

MiniDb – Indexing and Primary key

Επιμέλεια: Ιωακείμ Ελ-Χαττάμπ Μπριστογιάννης, Παρασκευή Παλιούρα, Μυρτώ – Μαρία Σπάθα, ΑΜ: Π19048, Π19129, Π19156

Όνομα Ομάδας: *Dream Team*

Στην παρούσα εργασία υλοποιήθηκαν οι παρακάτω λειτουργίες της miniDB:

1. Υποστήριξη πολλαπλών πρωτευόντων κλειδιών σε έναν πίνακα σχέσεων.
2. Υποστήριξη ευρετηρίου μέσω B+ δέντρου σε μία ή περισσότερες στήλες με διπλότυπα (ενέργειες δημιουργίας ευρετηρίου και επιλογής ευρετηριασμένων χαρακτηριστικών).
3. Υποστήριξη ευρετηρίου μέσω B+ δέντρου σε πολλαπλές στήλες ενός πίνακα (ενέργειες δημιουργίας ευρετηρίου και επιλογής ευρετηριασμένων χαρακτηριστικών).
4. Δημιουργία ευρετηρίου κατακερματισμού.

Παρακάτω, παρουσιάζονται αναλυτικότερα οι λειτουργίες που αναπτύχθηκαν για κάθε ενέργεια.

Ενέργεια δημιουργίας πίνακα με πολλαπλά πρωτεύοντα κλειδιά:

Για την υποστήριξη αυτής της ενέργειας έγιναν τροποποιήσεις στην επεξεργασία του query δημιουργίας πίνακα στη μέθοδο `create_query_plan` του αρχείου `mdb.py`, στη μέθοδο `__init__` της κλάσης `table` και στη μέθοδο `show` της κλάσης `table`.

Πιο συγκεκριμένα:

- Μέθοδος `create_query_plan`: Αρχικά, γίνεται αναγνώριση της μορφής του query και στη συνέχεια αποθηκεύεται το κάθε keyword σε ένα dictionary.
- Μέθοδος `__init__`: Το index του primary key είναι μια λίστα που αποθηκεύει το κάθε column του primary key.
- Μέθοδος `show`: Χρήση του καινούργιου primary key index για την εμφάνιση του κάθε column που ανήκει στο primary key.

Ενέργεια Δημιουργίας Ευρετηρίου:

Για την υποστήριξη αυτής της ενέργειας έγιναν τροποποιήσεις στην επεξεργασία του query δημιουργίας ευρετηρίου ('create index `index_name` on `table_name(columns_names)` using `btree/hash`') στη μέθοδο `create_query_plan` του αρχείου `mdb.py` και των μεθόδων δημιουργίας ευρετηρίου `create_index`, `construct_index` της κλάσης `database`.

Επιπλέον, για τη δημιουργία ευρετηρίου μέσω B+ δέντρου τροποποιήθηκαν οι μέθοδοι `__init__`, `insert` της κλάσης `btree`, ενώ για τη δημιουργία ευρετηρίου κατακερματισμού έχει δημιουργηθεί το αρχείο `hash.py`.

Πιο συγκεκριμένα:

- Μέθοδος `create_query_plan`: Έχει προστεθεί το action 'create_index', στο οποίο αναγνωρίζεται η μορφή του query και οι στήλες που βρίσκονται μετά το 'on' εισάγονται σε ένα dictionary.
- Μέθοδος `create_index`: Παίρνει ως είσοδο το όνομα του ευρετηρίου, το όνομα του πίνακα, τα ονόματα των στηλών πάνω στα οποία θα δημιουργηθεί το ευρετήριο και τον τύπο του ευρετηρίου.
Σε περίπτωση που ο τύπος του ευρετηρίου είναι `btree`, εισάγονται τα στοιχεία του ευρετηρίου (όνομα ευρετηρίου, όνομα πίνακα, όνομα ευρετηριασμένων στηλών) στον πίνακα `meta_indexes` σε μορφή string και γίνεται έλεγχος αν οι στήλες περιέχουν διπλότυπα. Σε περίπτωση που υπάρχουν διπλότυπα, στον αρχικό πίνακα προστίθεται μια κρυφή στήλη για κάθε στήλη που ευρετηριάζεται (η κρυφή στήλη περιέχει ακεραίους και είναι `auto incrementing`). Στη συνέχεια καλείται η μέθοδος `construct_index` και αποθηκεύεται η βάση.

