

Санкт-Петербургский государственный университет
Факультет Прикладной Математики - Процессов Управления
Кафедра технологии программирования

Буланина Екатерина Дмитриевна

Извлечение данных из учебного контента

Выпускная квалификационная работа бакалавра

Научный руководитель:
к. ф.-м. н., доцент Добрынин В. Ю.

Рецензент:
разработчик ЗАО «БИОКАД» Чуриков Н. С.

Заведующий кафедрой:
к. т. н., доцент Блеканов И. С.

Санкт-Петербург
2019

SAINT-PETERSBURG STATE UNIVERSITY

Applied Mathematics and Control Processes Faculty
Programming Technology Chair

Ekaterina Bulanina

Educational data mining

Graduation Thesis

Scientific supervisor:
Vladimir Dobrynin

Reviewer:
Nikita Churikov

Chairhead:
Ivan Blekanov

Saint-Petersburg
2019

Оглавление

Введение	4
1. Постановка задачи	6
2. Обзор литературы	7
3. Архитектура	9
4. Извлечение и обработка данных	10
4.1. Формат данных	10
4.2. Извлечение данных	11
4.3. Предобработка	12
5. Построение модели	14
5.1. Входные данные и необходимые метрики	14
5.2. Предсказание тем	15
6. Построение и расширение графа	19
6.1. Обозначения	19
6.2. Построение графа	19
6.3. Расширение графа	20
7. Апробация	22
7.1. Привязка заданий к темам	22
7.2. Построение графа тем	22
Заключение	24
Список литературы	25

Введение

За последние десятилетия электронная промышленность достигла колоссальных масштабов, а компьютеры и информационные технологии повсеместно вошли в быт людей. Человечество использует компьютеры для автоматизации всех сфер деятельности: производство и сельское хозяйство, логистика и транспорт, биология и медицина, а также образование.

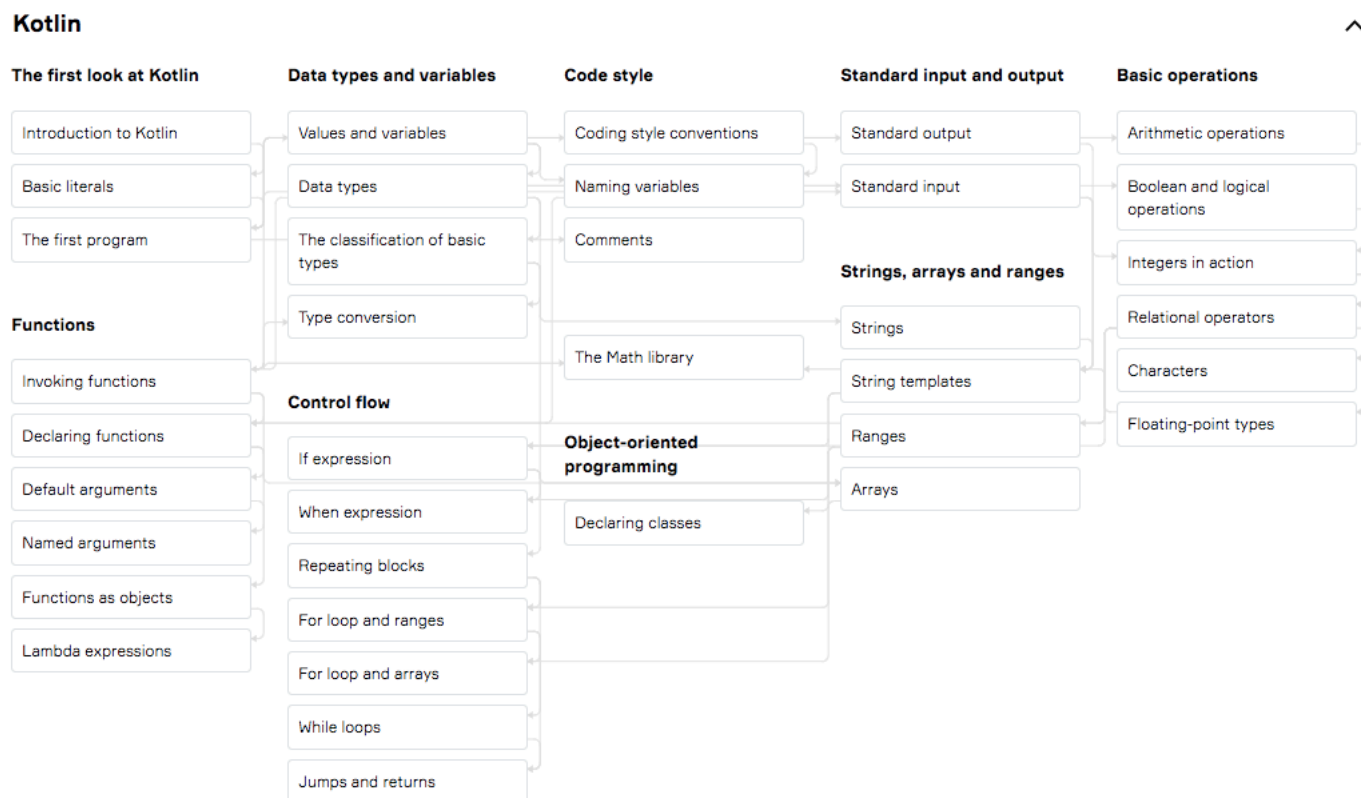
Если говорить об образовании, то в настоящее время существует множество систем управления обучением (Learning Management System), одним из примеров которых в частности является известный каждому студенту нашего университета BlackBoard [2]. Такие системы управления образованием предоставляют своим ученикам возможность удалённого обучения онлайн и контроля своих знаний при помощи интерактивных заданий.

