

Теория формальных языков. Лекции и практики. Заметки.

Семён Григорьев

29 октября 2020 г.

Содержание

1	План лекций	4
2	Лекция 1: Введение	4
3	Практика 1	7
3.1	Григорьев С.В.	7
4	Лекция 2: Регулярные языки	8
5	Практика 2	9
5.1	Григорьев С.В.	9
6	Лекция 3. Контекстно-свободные грамматики	9
7	Практика 3	9
7.1	Григорьев С.В.	9
8	Лекция 4	10
9	Практика 4	10
9.1	Григорьев С.В.	10
10	Лекция 5	10
10.1	Нормальная форма Хомского (НФХ)	10
10.2	Лемма о накачке для КС языков	11
11	Практика 5	12
11.1	Григорьев С.В.	12
12	Лекция 6	13
12.1	Свойства замкнутости КС языков	13
12.2	Алгоритм СЮК	13
13	Практика 6	13
13.1	Григорьев С.В.	13
14	Лекция 7	15
14.1	Алгоритм на основе матричного произведения	15
14.2	Алгоритм на основе тензорного произведения	15
15	Практика 7	15
15.1	Григорьев С.В.	15

16 Лекция 8	16
17 Практика 8	16
17.1 Григорьев С.В.	16
17.1.1 Лексический анализ	16
17.1.2 Абстрактный и конкретный синтаксис	17
18 Лекция 9	19
19 Практика 9	19
19.1 Григорьев С.В.	19
20 Лекция 10	25
21 Практика 10	25
21.1 Григорьев С.В.	25
22 Лекция 11	25
23 Практика 11	25
23.1 Григорьев С.В.	25
24 Лекция 12	26
25 Практика 12	26
25.1 Григорьев С.В.	26
26 Лекция 13	26
27 Практика 11	26
27.1 Григорьев С.В.	26

1 План лекций

1. Введение. Базовые определения. Обзор курса.
2. Регулярные языки, конечные автоматы (детерминированные, недетерминированные), регулярные выражения. Детерминизация, ε -замыкание, минимизация.
3. Взаимные преобразования способов задания.
4. Теоретико-языковые свойства регулярных языков. Лемма о накачке, замкнутость относительно операций.
5. Алгоритмы вычисления операций.
6. Лево(право)-линейные грамматики и регулярные языки.
7. Грамматики, переписывающие системы. КС-грамматики (обыкновенные грамматики). Вывод в грамматике, неоднозначные грамматики, существенно неоднозначные языки, дерево вывода.
8. Рекурсивные автоматы.
9. Лемма о накачке, замкнутость относительно операций, проверка пустоты.
10. Нормальная форма Хомского, СЮК.
11. LL
12. LR
13. !!!
14. !!!
15. !!!
16. !!!
17. !!!

2 Лекция 1: Введение

Алфавит, язык. Операции над строками. Операции над языками.

Какие вопросы можно задавать о языках: о пустоте, универсальности, о построении пересечения, о пустоте пересечения, о вложенности, об эквивалентности.

Базовые способы задания: перечисление, генератор, распознаватель.

Взаимосвязь теории формальных языков с другими областями, области её применения.

- Синтаксический анализ языков программирования: в компиляторах, интерпретаторах, средах разработки, других инструментах.
- Анализ естественных языков. Активность в этой области несколько спала, так как на передний план сейчас вышли различные методы машинного обучения. Однако и в этой области ведутся работы. Примеры конференций:

- International Conference on Parsing Technologies (IWPT-2020)
- FG: Formal Grammar (FG-2020)
- Статический анализ кода.
 - Различные задачи межпроцедурного анализа. Основной подход — language reachability. Основоположник — Томас Репс. Примеры работ.
 - * Thomas Reps. 1997. Program analysis via graph reachability. In Proceedings of the 1997 international symposium on Logic programming (ILPS '97). MIT Press, Cambridge, MA, USA, 5–19.
 - * Qirun Zhang and Zhendong Su. 2017. Context-sensitive data-dependence analysis via linear conjunctive language reachability. In Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages (POPL 2017). Association for Computing Machinery, New York, NY, USA, 344–358. DOI:<https://doi.org/10.1145/3037697.3037744>
 - * Kai Wang, Aftab Hussain, Zhiqiang Zuo, Guoqing Xu, and Ardalan Amiri Sani. 2017. Graspan: A Single-machine Disk-based Graph System for Interprocedural Static Analyses of Large-scale Systems Code. In Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '17). Association for Computing Machinery, New York, NY, USA, 389–404. DOI:<https://doi.org/10.1145/3037697.3037744>
 - * Lu Y., Shang L., Xie X., Xue J. (2013) An Incremental Points-to Analysis with CFL-Reachability. In: Jhala R., De Bosschere K. (eds) Compiler Construction. CC 2013. Lecture Notes in Computer Science, vol 7791. Springer, Berlin, Heidelberg
 - Интерливинг (или шафл) языков для верификации многопоточных программ.
 - * Approximating the Shuffle of Context-free Languages to Find Bugs in Concurrent Recursive Programs
 - * Flick N.E. (2015) Quotients of Unbounded Parallelism. In: Leucker M., Rueda C., Valencia F. (eds) Theoretical Aspects of Computing - ICTAC 2015. ICTAC 2015. Lecture Notes in Computer Science, vol 9399. Springer, Cham
 - Система типов Java: Radu Grigore, Java Generics are Turing Complete.
- Графовые базы данных. Поиск путей с ограничениями.
 - Maurizio Nol  and Carlo Sartiani. 2016. Regular Path Queries on Massive Graphs. In Proceedings of the 28th International Conference on Scientific and Statistical Database Management (SSDBM '16). Association for Computing Machinery, New York, NY, USA, Article 13, 1–12. DOI:<https://doi.org/10.1145/2949689.2949711>
 - Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019. An Experimental Study of Context-Free Path Query Evaluation Methods. In Proceedings of the 31st International Conference on Scientific and Statistical Database Management (SSDBM '19). Association for Computing Machinery, New York, NY, USA, 121–132. DOI:<https://doi.org/10.1145/3335783.3335791>
 - Jelle Hellings. Querying for Paths in Graphs using Context-Free Path Queries.
- Биоинформатика. В основном это анализ геномных и белковых последовательностей.
 - Witold Dyrka, Mateusz Pyzik, Francois Coste, and Hugo Talibart. Estimating probabilistic context-free grammars for proteins using contact map constraints.