Σημείωση: Στήλες με διπλότυπα χαρακτηρίζονται όλες οι στήλες που δεν είναι `primary key`.

Σε περίπτωση που ο τύπος του ευρετηρίου είναι `hash`, εισάγονται τα στοιχεία του ευρετηρίου (όνομα ευρετηρίου, όνομα πίνακα, όνομα ευρετηριασμένων στηλών) στον πίνακα `meta_indexes`, καλείται η μέθοδος `construct_index` και αποθηκεύεται η βάση.

- Μέθοδος `construct_index`: Παίρνει ως είσοδο το όνομα του ευρετηρίου, το όνομα του πίνακα, τα ονόματα των στηλών πάνω στα οποία θα δημιουργηθεί το ευρετήριο, τον τύπο του ευρετηρίου και μια λογική μεταβλητή για τον καθορισμό διπλοτύπων (`True` σε περίπτωση που υπάρχουν διπλότυπα).
Σε περίπτωση που ο τύπος του ευρετηρίου είναι `btree`, φτιάχνεται ένα αντικείμενο της κλάσης `btree` με `branching factor` 3. Στη συνέχεια, δημιουργούμε μια λίστα που περιέχει τις εγγραφές για κάθε στήλη που ευρετηριάζεται, καθώς και την τιμή της κρυφής στήλης για κάθε εγγραφή σε περίπτωση που υπάρχουν διπλότυπα. Τέλος, για κάθε στοιχείο αυτής της λίστας καλείται η μέθοδος `insert` της κλάσης `btree`.
Σε περίπτωση που ο τύπος του ευρετηρίου είναι `hash`, φτιάχνεται ένα αντικείμενο της κλάσης `hash_table` και για κάθε εγγραφή στο `primary key` του πίνακα, εισάγουμε την τιμή του και το δείκτη του στο αντικείμενο του `hash table`. Τέλος, αποθηκεύεται το ευρετήριο.

Σημείωση: Το ευρετήριο κατακερματισμού δημιουργείται μόνο στο μονοδιάστατο πρωτεύον κλειδί του πίνακα.

- Αρχείο `hash.py`: Αποτελείται από μια κλάση (τη `hash_table`), τα αντικείμενα της οποίας είναι οι πίνακες κατακερματισμού. Περιέχει τις παρακάτω μεθόδους:
 - `__init__`: Ορίζεται το μέγεθος του πίνακα που αντιστοιχεί στον μέγιστο αριθμό εγγραφών που αποθηκεύονται σε αυτόν (εδώ, έχουμε το πολύ 23 εγγραφές). Αρχικοποιείται ο πίνακας με τιμές `None` για κάθε εγγραφή. Ορίζεται επίσης ο αριθμός των γεμάτων

κάδων σε 0, όπως και το κατώφλι για το resize (0.75). Καλείται η συνάρτηση κατακερματισμού (crc32_hash) και το crc32_table. Τέλος, ορίζεται ένας τυχαίος αριθμός a, ο οποίος χρησιμοποιείται όταν το κλειδί είναι string.

- `__repr__`: Επιστρέφει τα ζευγάρια σε μορφή {key1: value1, key2: value2, ...}.
- `Load_Factor`: Υπολογίζει πόσο γεμάτος είναι ο πίνακας, ώστε να γίνει resize όταν φτάσουμε το κατώφλι που έχουμε ορίσει.
- `Encode`: Κωδικοποιεί το κλειδί που στέλνεται στην CRC32_hash, το οποίο υπολογίζεται από τη μαθηματική συνάρτηση polynomial rolling hash function (ορίζεται μέσα στον κώδικα της Encode).
- `Crc32_table`: Επιστρέφει έναν πίνακα από τιμές οι οποίες χρησιμοποιούν τον αλγόριθμο crc32 ως hash μέθοδο και θα χρησιμοποιηθούν στη crc32_hash.
- `Crc32_hash`: Πρόκειται για τη συνάρτηση κατακερματισμού. Κατακερματίζει το κλειδί που παράγεται από τη μέθοδο encode(). Στην ουσία, χρησιμοποιεί First Order Least Significant Bit (LSB), ορίζει το αρχικό CRC σε FFFFFFFF και, τέλος, συμπληρώνει το τελικό CRC.
- `Resize`: Αυξάνει το μέγεθος του πίνακα, μόλις φτάσει στο κατώφλι. Το μέγεθος του πίνακα αλλάζει στον μικρότερο πρώτο αριθμό μεγαλύτερο από το διπλάσιο του τωρινού μεγέθους. Για να βρούμε τον πρώτο αριθμό, εφαρμόζουμε το «Κόσκινο του Ερατοσθένη». Αφού αλλάξουμε το μέγεθος του πίνακα, όλες οι καταχωρίσεις κατακερματίζονται εκ νέου.
- `Insert`: Παίρνει σαν είσοδο το key, value και χειρίζεται την εισαγωγή στοιχείων στον πίνακα.

