

**Συστήματα Διαχείρισης και Ανάλυσης Δεδομένων**  
**Διδάσκων: Ιωάννης Κωτίδης**

Εαρινό εξάμηνο 2025-2026

**Εργασία**

Ανάθεση: 24-04-2026

Παράδοση: **06-05-2026 Ώρα (23:55)**

*Οδηγίες*

- *Η εργασία είναι ατομική και υποχρεωτική.*
- *Η υποβολή της εργασίας πρέπει να γίνει στο eclass.*
- *Το παραδοτέο σας θα πρέπει να είναι ένα αρχείο PDF με όνομα AM.pdf (όπου AM είναι ο αριθμός μητρώου σας. π.χ. "3230001.pdf").*
- *Πιθανή αντιγραφή θα τιμωρείται με μηδενισμό όλων των εμπλεκομένων.*

### **Βάση Δεδομένων RETAILDB**

Στόχος της εργασίας είναι η πρακτική εφαρμογή των γνώσεων που αποκομίσατε από τις διαλέξεις του μαθήματος σχετικά με τη δημιουργία ευρετηρίων και την βελτιστοποίηση των επερωτήσεων SQL. Για τον σκοπό της εργασίας θα χρησιμοποιήσετε την βάση δεδομένων **RETAILDB** η οποία περιέχει δεδομένα πωλήσεων μιας εικονικής (μη υπαρκτής) πολυεθνικής εταιρείας που εμπορεύεται έναν σημαντικό αριθμό προϊόντων. Οι βασικές οντότητες της βάσης αφορούν σε στοιχεία πελατών, προμηθευτών, προϊόντων και παραγγελιών. Τα δεδομένα των πινάκων δεν είναι πραγματικά αλλά έχουν δημιουργηθεί τυχαία για το σκοπό της συγκεκριμένης εργασίας.

Αρχικά θα δημιουργήσετε την βάση δεδομένων και θα φορτώσετε τα δεδομένα στους πίνακες, ακολουθώντας τις παρακάτω οδηγίες. Στη συνέχεια θα απαντήσετε στα ζητούμενα της εργασίας.

#### **1. Οδηγίες για την δημιουργία της βάσης δεδομένων RETAILDB**

Για να δημιουργήσετε την βάση δεδομένων και να φορτώσετε τις εγγραφές ακολουθήστε **ΠΡΟΣΕΚΤΙΚΑ** τα παρακάτω βήματα:

**Βήμα 1:** Κατεβάστε το αρχείο **retaildata.zip** (μέγεθος αρχείου 234 MB) από τον σύνδεσμο:  
<http://pages.aueb.gr/users/mkap/retaildata.zip>

**Βήμα 2:** Αποσυμπιέστε το αρχείο **retaildata.zip** στον φάκελο **C:\retaildata** (μέγεθος φακέλου 813MB).

**Βήμα 3:** Από το περιβάλλον του Microsoft Sql Server Management Studio εκτελέστε το SQL script "**CreateRetailDB.sql**" που δημιουργεί το λογικό σχήμα της βάσης.

**Βήμα 4:** Εκτελέστε το SQL script "**LoadRetailData.sql**" το οποίο θα φορτώσει δεδομένα στους πίνακες της βάσης. Το συγκεκριμένο script περιέχει εντολές της μορφής:

```
BULK INSERT customers           ! Πίνακας στον οποίο θα φορτωθούν τα δεδομένα  
FROM 'C:\retaildata\customers.txt' ! Αρχείο το οποίο περιέχει τα δεδομένα.  
WITH (FIRSTROW =2, FIELDTERMINATOR= '|', ROWTERMINATOR = '\n');
```

**Παράμετροι:**

- FIRSTROW=2 : Η πρώτη γραμμή του αρχείου περιέχει τα ονόματα των πεδίων και αγνοείται.
- FIELDTERMINATOR = '|' : Ο χαρακτήρας '|' δηλώνει το τέλος κάθε πεδίου της εγγραφής.
- ROWTERMINATOR='\n' : Ο χαρακτήρας αλλαγής γραμμής δηλώνει το τέλος κάθε εγγραφής του αρχείου.