LMS отслеживают различную информацию о своих пользователях: расписание и список курсов, которые посещает студент; оценки и баллы, полученные в результате выполнения электронных заданий. Всё это — колоссальные объёмы данных, которые потенциально можно использовать для улучшения образовательного процесса. Образовательные данные зачастую детализированы и точны, а их объём очень велик. Анализ данных в сфере образования включает в себя исследования, методы и инструменты, предназначенные для автоматического извлечения данных, связанных с учебной деятельностью людей и создаваемых в образовательных учреждениях.

К примеру, анализ данных в LMS может выявить связь между этапами обучения, к которым студент получил доступ в течение курса, и их итоговой оценкой курса. Аналогично, анализ данных стенограммы студента может выявить связь между оценкой студента по какому-то курсу и его решением изменить свою специализацию. Такая информация дает представление о процессе обучения, что позволяет учащимся и учителям принимать обоснованные решения о том, как взаимодействовать, предоставлять и управлять образовательными ресурсами.

Веб-сайт «Hyperskill» [9] представляет собой новую образовательную платформу для изучения языков программирования. В этой платформе структура образовательного материала представлена не линейно, как в других системах (например, «Stepik» [20] или «Coursera» [4]), а в виде графа. Пример такого графа для языка программирования Kotlin [13] изображен ниже на Рис. 1.

Рис. 1: Структура материала в платформе Hyperskill для Kotlin



Вершины этого графа соответствуют основным темам учебного курса (в дальнейшем употребляется термин «граф тем»), а ребра обозначают порядок их изучения. К вершинам графа привязываются различные элементы курса, такие как учебный текст или видеоматериалы, а также тренировочные задания. Выполняя задания в системе Hyperskill, студент может двигаться по ребру от одной вершины к другой или выбрать другую вершину, постепенно охватывая весь граф. Преимуществами платформы являются нелинейность образовательного процесса и возможность видеть взаимосвязь тем между собой.

1. Постановка задачи

Целью данной работы является создание алгоритма автоматической разметки существующего текстового учебного материала платформы «Stepik» для его дальнейшего переиспользования в новой платформе «Hyperskill». Результат разметки будет представлять собой направленный ациклический граф [3], вершины которого соответствуют основным темам учебного курса, а ребра обозначают порядок их изучения. При этом каждому элементу курса будет соответствовать некоторый путь в графе.