- James WJ Anderson, Paula Tataru, Joe Staines, Jotun Hein, and Rune Lyngso. Evolving stochastic context-free grammars for RNA secondary structure prediction.
- Ryan Zier-Vogel. Predicting RNA secondary structure using a stochastic conjunctive grammar.
- Машинное обучение.
 - Matt J. Kusner, Brooks Paige, José Miguel Hernández-Lobato. Grammar Variational Autoencoder. Опубликовано в 2017 году и уже больше 200 цитирований.
 - TAG Parsing with Neural Networks and Vector Representations of Supertags. К разговору об обработке естественных языков.
 - Jungo Kasai, Robert Frank, Pauli Xu, William Merrill, Owen Rambow. End-to-end Graph-based TAG Parsing with Neural Networks.
- Языки — это не только про строки.
 - Языки деревьев: Tree Automata Techniques and Applications.
 - Языки графов:
 - * Graph Grammars
 - * HYPEREDGE REPLACEMENT GRAPH GRAMMARS
 - * (Re)introducing Regular Graph Languages
 - * Hyperedge Replacement: Grammars and Languages
 - ...
- Теория групп. Как правило, это проблема слов группы или дополнение к ней.
 - Anisimov, A.V. Group languages. Cybern Syst Anal (1971) 7: 594.
 - David E. Muller, Paul E. Schupp, Groups, the Theory of ends, and context-free languages, Journal of Computer and System Sciences, Volume 26, Issue 3, 1983, Pages 295-310, ISSN 0022-0000
 - HOLT, D., REES, S., ROVER, C., & THOMAS, R. (2005). GROUPS WITH CONTEXT-FREE CO-WORD PROBLEM. Journal of the London Mathematical Society, 71(3), 643-657. doi:10.1112/S002461070500654X
 - Groups with Context-Free Co-Word Problem and Embeddings into Thompson’s Group V
 - Kropholler, R. & Spriano, D. (2019). Closure properties in the class of multiple context-free groups. Groups Complexity Cryptology, 11(1), pp. 1-15. Retrieved 13 Feb. 2020, from doi:10.1515/gcc-2019-2004
 - Word problems of groups, formal languages and decidability
- Прочая забавная математика.
 - Немного топологии в теории формальных языков: Salvati S. On is an n-MCFL. – 2018.
 - Salvati S. MIX is a 2-MCFL and the word problem in Z2 is captured by the IO and the OI hierarchies //Journal of Computer and System Sciences. – 2015. – Т. 81. – №. 7. – С. 1252-1277.

- О том, как задачи из теории графов связаны с теорией формальных языков: Abboud, Amir & Backurs, Arturs & Williams, Virginia. (2015). If the Current Clique Algorithms are Optimal, So is Valiant's Parser. 98-117. 10.1109/FOCS.2015.16.
- A context-free grammar for the Ramanujan-Shor polynomials

3 Практика 1

Детали о том, как будет проходить практика.

3.1 Григорьев С.В.

Немного про описания языков. Пописать языковые уравнения, грамматики. Посмотреть на операции над языками.

Постановка задачи на весь семестр.

Запросы к графовым базам данных. Контекст задачи, примеры графовых БД (RedisGraph, Neo4j, ...), задача о путях в принципе.

Ссылка на второй конспект.

Задача: реализовать свою "графовую миниБД".

Реализация: оформление, инструменты, языки.

- Ограничений на язык реализации нет.
- Ограничений на использование библиотек нет. Главное — не нарушать лицензии и чтобы можно было вносить изменения в библиотеку (при необходимости).
- Каждый создаёт под решение репозиторий на GitHub и снабжает его всем необходимым: readme, лицензия, CI-сборка с тестированием, инструкции по локальному развёртыванию.
- Разработка ведётся в отдельной ветке и когда очередная часть задачи готова к сдаче — делаем pull request в master и добавляем меня (gsvglit) в ревьюеры.

Задачи на дом.

1. Выбрать язык программирования, на котором будет вестись разработка.
2. Создать репозиторий на GitHub.
3. Настроить CI-сборку и тестирование.
4. Реализовать подгрузку графов из RDF используя готовые библиотеки.

4 Лекция 2: Регулярные языки

Иерархия Хомского. Проблемы с ней. Классы языков.

Граматики. Системы переписывания.

Регулярные множества. Регулярные языки. Регулярные выражения.

$$V^* = \bigcup_{i=0}^{\infty} V^i$$

Конечные автоматы. Система переходов.

Язык, задаваемый автоматом.

Понятие выводимости (\vdash^*).

Конфигурация: $\langle \text{Состояние}, \text{Остаток} \rangle$.

Полный автомат и вершина-сток.

Детерминизация, алгоритм Томпсона.