Πηγή για τον αλγόριθμο crc32: https://www.section.io/engineering-education/hashtables-implementation-using-crc32-algorithm/?fbclid=IwAR3EPcre8bxixTYCq3si6lFHJgzclq5gJDWT_2TCzFVwHhADtQyC07ADYik

Ενέργεια επιλογής ευρετηριασμένων χαρακτηριστικών:

Η επιλογή ευρετηριασμένων χαρακτηριστικών υποστηρίζεται για ευρετήριο μέσω B+ δέντρου.

Για την υποστήριξη αυτής της ενέργειας έγιναν τροποποιήσεις στις μεθόδους select της κλάσης database, select_with_btree της κλάσης table, split_condition του αρχείου misc.py, has_index της κλάσης database και find της κλάσης btree.

Πιο συγκεκριμένα:

- Μέθοδος select: Έχει προστεθεί ο έλεγχος των στηλών του πίνακα που γίνεται η επιλογή για το αν είναι ευρετηριασμένες (μέσω της μεθόδου has_index). Σε περίπτωση που είναι ευρετηριασμένες, εκτελείται η μέθοδος select_with_btree.

- Μέθοδος `select_with_btree`: Έχει προστεθεί η διάκριση του query επιλογής στις παρακάτω περιπτώσεις:
 - Σε περίπτωση που η επιλογή γίνεται σε δύο στήλες και η πρώτη στήλη κάνει ερώτηση ταυτότητας, εκτελείται η μέθοδος `find`.
 - Σε περίπτωση που η επιλογή γίνεται σε μία στήλη, εκτελείται η μέθοδος `find`.
 - Σε οποιαδήποτε άλλη περίπτωση, εκτελείται η αναζήτηση χωρίς ευρετήριο.
- Μέθοδος `split_condition`: Χωρίζει τη συνθήκη της επιλογής (πρόταση μετά το `where`) στα συνθετικά της (όνομα στήλης, τελεστής σύγκρισης, τιμή χαρακτηριστικού). Έχει προστεθεί υποστήριξη για επιλογή πολλών στηλών.
- Μέθοδος `has_index`: Ελέγχεται αν υπάρχει ευρετήριο για τις στήλες που παίρνει ως είσοδο και επιστρέφει `True` και το όνομά του σε περίπτωση που υπάρχει. Ο έλεγχος γίνεται εκτελώντας τη μέθοδο `select` για τον πίνακα `meta_indexes` αναζητώντας το όνομα του πίνακα και τα ονόματα των στηλών που θέλουμε να ευρετηριάζονται.
- Μέθοδος `find`: Κάνει αναζήτηση στο ευρετήριο και επιστρέφει τα αποτελέσματα. Διακρίνονται οι παρακάτω περιπτώσεις:
 - Σε περίπτωση που το ευρετήριο δεν έχει δημιουργηθεί πάνω σε στήλες με διπλότυπα (δηλαδή είναι στο `primary key`), γίνεται αναζήτηση της ιδανικής θέσης της επιθυμητής τιμής μέσα στο δέντρο. Στη συνέχεια, ανάλογα με τον τελεστή σύγκρισης που παρέχεται από το `query`, επιστρέφονται τα κατάλληλα αποτελέσματα.
 - Σε περίπτωση που το ευρετήριο έχει δημιουργηθεί πάνω σε στήλες με διπλότυπα, στην τιμή αναζήτησης προστίθεται μια τιμή για την κρυφή στήλη (0 ή 20000000, ανάλογα με τον τελεστή σύγκρισης) και γίνεται αναζήτηση της ιδανικής θέσης της επιθυμητής τιμής μέσα στο δέντρο. Στη συνέχεια, ανάλογα με τον τελεστή σύγκρισης που παρέχεται από το `query`, επιστρέφονται τα κατάλληλα αποτελέσματα, αφού πρώτα έχει γίνει ο ίδιος έλεγχος για τα δεξιά ή αριστερά 'αδέρφια' του κόμβου που βρέθηκε (για να προσθέσουμε στα αποτελέσματα και αυτά).