Для достижения этой цели были поставлены следующие задачи:

1. Провести анализ данных, предоставленных ООО «ИнтеллиДжей Лабс» [10]
 - Векторизовать решения заданий
 - Выбрать модель построения множества тем для каждого задания с учетом их решений
2. Реализовать алгоритм автоматического построения графа тем [12] по учебному контенту
 - Построение начальной версии графа по предоставленным данным и базе «Wikidata» [23]
 - Построение финальной версии графа по начальной в соответствии с решениями заданий
 - Провести апробацию алгоритма

2. Обзор литературы

В качестве основного источника для изучения методов машинного обучения и анализа данных я использовала книгу Андреаса Мюллера и Сары Гвидо «Введение в машинное обучение с Python» [14]. В этом учебнике рассматриваются как классические, так и современные подходы ко всем аспектам моделирования и внедрения систем сбора и обработки коллекций данных. В частности, в этой книге представлены все алгоритмы классификации, которые используются в данной работе. Все описания алгоритмов машинного обучения в этой книге сопряжены с наглядными иллюстрациями их работы.

Поскольку составной частью задачи является построение графа, для реализации этого было полезно использовать книгу Томаса Кормена «Алгоритмы. Построение и анализ» [3], получившую признание профессорами ведущих университетов за наиболее точное и доступное описание известных компьютерной науке алгоритмов, в том числе и алгоритмов на графах.

Для ознакомления с областью применения методов машинного обучения и интеллектуального анализа данных в образовательной сфере была использована книга «Handbook of Educational Data Mining» [8]. Данная книга содержит исчерпывающий обзор основных методов анализа образовательных данных, а также более двадцати реальных случаев применения этих методов на практике.

В большинстве случаев знания учеников в образовательных учреждениях оцениваются в соответствии с какой-то измерительной шкалой, например от 2 до 5 баллов за экзамен или от 0 до 100 баллов за тест. Авторы работы «Knowledge Spaces and Learning Spaces» [7] дают представление о том, как дать оценку умениям и навыкам студентов иным образом. Основная идея заключается в том, чтобы воспринимать полученные человеком знания как набор освоенных им «концепций», принимая при этом во внимание предыдущие результаты. Знания при этом представляются нелинейно в виде графа, в котором эти темы или концепции расположены с учётом трудности их понимания и взаимо-

связи друг с другом. Так, в вершинах графа находятся концепции, а ребра графа показывают порядок их изучения. Работа представляет большой интерес с теоретической точки зрения, потому что в ней формально описываются способы оценки полученных знаний, основанные на распределении вероятностей правильных ответов на тестовые вопросы. Подобная методика могла быть использована в адаптивных онлайн-курсах, один из которых используется в моей работе.

Вдохновляющим примером разметки учебного материала послужила статья «Toward the Automatic Labeling of Course Questions for Ensuring their Alignment with Learning Outcomes» [22], опубликованная в 2017 году на конференции, посвященной анализу данных в сфере образования. В статье была рассмотрена задача классификации экзаменационных вопросов в соответствии с таксономией Блума [21]. Данная работа интересна тем, что решает задачу разметки контента при помощи метода векторизации TF-IDF, который также будет применяться в ходе моей работы.

В работе «An Automatic Knowledge Graph Creation Framework from Natural Language Text» [11] описывается универсальный подход для расширения графа тем по произвольному тексту на естественном языке. Для этого из текста извлекаются «тройки», элементы которых представляют собой два объекта и отношение между ними. Затем авторами работы вводится векторная метрика оценки схожести таких «троек». Расширения графа новыми темами осуществляется с помощью алгоритма кластеризации таких троек и определения степени соответствия их какому-либо существующему элементу графа.