НКА: $\langle \Sigma, Q, s \in Q, T \in Q, \delta : Q \times \Sigma \rightarrow 2^Q \rangle$

ДКА: $\langle \Sigma, Q_d, s_d \in Q_d, T_d \in Q_d, \delta_d : Q_d \times \Sigma \rightarrow Q_d \rangle$, где:

- $Q_d = \{q_d \mid q_d \in 2^Q\}$,
- $s_d = \{s\}$,
- $T_d = \{q \in Q_d \mid \exists p \in T : p \in q\}$,
- $\delta_d(q, c) = \{\delta(a, c) \mid a \in q\}$.

ε -замыкание.

1. Транзитивное замыкание отношения ε -перехода.
2. Обработка финальных состояний
3. Добавление переходов: если $\delta(v_0, \varepsilon) = v_1, \delta(v_1, c) = v_2$, то добавим $\delta(v_0, c) = v_2$.
4. Удалим ε -переходы.

Эквивалентность автоматов. Эквивалентность состояний: состояния эквивалентны если нет различающей строки.

Минимизация.

Теорема Клини об эквивалентности автоматов и регулярных языков.

Построение автомата по регулярному выражению.

Построение регулярного выражения по автомату: устранение вершин.

5 Практика 2

5.1 Григорьев С.В.

Построение минимального ДКА по регулярному выражению.

Домашнее задание.

1. Реализовать функцию (можно с применением библиотек), которая принимает на вход регулярное выражение в виде строки и строит по нему минимальный ДКА.
2. Реализовать необходимые тесты на построение ДКА по регулярному выражению.
3. Реализовать (можно с применением библиотек) пересечение минимального ДКА и НКА без ϵ -переходов.
4. Реализовать необходимые тесты на пересечение ДКА и НКА.

6 Лекция3. Контекстно-свободные грамматики

Левосторонние и правосторонние грамматики и регулярные языки. Неразрешимость задачи проверки того, что грамматика задаёт регулярный язык. Статья на эту тему: Self-embedded context-free grammars with regular counterparts. Грамматика \rightarrow регулярка и регулярка \rightarrow грамматика.

Вывод цепочки в грамматике, левосторонний, правосторонний вывод, неоднозначные и однозначные грамматики. Примеры. Существенно неоднозначные языки.

Дерево вывода. Соотношение между деревьями и выводами. Примеры.

Расширенные контекстно-свободные грамматики.

7 Практика 3

7.1 Григорьев С.В.

Пересечение автоматов — это тензорное произведение матриц смежности. Пример.

Про коммутативность пересечения и некоммутативность тензорного произведения.

Домашнее задание.

1. Реализовать консольный клиент, позволяющий
 - (a) загрузить RDF-файл
 - (b) вывести список меток рёбер
 - (c) задать к загруженному графу регулярный запрос с возможностью указать представление результата: пустота ответа, автомат сдампит в файл в формате DOT (<https://www.graphviz.org/doc/info/lang.html>), пара (кол-во рёбер, кол-во вершин) в результирующем автомате
 - (d) выйти из клиента.
2. Подгрузку RDF и выполнение запросов реализовать на основе уже существующей функциональности.

3. Провести замеры производительности на графах из репозитория https://github.com/JetBrains-Research/CFPQ_Data. Графы брать из подпапки `data/graphs/RDF`. Так как в графах присутствуют одинаковые отношения, то можно один и тот же запрос выполнять на всех графах. Отчёт оформить в виде раздела в README репозитория в виде таблицы.

Эти эксперименты проводятся локально! Не надо таскать репозиторий с графами за собой. Для тестов клиента использовать маленькие синтетические RDF.

4. Реализовать необходимые тесты на работоспособность клиента.

8 Лекция 4

Рекурсивные автоматы. Построение, интерпретация.

9 Практика 4

9.1 Григорьев С.В.

Больше подробностей про рекурсивные автоматы: тотальная минимизация. Как их применять для КС запросов. Тензоры + транзитивное замыкание.

Домашнее задание.

1. Реализовать выполнение регулярных запросов через тензорное произведение. Для тензорного произведения использовать существующие библиотеки линейной алгебры. Обратите внимание на то, что матрицы должны быть разреженными. Скорее всего, удобно будет использовать представление в виде набора булевых матриц.
2. Интегрировать новую реализацию в клиент наравне со старой.
3. Провести замеры производительности на графах из репозитория https://github.com/JetBrains-Research/CFPQ_Data. Графы брать из подпапки `data/graphs/RDF`. Так как в графах присутствуют одинаковые отношения, то можно один и тот же запрос выполнять на всех графах. Отчёт оформить в виде раздела в README репозитория в виде таблицы. Сравнить с результатами предыдущей задачи.

Эти эксперименты проводятся локально! Не надо таскать репозиторий с графами за собой. Для тестов клиента использовать маленькие синтетические графы и запросы.

4. Реализовать необходимые тесты на работоспособность алгоритма через тензорное произведение.

10 Лекция 5

10.1 Нормальная форма Хомского (НФХ)

Определение 10.1. КС грамматика находится в нормальной форме Хомского если любое правило имеет один из трёх видов:

1. $S \rightarrow \varepsilon$
2. $N_i \rightarrow t_j$
3. $N_i \rightarrow N_j N_k, N_j \neq S, N_k \neq S$

Важно: стартовый нетерминал не встречается в правых частях правил, ε -продукция только для стартового нетерминала.

Note. Любую КС грамматику можно преобразовать к нормальной форме Хомского.

Преобразование в НФХ. Шаги.

1. Устранение длинных правил.
2. Устранение ε -правил.
3. Устранение цепных правил.
4. Устранение бесполезных нетерминалов
 - (a) Удаление непорождающих нетерминалов
 - (b) Удаление недостижимых нетерминалов
5. Устранение продукций с правой частью длины 2, содержащей терминалы.

Надо не забыть добавить новый стартовый нетерминал, если нужно: чтобы вывести из него ε и чтобы не встречался в правых частях правил.

