν ισχύει, αυξάνονται όλοι οι δείκτες του πίνακα  $COUNTER(W, q_i) \forall q_i \in M$  κατά 1. Αν δεν ισχύει τίποτα από αυτά, ελέγχεται η σχέση υπερίσχυσης για κάθε κελί της δομής  $GridQ$  και κάθε κόμβο της δομής  $RTreeW$ . Για τις περιπτώσεις που διαπιστωθεί υπερίσχυση του σημείου  $p_{in}$  έναντι μιας ομάδας  $M$  σημείων αναφοράς με χρήση μιας ομάδας  $V$  προτιμήσεων, αυξάνονται κατά 1 οι τιμές του  $COUNTER$  που αφορούν τα  $W_a \in V, q_i \in M$ . Κάθε φορά που αυξάνεται κάποια τιμή του πίνακα  $COUNTER$  κατά 1, αποθηκεύονται οι τιμές  $p_{in}, q_i, W_a$  για τις οποίες πραγματοποιήθηκε η αύξηση του  $COUNTER$ .

**Αλγόριθμος:** Αλγόριθμος 4.8

```
function increaseCounter(a, i) {
    COUNTER(W_a, q_i)++;
    W_inf(p_in, q_i).add(W_a);
}

if (MDR  $\succ^{pto}$  p_in) {
    BUFFER.remove(p_in);
    exit;
}

else if (p_in  $\succ^{pto}$  MAR) {
    for (a = 0; a < m; a++) {
        for (i = 0; i < n; i++) {
            increaseCounter(a, i);
        }
    }
    exit;
}

foreach (Cell M in GridQ) {
    if (M  $\succ^{pto}$  p_in) {
        continue;
    }
}
```

```

else if ( $p_{in} \succ^{pto} M$ ) {
  for ( $a = 0; a < m; a++$ ) {
    foreach ( $q_i$  in  $M$ ) {
      increaseCounter( $a, i$ );
    }
  }
}
V = RTreeW.getRoot();
heapW.enqueue(V);
while (!heapW.isEmpty()) {
  V = heapW.dequeue();
  if ( $u_V(M) < \ell_V(p_{in})$ ) {
    continue; //  $p_{in}$  does not affect  $RTOP_k(M)$  by  $V$  preferences
  }
  else if ( $u_V(p_{in}) < \ell_V(M)$ ) {
    for each ( $W_a$  in  $V$ ) {
      for each ( $q_i$  in  $M$ ) {
        increaseCounter( $a, i$ ); //  $p_{in}$  affects all  $RTOP_k(M)$ 
      }
    }
  }
  else if (!( $V$  is point)) {
    heapW.enqueue( $V.expand$ );
  }
  else { //  $V$  is a point
    score( $p_{in}$ ) =  $f_V(p_{in})$ 
    for each ( $q_i$  in  $M$ ) {
      if ( $score(p_{in}) < SCORE(V, q_i)$ ) {
        increaseCounter( $a, i$ );
      }
    }
  }
}
}

```

Αλγόριθμος 4.8: Αλγόριθμος μονάδας εισόδου

### Εκκαθαριστική Μονάδα (exrunge module)

**Είσοδος:**  $p_{exp}, W_{inf}(p_{exp})$

**Έξοδος:**  $COUNTER(W, Q)$

**Περιγραφή:** Για κάθε προτίμηση  $W_a$  που υπάρχει σε όλες τις λίστες προτιμήσεων των κελιών  $W_{inf}(p_{exp}, q_i) \forall q_i \in Q$ , μειώνεται η τιμή του δείκτη  $COUNTER(W_a, q_i)$  κατά 1.

**Αλγόριθμος:** Αλγόριθμος 4.9

```

for ( $i = 0; i < n; i++$ ) {
  for each ( $W_a$  in  $W_{inf}(p_{exp}, q_i)$ ) {
     $COUNTER(W_a, q_i)--$ ;
  }
}

```

Αλγόριθμος 4.9: Αλγόριθμος εκκαθαριστικής μονάδας

Πανεπιστήμιο Πειραιώς

## Κεφάλαιο 5

# Πειραματική αξιολόγηση

Στο προηγούμενο κεφάλαιο παρουσιάστηκε και αναλύθηκε το προτεινόμενο μοντέλο επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Παρατέθηκαν επίσης τρεις αλγόριθμοι επεξεργασίας αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων, οι οποίοι μπορούν να εκτελούνται ακολουθώντας το προτεινόμενο μοντέλο.

Για την επιβεβαίωση της ορθότητας των προτεινόμενων αλγορίθμων, απαραίτητο είναι να υλοποιηθούν και να ελεγχθούν πειραματικά σε εργαστηριακό επίπεδο. Εκτός όμως της ορθότητας, εξίσου σημαντική είναι και η επιβεβαίωση των πλεονεκτημάτων (σε όρους αποδοτικότητας) των Αλγορίθμων 2 και 3. Στο κεφάλαιο αυτό παρατίθενται τα αποτελέσματα από την πειραματική αξιολόγηση που πραγματοποιήθηκε στους αλγορίθμους που παρουσιάστηκαν στο προηγούμενο κεφάλαιο. Τα αποτελέσματα που συλλέγησαν αναλύονται ώστε να δειχθεί ή μη η δυνατότητα των αλγορίθμων αυτών να υποστηρίξουν συστήματα λήψης αποφάσεων πραγματικού χρόνου.

Στόχοι του κεφαλαίου αυτού είναι:

- Η επιλογή των κατάλληλων παραμέτρων αξιολόγησης του προτεινόμενου μοντέλου και αλγορίθμων.
- Η παρουσίαση της οργάνωσης που ακολουθήθηκε για την εκτέλεση των εργαστηριακών πειραμάτων.
- Η παράθεση των αποτελεσμάτων υπό μορφή γραφημάτων για καλύτερη κατανόηση των μετρήσεων που πραγματοποιήθηκαν.
- Η εξαγωγή συμπερασμάτων για την αποδοτικότητα και την ορθότητα των προτεινόμενων μεθόδων επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων.

## 5.1 Παράμετροι αξιολόγησης

Τα συστήματα πραγματικού χρόνου, λόγω της φύσης τους, απαιτούν ιδιαίτερη προσοχή στην επίτευξη καλών επιδόσεων. Ένας από τους πιο σημαντικούς παράγοντες στη λειτουργία των συστημάτων πραγματικού χρόνου, αποτελεί η ταχύτητα επεξεργασίας των εισερχόμενων δεδομένων: το σύστημα πρέπει να είναι ικανό να επεξεργάζεται την εισερχόμενη πληροφορία με ρυθμό τουλάχιστον ίδιο με το ρυθμό εισόδου της πληροφορίας στο σύστημα. Για το σκοπό αυτό, και όπως αναφέρθηκε στο προηγούμενο κεφάλαιο, όλες οι δομές που χρησιμοποιούνται, αποθηκεύονται στην κύρια μνήμη του υπολογιστικού συστήματος, για την όσο το δυνατόν ταχύτερη επεξεργασία των πληροφοριών.

Από τις ιδιαιτερότητες αυτές των συστημάτων πραγματικού χρόνου, προκύπτει η ανάγκη αξιολόγησης της λειτουργίας του προτεινόμενου μοντέλου, λαμβάνοντας υπόψη το ρυθμό επεξεργασίας των δεδομένων και το μέγεθος της χρησιμοποιούμενης κύριας μνήμης. Ωστόσο στην παρούσα ερευνητική εργασία, στόχος αποτελεί και η σύγκριση των τριών προτεινόμενων αλγορίθμων, σε όρους υπολογιστικής πολυπλοκότητας. Έτσι οι παράμετροι αξιολόγησης που λαμβάνονται υπόψη στη συνέχεια, συνοψίζονται ως εξής:

- Μέσος ρυθμός επεξεργασίας της εισερχόμενης πληροφορίας.
- Μέγιστες απαιτήσεις σε χώρο κύριας μνήμης.
- Σύνολο πλήθους συγκρίσεων βαθμολόγησης (και συγκρίσεων κυριαρχίας) μεταξύ δύο σημείων του χώρου δεδομένων.

## 5.2 Οργάνωση πειραμάτων

Τα εργαστηριακά πειράματα εκτελέστηκαν χρησιμοποιώντας ένα ηλεκτρονικό υπολογιστή με Intel Core i7-2600 CPU (4 cores, 8 threads, 3.4 GHz) και 11 GB RAM σε περιβάλλον Windows 8.1. Όλοι οι αλγόριθμοι υλοποιήθηκαν σε γλώσσα Java και εκτελέστηκαν σε Java Runtime Environment 1.8.0\_20.

Οι μετρήσεις που παρουσιάζονται στη συνέχεια, πραγματοποιήθηκαν λαμβάνοντας υπόψη τις επιδόσεις του συστήματος μόνο αφού το *BUFFER* έχει γεμίσει με *T* πλήθους πλειάδες. Πραγματοποιήθηκε δηλαδή μέτρηση της επίδοσης του συστήματος, κατά την διάρκεια της πραγματικής του λειτουργίας και όχι κατά τη διάρκεια της αρχικοποίησής του.

Για τη μέτρηση του ρυθμού επεξεργασίας της εισερχόμενης πληροφορίας, μετρήθηκε ο απαιτούμενος χρόνος επεξεργασίας κάθε εισερχόμενης και εξερχόμενης πλειάδας. Η μέτρηση αυτή, πραγματοποιήθηκε χρησιμοποιώντας την εντολή της Java `System.nanoTime()`. Συγκεκριμένα μετρήθηκε η χρονική στιγμή εισόδου μιας πλειάδας στη Μονάδα Εισόδου και η χρονική στιγμή που η Μονάδα Εισόδου ολοκλήρωσε την επεξεργασία της πλειάδας αυτής. Οι δύο χρονικές στιγμές αφαιρούνται μεταξύ τους και προκύπτει το κόστος της πλειάδας σε νανοδευτερόλεπτα για τη Μονάδα Εισόδου. Αντίστοιχα εκτελείται και για την Εκκαθαριστική Μονάδα. Τα κόστη αυτά συγκεντρώνονται, ανάγονται σε δευτερόλεπτα και στη συνέχεια από αυτά εξάγεται το μέσο πλήθος των πλειάδων που επεξεργάζεται το σύστημα ανά δευτερόλεπτο. Η τιμή αυτή αποτελεί το ρυθμό επεξεργασίας που επιτυγχάνει το σύστημα, κατά μέσο όρο.