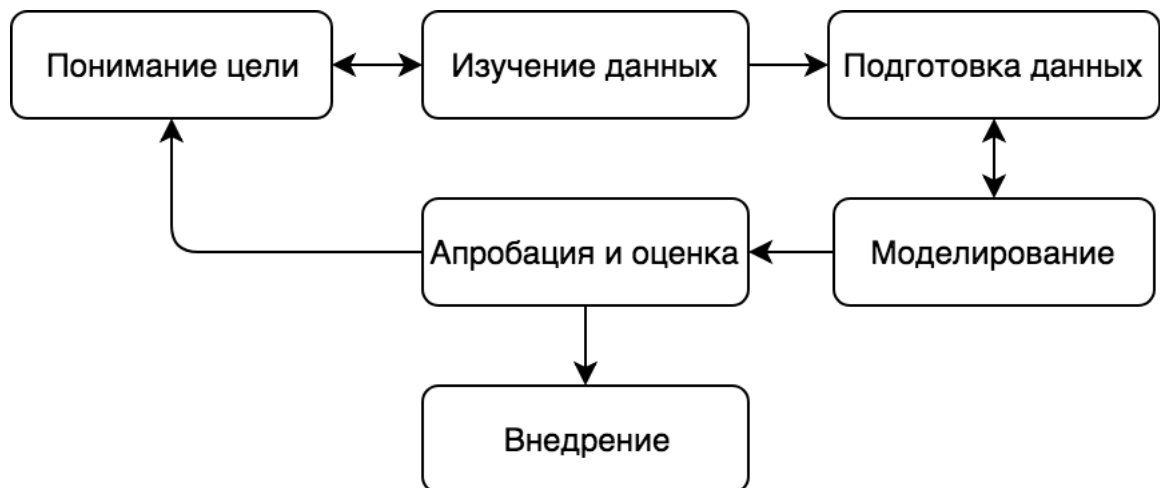
3. Архитектура

Работа над задачей была организована согласно принципам методологии CRISP-DM (CRoss Industry Standard Process for Data Mining) [19] и представляла собой три этапа: сбор и исследование данных, подготовка данных и моделирование. Моделирование включает в себя построение модели машинного обучения для привязки задач к темам и создание графа тем по результатам. Ниже на Рис. 2 представлена схема работы с проектом согласно методологии.

Для реализации поставленной задачи использовался язык программирования Python и библиотека алгоритмов машинного обучения Scikit-Learn [18], а также язык описания графов DOT [5].

Исходный код разделён на три независимых части: обработка данных, моделирование и построение графа. Обработка данных происходит в два этапа: приведение к единому виду и экстракция необходимых для моделирования данных, затем подготовка данных. Для моделирования использовалось несколько алгоритма машинного обучения для сравнения полученных результатов и определения наиболее точного метода.

Рис. 2: Этапы цикла исследования данных согласно CRISP-DM



4. Извлечение и обработка данных

4.1. Формат данных

Каждый курс на платформе Stepik состоит из *степов* — шагов, которые могут представлять собой просмотр видеоматериалов, решение проверочных заданий или работу с текстом. *Адаптивные* курсы подбирают задания индивидуально для каждого пользователя платформы, основываясь на уровне знаний, личной оценке той или иной задачи пользователем, а также затраченном времени. Курс «Adaptive Python» [1], для которого необходимо построить граф тем, представляет собой адаптивный тренажёр для обучению языку программирования Python [15]. Степами этого курса являются только задачи, которые необходимо решать при помощи Python — пример степа изображен ниже.