**ΠΡΟΣΟΧΗ:** Αν τοποθετήσετε τα δεδομένα σε φάκελο διαφορετικό από τον '**C:\retaildata**' θα πρέπει να τροποποιήσετε ανάλογα το path. Για παράδειγμα αν τοποθετήσετε τα δεδομένα στον φάκελο '**C:\DATA**' η παραπάνω εντολή πρέπει να αλλάξει ως εξής:

```
BULK INSERT customers  
FROM 'C:\DATA\customers.txt'  
WITH (FIRSTROW =2, FIELDTERMINATOR= '|', ROWTERMINATOR = '\n');
```

**Σημείωση:** Για την διαδικασία της μαζικής εισαγωγής των δεδομένων απαιτούνται περίπου 60 sec σε ένα υπολογιστή με δίσκο SSD. Το μέγεθος της βάσης είναι 2.1 GB (data & log files).

## 2. Περιγραφή των πινάκων της βάσης

Ακολουθεί η περιγραφή των πινάκων και των δεδομένων της βάσης.

<b>REGIONS: Πίνακας με τις γεωγραφικές περιοχές που δραστηριοποιείται η εταιρεία.</b>	
<b>Αριθμός εγγραφών=5</b>	
<b>regionkey</b>	Κωδικός περιοχής
<b>region</b>	Γεωγραφική περιοχή

<b>NATIONS: Πίνακας με τα κράτη που δραστηριοποιείται η εταιρεία.</b>	
<b>Αριθμός εγγραφών=25</b>	
<b>nationkey</b>	Κωδικός κράτους
<b>nation</b>	Κράτος

<b>CUSTOMERS: Πίνακας με τα στοιχεία των πελατών. Αριθμός εγγραφών=150.000</b>	
<b>custkey</b>	Κωδικός πελάτη
<b>cname</b>	Όνομα πελάτη
<b>cphone</b>	Τηλέφωνο πελάτη
<b>c_acctbal</b>	Υπόλοιπο λογαριασμού πελάτη
<b>market_segment</b>	Εμπορικός τομέας (π.χ. BUILDING, MACHINERY, FURNITURE κ.λπ.)
<b>nationkey</b>	Κωδικός κράτους πελάτη.
<b>c_comment</b>	Σχόλια υπό τη μορφή ελεύθερου κειμένου.

<b>SUPPLIERS: Πίνακας με τα στοιχεία των προμηθευτών. Αριθμός εγγραφών=10.000</b>	
<b>suppkey</b>	Κωδικός προμηθευτή
<b>sname</b>	Όνομα προμηθευτή
<b>nationkey</b>	Κωδικός κράτους προμηθευτή
<b>s_acctbal</b>	Υπόλοιπο λογαριασμού προμηθευτή
<b>s_comment</b>	Σχόλια υπό τη μορφή ελεύθερου κειμένου.

<b>PARTS: Πίνακας με τα προϊόντα που εμπορεύεται η εταιρεία. Αριθμός εγγραφών=200.000</b>	
<b>partkey</b>	Κωδικός προϊόντος
<b>ptype</b>	Τύπος προϊόντος.
<b>psize</b>	Μέγεθος προϊόντος.
<b>brand</b>	Μάρκα προϊόντος
<b>pname</b>	Ονομασία προϊόντος
<b>container</b>	Είδος συσκευασίας.
<b>manufacturer</b>	Κατασκευαστής προϊόντος.
<b>retailprice</b>	Τιμή προϊόντος.

<b>PARTSUPP: Πίνακας που συνδέει τα προϊόντα με τους προμηθευτές. Αριθμός εγγραφών=800.000</b>	
<b>partkey</b>	Κωδικός προϊόντος
<b>suppkey</b>	Κωδικός προμηθευτή.
<b>supplycost</b>	Κόστος προμήθειας του συγκεκριμένου προϊόντος από τον συγκεκριμένο προμηθευτή.
<b>availqty</b>	Διαθέσιμη ποσότητα.
<b>ps_comment</b>	Σχόλιο υπό την μορφή ελεύθερου κειμένου.

