



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет
имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ _____ «Информатика и системы управления»

КАФЕДРА _____ «Программное обеспечение ЭВМ и информационные технологии»

РАСЧЕТНО-ПОЯСНИТЕЛЬНАЯ ЗАПИСКА

К К ВЫПУСКНОЙ КВАЛИФИКАЦИОННОЙ РАБОТЕ

НА ТЕМУ:

«Метод построения поисковых индексов в реляционной
базе данных на основе глубоких нейронных сетей»

| | | | |
|-----------------|----------------------------|--------------------------|--|
| Студент: | <u>ИУ7-83Б</u> (группа) | _____ (подпись, дата) | <u>М. Д. Маслова</u> (И. О. Фамилия) |
| Руководитель: | | _____ (подпись, дата) | <u>А. А. Оленев</u> (И. О. Фамилия) |
| Нормоконтролер: | | _____ (подпись, дата) | <u>Д. Ю. Мальцева</u> (И. О. Фамилия) |

2023 г.

1 Технологическая часть

1.1 Выбор средств программной реализации

Для реализации метода построения индекса в реляционной базе данных на основе глубоких нейронных сетей в качестве языка программирования выбран Python 3.10 [**python**], так как предоставляет широкий выбор библиотек для глубокого обучения и визуализации его результатов. В качестве библиотеки глубокого обучения выбран TensorFlow 2.11.0 [**tf**] и работающий поверх нее высокоуровневый программный интерфейс Keras 2.11.0 [**keras**]. Для работы с массивами данных выбрана библиотека numpy [**numpy**].

В качестве реляционной системы управления базами данных выбрана SQLite [**sqlite**], предоставляющая программный интерфейс виртуальных таблиц, позволяющих релизовать пользовательский поисковый индекс. Виртуальные таблицы являются одним из видов расширений SQLite, программный интерфейс которых предоставляется на языке C [**c**], который и выбран в качестве языка программирования для взаимодействия с реляционной базой данных.

Для обеспечения взаимодействия между компонентом работы с базой данных и компонентом, непосредственно реализующий индекс, используются библиотеки языка C `Python.h` для работы с объектами языка Python и `numpy/arrayobject.h`, предоставляющая программный интерфейс для работы с numpy-массивами, которые являются основным типом данных, через который происходит взаимодействие модулей.

1.2 Реализация программного обеспечения

1.2.1 Форматы входных и выходных данных

Основой для построения индекса в качестве входных выступают данные из таблицы реляционной базы данных SQLite. Требованием к таблице является наличие атрибута целочисленного типа (INTEGER) с уникальными значениями (UNIQUE).

Для создания индекса в аргументах соответствующего запроса должен присутствовать идентификатор столбца, удовлетворяющего требованию описанному выше, и строковое имя модели индекса, на основе которой он строится (FCNN2 — для модели с двумя скрытыми слоями, FCNN3 — с тремя).

В качестве входных данных компонента, реализующего индекса, является набор значений столбца, идентификатор которого был указан в аргументах запроса на создание индекса, и идентификаторы строк ROWID индексируемой таблицы.

Выходным данных является указатель на объект, представляющий индекс на основе глубокой нейронной сети, обученный на предоставленных входных данных.

Формат входных и выходных данных операций с индексом (поиска и вставки) описан в следующем пункте.

1.2.2 Поддерживаемые виды запросов

Разработанное программное обеспечения предоставляет возможность работы только с запросами фильтрации, представленными оператором WHERE следующими со следующими условиями:

- `column operator value`,
где `column` — имя проиндексированного столбца,
`operator` — одна из операций сравнения: `=`, `<`, `>`, `<=`, `>=`,
`value` — некоторое целочисленное значение.
- `column BETWEEN value1 AND value2`,
где `column` — имя проиндексированного столбца,
`value1`, `value2` — целочисленные значения, представляющие нижнюю и верхнюю границы диапазона.

Выходным значением из модуля на языке Python, реализующего индекс, по данным запросам является `numpy`-массив с соответствующими запросу значениями ROWID, а результатом работы программного обеспечения — набор записей таблицы с ROWID из представленного массива.

Для вставки поддерживается стандартный запрос `INSERT`, результатом которого является индекс с переобученной на новых данных моделью. Запросы удаления и изменения не поддерживаются, так как являются вторичными для оценки работы метода.