Рис. 3: Степ (шаг) в курсе «Adaptive Python»

Вы в адаптивном режиме [Следующий рекомендованный урок](#) [Урок слишком сложный](#) [Урок слишком простой](#)

Squirrels and nuts - 1
N squirrels found K nuts and decided to divide them equally. Determine how many nuts each squirrel will get.
Input data format
There are two positive numbers N and K, each of them is not greater than 10000.
Sample Input:
3
14
Sample Output:
4

Программирование — Напишите программу. Тестируется через stdin → stdout

Time Limit: 1 seconds
Memory Limit: 256 MB

```
1 # put your python code here
2
3
4
5
6
```

Отправить [Запустить код](#)

Для извлечения необходимых данных курса «Adaptive Python» было использовано Stepik REST API ¹ — программный интерфейс платформы Stepik, позволяющий автоматически обращаться к её ресурсам. Компанией JetBrains были предоставлены данные о степах и попытках (как удачных, так и неудачных) решения задач этих степов пользователями платформы.

В предоставленном наборе данных каждому степу соответствовал набор *тем* — некоторых «тэгов», описывающих содержательную часть задания (к примеру, «Целые числа» или «ООП»). Помимо этого, для каждой темы был указан её идентификатор в базе знаний «Wikidata» [23]. В данных о попытках решения для каждого степа помимо кода решения на языке Python было также указано, являлось ли решение верным. Данные были представлены в формате рабочей книги Microsoft Excel (.xlsx) ² и текстовом формате представления таблицы Comma-Separated Values (.csv) ³. Для обработки этих форматов были использованы библиотеки `openpyxl` [25] и `pandas` [26].

4.2. Извлечение данных

Полученные данные содержали несущественные для решения поставленной задачи значения и были представлены в разных форматах, поэтому для удобства было принято решение в дальнейшем использовать формат представления JSON ⁴.

Формулировки заданий курса было решено исключить из рассмотрения, поскольку они зачастую слишком абстрактны и от студента курса требуется самостоятельно составить требования к решению. Более полезным для привязки того или иного степа курса к набору тем будет анализ кода решения, полученного платформой от пользователей. Поэтому из предоставленных данных были извлечены попытки решения пользователей того или иного задания, а также существующий список

¹REpresentational State Transfer Application Programming Interface

²формат данных программы для работы с электронными таблицами

³значения колонок таблицы разделены запятой

⁴JavaScript Object Notation — основанный на языке программирования JavaScript формат обмена данными, преимуществом которого является удобство чтения как человеком, так и компьютером

тем, соответствующих заданию.

4.3. Предобработка

Решения пользователей представляют собой тексты кода на языке Python. Для построения модели было решено взять только верные попытки решения задач. Коды этих попыток были протестированы внутренней системой Stepik и удовлетворяют условиям заданий. К тому же, среди неверных попыток оказывается слишком много пустых или синтаксически некорректных блоков кода.

Поскольку степы адаптивного курса представляют собой классические задачи по программированию, коды правильных решений в большинстве своём по структуре оказываются очень похожи. Однако они всё же отличаются названиями переменных, а также наличием или отсутствием в них комментариев, которые могут быть самыми разнообразными. В ходе экспериментов было обнаружено, что различное именование переменных отрицательно влияет на результат моделирования: одну и ту же целочисленную переменную-счетчик разные студенты могут назвать x , a , i , k , n ; переменную для размера массива называют *size*, *length*, *count*.

Поэтому было принято решение предварительно обработать коды таким образом, чтобы очистить их от комментариев, а также переименовать все переменные единым образом. Процесс предобработки состоит из трех этапов: синтаксический анализ, унификация переменных и pretty-printing.

Синтаксический анализ

На первом этапе для исходного кода строится его синтаксическое дерево с помощью парсера из библиотеки ast [24].