<b>ORDERS : Πίνακας με τις παραγγελίες των πελατών. Αριθμός εγγραφών=1.500.000</b>	
<b>orderkey</b>	Κωδικός παραγγελίας
<b>orderdate</b>	Ημερομηνία παραγγελίας
<b>custkey</b>	Κωδικός πελάτη
<b>orderpriority</b>	Προτεραιότητα παραγγελίας
<b>totalprice</b>	Συνολική αξία παραγγελίας
<b>o_comment</b>	Σχόλια υπό την μορφή ελεύθερου κειμένου

<b>LINEITEM: Πίνακας με τα προϊόντα των παραγγελιών. Αριθμός εγγραφών=4.423.659</b>	
<b>orderkey</b>	Κωδικός παραγγελίας.
<b>linenumber</b>	Γραμμή παραγγελίας. Μία γραμμή παραγγελίας περιέχει μια συγκεκριμένη ποσότητα ενός συγκεκριμένου προϊόντος (partkey) το οποίο έχει προμηθευτεί από συγκεκριμένο προμηθευτή (suppkey).
<b>discount</b>	Συντελεστής έκπτωσης γραμμής παραγγελίας.
<b>price</b>	Τιμή γραμμής παραγγελίας (price = quantity*suppkeycost).
<b>suppkey</b>	Κωδικός προμηθευτή.
<b>quantity</b>	Ποσότητα γραμμής παραγγελίας.
<b>returnflag</b>	Ένδειξη επιστροφής (R=το προϊόν της συγκεκριμένης γραμμής παραγγελίας επεστράφη).
<b>partkey</b>	Κωδικός προϊόντος.
<b>tax</b>	Συντελεστής φορολογίας γραμμής παραγγελίας.
<b>shipdate</b>	Προγραμματισμένη ημερομηνία αποστολής.
<b>receiptdate</b>	Ημερομηνία παραλαβής.
<b>commitdate</b>	Καταληκτική ημερομηνία παράδοσης: ημερομηνία μέχρι την οποία το προϊόν της γραμμής παραγγελίας πρέπει να έχει παραδοθεί στον πελάτη. Συνήθως η καταληκτική ημερομηνία παράδοσης ορίζεται σε κάποια σύμβαση ή συμφωνία με τον πελάτη.
<b>shipmode</b>	Τρόπος αποστολής.
<b>shipinstruct</b>	Οδηγίες αποστολής
<b>l_comment</b>	Σχόλια υπό την μορφή ελεύθερου κειμένου.

### 3. Ζητούμενα εργασίας

Ακολουθούν τα ζητούμενα της εργασίας. Για την απάντηση των ζητημάτων **δεν επιτρέπεται καμία απολύτως τροποποίηση του σχήματος** εκτός φυσικά από την δημιουργία των ζητούμενων ευρετηρίων. Επίσης **απαγορεύεται** η δημιουργία και η χρήση όψεων (views).

Σε κάθε ζήτημα δεν αρκεί μόνο να παραθέσετε τα ερωτήματα σε γλώσσα SQL ή/και τις εντολές δημιουργίας των ευρετηρίων που ζητούνται. Σε κάθε περίπτωση **πρέπει να τεκμηριώσετε τις απαντήσεις σας και να παραθέσετε στοιχεία που επιβεβαιώνουν τους ισχυρισμούς σας**. Για παράδειγμα:

- Σε περιπτώσεις που ζητείται να αποδείξετε ότι ένα ευρετήριο επιταχύνει ένα ερώτημα, εκτελέστε το ερωτήμα δίχως το ευρετήριο και εξετάστε το πλάνο εκτέλεσης. Αφού δημιουργήσετε το ευρετήριο εκτελέστε εκ νέου το ερωτήμα και επανεξετάστε το πλάνο εκτέλεσης. Συγκρίνοντας τα δύο πλάνα μπορείτε να καταλήξετε σε συμπεράσματα σχετικά με την καταλληλότητα του ευρετηρίου.
- Σε περιπτώσεις που πρέπει να συγκρίνετε ένα η περισσότερα ερωτήματα, εκτελέστε τα όλα μαζί σε δέσμη και εξετάστε τα πλάνα εκτέλεσης. Ο SQL server δείχνει το κόστος κάθε ερωτήματος ως ποσοστό επί του συνολικού κόστους εκτέλεσης της δέσμης.
- Ενεργοποιήστε τα στατιστικά στοιχεία I/O με την εντολή: **set statistics io on**. Με τον τρόπο αυτό μπορείτε να βλέπετε κάθε φορά που εκτελείτε ένα ερωτήμα πόσες σελίδες διαβάζονται από τον δίσκο ή/και από την μνήμη (buffer).
- Μπορείτε να ενεργοποιήσετε τα στατιστικά στοιχεία σχετικά με τον χρόνο εκτέλεσης του ερωτήματος με την εντολή **set statistics time on**.
- Κάθε φορά πριν την εκτέλεση ενός ερωτήματος, εκτελέστε τις παρακάτω εντολές που "καθαρίζουν" τους buffers που χρησιμοποιεί ο SQL server για την αποθήκευση των δεδομένων και των πλάνων εκτέλεσης:  
**checkpoint**  
**dbcc dropcleanbuffers**

Με τον τρόπο αυτό διασφαλίζετε ότι, το ερωτήμα που θα εκτελέσετε δεν θα χρησιμοποιήσει τυχόν σελίδες που υπάρχουν στην μνήμη από προηγούμενες εκτελέσεις του ιδίου ή/και άλλων ερωτημάτων. Σε αντίθετη περίπτωση μπορεί να οδηγηθείτε σε λάθος συμπεράσματα.

- Κάθε ζήτημα πρέπει να το αντιμετωπίσετε **ανεξάρτητα από τα υπόλοιπα** και να το υλοποιήσετε στο αρχικό στιγμιότυπο της βάσης. Για παράδειγμα αν θέλετε να εξετάσετε κατά πόσο ένα ευρετήριο κάνει πιο αποδοτικό ένα ερώτημα, βεβαιωθείτε ότι έχετε διαγράψει (drop index) τα ευρετήρια που έχετε δημιουργήσει για την βελτιστοποίηση άλλων ερωτημάτων.
- Να λαμβάνετε υπόψη ότι Ο SQL Server δημιουργεί εξορισμού ευρετήριο συστάδων (clustered index) στο πρωτεύων κλειδί κάθε πίνακα της βάσης

### Ζήτημα 1 [ 25 μονάδες]

Το παρακάτω ερωτήμα εμφανίζει τον κωδικό παραγγελίας και την ημερομηνία παραγγελίας για παραγγελίες που περιλαμβάνουν τουλάχιστον ένα προϊόν μάρκας «Origin». Κάθε παραγγελία εμφανίζεται μία μόνο φορά.

```
SELECT DISTINCT orders.orderkey, orders.orderdate
  FROM orders JOIN lineitem ON orders.orderkey = lineitem.orderkey
        JOIN parts ON parts.partkey = lineitem.partkey
 WHERE brand = 'Origin';
```

Ζητείται να δημιουργήσετε ένα ή περισσότερα ευρετήρια που επιταχύνουν την εκτέλεση του ερωτήματος. Να παραθέσετε τις εντολές δημιουργίας των ευρετηρίων καθώς και στοιχεία που τεκμηριώνουν ότι το ευρετήριο (ή τα ευρετήρια) που δημιουργήσατε βελτιώνει την απόδοση του ερωτήματος.