1.2.3 Программный интерфейс виртуальных таблиц

Виртуальные таблицы SQLite [`vtable`] — это объект базы данных, представляющий с точки зрения инструкций SQL обычную таблицу или представле-

ние, но обрабатывающий запросы посредством вызова функций программного интерфейса, которые реализуются пользователем.

Виртуальные таблицы являются расширением SQLite, регистрация которых происходит с использованием макросов и функции инициализации, представленных на линстинге 1.1.

Листинг 1.1 — Инициализация расширения

```
1  #include <sqlite3ext.h>
2
3  SQLITE_EXTENSION_INIT1
4
5  int sqlite3_extension_init(sqlite3 *db,
6                           char **pzErrMsg,
7                           const sqlite3_api_routines *pApi)
8  {
9      int rc = SQLITE_OK;
10     SQLITE_EXTENSION_INIT2(pApi);
11
12     /* инициализация расширения */
13
14     return rc;
15 }
```

При инициализации расширения виртуальной таблицы должен быть зарегистрирован модуль, описывающийся структурой `sqlite3_module`, с помощью функции регистрации `sqlite3_create_module`, подробное описание которых представлено на листинге 1.2

Листинг 1.2 — Структура и функции для регистрации модуля виртуальной таблицы

```
1  struct sqlite3_module {
2      int iVersion;
3      int (*xCreate)(sqlite3*, void *pAux,
4                   int argc, char *const*argv,
5                   sqlite3_vtab **ppVTab,
6                   char **pzErr);
7      int (*xConnect)(sqlite3*, void *pAux,
8                    int argc, char *const*argv,
9                    sqlite3_vtab **ppVTab,
10                   char **pzErr);
11     int (*xBestIndex)(sqlite3_vtab *pVTab, sqlite3_index_info*);
12     int (*xDisconnect)(sqlite3_vtab *pVTab);
```

Продолжение листинга 1.2

```
13     int (*xDestroy)(sqlite3_vtab *pVTab);
14     int (*xOpen)(sqlite3_vtab *pVTab,
15                 sqlite3_vtab_cursor **ppCursor);
16     int (*xClose)(sqlite3_vtab_cursor*);
17     int (*xFilter)(sqlite3_vtab_cursor*,
18                   int idxNum, const char *idxStr,
19                   int argc, sqlite3_value **argv);
20     int (*xNext)(sqlite3_vtab_cursor*);
21     int (*xEof)(sqlite3_vtab_cursor*);
22     int (*xColumn)(sqlite3_vtab_cursor*, sqlite3_context*, int);
23     int (*xRowid)(sqlite3_vtab_cursor*, sqlite_int64 *pRowid);
24     int (*xUpdate)(sqlite3_vtab *, int,
25                   sqlite3_value **, sqlite_int64 *);
26     /* представлены указатели на реализующиеся методы */
27     /* при инициализации структуры, */
28     /* остальные поля принимают значение 0 */
29 };
30
31 int sqlite3_create_module(
32     sqlite3 *db,          /* соединение для регистрации модуля */
33     const char *zName,    /* имя модуля */
34     const sqlite3_module *, /* ссылка на структуру модуля */
35     void *,              /* данные для xCreate/xConnect */
36 );
```

Методы, сигнатуры которых представлены на листинге 1.2, можно разделить на две группы.

Методы для взаимодействие с таблицей, как с некоторым объектом, к которым относятся:

- `xCreate` — создание виртуальной таблицы, при выполнении соответствующего запроса, представленного на листинге 1.3;
- `xConnect` — подключение к виртуальной таблице, вызывющийся при выполнении любого запроса к таблице, который является первым при повторном подключении к базе данных;
- `xDestroy` — удаление виртуальной таблицы при выполнении запроса, представленного на листинге 1.4;
- `xDisconnect` — удаление подключения к виртуальной таблице.

Листинг 1.3 — Запрос на создание виртуальной таблицы

```
1 CREATE VIRTUAL TABLE <имя_таблицы> USING <имя_модуля>(arg1, ...);
```

Листинг 1.4 — Запрос на удаление виртуальной таблицы

```
1 DROP TABLE <имя_таблицы>;
```

Данные методы работают со структурой `sqlite3_vtab`, представленной на листинге `??`. Для реализации нужных функциональностей указатель на данную структуру включается в пользовательскую, с которой уже работают представленные методы посредством преобразования типов. Это дает возможность передавать между методами виртуальной таблицы нужные данные.