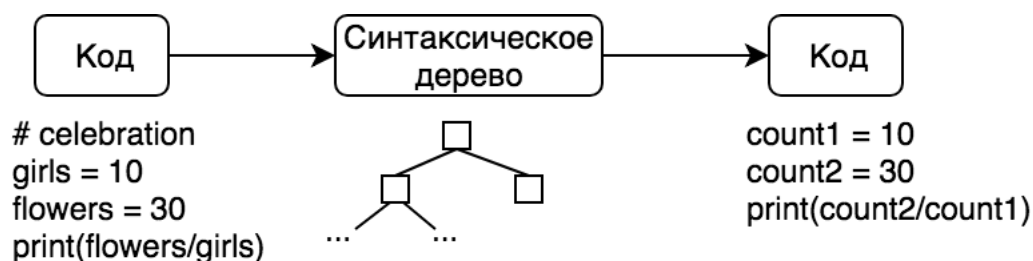
Унификация переменных

Далее для каждого синтаксического дерева запускается процесс унификации переменных. Он состоит из двух рекурсивных обходов дерева, которые представляют собой реализации стандартного для деревьев шаблона проектирования «Посетитель» [6]. Сначала первый посетитель `VariablesCollector` обходит дерево, собирает все имена переменных и назначает каждому имени свой номер. Затем второй посетитель `VariablesRenamer` переименовывает все переменные, назначая им имена i_1 , i_2 и т.д. в соответствии с номером.

Pretty-printing

Pretty-printing [16] — это форматирование исходного кода в соответствии с некоторым заданным стилем. Модели, использованные в работе, оперируют с текстами, а не деревьями, поэтому по итогу предобработки мы должны снова получить текст программы. Как и в случае с разными названиями переменных, модели показывают худшие результаты на текстах с разными стилями оформления (например, с разной табуляцией, отступами и переносами или отсутствием переносов строк). Поэтому каждое синтаксическое дерево печатается в некотором едином стиле. Эта печать реализуется с помощью `pretty-printer` — ещё одного посетителя, который обходит все вершины и печатает текст для каждой из них.

Рис. 4: Процесс обработки кода



5. Построение модели

5.1. Входные данные и необходимые метрики

Для обучения модели классификации по нескольким меткам (Multi-label classification ⁵) извлечённые данные были разделены на тренировочные и тестовые. Для этого данные были перемешаны с помощью генератора случайных чисел, затем 70% данных были отмечены как тренировочные, а 30% будут затем использоваться для проверки обученной модели.

Рассмотрим метрики для определения точности классификации.

Accuracy

Простейшая метрика, представляющая собой долю верно классифицированных объектов среди всего размера обучающей выборки.

$$\text{Accuracy} = \frac{|\{\text{correct}\}|}{|\{\text{train sample}\}|}$$

Precision and recall

Точность классификатора определяется как доля верно классифицированных объектов среди всех объектов, отнесённых классификатором к классу k :

$$\text{Precision}_k = \frac{|\{\text{relevant}\} \cap \{\text{retrieved}\}|}{|\{\text{retrieved}\}|}$$

Точность отвечает на вопрос: сколько из всех выбранных объектов действительно принадлежат классу k ?

Полнота определяется как доля верно классифицированных объектов среди всех объектов класса k в тестовой выборке:

$$\text{Recall}_k = \frac{|\{\text{relevant}\} \cap \{\text{retrieved}\}|}{|\{\text{relevant}\}|}$$

⁵задача классификации, при которой каждому объекту сопоставляется один или несколько классов

Полнота отвечает на вопрос: сколько из всех действительно принадлежащих классу k объектов были выбраны?

F1-score

F1-мера определяется как среднее гармоническое между точностью и полнотой класса k :

$$\text{F1-score}_k = \frac{2 \cdot \text{Precision}_k \cdot \text{Recall}_k}{\text{Precision}_k + \text{Recall}_k}$$

Как видно, F1-мера стремится к нулю, если точность или полнота стремится к нулю. При этом одинаковый вес придаётся обоим этим метрикам.

Общая формула для оценки $F\beta$ -score, которая придаёт больший вес полноте на порядок β :

$$F\beta_k = \frac{(1 + \beta^2) \cdot \text{Precision}_k \cdot \text{Recall}_k}{\beta^2 \cdot \text{Precision}_k + \text{Recall}_k}$$

Support

Метрика поддержки показывает для класса k количество объектов, ему принадлежащих. Данное значение служит коэффициентом для вычисления средних значений Precision и Recall по всем классам.