### Ζήτημα 2 [ 25 μονάδες]

Το παρακάτω ερωτήμα εμφανίζει τον κωδικό προϊόντος, την ονομασία του προϊόντος, το όνομα του προμηθευτή και το κόστος προμήθειας για προϊόντα που προμηθεύονται από Ευρωπαίους προμηθευτές και για τα οποία η διαθέσιμη ποσότητα είναι μεγαλύτερη από 5000 τεμάχια.

```
SELECT parts.partkey, pname, sname, supplycost
 FROM parts
  JOIN partsupp ON parts.partkey = partsupp.partkey
  JOIN suppliers ON suppliers.supplykey = partsupp.supplykey
  JOIN nations ON nations.nationkey = suppliers.nationkey
  JOIN regions ON regions.regionkey = nations.regionkey
 WHERE region = 'EUROPE'
    and partsupp.availqty > 5000;
```

Ζητείται να δημιουργήσετε ένα ή περισσότερα ευρετήρια που επιταχύνουν την εκτέλεση του ερωτήματος. Να παραθέσετε τις εντολές δημιουργίας των ευρετηρίων καθώς και στοιχεία που τεκμηριώνουν ότι το ευρετήριο (ή τα ευρετήρια) που δημιουργήσατε βελτιώνει την απόδοση του ερωτήματος.

### Ζήτημα 3 [ 25 μονάδες]

Ζητήθηκε από έναν προγραμματιστή να γράψει ένα ερωτήμα που εμφανίζει τη συνολική αξία των επείγουσών παραγγελιών (1-URGENT) ανά ημέρα για τον Ιανουάριο του 1994. Ο προγραμματιστής έγραψε το ακόλουθο ερωτήμα:

```
SELECT orderdate, sum(totalprice)
  FROM orders
 WHERE (MONTH(orderdate) = 1 AND YEAR (orderdate) = 1994) AND
        (orderpriority='1-URGENT')
 GROUP BY orderdate
 ORDER BY orderdate;
```

Ο υπεύθυνος του τμήματος πωλήσεων δεν είναι ικανοποιημένος από την απόδοση του ερωτήματος. Καλείστε να προτείνετε τρόπους βελτιστοποίησης του ερωτήματος. Μπορείτε να πειραματιστείτε με τη δημιουργία ευρετηρίων, σε συνδυασμό με την συγγραφή εναλλακτικών ερωτήσεων που παράγουν το επιθυμητό αποτέλεσμα. Να παραθέσετε στοιχεία που τεκμηριώνουν την απάντησή σας.

#### Ζήτημα 4 [25 μονάδες]

Τα παρακάτω δύο ερωτήματα εμφανίζουν, αντίστοιχα:

1. όλα τα στοιχεία των προμηθευτών που προμηθεύουν προϊόντα τύπου «LARGE PLATED COPPER»,
2. όλα τα στοιχεία των προϊόντων που προμηθεύει ο προμηθευτής με όνομα «Runolfsson LLC».

```
1.  
SELECT suppliers.*  
FROM parts JOIN partsupp ON parts.partkey = partsupp.partkey  
JOIN suppliers ON suppliers.supkey = partsupp.supkey  
WHERE ptype = 'LARGE PLATED COPPER';
```

```
2.  
SELECT parts.*  
FROM suppliers JOIN partsupp ON partsupp.supkey = suppliers.supkey  
JOIN parts ON parts.partkey = partsupp.partkey  
WHERE sname = 'Runolfsson LLC';
```

Να δημιουργήσετε κατάλληλα ευρετήρια που επιταχύνουν την εκτέλεση των παραπάνω ερωτημάτων. Να παραθέσετε τις εντολές δημιουργίας των ευρετηρίων, καθώς επίσης και στοιχεία που αποδεικνύουν ότι τα ευρετήρια που δημιουργήσατε επιταχύνουν την εκτέλεση των ερωτημάτων.

**ΠΡΟΣΟΧΗ:** Οι απαντήσεις σας στα παραπάνω ζητήματα πρέπει να περιλαμβάνουν τόσο τα σχετικά πλάνα εκτέλεσης και στατιστικά I/O όσο και σύντομη αλλά σαφή αιτιολόγηση της επιλογής των ευρετηρίων. **Δεν αρκεί η απλή παράθεση εικόνων και μετρήσεων χωρίς σχολιασμό.**

ΔΙΑΓΡΑΜΜΑΤΙΚΗ ΑΝΑΠΑΡΑΣΤΑΣΗ ΛΟΓΙΚΟΥ ΣΧΗΜΑΤΟΣ ΤΗΣ ΒΑΣΗΣ ΔΕΔΟΜΕΝΩΝ RETAILDB