Важно. Порядок применения шагов преобразования важен.

1. Второй шаг можно поднять вверх, но это приведёт к более существенному разрастанию результирующей грамматики.
2. Подшаги шага 4 нельзя менять местами. Попробуйте поприменять их к грамматике:

$$\begin{aligned} S &\rightarrow AB \mid a \\ A &\rightarrow b \end{aligned}$$

Материалы по преобразованию в НФХ.

10.2 Лемма о накачке для КС языков

Теорема 10.1. Пусть L — контекстно-свободный язык над алфавитом Σ , тогда существует такое n , что для любого слова $\omega \in L$, $|\omega| \geq n$ найдутся слова $u, v, x, y, z \in \Sigma^*$, для которых верно: $uvxyz = \omega$, $vy \neq \varepsilon$, $|vxy| \leq n$ и для любого $k \geq 0$ $uv^kxy^kz \in L$.

Идея доказательства леммы о накачке.

1. Для любого КС языка можно найти грамматику в нормальной форме Хомского.

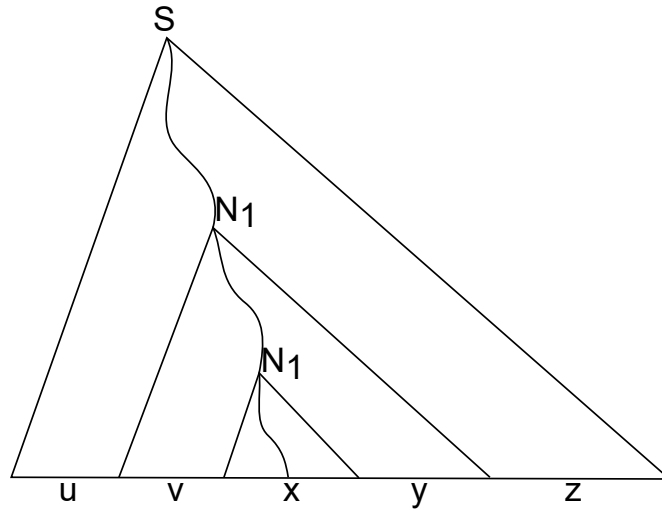


Рис. 1: Разбиение цепочки для леммы о накачке

2. Очевидно, что если брать достаточно длинные цепочки, то в дереве вывода этих цепочек, на пути от корня к какому-то листу обязательно будет нетерминал, встречающийся минимум два раза. Если m — количество нетерминалов в НФХ, то длины 2^{m+1} должно хватить. Это и будет n из леммы.
3. Возьмём путь, на котором есть хотя бы дважды повторяется некоторый нетерминал. Скажем, это нетерминал N_1 . Пойдём от листа по этому пути. Найдём первое появление N_1 . Цепочка, задаваемая поддеревом для этого узла — это x из леммы.
4. Пойдём дальше и найдём второе появление N_1 . Цепочка, задаваемая поддеревом для этого узла — это uxy из леммы.
5. Теперь мы можем копировать кусок дерева между этими повторениями N_1 и таким образом накачивать исходную цепочку.

Надо только проверить выполнение ограничений на длины.

Материалы по лемме о накачке для КС языков.

Проверить неконтекстно-свободность языка $L = \{a^n b^n c^n \mid n > 0\}$.

11 Практика 5

11.1 Григорьев С.В.

Преобразование в нормальную форму Хомского.

Формат входа:

1. Одна продукция на строку.
2. Продукция — это список терминалов и нетерминалов через пробел, начинающийся с нетерминала (левая часть продукции).
3. Нетерминалы — заглавные буквы с опциональным числовым суффиксом.
4. Терминалы — строчные буквы с опциональным числовым суффиксом.

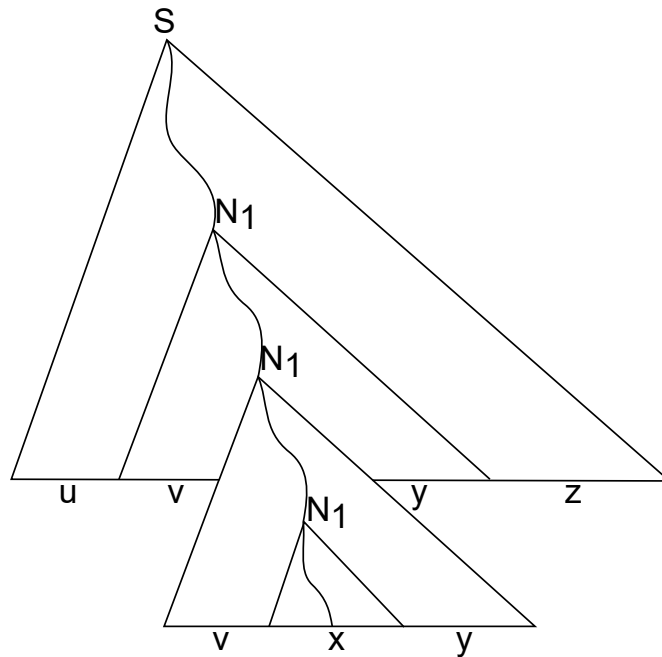


Рис. 2: Пример накачки цепочки с рисунка 1

5. Специальный символ `eps` для обозначения ε .

Пример входа, описывающего грамматику $S \rightarrow aSbS \mid \varepsilon$:

```
S a S b S
S eps
```

Домашнее задание.

1. Реализовать преобразование в нормальную форму Хомского. На входе файл с грамматикой, на выходе — файл с грамматикой в НФХ в том же формате, что и вход.

12 Лекция 6

12.1 Свойства замкнутости КС языков

Глава 2.6 конспекта.

12.2 Алгоритм СҮК

Глава 4.1 “Алгоритм СҮК” конспекта.