Методы прохода по записям таблицы, использующие для этого структуру курсора `sqlite3_vtab_cursor`, представленную на листинге `??`, над которой также реализуют обертку для хранения необходимых для обработки переменных.

Данная группа представлена методом обработки вставки, удаления и изменения записи (последний) и методами для прохода по записям таблицы при поиске:

- `xOpen` — создание и инициализации структуры курсора;
- `xBestIndex` — получение параметров фильтрации и выбор лучшего индекса для обработки запроса;
- `xFilter` — получение соответствующих параметрам фильтрации записей, установка курсора на первую из них;
- `xEof` — проверка окончания списка выбранных записей;
- `xNext` — переход к следующей записи;
- `xColumn` — обработка столбца записи;
- `xClose` — удаление структуры курсора;
- `xUpdate` — реализация запросов вставки, удаления и изменения.

Для реализации метода построения индекса используются оберточные структуры для виртуальной таблицы и курсора, представленные на листинге `??`.

На листинге `??` представлена реализация метода создания индекса, реализованного в качестве `xCreate`. Код инициализации и запуска обучения индекса в Python через программный интерфейс приведен на листинге `??`.

Обработка параметров фильтрации приведена на листинге ???. Получение массива подходящих строк и проход для их вывода приведены на листингах ??, ?? соответственно.

1.2.4 Реализация индекса

...

1.3 Сборка программного обеспечения

При сборке расширения SQLite компиляция и линковка происходит с флагами, обычно используемыми при сборке динамических библиотек: флаг `-fPIC` при компиляции в объектные файлы для создания позиционно-независимого кода и флаг `-shared` для получения файла динамической библиотеки. Дополнительными флагами при линковке являются флаги подключения библиотек `-lsqlite3` и `-lpthread`. Также при компиляции требуется указание путей к заголовочным файлам `Python.h` и `numpy/arrayobject.h`. Их автоматическое получение, а также ключевые моменты сборки приведены на листинге ??.

Для работы компонента индекса, реализованного на Python, требуется установка зависимостей из уже сформированного файла `requirements.txt` путем выполнения команды, представленной на листинге ??. Также для штатной работы программного обеспечения требуется прописать путь к модулям, реализованным на языке Python, что приведено на том же листинге.

1.4 Взаимодействие с программным обеспечением

Взаимодействие с программным обеспечением происходит через командную строку `sqlite3`. Пример работы представлен на листинге ??.

1.5 Результаты тестирования

После разработки программного обеспечения было проведено автоматическое тестирование модуля индекса на основе глубоких нейронных сетей. Были использованы следующие классы эквивалентности:

- поиск существующего единичного ключа;
- поиск наименьшего и наибольшего в наборе ключей;
- поиск несуществующего ключа;

- поиск ключей по условию $<$ с существующим ключом в качестве границы;
- поиск ключей по условию $<$ с несуществующим ключом в качестве границы;
- два предыдущих класса по условиям $>$, $<=$, $>=$;
- поиск по диапазону с двумя границами;
- поиск по диапазону с двумя границами, в который попадают все ключи;
- поиск по диапазону с двумя границами, в который не попадает ни один ключ.

Результаты тестирования представлены на листинге 1.5

Листинг 1.5 — Результаты автоматического тестирования

```

1  ===== short test summary info =====
2  PASSED tests.py::TestLindex::test_train
3  PASSED tests.py::TestLindex::test_middle
4  PASSED tests.py::TestLindex::test_first
5  PASSED tests.py::TestLindex::test_last
6  PASSED tests.py::TestLindex::test_not_exist
7  PASSED tests.py::TestLindex::test_range_lw_in
8  PASSED tests.py::TestLindex::test_range_lw_out
9  PASSED tests.py::TestLindex::test_range_gr_in
10 PASSED tests.py::TestLindex::test_range_gr_out
11 PASSED tests.py::TestLindex::test_range_le_in
12 PASSED tests.py::TestLindex::test_range_le_out
13 PASSED tests.py::TestLindex::test_range_ge_in
14 PASSED tests.py::TestLindex::test_range_ge_out
15 PASSED tests.py::TestLindex::test_range_gle
16 PASSED tests.py::TestLindex::test_range_gle_all
17 PASSED tests.py::TestLindex::test_range_gle_none
18 PASSED tests.py::TestLindex::test_insert
19 ===== 17 passed, 1 warning in 4.36s =====

```

Также было проведено ручное интеграционное тестирование программного обеспечения.