Ακολουθούν Screenshots από την εκτέλεση της εφαρμογής:

PyCharm Professional Edition

File Edit View Navigate Code Refactor Run Tools Git Window Help

miniDB - database.py

Project: miniDB

Structure:

- miniDB
- __init__.py
- dbdata
- instructor_db
- instructor
- advisor.py

Code:

```
def _construct_index(table_name, index_name, index_type, selected_columns_names):
    # For each record in the primary key of the table, insert its value and index to the hash
    if selected_columns_names[0] == self.tables[table_name].pk[0]:
        for idx, key in enumerate(self.tables[table_name].columns[selected_columns_names[1]]):
            h.insert(key, idx)
    # save the hash
    Database._construct_index = h if index_type == 'hash'
```

Terminal:

```
(miniDB) create index inst_multi on instructor(dept_name, salary) using btree
btree index created
(miniDB) select * from instructor
id (str) apno name (str) dept_name (str) salary (int)
-----
10101 sriniwasan comp. sci. 45000
12123 wu finance 90000
13121 mearns music 40000
20222 elnsteth physics 90000
22345 el nate history 80000
23456 gold physics 87000
45565 kate comp. sci. 75000
50505 califfari history 62000
76543 slugh finance 80000
76764 erick biology 72000
83521 brawet comp. sci. 92000
85505 kin elec. eng. 80000

(miniDB) select * from instructor where dept_name = 'finance' and salary = 70000
Traceback (most recent call last):
```

PyCharm Professional Edition

File Edit View Navigate Code Refactor Run Tools Git Window Help

miniDB - btree.py

Project: miniDB

Structure:

- __init__.py
- dbdata
- instructor_db
- instructor
- advisor.py
- btree.py

Code:

```
target_node = self.nodes[target_node.right.stabbing]
for data in enumerate(target_node.values):
    results.extend(target_node.stra)
return results

elif operator[0] == 'lt':
    if not self.is_duplicate(self.operator[0], self.operator[1]):
```

Terminal:

```
(miniDB) create index inst_id on instructor(id) using btree
btree index created
(miniDB) select * from instructor
id (str) apno name (str) dept_name (str) salary (int)
-----
10101 sriniwasan comp. sci. 45000
12123 wu finance 90000
13121 mearns music 40000
20222 elnsteth physics 90000
22345 el nate history 80000
23456 gold physics 87000
45565 kate comp. sci. 75000
50505 califfari history 62000
76543 slugh finance 80000
76764 erick biology 72000
83521 brawet comp. sci. 92000
85505 kin elec. eng. 80000

(miniDB) select * from instructor where id = 10101
Found
id (str) apno name (str) dept_name (str) salary (int)
-----
10101 sriniwasan comp. sci. 45000
```

PyCharm Professional Edition

File Edit View Navigate Code Refactor Run Tools Git Window Help

miniDB - btree.py

Project: miniDB

Structure:

- __init__.py
- dbdata
- instructor_db
- instructor
- advisor.py
- btree.py

Code:

```
def _select(self, from_instructor where id <= 76764):
    id (str) apno name (str) dept_name (str) salary (int)
    -----
    76764 erick biology 72000
    76543 slugh finance 80000
    50505 califfari history 62000
    45565 kate comp. sci. 75000
    23456 gold physics 87000
    22345 el nate history 80000
    20222 elnsteth physics 90000
    13121 mearns music 40000
    12123 wu finance 90000
    10101 sriniwasan comp. sci. 45000

def _select(self, from_instructor where id <= 76764):
    id (str) apno name (str) dept_name (str) salary (int)
    -----
    76764 erick biology 72000
    83521 brawet comp. sci. 92000
    85505 kin elec. eng. 80000
    83521 brawet comp. sci. 92000
    85505 kin elec. eng. 80000
```