13 Практика 6

13.1 Григорьев С.В.

Алгоритм СҮК и алгоритм Хеллингса.

Формат входа для СҮК:

1. Грамматика: смотри предыдущее ДЗ.
2. Входная строка: терминалы разделены пробелами.

Пример входа, описывающего грамматику $S \rightarrow aSbS \mid \varepsilon$:

```
S a S b S
S eps
```

Пример входной строки:

```
a a b a a b b b
```

Формат входа алгоритма Хеллингса:

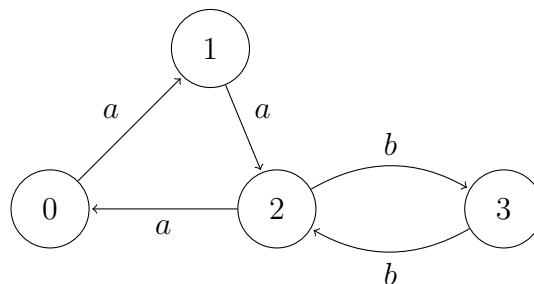
1. Грамматика: тот же формат, что и для СУК. НО! ИСпользуем преобразование а ослабленную НФХ.
2. Входной граф: файл в котором на каждой строке записано ребро в виде тройки

$\langle \text{вершина} \text{ метка_ребра } \text{вершина} \rangle$.

Элементы тройки разделены пробелами.

3. Можно считать, что все вершины графа — числа от нуля, идущие подряд.

Пример входного графа:



Пример описания входного графа:

```
0 a 1
1 a 2
2 a 0
2 b 3
3 b 2
```

Домашнее задание.

1. Реализовать алгоритм СУК для линейного входа. На вход принимаются два файла: с грамматикой и входной строкой. Результат (выводится ли входная цепочка в грамматике) печатается в консоль.
2. Реализовать алгоритм Хеллингса. На вход принимается файл с грамматикой и файл с графом. В результирующий файл печатается грамматика в ослабленной НФХ (с которой непосредственно работал алгоритм) и множество пар достижимых вершин для стартового нетерминала (одна пара на строку, две вершины через пробел)

14 Лекция 7

Алгоритмы решения задачи контекстно-свободной достижимости, основанные на операциях линейной алгебры.

14.1 Алгоритм на основе матричного произведения

Глава 5.1 конспекта.

14.2 Алгоритм на основе тензорного произведения

Глава 6 конспекта.

15 Практика 7

15.1 Григорьев С.В.

Алгоритмы на основе линейной алгебры.

Для реализации предлагается использовать следующие библиотеки. Так как с булевыми не везде хорошо, то будем использовать те типы, которые поддерживаются: `Int`, `float` и т.д.

- Для языка Python — разреженные матрицы в `scipy` и соответствующие операции работы с ними: `scipy.sparse.kron` и обычное матричное произведение. Предпочтительный формат разреженных матриц — CSR.
- Для языка Kotlin — `la4j`. Операции: кронекер и обычное умножение.

Поэлементное сложение есть и там и там.

Для алгоритма на матричном умножении всё точно так же, как и в предыдущей ДЗ для Хеллингса.

Для тензорного произведения расширим формат представления входной грамматики. Одна строка на нетерминал. Терминалы, нетерминалы, ε обозначаются как и раньше. Как и раньше, левая часть от правой отделена пробелом. В правой части можно использовать конструкции регулярных выражений: альтернатива, звезда клин, групперирующие скобки. Этот набор можно расширять по своему усмотрению.

Пример входа, описывающего грамматику $S \rightarrow (aSb)^* \mid \varepsilon$:

`S (a S b)* | eps`

Домашнее задание. Время на выполнение — две недели. Один из алгоритмов — на первую, оставшийся и эксперименты — на вторую.

1. Реализовать алгоритм, основанный на матричном умножении. На вход принимаются два файла: с грамматикой и входным графом. В результирующий файл печатается грамматика в ослабленной НФХ (с которой непосредственно работал алгоритм) и множество пар достижимых вершин для стартового нетерминала (одна пара на строку, две вершины через пробел).

2. Реализовать алгоритм, основанный на тензорном произведении. На вход принимается файл с грамматикой и файл с графом. В результирующий файл печатается матрица смежности рекурсивного автомата (с которым непосредственно работал алгоритм, построчно, элементы разделены пробелом, пустая ячейка обозначается символом '.') и множество пар достижимых вершин для стартового нетерминала (одна пара на строку, две вершины через пробел).
3. Сравнить производительность трёх реализованных алгоритмов (Хеллингс, матричное произведение, тензорное произведение). Результат — описание эксперимента и таблица сравнения в readme.

16 Лекция 8

Нисходящий синтаксический анализ: рекурсивный спуск и LL(k).

Глава 7, 8.1 и 8.2 конспекта.

17 Практика 8

17.1 Григорьев С.В.

Разработка синтаксических анализаторов: лексический и синтаксический анализы, абстрактный и конкретный синтаксис.

17.1.1 Лексический анализ

$$E \rightarrow n \mid E + E \mid E * E$$

Что такое n (число)? Это абстракция.

$$n = [1 - 9][0 - 9]^*$$

$$n = 0 \mid ((-)?[1 - 9][0 - 9]^*)$$

$$n = 0 \mid ((-)?[1 - 9][0 - 9]^*(.)[0 - 9]^*[1 - 9])$$

То же самое верно и для других терминалов.

$$E \rightarrow n \mid E \text{ op_plus } E \mid E \text{ op_pow } E \mid E \text{ op_mult } E$$

$$\text{op_pow} = \wedge$$

$$\text{op_pow} = **$$

Для введения этой абстракции анализ языков разделяют на две стадии:

1. Лексический анализ: переводим последовательность “символов с клавиатуры” в последовательность токенов/терминальных символов. Работает на основе регулярных выражений (как правило).
2. Синтаксический анализ: переводим последовательность терминальных символов в структурное представление (дерево разбора). Работает на основе КС грамматик.