Για τη μέτρηση της μέγιστης χρήσης της κύριας μνήμης του συστήματος, υλοποιήθηκε κώδικας σε Java, ο οποίος πραγματοποιεί μέτρηση της χρησιμοποιούμενης κύριας μνήμης, εξετάζοντας τις δομές που βρίσκονται σε χρήση. Για την πραγματοποίηση της μέτρησης αυτής, θεωρείται ότι οι τιμές που μπορεί να λαμβάνει κάθε πλειάδα σε οποιαδήποτε διάστασή της, είναι μεγέθους 4 Bytes (Integer ή Float). Η μέτρηση αυτή, πραγματοποιείται κάθε φορά που εισέρχεται ή εξέρχεται κάποια πλειάδα από το σύστημα. Τελικά συγκεντρώνονται όλες οι διαθέσιμες μετρήσεις και επιλέγεται η μέγιστη από αυτές. Η τιμή αυτή, αποτελεί τη μέγιστη απαίτηση που έχει το σύστημα σε χώρο κύριας μνήμης.

Για τη μέτρηση της υπολογιστικής πολυπλοκότητας κάθε αλγορίθμου, χρησιμοποιήθηκε ένας μετρητής που αυξάνει κατά 1 κάθε φορά που πραγματοποιείται μια σύγκριση μεταξύ της βαθμολόγησης δύο σημείων του χώρου δεδομένων ή/και σύγκριση ύπαρξης κυριαρχίας μεταξύ αυτών. Η τιμή αυτή αποτελεί τη συγκριτική βαθμολόγηση της υπολογιστικής πολυπλοκότητας κάθε αλγορίθμου.

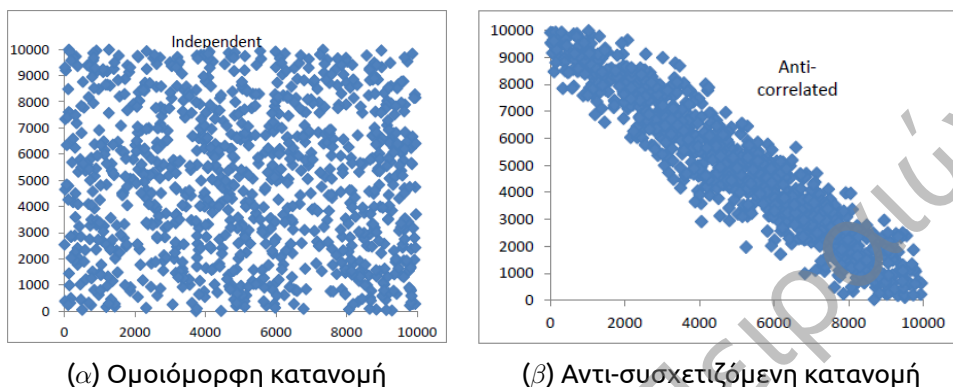
Τα σύνολα δεδομένων  $S$  και  $Q$  που μελετήθηκαν, παράγονται τυχαία χρησιμοποιώντας κάθε φορά μια εκ των δύο συνθετικών κατανομών: ομοιόμορφη (uniform - independent) και αντι-συσχετιζόμενη (anti-correlated). Οι λεπτομέρειες των κατανομών αυτών, παρουσιάζονται στη συνέχεια:

**Uniform - Independent:** Γι' αυτό το είδος της κατανομής, όλες οι τιμές των γνωρισμάτων παράγονται ανεξάρτητα χρησιμοποιώντας μια ομοιόμορφη κατανομή.

**Anti-correlated:** Μια τέτοιου είδους κατανομή αναπαριστά ένα περιβάλλον στο οποίο τα σημεία που είναι καλά σε μια διάσταση, δεν είναι καλά, σε μία ή περισσότερες από τις άλλες διαστάσεις. Τα σημεία παράγονται για παράδειγμα ως εξής: Επιλέγουμε ένα επίπεδο κάθετο στη γραμμή από  $(0, \dots, 0)$  έως  $(10.000, \dots, 10.000)$  χρησιμοποιώντας μια κανονική κατανομή. Χρησιμοποιούμε

μια κανονική κατανομή με πολύ μικρή απόκλιση ώστε όλα τα σημεία να τοποθετούνται σε επίπεδα που είναι κοντά στο επίπεδο μέσω του σημείου (5.000,...,5.000). Μέσα στο επίπεδο, οι διαφορετικές τιμές των γνωρισμάτων δημιουργούνται χρησιμοποιώντας μια ομοιόμορφη κατανομή.

Στο Σχήμα 5.1 παρουσιάζονται ενδεικτικά οι γραφικές απεικονίσεις για τις παραπάνω κατανομές δεδομένων, για 1000 σημεία στο διδιάστατο χώρο.



Σχήμα 5.1: Κατανομές δεδομένων

Για την πραγματοποίηση των πειραμάτων, επελέγησαν τέσσερις παράμετροι, ώστε μεταβάλλοντας τις τιμές αυτών, να μελετάται η επίδρασή των παραμέτρων στην απόδοση του συστήματος. Οι παράμετροι αυτές είναι οι εξής: (i) πλήθος διαστάσεων των συνόλων  $S$  και  $Q$ , (ii) μέγεθος του συνόλου  $Q$ , (iii) μέγεθος του συρόμενου χρονικού παραθύρου και (iv) μέγεθος του συνόλου  $W$ . Για κάθε μια από τις παραπάνω παραμέτρους, έγινε επιλογή μιας καθολικής τιμής η οποία παραμένει ίδια για κάθε πείραμα, μεταβάλλοντας κάθε φορά μόνο μια εκ των παραπάνω παραμέτρων. Ο Πίνακας 5.1 παρουσιάζει τονισμένες τις καθολικές τιμές των παραπάνω παραμέτρων, και μη τονισμένες τις υπόλοιπες τιμές που λαμβάνουν οι μεταβλητές στα εκάστοτε πειράματα. Πρέπει να σημειωθεί ότι κάθε σύνολο τιμών των παραπάνω παραμέτρων, ελέγχθηκε και στις δύο κατανομές που παρουσιάστηκαν νωρίτερα.

Παράμετρος	Τιμές
Dimensionality	2, 4, 6
$ Q $	50, 100, 200
Sliding Window	10K, 100K, 1000K
$ W $	10K, 20K, 40K

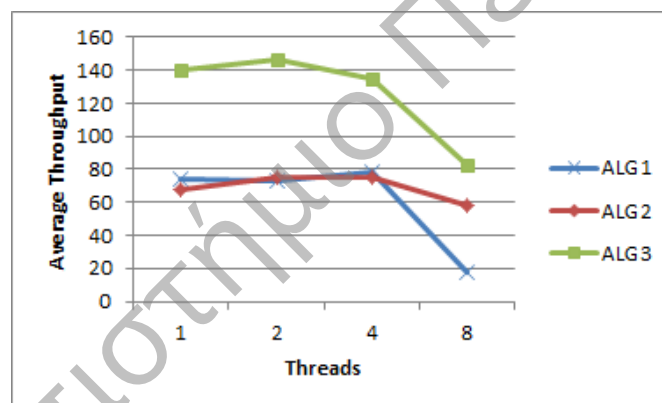
Πίνακας 5.1: Παράμετροι εκτέλεσης πειραμάτων

Σε αυτό το σημείο, πρέπει να σημειωθεί ότι οι αλγόριθμοι που αναπτύχθηκαν και για τις τρεις περιπτώσεις, εκμεταλλεύονται την ύπαρξη πολλών φυσικών πυρήνων στη

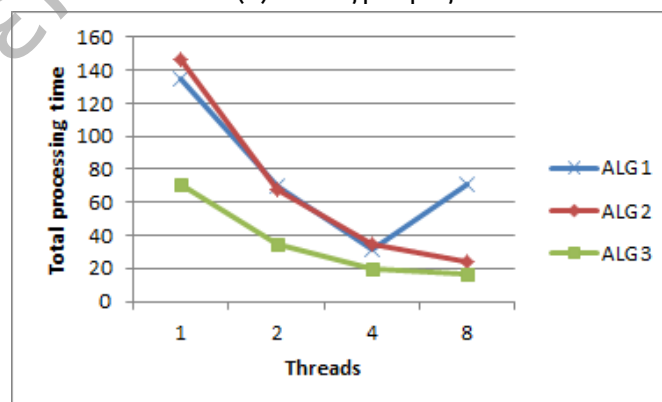


CPU του συστήματος, με τη χρήση νημάτων (threads). Συγκεκριμένα κάθε νεοεισερχόμενη (ή εξερχόμενη) πλειάδα, επεξεργάζεται από διαφορετικό νήμα, ώστε να χρησιμοποιούνται όλοι οι διαθέσιμοι πυρήνες του συστήματος. Η ύπαρξη νημάτων, στη λειτουργία του συστήματος, μειώνει το συνολικό απαιτούμενο χρόνο επεξεργασίας των πλειάδων. Λόγω όμως της αναγκαιότητας συγχρονισμού μεταξύ των νημάτων που εκτελούνται ταυτόχρονα στο σύστημα, η εκάστοτε πλειάδα απαιτεί παραπάνω χρόνο επεξεργασίας (συμπεριλαμβανομένου του χρόνου συγχρονισμού των νημάτων), μειώνοντας με αυτόν τον τρόπο το ρυθμό επεξεργασίας της εκάστοτε πλειάδας.

Η μείωση του ρυθμού, όπως φαίνεται και στο Σχήμα 5.2 δεν επηρεάζει αρνητικά το συνολικό απαιτούμενο χρόνο εκτέλεσης του εκάστοτε αλγορίθμου. Έπειτα από αυτήν την παρατήρηση, αποφασίστηκε όλα τα πειράματα της παρούσας ερευνητικής εργασίας, να εκτελούνται χρησιμοποιώντας σταθερά 8 νήματα. Οι ρυθμοί εκτέλεσης που θα παρουσιάζονται στη συνέχεια θα είναι ελαφρώς χειρότεροι, λόγω του φαινομένου που περιγράφηκε προηγουμένως, ωστόσο πρέπει να λαμβάνεται υπόψη ότι οι ρυθμοί αυτοί αφορούν μόνο ένα νήμα κάθε φορά.



