

Умный поиск курсов на графах

Ахтямов Алексей (mralexheimk@yandex.ru)

October 26, 2022

Введение

В текущем документе описывается построение эффективной конфигурации базы данных PostgreSQL с таблицами курсов и тэгов, создание умной системы поиска курсов, основанной на графах и извлечение тэгов из пользовательского ввода (в поле поиска курсов) с помощью расстояния Левенштейна.

Реализация на Python представлена тут [GitHub](#).

Структура базы данных

База данных будет содержать 2 таблицы ('coursers' и 'tags')

Coursers:

Тип: *динамически обновляемая*

Пример создания 'coursers' таблицы с помощью SQL запроса (нам важны только столбцы **id**, **title** (название курса) и **tags** (тэги)):

```
1 CREATE TABLE courserts (  
2     id SERIAL PRIMARY KEY,  
3     title VARCHAR (255) NOT NULL,  
4     author VARCHAR (255) NOT NULL,  
5     views INTEGER NOT NULL,  
6     likes INTEGER NOT NULL,  
7     tags TEXT NOT NULL)
```

Таблица будет содержать список добавленных курсов.
Столбец **tags** содержит набор строк (тэгов), разделенных специальным символом (в нашей реализации это символ '@').

Tags:

Тип: *периодически обновляемая*

Заполнение происходит периодически из таблицы '**coursers**'

Создание '**tags**' таблицы с помощью SQL запроса:

```
1 CREATE TABLE tags (  
2     id SERIAL PRIMARY KEY,  
3     tag VARCHAR(255) UNIQUE NOT NULL,  
4     coursers_count INTEGER NOT NULL,  
5     related_tags TEXT)
```

Таблица будет хранить список всех тэгов.

Столбец **coursers_count** хранит число курсов, использующих тэг.

Столбец **related_tags** хранит индексы тэгов, связанных с текущим.

Для ускорения обновления таблицы и ограничения числа смежных вершин в графе введем ограничение на количество тэгов, связанных с текущим.

Обозначим это ограничение сверху как RT_{max} . При периодическом обновлении таблицы в столбец **related_tags** войдут индексы тех тэгов, **coursers_count** которых в таблице наибольший.

Построение графа

Визуальный пример графа с тэгами и курсами (синие - тэги, зеленые - курсы).

Обозначим число курсов как **C**, а число тэгов как **T**.

Граф - множество вершин, соединенных между собой каким-то образом.

Ориентированный граф - граф, с ориентированным соединением между вершинами.

Вершинами нашего графа будут все тэги и все курсы (**C + T** вершин).

Каждая вершина имеет максимум RT_{max} выходящих из нее ребер, ведущих

к вершинам-тэгам и максимум C выходящих ребер, ведущих к вершинам-курсам.

В худшем случае, когда каждый тэг есть в каждом курсе, используемая для хранения графа память равна $O('число вершин' + 'число связей') = O((T + C) + (T + C) * (RT_{max} + C))$

Перебирая значения двух переменных при $RT_{max} = 5$, получаем:

	coursersCount	tagsCount	Memory
0	1.000e+03	2.000e+04	2.113e+07
1	2.000e+03	4.000e+04	8.425e+07
2	5.000e+03	1.000e+05	5.256e+08
3	1.000e+04	2.000e+05	2.101e+09
4	5.000e+04	1.000e+06	5.251e+10
5	2.000e+05	4.000e+06	8.400e+11
6	5.000e+05	1.000e+07	5.250e+12
7	1.000e+06	2.000e+07	2.100e+13

В худшем случае, при отсутствии модерации, оптимально хранить ~ 10.000 курсов и ~ 200.000 тэгов с затратами меньше 1 ГБ.

Предположим, что модерация есть, тогда, пусть M - среднее число курсов, содержащих определенный тэг (число курсов на тэг).

В среднем случае, требуемая память составляет $O('число вершин' + 'число связей') = O((T + C) + (T + C) * (RT_{max} + M))$

Перебирая значения трех переменных при $RT_{max} = 5$, получаем:

	coursersCount	tagsCount	averageCoursersOnTag	Memory
0	1.000e+03	2.000e+04	20	5.460e+05
1	2.000e+03	6.000e+04	30	2.232e+06
2	5.000e+03	2.000e+05	40	9.430e+06
3	1.000e+04	5.000e+05	50	2.856e+07
4	5.000e+04	3.000e+06	60	2.013e+08
5	2.000e+05	1.400e+07	70	1.079e+09
6	5.000e+05	4.000e+07	80	3.483e+09
7	1.000e+06	9.000e+07	90	8.736e+09

Для хранения графа с миллионом курсов и 90-та миллионами тэгов требуемая память составляет $\sim 9 \cdot 10^9$ бит (4 ГБ).

Инициализация графа происходит периодически после обновления таблицы 'tags'.

Обработка пользовательского ввода

Для поиска тэгов в пользовательской строке разобьем строку на слова (разделенные пробелом) переберем все *размещения* из количества слов по i , $\forall i \leq N$ (N - максимальное число слов в тэге, в нашей реализации оно равно 3).

В каждом размещении объединим слова в строку и с помощью бинарного поиска найдем окрестность похожих на нее тэгов. Выберем из них тот, расстояние по Левенштейну до которого минимально и, если это расстояние меньше какой-то верхней границы - берем этот тэг для дальнейшей рекомендации по нему.

Расстояние Левенштейна - минимальное число операций (удаление, вставка и замена символа), необходимых для превращения одной строки в другую. Чем оно меньше - тем более похожи строки. В нашей реализации, для избежания эквивалентности коротких тэгов, операция замены стоит в 2 раза больше других операций.

После обработки мы получим список тэгов, по которым мы будем рекомендовать пользователю курсы.

Поиск курсов

После получения списка тэгов, добавим в этот список семантически похожие тэги.

В нашем графе, если вершина-тэг A идет в вершину-тэг B означает, что B - более высокая абстракция над A или, по-другому, A - частный случай B .

Это обеспечивается за счет ограничения RT_{\max} , которое заставляет граф строить связи с теми вершинами-тэгами, количество курсов с которыми наибольшее.

Данный факт хорошо прослеживается в данном выше визуальном примере, где такие общие тэги, как 'программирование', 'для новичков' имеют сравнительно большее число вершин, входящих в них.

Для добавления новых тэгов, переберем все смежные вершины-тэги в графе у списка тэгов и добавим в список. Прделаем эту операцию столько раз, какую глубину мы указали в конфигурации.

В итоге, мы имеем список тэгов и информацию о том, на какой глубине получен этот тэг. Переберем все эти тэги в графе и найдем все смежные вершины-курсы и добавим в нашу рекомендацию. В зависимости от глубины, ключевых слов в названии курса или рейтинга сформируем и отсортируем список рекомендуемых курсов.

Полезные ссылки

1. [Python - psycorg2 библиотека \(работа с базой данных\)](#)
2. [Python - pyvis библиотека \(работа с графами\)](#)
3. [Расстояние Левенштейна](#)
4. [PyVis Network arxiv.org article](#)