Да, можно и не разделять: scannerless parsing.

17.1.2 Абстрактный и конкретный синтаксис

Структурное представление текста, удобное для решения тех или иных задач, не содержит многие детали исходного текста.

Абстрактное описание синтаксиса:

```
type stmt =
  ...
  | IfStmt of expr*List<stmt>*List<stmt>
  ...
```

Примеры конкретного синтаксиса:

```
...
if ( a + b > 0 )
then
  ...
else
  ...
...
```

```
...
if ( a + b > 0 )
{
  ...
}
else
{
  ...
}
...
```

```
...
cond_branch ( a + b > 0 )( ... )( ... )
...
```

Ещё пример абстрактного синтаксиса:

```
type arithExpr =
  | Num of int
  | OpPlus of arithExpr * arithExpr
  | OpMult of arithExpr * arithExpr
  | OpPow of arithExpr * arithExpr
```

Важно то, что здесь нет ни слова про скобки, приоритеты, ассоциотивность.

$12 + 3 * 4 \sim 5 + 2$

$(12 + 3) * 4 \sim 5 + 2$

$12 + 3 * 4 \sim (5 + 2)$

$(12 + 3 * 4) \sim (5 + 2)$

Потому что многие свойства естественным образом выражается структурой дерева и, соответственно, детали исходного текста теряются.

$12 + 3 * 4 \sim 5 + 2$

```
OpPlus
  (Num 12)
  (OpPlus
    (OpMult
      (Num 3)
      (OpPow
        (Num 4)
        (Num 5))))
  (Num 2)
```

$(12 + 3) * 4 \sim 5 + 2$

```
OpPlus
  (OpMult
    (OpPlus
      (Num 12)
      (Num 3))
    (OpPow
      (Num 4)
      (Num 5)))
  (Num 2))
```

$12 + 3 * 4 \sim (5 + 2)$

```
OpPlus
  (Num 12)
  (OpMult
    (Num 3)
    (OpPow
      (Num 4)
      (OpPlus
        (Num 5)
        (Num 2))))
```

$(12 + 3 * 4) \wedge (5 + 2)$

```
OpPow
  (OpPlus
    (Num 12)
    (OpMult
      (Num 3)
      (Num 4))
  )
  (OpPlus
    (Num 5)
    (Num 2))
```

Но важно помнить, что для разных задач нужна разная информация.

Можно начинать знакомство с ANTLR, так как он будет использоваться в следующих домашних работах.

18 Лекция 9

Нисходящий синтаксический анализ: GLL и его применение для поиска путей с КС ограничениями.

Глава 8.3 и 8.4 конспекта.

19 Практика 9

19.1 Григорьев С.В.

Абстрактный и конкретный синтаксис языка запросов к графам, семантика языка запросов к графам.

Абстрактный синтаксис скрывает детали синтаксического анализа и конкретного “текстового” представления синтаксиса, оставляя только важную для дальнейшего анализа информацию о структуре программы.

Первая, минималистичная, версия языка запросов к графам.

```
type script: Seq of List<stmt>
```

```
type stmt:
  | Connect of String
  | Select of obj_expr * graph_expr
```

```
type graph_expr =
  | Intersect of graph_expr * graph_expr
  | Query of pattern
  | GraphName of String
```

```

type obj_expr:
  | Edges
  | Conut of obj_expr

type pattern:
  | smb of String
  | star of pattern
  | plus of pattern
  | alt of pattern * pattern
  | seq of pattern * pattern
  | option of pattern

```

Теперь поговорим о семантике данного языка.

Семантика — отображение из программ в семантический домен. Программы обычно представляются как абстрактные синтаксические деревья. Семантический домен — то, что описывает смысл программы в контексте решаемой задачи.

Семантика бывает разной. Классические примеры: динамическая семантика (вычисление программ) и статическая семантика (вывод типов).

Доменом динамической семантики арифметических выражений является множество функций из функции, сопоставляющей переменным целочисленные значения, в целые числа.

Семантика часто обозначается скобками $\llbracket \cdot \rrbracket$ и $\llbracket \cdot \rrbracket$.

Семантика арифметических выражений:

$$\llbracket \cdot \rrbracket : (X \rightarrow \mathbb{N}) \rightarrow \mathbb{N}$$

$$\llbracket n \rrbracket(\sigma) = n, n \in \mathbb{N}$$

$$\llbracket v \rrbracket(\sigma) = \sigma(v)$$

$$\llbracket e_1 \oplus e_2 \rrbracket(\sigma) = \llbracket e_1 \rrbracket(\sigma) \otimes \llbracket e_2 \rrbracket(\sigma)$$

$$\llbracket -e \rrbracket(\sigma) = \llbracket 0 - e \rrbracket(\sigma)$$

Здесь \oplus — конструкторы бинарных операций, \otimes — соответствующие арифметические функции.

В нашем случае, исполнение скрипта на получившемся языке достаточно естественно описывается состоянием из текущей базы данных, текущего графа, и текущего результата вычислений. Значение этой тройки и будут семантическим доменом нашего языка.