(α) Μέσος ρυθμός



(β) Συνολικός χρόνος

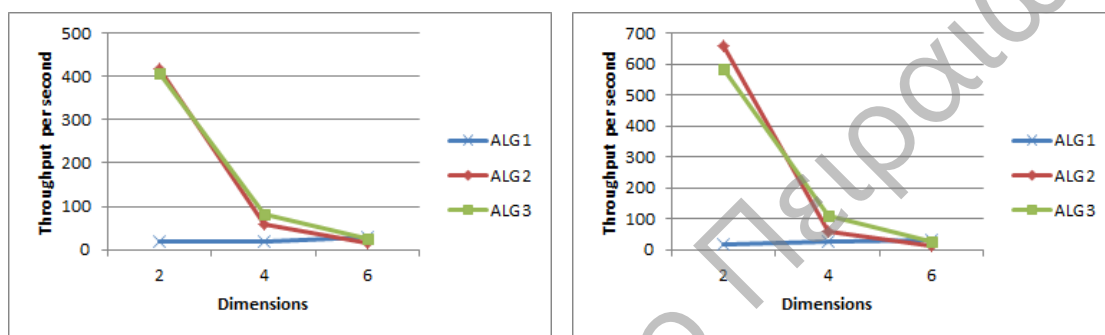
Σχήμα 5.2: Απόδοση του συστήματος μεταβάλλοντας τον αριθμό των νημάτων

Κατά τη διάρκεια λειτουργίας του κάθε πειράματος, το σύστημα, αρχικά δέχεται ως είσοδο  $T$  πλήθους πλειάδες, ώστε να γεμίσει ο *BUFFER*. Κατά την διαδικασία αυτή,

καμία μέτρηση δεν πραγματοποιείται στις επιδόσεις του συστήματος. Στη συνέχεια, δέχεται για 10 φορές (timestamps) δεδομένα εισόδου, ίσα με 1000 πλειάδες. Σύνολο δηλαδή κάθε πείραμα, δέχεται ως είσοδο 10K πλειάδες κατά τη διάρκεια της πραγματικής λειτουργίας του. Στη συνέχεια παρουσιάζονται τα αποτελέσματα των πειραμάτων που εκτελέστηκαν στα πλαίσια της παρούσας εργασίας.

## 5.3 Αποτελέσματα

### 5.3.1 Μεταβάλλοντας το πλήθος των διαστάσεων

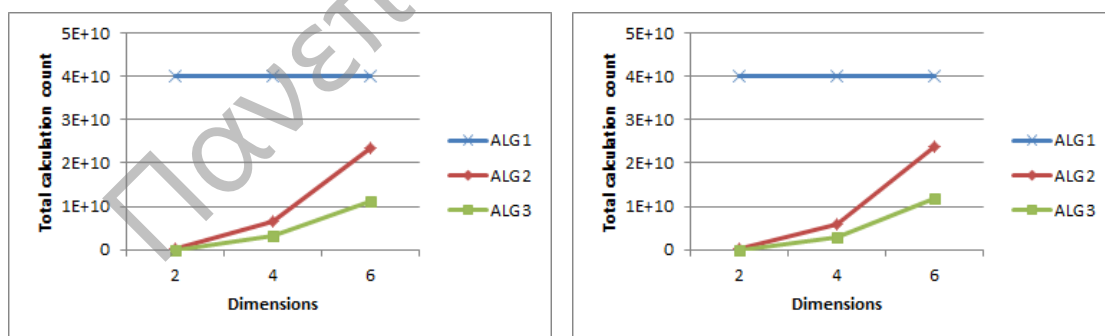


(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.3: Μέσος ρυθμός επεξεργασίας

Στο Σχήμα 5.3 φαίνεται καθαρά πως με την αύξηση των διαστάσεων των συνόλων δεδομένων, μειώνεται αισθητά η αποδοτικότητα του συστήματος. Ενδιαφέρον έχει το γεγονός, ότι στις 6 διαστάσεις και οι τρεις αλγόριθμοι αποδίδουν το ίδιο καλά.

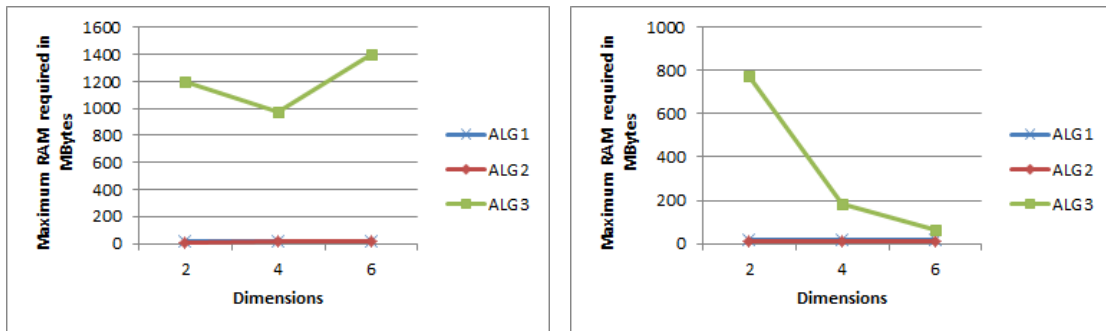


(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.4: Πολυπλοκότητα υπολογιστικής επεξεργασίας

Στο Σχήμα 5.4 φαίνεται πως η πολυπλοκότητα του Αλγορίθμου 1 δεν επηρεάζεται από το πλήθος των διαστάσεων, διατηρώντας όμως σταθερά το μέγιστο πλήθος των απαιτούμενων υπολογισμών. Ο Αλγόριθμος 3 παρουσιάζει τη μικρότερη υπολογιστική πολυπλοκότητα.



(α) Ομοιόμορφη κατανομή

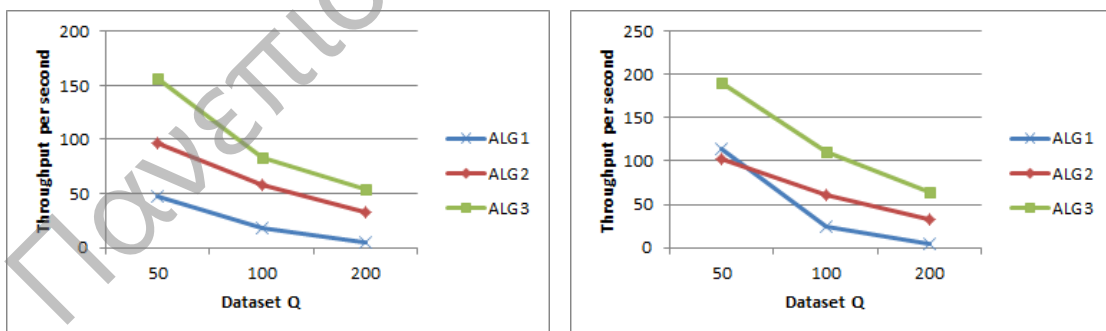
(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.5: Μέγιστος απαιτούμενος χώρος κύριας μνήμης

Στο Σχήμα 5.5 φαίνεται πως οι Αλγόριθμοι 2 και 3, έχουν περίπου τις ίδιες απαιτήσεις σε κύρια μνήμη ενώ ο αλγόριθμος 3 έχει σημαντικά μεγαλύτερες απαιτήσεις. Αυτό συμβαίνει λόγω της λίστας  $W_{inf}$  που συντηρεί ο Αλγόριθμος 3 και η οποία μπορεί να χρησιμοποιεί ένα πολύ μεγάλο μέγεθος κύριας μνήμης.

### 5.3.2 Μεταβάλλοντας το μέγεθος του $Q$

Στο Σχήμα 5.6 φαίνεται ότι η αύξηση του πλήθους των σημείων  $q_i \in Q$  μειώνει σχεδόν γραμμικά το ρυθμό επεξεργασίας του συστήματος. Εξάριση αποτελεί ο Αλγόριθμος 1 για την αντι-συσχετιζόμενη κατανομή, όπου για 50 query points, παρουσιάζει ίδιο ρυθμό επεξεργασίας με τον Αλγόριθμο 2. Συνολικά ο Αλγόριθμος 3 είναι καλύτερος των άλλων δύο, για οποιοδήποτε μέγεθος του συνόλου  $Q$ .

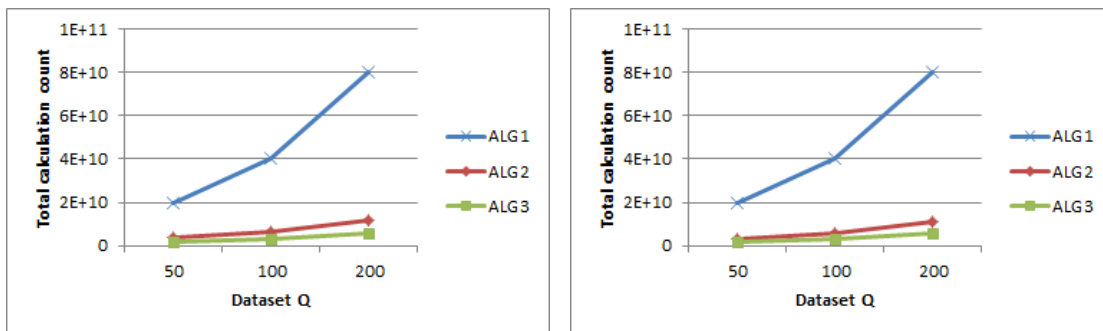


(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.6: Μέσος ρυθμός επεξεργασίας

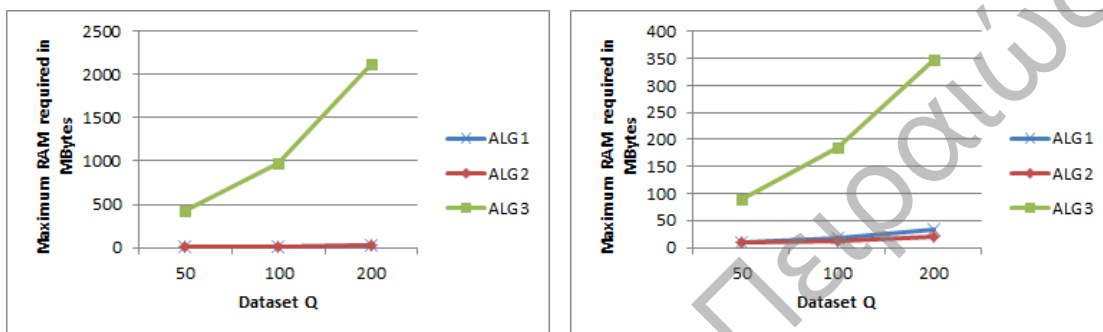
Στο Σχήμα 5.7 φαίνεται καθαρά ότι ο Αλγόριθμος 1, πραγματοποιεί πολύ μεγαλύτερο πλήθος υπολογισμών, όσο αυξάνεται το μέγεθος του συνόλου  $Q$ . Οι αλγόριθμοι 2 και 3 επηρεάζονται λιγότερο από την αύξηση του συνόλου αυτού, ενώ πάλι ο Αλγόριθμος 3 είναι ο καλύτερος, σε όλες τις περιπτώσεις.



(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.7: Πολυπλοκότητα υπολογιστικής επεξεργασίας



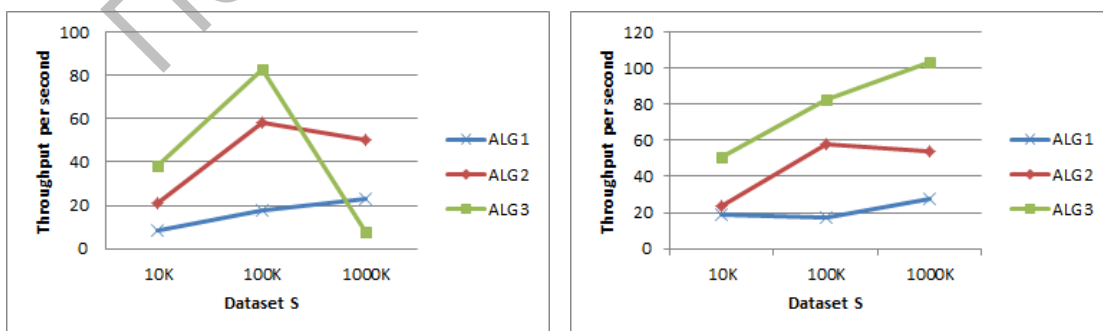
(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.8: Μέγιστος απαιτούμενος χώρος κύριας μνήμης

Στο Σχήμα 5.8 φαίνεται ότι όσο το σύνολο  $Q$  μεγαλώνει, ο Αλγόριθμος 3 αυξάνει σημαντικά τις απαιτήσεις του όσον αφορά την κύρια μνήμη που χρησιμοποιεί. Οι αλγόριθμοι 1 και 2 έχουν σχεδόν τις ίδιες απαιτήσεις σε κύρια μνήμη. Η αυξημένη χρήση κύριας μνήμης του Αλγορίθμου 3, οφείλεται στη χρήση της δομής  $W_{inf}$  που περιγράφηκε νωρίτερα.

### 5.3.3 Μεταβάλλοντας το μέγεθος του συρόμενου χρονικού παραθύρου

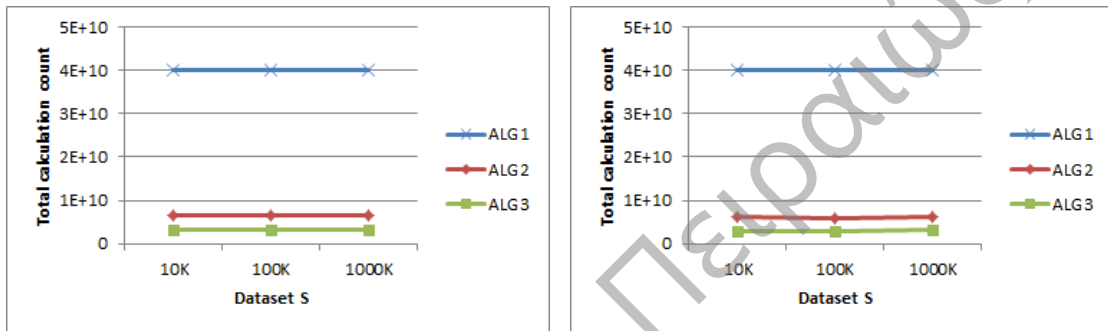


(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.9: Μέσος ρυθμός επεξεργασίας

Στο Σχήμα 5.9 φαίνεται ότι αυξάνοντας το μέγεθος του συρόμενου χρονικού παραθύρου, οι επιδόσεις των αλγορίθμων βελτιώνονται. Αυτό πιθανότατα οφείλεται στο πως το Λειτουργικό Σύστημα χειρίζεται το caching στη CPU. Ενδιαφέρον παρουσιάζει το γεγονός ότι για 1000K μέγεθος συρόμενου χρονικού παραθύρου, και ομοιόμορφη κατανομή, ο Αλγόριθμος 3 παρουσιάζει πολύ χειρότερες επιδόσεις απο οποιοδήποτε αλγόριθμο. Αυτό οφείλεται, όπως θα δουμε και στη συνέχεια, στις πολύ αυξημένες απαιτήσεις κύριας μνήμης που έχει ο Αλγόριθμος 3, λόγω της  $W_{inf}$ . Η δραματική αύξηση της ζήτησης της κύριας μνήμης, δημιουργεί καθυστέρηση στην εκτέλεση του Αλγορίθμου 3, λόγω του αυξημένου κόστους εκτέλεσης του garbage collector του Java Runtime Environment.

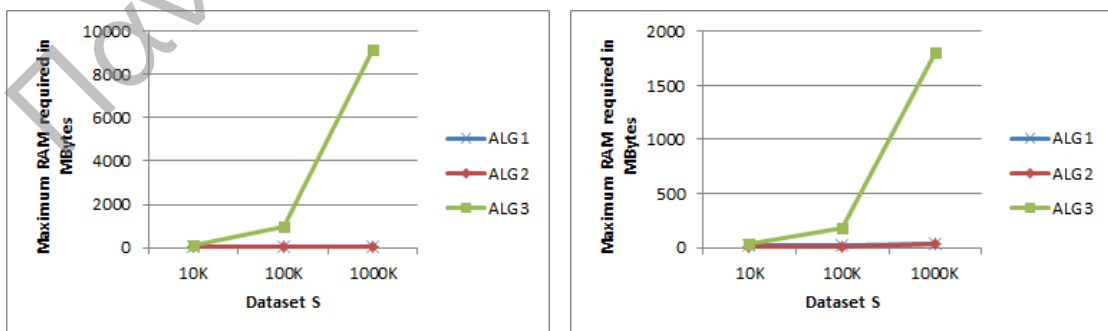


(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.10: Πολυπλοκότητα υπολογιστικής επεξεργασίας

Στο Σχήμα 5.10 φαίνεται, όπως αναμενόταν, ότι το πλήθος των υπολογισμών δεν επηρεάζεται από τη μεταβολή του μεγέθους του συρόμενου χρονικού παραθύρου. Αυτό συμβαίνει γιατί το πλήθος των υπολογισμών που απαιτούνται, εξαρτάται κάθε φορά, από το πλήθος των εισερχόμενων πλειάδων στο σύστημα και όχι από το μέγεθος του συρόμενου χρονικού παραθύρου. Ωστόσο φαίνεται πως ο Αλγόριθμος 1 απαιτεί σταθερά πολύ μεγαλύτερο πλήθος υπολογισμών από τους άλλους δύο αλγορίθμους.



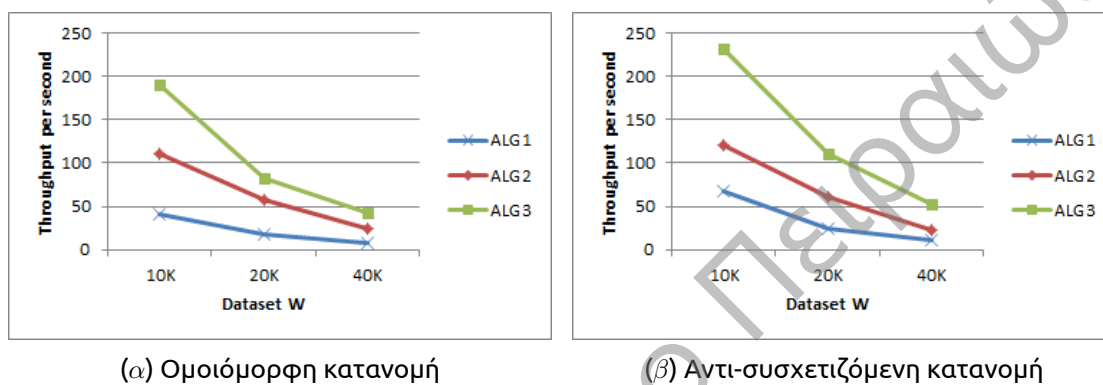
(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.11: Μέγιστος απαιτούμενος χώρος κύριας μνήμης

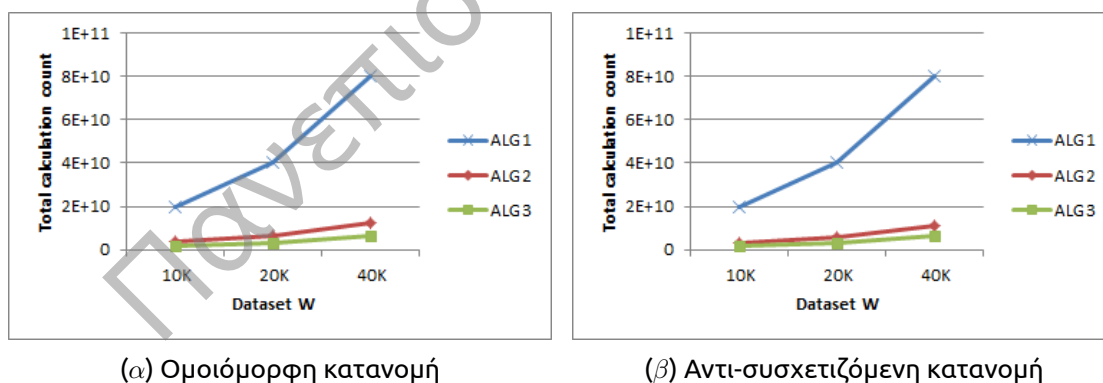
Στο Σχήμα 5.11 φαίνεται αυτό που περιγράφηκε και νωρίτερα. Ο Αλγόριθμος 3, έχει σημαντικά αυξημένες απαιτήσεις σε κύρια μνήμη από τους άλλους δύο. Στην περίπτωση της ομοιόμορφης κατανομής μάλιστα, οι απαιτήσεις του Αλγορίθμου 3 σε κύρια μνήμη, φτάνουν μέχρι και τα 9 GB. Ο Αλγόριθμος 3 με χρήση της  $W_{inf}$  δεν ενδείκνυται για χρήση σε περιβάλλοντα που απαιτούν μεγάλο συρόμενο χρονικό παράθυρο. Οι Αλγόριθμοι 2 και 3, έχουν περίπου ίδιες απαιτήσεις σε κύρια μνήμη, ανεξάρτητα από το μέγεθος του συρόμενου χρονικού παραθύρου.