При обработке непустой последовательности инструкций сначала вычисляем результат выполнения первой инструкции, затем на нем — всех остальных.

$$\frac{c \xrightarrow{hd} c', \quad c' \xrightarrow{Seq \ tl} c''}{c \xrightarrow{Seq \ (hd::tl)} c''}$$

При обработке команды **Connect** сохраняем её аргумент в конфигурации как путь к текущей базе данных.

$$\frac{}{\langle db, graph, res \rangle \xrightarrow{Connect\ db'} \langle db', graph, res \rangle}$$

При обработке инструкции **Select** сперва вычисляется граф, из которого будут выбираться объекты, затем вычисляются сами объекты.

$$\frac{c \xrightarrow{graph} c', \quad c' \xrightarrow{obj} c''}{c \xrightarrow{Select\ obj\ graph} c''}$$

При обработке инструкции **Graph** загружаем в окружение граф из файла с именем, соответствующим аргументу этой инструкции, лежащем в базе (директории), соответствующей значению `cur_db`.

$$\frac{}{\langle db, graph, res \rangle \xrightarrow{GraphName\ name} \langle db, loadGraph(String.concat(db, name)), res \rangle}$$

При обработке инструкции **Intersect** сперва вычисляем оба подвыражения, затем перескаем получившиеся графы как конечные автоматы. Считаем все вершины стартовыми и финальными одновременно.

$$\frac{\langle db, graph, res \rangle \xrightarrow{graph_1} \langle db, graph', res \rangle, \quad \langle db, graph, res \rangle \xrightarrow{graph_2} \langle db, graph'', res \rangle}{\langle db, graph, res \rangle \xrightarrow{Intersect\ graph_1\ graph_2} \langle db, intersectFA(graph', graph''), res \rangle}$$

При обработке инструкции **Query** необходимо построить минимальный детерминированный автомат по регулярному выражению, задаваемому шаблоном.

$$\frac{}{\langle db, graph, res \rangle \xrightarrow{Query\ pattern} \langle db, buildMinDFA(pattern), res \rangle}$$

При обработке инструкции **Edges** в текущий результат вычислений записывается множество троек, соответствующих рёбрам текущего графа.

$$\frac{}{\langle db, graph, res \rangle \xrightarrow{Edges} \langle db, graph, \{(v_i, e, v_j) \mid (v_i, e, v_j) \in graph.Edges\} \rangle}$$

При обработке инструкции **Count** сперва вычисляется подвыражение, а затем в текущий результат вычислений записывается размер результата вычисления подвыражения, если результат был множеством.

$$\frac{\langle db, graph, res \rangle \xrightarrow{expr} \langle db, graph, res' \rangle}{\langle db, graph, res \rangle \xrightarrow{Count\ expr} \langle db, graph, Set.count(res') \rangle} \quad type_of(res') \equiv Set$$

$$\begin{aligned}
&script \rightarrow \varepsilon \\
&\quad | stmt \ SEMI \ script \\
&stmt \rightarrow KW_CONNECT \ KW_TO \ STRING \\
&\quad | KW_LIST \ KW_ALL \ KW_GRAPHS \\
&\quad | select_stmt \\
&\quad | named_pattern_stmt \\
named_pattern &\rightarrow NT_NAME \ OP_EQ \ pattern \\
select_stmt &\rightarrow KW_SELECT \ obj_expr \ KW_FROM \ STRING \ KW_WHERE \ where_expr \\
obj_expr &\rightarrow vs_info \\
&\quad | KW_COUNT \ vs_info \\
&\quad | KW_EXISTS \ vs_info \\
vs_info &\rightarrow LBR \ IDENT \ COMMA \ IDENT \ RBR \\
&\quad | IDENT \\
where_expr &\rightarrow LBR \ v_expr \ RBR \\
&\quad OP_MINUS \ pattern \ OP_MINUS \ OP_GR \\
&\quad LBR \ v_expr \ RBR \\
v_expr &\rightarrow IDENT \mid UNDERSCORE \mid IDENT \ DOT \ KW_ID \ OP_EQ \ INT \\
pattern &\rightarrow alt_elem \\
&\quad | alt_elem \ MID \ pattern \\
alt_elem &\rightarrow seq \\
&\quad | LBR \ RBR \\
seq &\rightarrow seq_elem \\
&\quad | seq_elem \ seq \\
seq_elem &\rightarrow prim_pattern \\
&\quad | prim_pattern \ OP_STAR \\
&\quad | prim_pattern \ OP_PLUS \\
&\quad | prim_pattern \ OP_Q \\
prim_pattern &\rightarrow IDENT \\
&\quad | NT_NAME \\
&\quad | LBR \ pattern \ RBR
\end{aligned}$$

Токены (терминальный алфавит).

$LBR = '('$
 $RBR = ')'$
 $COMMA = ','$
 $SEMI = ';'$
 $MID = '|'$
 $DOT = '.'$
 $OP_STAR = '*'$
 $OP_PLUS = '+'$
 $OP_Q = '?'$
 $OP_MINUS = '-'$
 $OP_GR = '>'$
 $OP_EQ = '='$
 $KW_ID = 'ID'$
 $KW_COUNT = 'COUNT'$
 $KW_EXISTS = 'EXISTS'$
 $KW_FROM = 'FROM'$
 $KW_WHERE = 'WHERE'$
 $KW_LIST = 'LIST'$
 $KW_ALL = 'ALL'$
 $KW_GRAPHS = 'GRAPHS'$
 $KW_CONNECT = 'CONNECT'$
 $KW_TO = 'TO'$
 $IDENT = [a - z][a - z]^*$
 $INT = 0 \mid [1 - 9][0 - 9]^*$
 $NT_NAME = [A - Z][a - z]^*$
 $STRING = '([aA - zZ] \mid [0 - 9] \mid (' - ' \mid ' ' \mid ' _ ' \mid ' / ' \mid ' . '))^*'$

Пример скрипта:

```
CONNECT TO [\home\user\graph_db];
S = a S b S | ( ) ;
SELECT COUNT(u) FROM [g1.txt] where (v.id = 10) - S -> (u);
```

Задача на дом.

1. Описать конкретный синтаксис языка запросов к графам, согласованный с приведённым абстрактным синтаксисом. Результат: грамматика и лексика в маркдаун (*.md) файле в репозитории (например, описание языка в readme.md). С примерами выражений на языке и текстовым описанием его особенностей.
2. Синтаксический анализатор предложенного языка на основе имеющейся реализации СУК. На вход принимается строка (с консоли или из файла), в которой токены разделены пробелами и переносами строк. Для приведённого выше скрипта вход выглядит следующим образом:

```

KW_CONNECT KW_TO STRING SEMI
NT_NAME OP_EQ IDENT NT_NAME IDENT NT_NAME MID LBR RBR SEMI
KW_SELECT KW_COUNT LBR IDENT RBR
KW_FROM STRING
KW_WHERE LBR IDENT DOT KW_ID OP_EQ INT RBR
OP_MINUS NT_NAME OP_MINUS OP_GR LBR IDENT RBR SEMI

```

Расширенная версия.

```

type script: Seq of List<stmt>

type stmt:
  | Connect of String
  | Select of obj_expr * graph_expr
  | NamedPattern of String * regexp

type graph_expr =
  | ReachabilityInfo of graph_expr * path_pattern
  | Intersect of graph_expr * graph_expr
  | GraphName of String

type obj_expr:
  | Vertices of String
  | Conut of String
  | Exists of String
  | Pair of obj_expr * obj_expr

vs_info:
  | pair of v_info * v_info
  | single of v_info

v_info:
  | name of String
  | any

path_pattern:
  | path of v_expr * pattern * v_expr

v_expr:
  | name of String
  | any
  | v_pattern

v_pattern:
  | id_pattern of Int

pattern:
  | smb of String
  | star of regexp

```



```
| plus of regexp
| alt of regexp * regexp
| seq of regexp * regexp
| option of regexp
```

20 Лекция 10

Восходящий синтаксический анализ: алгоритмы семейства LR.

Глава 9.1 конспекта.

21 Практика 10

21.1 Григорьев С.В.

Синтаксический анализ с использованием генератора синтаксических анализаторов ANTLR.

Домашняя страница ANTLR.

ANTLR для Python.

ANTLR для Java.

Задача на дом.

1. Реализовать синтаксический анализатор языка запросов к графовым БД (разработанный в рамках предыдущей ДЗ) с использованием ANTLR (или другого инструмента создания синтаксических анализаторов). Ожидаемая функциональность следующая.
 - Чтение скрипта из файла.
 - Чтение скрипта с консоли.
 - Вывод в консоль сообщения о корректности/некорректности скрипта.
 - Вывод в файл дерева разбора скрипта в формате GraphViz/DOT.

22 Лекция 11

Восходящий синтаксический анализ: алгоритмы семейства LR.

Глава 9.2 конспекта.

23 Практика 11

23.1 Григорьев С.В.

Исполнение скриптов запросов.

ANTLR-пакет для Ubuntu.

Пример калькулятора.

Задача на дом.

1. Автоматизированная сборка с генерацией файлов по грамматике. Сгенерированные файлы удалены из репозитория.
2. Поддержка выполнения части команд языка запросов. На данном этапе считаем, что скрипт полностью корректен.
 - `connect` полностью
 - `list` полностью
 - `named_pattern` можно без конструкций регулярных выражений (звезда Клини, альтернатива)
 - `select` можно только `exists` и `v_expr` только имя (`name of String`)
3. Тесты. Очень много тестов. Каждая смысловая конструкция должна быть проверена.

24 Лекция 12

Контрольная работа по КС языкам.

25 Практика 12

25.1 Григорьев С.В.

Исполнение скриптов запросов.

ANTLR-пакет для Ubuntu.

Пример калькулятора.

Задача на дом.

1. Полная поддержка языка запросов.
2. Тесты. Очень много тестов. Каждая смысловая конструкция должна быть проверена.

26 Лекция 13

Сиинтаксически управляемая трансляция, атрибутивные грамматики.

Презентация.

27 Практика 11

27.1 Григорьев С.В.

Дополнительные задачи на зачёт.

Для получения зачёта необходимо выполнить следующие условия.

1. Сдать все семестровые задачи.

2. Если какие-то из задач в семестре были сданы с нарушением условий, установленных в начале семестра, то необходимо решить дополнительные задачи. Какие именно, указано в таблице с результатами, в колонке “Доп. задача”.

Дополнительные задачи. В таблице указаны их номера. Красным выделены те, которые нужно решить. Для каждой задачи должно быть представлено:

- Расширение абстрактного синтаксиса (добавлено в документацию в репозитории)
 - Расширение конкретного синтаксиса (расширена документация в репозитории)
 - Примеры (добавлено в документацию в репозитории)
 - Реализация
 - Тесты
1. Расширить команду **list** так, чтобы можно было опционально указывать путь к БД, графы в которой хочется вывести. По умолчанию всё так же выводятся графы из текущей БД.
 2. Расширить команду **list** так, чтобы с её помощью можно было вывести множество различных меток рёбер в указанном графе.
 3. Расширить команду **select** возможностью опционально указывать алгоритм, с помощью которого выполнять текущий запрос. Воспроизвести эксперимент из ДЗ 5. Повлияло ли использование языка запросов на результаты экспериментов?
 4. Расширить команду **select** возможностью в качестве графа-источника использовать результат операций объединения, пересечения и дополнения графов из текущей базы данных. То есть во **from** можно писать выражение над графами типа **from** $(g_1 \cup g_2) \cap g_3$. Данные операции должны трактоваться как операции над языками, задаваемыми автоматом, где переходы определяются графом, а стартовые и финальные состояния — все вершины графа.