### 5.3.4 Μεταβάλλοντας το μέγεθος του $W$



Σχήμα 5.12: Μέσος ρυθμός επεξεργασίας

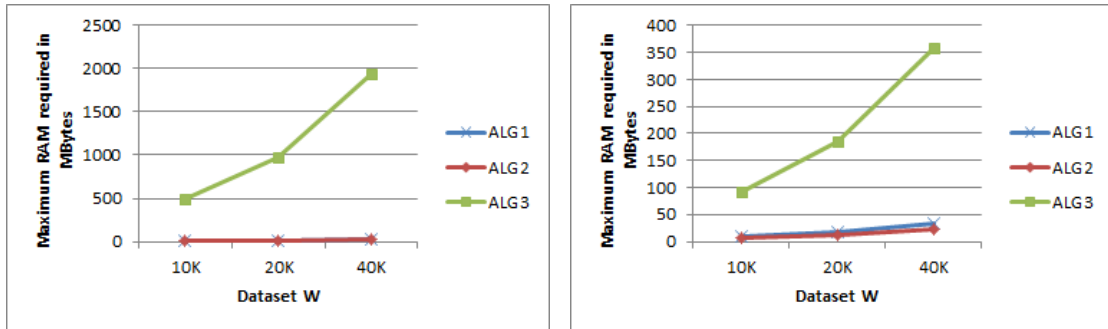
Στο Σχήμα 5.12 φαίνεται ότι η γραμμική μεταβολή του μεγέθους του συνόλου  $W$ , επηρεάζει σχεδόν γραμμικά το ρυθμό επεξεργασίας της εισερχόμενης πληροφορίας.



Σχήμα 5.13: Πολυπλοκότητα υπολογιστικής επεξεργασίας

Στο Σχήμα 5.13 φαίνεται ότι ο Αλγόριθμος 1 έχει σημαντικά μεγαλύτερη πολυπλοκότητα όσο αυξάνονται οι προτιμήσεις των χρηστών. Οι Αλγόριθμοι 2 και 3 επηρεάζονται πολύ λιγότερο.

Στο Σχήμα 5.14 φαίνονται για μια ακόμη φορά οι σημαντικά μεγαλύτερες ανάγκες του Αλγορίθμου 3 σε κύρια μνήμη. Ειδικά στην ομοιόμορφη κατανομή, η απαιτούμενη



(α) Ομοιόμορφη κατανομή

(β) Αντι-συσχετιζόμενη κατανομή

Σχήμα 5.14: Μέγιστος απαιτούμενος χώρος κύριας μνήμης

μνήμη αγγίζει τα 2 GB για 40K προτιμήσεις χρηστών. Οι αλγόριθμοι 2 και 3 έχουν παρόμοιες (και σημαντικά χαμηλότερες) απαιτήσεις σε κύρια μνήμη.

## 5.4 Συμπεράσματα αξιολόγησης

Από την εκτέλεση των πειραμάτων, γίνεται σαφές ότι ο Αλγόριθμος 3, είναι σε γενικές γραμμές ο πιο γρήγορος αλγόριθμος, έχοντας όμως πολύ μεγάλες απαιτήσεις σε κύρια μνήμη. Ο αλγόριθμος 2 επιτυγχάνει σημαντικά καλύτερες επιδόσεις από τον Αλγόριθμο 1, ενώ καταφέρνει να διατηρεί τις ίδιες απαιτήσεις σε κύρια μνήμη με εκείνον. Ο Αλγόριθμος 1, είναι όπως αναμενόταν ο πιο ακριβός αλγοριθμικός σε όρους υπολογιστικής πολυπλοκότητας και ρυθμού επεξεργασίας, δικαιολογώντας τον τίτλο του αφελή αλγορίθμου.

Η μεταβολή των διαστάσεων των συνόλων, έχει ιδιαίτερο ενδιαφέρον καθώς στις 2 διαστάσεις, οι Αλγόριθμοι 2 και 3 παρουσιάζουν πολύ υψηλές επιδόσεις σε σχέση με τις 4 και 6 διαστάσεις. Αντίθετα, στις 6 διαστάσεις, όλοι οι αλγόριθμοι αποδίδουν περίπου το ίδιο. Αυτό συμβαίνει γιατί στις 6 διαστάσεις η υπολογιστική πολυπλοκότητα των αλγορίθμων συγκλίνει σημαντικά, επιτρέποντας στο ρυθμό επεξεργασίας να μειωθεί. Ενδιαφέρον παρουσιάζει ακόμη το γεγονός ότι η αντι-συσχετιζόμενη κατανομή έχει καλύτερες επιδόσεις από την ομοιόμορφη, τόσο όσον αφορά το ρυθμό επεξεργασίας, όσο και για την απαιτούμενη κύρια μνήμη. Οι Αλγόριθμοι 1 και 2, έχουν ίδιες απαιτήσεις σε κύρια μνήμη, οι οποίες δεν αλλάζουν με τη μεταβολή του πλήθους των διαστάσεων. Ο Αλγόριθμος 3 αντίθετα επηρεάζεται από τις διαστάσεις με τρόπο που δεν μπορεί να προβλεφθεί, λόγω της δυναμικής φύσης της δομής  $W_{inf}$ , και στην αλλαγή της συμπεριφοράς της, όταν μεταβάλλεται το πλήθος των διαστάσεων.

Η μεταβολή του πλήθους των δεδομένων στα σύνολα  $Q$  και  $W$  παρουσιάζει παρόμοια αποτελέσματα. Ο Αλγόριθμος 3 επιτυγχάνει τους υψηλότερους ρυθμούς επεξεργασίας και τη χαμηλότερη υπολογιστική πολυπλοκότητα έχοντας όμως σημαντικά αυξημένες απαιτήσεις σε κύρια μνήμη. Ο Αλγόριθμος 2 καταφέρνει να αυξήσει δραματικά τους ρυθμούς επεξεργασίας και να μειώσει την υπολογιστική πολυπλοκότητα σε σχέση με τον Αλγόριθμο 1, ενώ παράλληλα επιτυγχάνει την ίδια χρήση κύριας μνήμης. Ο Αλγόριθμος 1, καταφέρνει να υπερισχύσει έναντι του Αλγορίθμου 2 μόνο σε μία περίπτωση, όπου τα δεδομένα εισόδου είναι πολύ μικρά. Ενδιαφέρον παρουσιάζει επίσης το γεγονός ότι οι επιδόσεις των αλγορίθμων σε αντι-συσχετιζόμενες κατανομές, είναι πολύ καλύτερες από τις επιδόσεις σε ομοιόμορφες κατανομές.

Τέλος ενδιαφέρον παρουσιάζει και η συμπεριφορά των επιδόσεων με τη μεταβολή του μεγέθους του συρόμενου χρονικού παραθύρου. Οι ρυθμοί επεξεργασίας ακολουθούν αυξητική τάση, αν εξαιρεθεί το γεγονός ότι για 1000K μέγεθος συρόμενου παραθύρου, ο Αλγόριθμος 3 παρουσιάζει πολύ μεγάλη πτώση. Ωστόσο αυτό συμβαίνει λόγω της ιδιαίτερα υψηλής απαίτησης σε κύρια μνήμη, και του χρόνου εκτέλεσης που απαιτεί ο garbage collector της Java για να εκτελεστεί. Η αυξητική τάση που παρατηρείται, μπορεί να δικαιολογηθεί, αν αναλογιστεί κανείς την ίδια υπολογιστική πολυπλοκότητα που παρουσιάζουν οι Αλγόριθμοι, ανεξάρτητα από το μέγεθος του συρόμενου παραθύρου. Το γεγονός αυτό, σε συνδυασμό με το καλύτερο caching που πραγματοποιεί το λειτουργικό σύστημα σε μια διεργασία που εκτελείται για μεγαλύτερο χρονικό διάστημα, δικαιολογούν τις αυξημένες επιδόσεις του συστήματος. Οι απαιτήσεις σε κύρια μνήμη για τον Αλγόριθμο 3, γίνονται απαγορευτικές για μεγέθη 1000K, επιτρέποντας μόνο τη χρήση των άλλων δύο αλγορίθμων σε συστήματα πραγματικού χρόνου. Για μια ακόμη φορά, οι ομοιόμορφες κατανομές, παρουσιάζουν χειρότερες επιδόσεις από τις αντι-συσχετιζόμενες.



## Κεφάλαιο 6

# Συμπεράσματα και μελλοντική έρευνα

### 6.1 Σύνοψη

Στην εποχή του Διαδικτύου των Πραγμάτων (Internet of Things) το πλήθος των διασυνδεδεμένων συσκευών αυξάνεται συνεχώς. Συσκευές συλλέγουν διαρκώς δεδομένα, χωρίς να απαιτείται η ανθρώπινη παρέμβαση, και τα αποστέλλουν προς κεντρικούς εξυπηρετητές για επεξεργασία ή και αποθήκευση. Αυτό έχει σαν αποτέλεσμα τη διαρκή αύξηση του όγκου των διακινούμενων δεδομένων και της ανάγκης για αποδοτική διαχείρισή τους, τόσο σε όρους ταχύτητας επεξεργασίας τους, όσο και σε όρους αποδοτικής αποθήκευσής τους. Οι υποδομές υπολογιστικών νεφών μπορούν να βοηθήσουν αποτελεσματικά στη διαχείριση του τεράστιου όγκου της διακινούμενης πληροφορίας.

Ο χώρος των βάσεων δεδομένων δεν μπορεί να μείνει ανεπηρέαστος από την ολοένα αυξανόμενη χρήση των υπολογιστικών νεφών. Νέες τεχνικές για την αποδοτική διαχείριση του πολύ μεγάλου όγκου πληροφορίας κάνουν την εμφάνισή τους. Οι τεχνικές αυτές δημιουργούν ένα νέο ερευνητικό πεδίο, που ονομάζεται Ανάλυση Μεγάλων Δεδομένων (Big Data Analytics). Ένα πολύ σύνηθες χαρακτηριστικό των Μεγάλων Δεδομένων αποτελεί η ύπαρξη πολλών διαστάσεων δεδομένων, οι οποίες μπορεί να μην είναι εξ αρχής γνωστές ή σαφώς καθορισμένες.

Η πολυκριτηριακή ανάλυση αποφάσεων, αποτελεί κλάδο της επιχειρησιακής έρευνας, που δεν μπορεί να μείνει ανεπηρέαστος από τις αλλαγές στη φύση των διακινούμενων δεδομένων. Το φαινόμενο της αύξησης των διαστάσεων των δεδομένων, επηρεάζει σημαντικά την αποδοτικότητα των υπάρχουσών τεχνικών υποστήριξης της

διαδικασίας λήψης απόφασης. Έτσι εμφανίστηκε η ανάγκη ανάπτυξης ενός νέου συνόλου αποδοτικών τεχνικών, οι οποίες θα εκμεταλλεύονται τη φύση των υποδομών υπολογιστικών νεφών. Στην παρούσα διπλωματική εργασία, μελετήθηκε η υποστήριξη της διαδικασίας λήψης αποφάσεων, μέσω των ερωτημάτων k κορυφαίων σημείων.

Στην υπάρχουσα βιβλιογραφία των ερωτημάτων k κορυφαίων σημείων, οι τεχνικές που μελετώνται, θέτουν στο επίκεντρό τους τον πελάτη ενός προϊόντος ή μιας υπηρεσίας: εξετάζονται τρόποι αποδοτικής εξυπηρέτησης των αναγκών του, παρέχοντας σε κάθε πελάτη εξατομικευμένη αναζήτηση προϊόντων ή και υπηρεσιών που ανταποκρίνονται στις δικές του προτιμήσεις. Στην παρούσα εργασία ωστόσο, τα ερωτήματα k κορυφαίων σημείων εξετάζονται από τη σκοπιά του παρόχου του προϊόντος ή της υπηρεσίας: εξετάζονται αποδοτικές μέθοδοι αξιολόγησης ενός προϊόντος ή μιας υπηρεσίας, σύμφωνα με τις προτιμήσεις υποψηφίων πελατών (χρηστών). Πιο συγκεκριμένα, στα αντίστροφα ερωτήματα k κορυφαίων σημείων, παρέχεται η δυνατότητα αναζήτησης των χρηστών εκείνων, οι οποίοι θα προτιμούσαν το προϊόν ενός παρόχου έναντι ενός άλλου. Στη γενική τους περίπτωση τα αντίστροφα ερωτήματα k κορυφαίων σημείων, παρουσιάζουν τις προτιμήσεις εκείνες, για τις οποίες ένα προϊόν ή μια υπηρεσία ανήκει στα k πλήθους καλύτερα.

Σε ένα περιβάλλον όμως που παρατηρούνται διαρκώς αλλαγές, νέα προϊόντα ή υπηρεσίες κάνουν διαρκώς την εμφάνισή τους. Προς το σκοπό αυτό, είναι σημαντική η δημιουργία ενός συστήματος που παρακολουθεί τις εξελίξεις και ενημερώνει τον πάροχο, οποιαδήποτε στιγμή το θελήσει, για την προτίμηση ή μη ενός προϊόντος του. Το σύστημα αυτό, ακολουθεί τους κανόνες των συστημάτων πραγματικού χρόνου, όπου ο απαιτούμενος χρόνος επεξεργασίας της εισερχόμενης πληροφορίας, είναι εξίσου σημαντικός με την ορθότητα των αποτελεσμάτων που παράγονται. Η παρούσα διπλωματική εργασία, ακολουθώντας τις επιταγές της σύγχρονης εποχής, πραγματεύεται την αποδοτική διαχείριση των αντίστροφων ερωτημάτων k κορυφαίων σημείων σε συστήματα πραγματικού χρόνου, μέσω ροών δεδομένων.

Στο Κεφάλαιο 3 παρουσιάστηκαν ένα σύνολο από τεχνικές αντιμετώπισης παρόμοιων προβλημάτων. Πρέπει ωστόσο να αναφερθεί ότι από την έρευνα που πραγματοποιήθηκε δεν εντοπίστηκε η ύπαρξη μεθόδου επίλυσης αντίστροφων ερωτημάτων k κορυφαίων σημείων σε ροές δεδομένων. Έτσι στο Κεφάλαιο 3 παρουσιάστηκαν μέθοδοι που αφορούν στην επίλυση παρόμοιων ερωτημάτων σε συστήματα ροών δεδομένων. Ορισμένα παραδείγματα αποτελούν τα ερωτήματα κορυφογραμμής και τα ερωτήματα k κορυφαίων σημείων. Από την διαδικασία αυτή, εξήχθησαν πολύ σημαντικά συμπεράσματα για τις υπάρχουσες τεχνικές προσέγγισης παρόμοιων ερωτημάτων.

Στο Κεφάλαιο 4 αναλύθηκε η προτεινόμενη αρχιτεκτονική του συστήματος επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Η αρχιτεκτονική αυτή αποτελείται από τρεις διακριτές και σαφώς ορισμένες μονάδες επεξεργασίας δεδομένων. Περιγράφεται επίσης ο τρόπος με τον οποίο τα δεδομένα αποθηκεύονται στην κύρια μνήμη, για την αύξηση της αποδοτικότητας του προτεινόμενου μοντέλου. Στη συνέχεια με την παρατήρηση ενός συνόλου πολύ ενδιαφέρουσων ιδιοτήτων, προτείνονται τρεις αλγόριθμοι επεξεργασίας ερωτημάτων  $k$  κορυφαίων σημείων. Ο πρώτος αλγόριθμος, αποτελεί τη βασική μέθοδο επεξεργασίας τέτοιου είδους ερωτημάτων, ενώ οι άλλοι δύο βελτιώνουν την αποδοτικότητα του μοντέλου εκμεταλλευόμενοι την ύπαρξη των δοθέντων ιδιοτήτων.

Στο Κεφάλαιο 5 εξετάζονται οι προτεινόμενες τεχνικές ως προς την αποδοτικότητά τους. Με την υλοποίηση των προτεινόμενων αλγόριθμων, κατέστη δυνατή η σύγκριση των επιδόσεών τους σε όρους ρυθμού επεξεργασίας, υπολογιστικής πολυπλοκότητας και απαιτούμενου χώρου κύριας μνήμης. Μέσα από τη διαδικασία αυτή, έγινε σαφές ότι ο Αλγόριθμος 3 είναι ο πιο αποδοτικός σε όρους υπολογιστικής πολυπλοκότητας και ρυθμού επεξεργασίας, έχοντας όμως πολύ μεγάλες απαιτήσεις σε κύρια μνήμη. Αυτό έχει σαν αποτέλεσμα τη μη ικανοποιητική απόδοσή του σε πολύ μεγάλο όγκο δεδομένων. Ο Αλγόριθμος 2 από την άλλη, καταφέρνει μια εντυπωσιακή αύξηση του ρυθμού επεξεργασίας σε σχέση με τον Αλγόριθμο 1, διατηρώντας παράλληλα ίδιες απαιτήσεις σε κύρια μνήμη. Έτσι ανάλογα με τις ανάγκες του εκάστοτε συστήματος, μπορεί κάθε φορά να επιλέγεται ο κατάλληλος αλγόριθμος επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων.

## 6.2 Συμπεράσματα

Μέσα από την εκπόνηση της παρούσας ερευνητικής διπλωματικής εργασίας και κυρίως κατά τη διαδικασία αξιολόγησης του προτεινόμενου μοντέλου επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων, προέκυψαν οι εξής παρατηρήσεις:

1. Στο χώρο της Ανάλυσης Μεγάλων Δεδομένων απαιτείται η δημιουργία νέων, πιο αποδοτικών μεθόδων, για την υποστήριξη της διαδικασίας λήψης απόφασης. Οι παραδοσιακές μέθοδοι δεν μπορούν να έχουν εφαρμογή σε συστήματα υπολογιστικών νεφών.
2. Τα αντίστροφα ερωτήματα  $k$  κορυφαίων σημείων, μπορούν να αποδειχθούν ένα πολύ χρήσιμο εργαλείο στη διαδικασία λήψης απόφασης για το σχεδιασμό ενός νέου προϊόντος ή μιας υπηρεσίας.

3. Η ύπαρξη πολλών διαστάσεων στα σύνολα δεδομένων, μπορεί να μειώσει δραματικά το ρυθμό επεξεργασίας της εισερχόμενης πληροφορίας, χωρίς κατ' ανάγκη να αυξάνεται η υπολογιστική πολυπλοκότητα των αλγορίθμων.
4. Τα μοντέλα επεξεργασίας ερωτημάτων ροών δεδομένων, απαιτούν πολύ καλούς ρυθμούς επεξεργασίας της εισερχόμενης πληροφορίας, ενώ παράλληλα πρέπει να διατηρείται καλή διαχείριση της διαθέσιμης κύριας μνήμης.
5. Το λειτουργικό σύστημα καθώς και η γλώσσα προγραμματισμού στην οποία υλοποιείται ένα σύστημα πραγματικού χρόνου, μπορεί να έχουν ισχυρή επίδραση στη συνολική αποδοτικότητά του.
6. Η χρήση ευρετηριακών δομών στα σύνολα δεδομένων (όπως R Tree και Grid), μπορεί να αυξήσει σε σημαντικό βαθμό τη συνολική απόδοση του συστήματος, εφόσον τα σύνολα δεδομένων είναι πολύ μεγάλα. Όπως παρατηρήθηκε από τα πειράματα που εκτελέστηκαν, η μείωση του ρυθμού επεξεργασίας της εισερχόμενης πληροφορίας, στους αλγορίθμους που χρησιμοποιούν ευρετηριακές δομές, δεν ακολουθεί τη γραμμικότητα της αύξησης του πλήθους των δεδομένων.
7. Η χρήση ευρετηριακών δομών στα σύνολα δεδομένων (όπως R Tree και Grid), μπορεί να μειώσει την αποδοτικότητα του προτεινόμενου μοντέλου, εφόσον τα σύνολα δεδομένων είναι πολύ μικρά. Παράδειγμα αποτελεί το πείραμα που διεξήχθη με χρήση 50 query points σε αντι-συσχετιζόμενη κατανομή.
8. Οι πολύ μεγάλες απαιτήσεις σε κύρια μνήμη, μπορούν να μειώσουν δραματικά το ρυθμό επεξεργασίας των εισερχόμενων δεδομένων, είτε λόγω της ύπαρξης του garbage collector, είτε λόγω της μη διαθεσιμότητας αρκετής κύριας μνήμης στο υπολογιστικό σύστημα.

### 6.3 Μελλοντική έρευνα

Το μοντέλο που παρουσιάστηκε στην παρούσα διπλωματική εργασία, αποτελεί ένα πολύ αποδοτικό μοντέλο επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Ωστόσο ορισμένα ενδιαφέροντα ζητήματα παραμένουν ανοικτά προς μελλοντική έρευνα:

Αρκετό ενδιαφέρον θα παρουσιάζε, η διεξαγωγή πειραματικών μετρήσεων σε σύνολα πραγματικών δεδομένων και η συγκριτική μελέτη των αποτελεσμάτων με εκείνα που παρατηρήθηκαν στα σύνολα συνθετικών δεδομένων. Θα μπορούσαν επίσης να χρησιμοποιηθούν ανάμεικτα σύνολα δεδομένων, όπου τα σημεία αλλού κατανέμονται στο χώρο ως anti-correlated και αλλού ως independent.

Ακόμη θα μπορούσαν να μελετηθούν αντίστοιχες μέθοδοι επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων σε κατανεμημένες ροές δεδομένων. Το μοντέλο αυτό, είναι αρκετά σύνηθες ειδικά στην εποχή της νεφοϋπολογιστικής, όπου διασυνδεδεμένες συσκευές από όλο τον κόσμο επικοινωνούν μεταξύ τους, παράγοντας πολύ μεγάλο όγκο διακινούμενης πληροφορίας.

Ενδιαφέρον θα παρουσίαζε επίσης η προσπάθεια μιας πιο αποδοτικής προσέγγισης στη δημιουργία αλγορίθμου επίλυσης αντίστροφων ερωτημάτων  $k$  κορυφαίων σημείων. Τα προβλήματα που παρατηρήθηκαν στους προτεινόμενους αλγορίθμους, αφορούν κυρίως στην πολύ μεγάλη απαίτηση σε κύρια μνήμη (ειδικά για τον Αλγόριθμο 3), καθώς και το πρόβλημα της δραματικής μείωσης της απόδοσης σε δεδομένα πολλών διαστάσεων.

Πανεπιστήμιο Πειραιώς

## Παράρτημα Α

### Δείγμα πειραματικών μετρήσεων

W10K ALG1 (NS)	W10K ALG2 (NS)	W10K ALG3 (NS)
18209738	6340	199832
23740125	0	604
20020599	146585898	12420052
25332439	27344339	31996
25457711	7641298	35620
25409413	15649960	5411756
25119326	17481951	6820841
25289575	5419907	6522603
25541930	11770750	23244
25178188	302	604
25337268	12319231	17585188
24779431	302	302
25007336	12151698	13049733
25127475	12258556	14484778
25904163	6224969	6918644
25947027	17150810	7924142
24887799	11011873	6296509
25879108	22719534	32508569
25272671	13601535	14008140
25145889	5075483	5769763
24567524	302	302
24988922	302	0
25285349	13007775	10733558
24877534	6272662	6089131
25131400	12530232	905
24973226	5839190	302
25314027	12841751	13892226
24837690	12233804	11636724

Πανεπιστήμιο Πειραιώς



# Βιβλιογραφία

- [1] Cisco visual networking index: Forecast and methodology, 2013–2018. Technical report, 2014.
- [2] Hermann Kopetz. Internet of things. In *Real-Time Systems*, pages 307–323. Springer, 2011.
- [3] Steven Miller. Big data analytics. 2013.
- [4] R Bellman. Curse of dimensionality. *Adaptive control processes: a guided tour*. Princeton, NJ, 1961.
- [5] Bernard Roy. *Multicriteria methodology for decision aiding*, volume 12. Springer, 1996.
- [6] Peter Mell and Tim Grance. The nist definition of cloud computing. 2011.
- [7] Divyakant Agrawal, Sudipto Das, and Amr El Abbadi. Big data and cloud computing: current state and future opportunities. In *Proceedings of the 14th International Conference on Extending Database Technology*, pages 530–533. ACM, 2011.
- [8] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [9] Harold Lim, Yuzhang Han, and Shivnath Babu. How to fit when no one size fits. In *CIDR*, 2013.
- [10] Albert Bifet. Mining big data in real time. *Informatica (Slovenia)*, 37(1):15–20, 2013.
- [11] Chengwei Wang, Infantdani Abel Rayan, and Karsten Schwan. Faster, larger, easier: Reining real-time big data processing in cloud. In *Proceedings of the Posters and Demo Track, Middleware '12*, pages 4:1–4:2, New York, NY, USA, 2012. ACM. ISBN 978-1-4503-1612-5. doi: 10.1145/2405153.2405157. URL <http://doi.acm.org/10.1145/2405153.2405157>.

- [12] Rayman Preet Singh, S. Keshav, and Tim Brecht. A cloud-based consumer-centric architecture for energy data analytics. In *Proceedings of the Fourth International Conference on Future Energy Systems, e-Energy '13*, pages 63–74, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-2052-8. doi: 10.1145/2487166.2487174. URL <http://doi.acm.org/10.1145/2487166.2487174>.
- [13] Kyriakos Mouratidis, Spiridon Bakiras, and Dimitris Papadias. Continuous monitoring of top-k queries over sliding windows. In *Proceedings of the 2006 ACM SIGMOD International Conference on Management of Data, SIGMOD '06*, pages 635–646, New York, NY, USA, 2006. ACM. ISBN 1-59593-434-0. doi: 10.1145/1142473.1142544. URL <http://doi.acm.org/10.1145/1142473.1142544>.
- [14] Brian Babcock and Chris Olston. Distributed top-k monitoring. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data, SIGMOD '03*, pages 28–39, New York, NY, USA, 2003. ACM. ISBN 1-58113-634-X. doi: 10.1145/872757.872764. URL <http://doi.acm.org/10.1145/872757.872764>.
- [15] Akrivi Vlachou, Christos Doulkeridis, and Kjetil Nørkvåg. Distributed top-k query processing by exploiting skyline summaries. *Distributed and Parallel Databases*, 30(3-4):239–271, 2012.
- [16] Parisa Haghani, Sebastian Michel, and Karl Aberer. Evaluating top-k queries over incomplete data streams. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pages 877–886. ACM, 2009.
- [17] Martin Burkhart and Xenofontas Dimitropoulos. Fast privacy-preserving top-k queries using secret sharing. In *Computer Communications and Networks (ICCCN), 2010 Proceedings of 19th International Conference on*, pages 1–7. IEEE, 2010.
- [18] Man Lung Yiu and Nikos Mamoulis. Multi-dimensional top-k dominating queries. *The VLDB Journal*, 18(3):695–718, 2009.
- [19] Xi Zhang and Jan Chomicki. Semantics and evaluation of top-k queries in probabilistic databases. *Distributed and parallel databases*, 26(1):67–126, 2009.
- [20] Akrivi Vlachou, Christos Doulkeridis, Kjetil Nørkvåg, and Yannis Kotidis. Top-k queries. In *Peer-to-Peer Query Processing over Multidimensional Data*, pages 63–72. Springer, 2012.
- [21] Tingjian Ge, Stan Zdonik, and Samuel Madden. Top-k queries on uncertain data: on score distribution and typical answers. In *Proceedings of the 2009 ACM*

- SIGMOD International Conference on Management of data*, pages 375–388. ACM, 2009.
- [22] M.A Soliman, IF. Ilyas, and K. Chen-Chuan Chang. Top-k query processing in uncertain databases. In *Data Engineering, 2007. ICDE 2007. IEEE 23rd International Conference on*, pages 896–905, April 2007. doi: 10.1109/ICDE.2007.367935.
- [23] Akrivi Vlachou, Christos Doulkeridis, Yannis Kotidis, and Kjetil Norvag. Reverse top-k queries. In *Data Engineering (ICDE), 2010 IEEE 26th International Conference on*, pages 365–376. IEEE, 2010.
- [24] Ian Foster, Yong Zhao, Ioan Raicu, and Shiyong Lu. Cloud computing and grid computing 360-degree compared. In *Grid Computing Environments Workshop, 2008. GCE'08*, pages 1–10. Ieee, 2008.
- [25] Rajkumar Buyya, Chee Shin Yeo, Srikumar Venugopal, James Broberg, and Ivona Brandic. Cloud computing and emerging it platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation computer systems*, 25(6):599–616, 2009.
- [26] Michael Armbrust, Armando Fox, Rean Griffith, Anthony D Joseph, Randy Katz, Andy Konwinski, Gunho Lee, David Patterson, Ariel Rabkin, Ion Stoica, et al. A view of cloud computing. *Communications of the ACM*, 53(4):50–58, 2010.
- [27] Sikder Sunbeam Islam, Muhammad Baqer Mollah, Md Imanul Huq, and M Aman Ullah. Cloud computing for future generation of computing technology. In *Cyber Technology in Automation, Control, and Intelligent Systems (CYBER), 2012 IEEE International Conference on*, pages 129–134. IEEE, 2012.
- [28] Shyam Kumar Doddavula and Amit Wasudeo Gawande. Adopting cloud computing: enterprise private clouds. *SETLabs Briefings*, 7(7):11–18, 2009.
- [29] Daniel J Abadi. Data management in the cloud: Limitations and opportunities. *IEEE Data Eng. Bull.*, 32(1):3–12, 2009.
- [30] Jim Gray et al. The transaction concept: Virtues and limitations. In *VLDB*, volume 81, pages 144–154, 1981.
- [31] Seth Gilbert and Nancy Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *ACM SIGACT News*, 33(2): 51–59, 2002.

- [32] Brian F Cooper, Eric Baldeschwieler, Rodrigo Fonseca, James J Kistler, PPS Narayan, Chuck Neerdaels, Toby Negrin, Raghu Ramakrishnan, Adam Silberstein, Utkarsh Srivastava, et al. Building a cloud for yahoo! *IEEE Data Eng. Bull.*, 32(1):36–43, 2009.
- [33] Azza Abouzeid, Kamil Bajda-Pawlikowski, Daniel Abadi, Avi Silberschatz, and Alexander Rasin. Hadoopdb: an architectural hybrid of mapreduce and dbms technologies for analytical workloads. *Proceedings of the VLDB Endowment*, 2(1):922–933, 2009.
- [34] S. Sakr, A. Liu, D.M. Batista, and M. Alomari. A survey of large scale data management approaches in cloud environments. *Communications Surveys Tutorials, IEEE*, 13(3):311–336, Third 2011. ISSN 1553-877X. doi: 10.1109/SURV.2011.032211.00087.
- [35] Hermann Kopetz. *Real-time systems: design principles for distributed embedded applications*. Springer, 2011.
- [36] Rajib Mall. *Real-Time Systems: Theory and Practice*. Pearson Education India, 2009.
- [37] Fan Liu, Ajit Narayanan, and Quan Bai. *Real-time systems*. 2000.
- [38] Zdzisaw Pawlak and Roman Sowinski. Rough set approach to multi-attribute decision analysis. *European Journal of Operational Research*, 72(3):443–459, 1994.
- [39] Panos M Pardalos, Yannis Siskos, and Constantin Zopounidis. *Advances in multicriteria analysis*. Springer, 1995.
- [40] Ralph E Steuer. *Multiple criteria optimization: theory, computation, and application*. Krieger Malabar, 1989.
- [41] Theo J Stewart. A critical survey on the status of multiple criteria decision making theory and practice. *Omega*, 20(5):569–586, 1992.
- [42] Jan Chomicki, Paolo Ciaccia, and Niccolo' Meneghetti. Skyline queries, front and back. *ACM SIGMOD Record*, 42(3):6–18, 2013.
- [43] S Borzsony, Donald Kossmann, and Konrad Stocker. The skyline operator. In *Data Engineering, 2001. Proceedings. 17th International Conference on*, pages 421–430. IEEE, 2001.
- [44] Hsiang-Tsung Kung, Fabrizio Luccio, and Franco P Preparata. On finding the maxima of a set of vectors. *J. acm*, 22(4):469–476, 1975.

- [45] Franco P Preparat and Michael Ian Shamos. Computational geometry: an introduction. 1985.
- [46] Yuan-Chi Chang, Lawrence Bergman, Vittorio Castelli, Chung-Sheng Li, Ming-Ling Lo, and John R. Smith. The onion technique: Indexing for linear optimization queries. *SIGMOD Rec.*, 29(2):391–402, May 2000. ISSN 0163-5808. doi: 10.1145/335191.335433. URL <http://doi.acm.org/10.1145/335191.335433>.
- [47] Vagelis Hristidis, Nick Koudas, and Yannis Papakonstantinou. Prefer: A system for the efficient execution of multi-parametric ranked queries. *SIGMOD Rec.*, 30(2):259–270, May 2001. ISSN 0163-5808. doi: 10.1145/376284.375690. URL <http://doi.acm.org/10.1145/376284.375690>.
- [48] Dong Xin, Chen Chen, and Jiawei Han. Towards robust indexing for ranked queries. In *Proceedings of the 32nd International Conference on Very Large Data Bases, VLDB '06*, pages 235–246. VLDB Endowment, 2006. URL <http://dl.acm.org/citation.cfm?id=1182635.1164149>.
- [49] Ihab F. Ilyas, George Beskales, and Mohamed A. Soliman. A survey of top-k query processing techniques in relational database systems. *ACM Comput. Surv.*, 40(4):11:1–11:58, October 2008. ISSN 0360-0300. doi: 10.1145/1391729.1391730. URL <http://doi.acm.org/10.1145/1391729.1391730>.
- [50] Flip Korn and S. Muthukrishnan. Influence sets based on reverse nearest neighbor queries. *SIGMOD Rec.*, 29(2):201–212, May 2000. ISSN 0163-5808. doi: 10.1145/335191.335415. URL <http://doi.acm.org/10.1145/335191.335415>.
- [51] Evangelos Dellis and Bernhard Seeger. Efficient computation of reverse skyline queries. In *Proceedings of the 33rd International Conference on Very Large Data Bases, VLDB '07*, pages 291–302. VLDB Endowment, 2007. ISBN 978-1-59593-649-3. URL <http://dl.acm.org/citation.cfm?id=1325851.1325887>.
- [52] Brian Babcock, Shivnath Babu, Mayur Datar, Rajeev Motwani, and Jennifer Widom. Models and issues in data stream systems. In *Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 1–16. ACM, 2002.
- [53] Katja Hose and Akrivi Vlachou. A survey of skyline processing in highly distributed environments. *The VLDB Journal—The International Journal on Very Large Data Bases*, 21(3):359–384, 2012.

- [54] Xuemin Lin, Yidong Yuan, Wei Wang, and Hongjun Lu. Stabbing the sky: Efficient skyline computation over sliding windows. In *Data Engineering, 2005. ICDE 2005. Proceedings. 21st International Conference on*, pages 502–513. IEEE, 2005.
- [55] Yufei Tao and Dimitris Papadias. Maintaining sliding window skylines on data streams. *Knowledge and Data Engineering, IEEE Transactions on*, 18(3):377–391, 2006.
- [56] Shengli Sun, Zhenghua Huang, Hao Zhong, Dongbo Dai, Hongbin Liu, and Jinjiu Li. Efficient monitoring of skyline queries over distributed data streams. *Knowledge and information systems*, 25(3):575–606, 2010.
- [57] Nikos Sarkas, Gautam Das, Nick Koudas, and Anthony KH Tung. Categorical skylines for streaming data. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 239–250. ACM, 2008.
- [58] Wenlin He, Cuiping Li, and Hong Chen. Maintaining the dominant representatives on data streams. In *Database and Expert Systems Applications*, pages 704–718. Springer, 2009.
- [59] Yuan Fang and Chee-Yong Chan. Efficient skyline maintenance for streaming data with partially-ordered domains. In *Database Systems for Advanced Applications*, pages 322–336. Springer, 2010.
- [60] Anna C Gilbert, Yannis Kotidis, S Muthukrishnan, and Martin Strauss. Surfing wavelets on streams: One-pass summaries for approximate aggregate queries. In *VLDB*, volume 1, pages 79–88, 2001.
- [61] Cheqing Jin, Ke Yi, Lei Chen, Jeffrey Xu Yu, and Xuemin Lin. Sliding-window top-k queries on uncertain streams. *Proc. VLDB Endow.*, 1(1):301–312, August 2008. ISSN 2150-8097. doi: 10.14778/1453856.1453892. URL <http://dx.doi.org/10.14778/1453856.1453892>.
- [62] Ming Hua, Jian Pei, Wenjie Zhang, and Xuemin Lin. Efficiently answering probabilistic threshold top-k queries on uncertain data. In *ICDE*, volume 8, pages 1403–1405, 2008.
- [63] Akrivi Vlachou, Christos Doulkeridis, Yannis Kotidis, and Kjetil Norvag. Monochromatic and bichromatic reverse top-k queries. *Knowledge and Data Engineering, IEEE Transactions on*, 23(8):1215–1229, 2011.
- [64] Akrivi Vlachou, Christos Doulkeridis, Kjetil Nørnvåg, and Yannis Kotidis. Branch-and-bound algorithm for reverse top-k queries. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, SIGMOD

- '13, pages 481–492, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-2037-5. doi: 10.1145/2463676.2465278. URL <http://doi.acm.org/10.1145/2463676.2465278>.
- [65] Dimitris Papadias, Panos Kalnis, Jun Zhang, and Yufei Tao. Efficient olap operations in spatial data warehouses. In ChristianS. Jensen, Markus Schneider, Bernhard Seeger, and VassilisJ. Tsotras, editors, *Advances in Spatial and Temporal Databases*, volume 2121 of *Lecture Notes in Computer Science*, pages 443–459. Springer Berlin Heidelberg, 2001. ISBN 978-3-540-42301-0. doi: 10.1007/3-540-47724-1\_23. URL [http://dx.doi.org/10.1007/3-540-47724-1\\_23](http://dx.doi.org/10.1007/3-540-47724-1_23).
- [66] Akrivi Vlachou, Christos Doulkeridis, and Kjetil Nørnvåg. Monitoring reverse top-k queries over mobile devices. In *Proceedings of the 10th ACM International Workshop on Data Engineering for Wireless and Mobile Access*, pages 17–24. ACM, 2011.
- [67] Vagelis Hristidis, Nick Koudas, and Yannis Papakonstantinou. Prefer: A system for the efficient execution of multi-parametric ranked queries. In *ACM SIGMOD Record*, volume 30, pages 259–270. ACM, 2001.
- [68] Dong Xin, Chen Chen, and Jiawei Han. Towards robust indexing for ranked queries. In *Proceedings of the 32nd international conference on Very large data bases*, pages 235–246. VLDB Endowment, 2006.