5.2. Предсказание тем

Для того, чтобы привязать уроки онлайн-курса к темам, за основу был взят принцип «мешка слов» [14], в котором используется векторная модель представления решений задач. Предсказанные на данном этапе темы впоследствии будут соответствовать некоторому пути в графе тем.

TF-IDF представление

Для векторизации решений заданий курса используется TF-IDF мера. Введем необходимые обозначения и определения:

- *Документ* — текстовые данные, состоящие из набора лексем. В нашем случае документом является текстовое представление решения пользователя.
- *TF (term frequency)* — отношение числа вхождения определенной лексемы n_i к общему количеству лексем в документе:

$$TF = \frac{n_i}{\sum_k n_k}$$

- *IDF (inverse document frequency)* — обратная документная частота, определяемая следующим образом:

$$IDF = \log \frac{N}{df}$$

N — количество документов в коллекции

df — количество документов, содержащих n_i .

- *TF-IDF* — мера, определяющая вес каждой лексемы в каждом документе:

$$TF-IDF = TF \cdot IDF$$

Мера *TF-IDF* уменьшается, когда лексема встречается в большом количестве документов, и увеличивается, если лексема встречается часто в небольшом количестве документов.

Задача привязки заданий к темам рассматривается в работе как задача классификации по k категориям. Для этого используются представленные ниже классификаторы, результаты работы которых будут сопоставлены в соответствии с приведенными выше метриками.

One vs Rest

Задача классификации данных (в нашем случае заданий) по k категориям может быть сведена к k задачам бинарной классификации. Для каждого из классов обучается свой классификатор, определяющий принадлежность либо не принадлежность классу. Так работает классификатор One vs Rest. Помимо вычислительной эффективности и простоты реализации, преимуществом данной классификации является легкость интерпретации полученных данных.

K Neighbors

Алгоритм k -ближайших соседей — простейший алгоритм классификации и регрессии для задач обучения с учителем. Классификатор использует метрику (например, стандартную в евклидовом пространстве) для вычисления расстояния между объектами выборки, и определяет класс каждого из объектов как наиболее часто встречающийся среди k его ближайших по метрике соседей. Несмотря на недостатки алгоритма, например необходимость настройки параметров, было принято решение попробовать его как классический и наиболее простой вариант.

Multinomial Naive Bayesian Classifier

Мультиномиальный наивный байесовский классификатор представляет собой вероятностную модель с применением теоремы Байеса в предположении о независимости. Классификатор определяет искомый класс объекта как класс, имеющий максимальную *апостериорную* вероятность. Для избежания потери значимых разрядов при умножении зачастую используют сумму логарифмов вероятностей вследствие монотонности логарифмической функции и равенства $\log(x \cdot y) = \log(x) + \log(y)$.

Decision Tree

Алгоритм решающих деревьев — логический алгоритм классификации, задающийся бинарным деревом, используется во множестве задач регрессии и классификации. Каждому листу дерева соответствует класс, и при классификации того или иного объекта алгоритм проходит весь путь от корня до некоторого листа, определяя тем самым класс. Для построения алгоритма необходимо выбрать признак разделения и минимальный порог (если необходим).

Random Forest

Случайный лес — ансамбль решающих деревьев, эффективно используемый для обработки данных с большой размерностью. Основными преимуществами данного метода машинного обучения являются высокая масштабируемость и параллелизуемость. Один из его недостатков — склонность к переобучению на зашумлённых данных.

6. Построение и расширение графа

6.1. Обозначения

Введем необходимые обозначения и определения:

- \mathcal{T} — множество всех тем онлайн-курса.
- *Урок* — множество тем $\{t_i\}$ одного занятия онлайн-курса.
- \mathcal{L} — множество всех уроков онлайн-курса.
- G — направленный граф тем; его вершинами являются элементы множества \mathcal{T} , а дуги соответствуют порядку изучения тем
- $count(t_i) = |\{l \in \mathcal{L} : t_i \in l\}|$ — количество уроков с этой темой
- $neighbors(t_i)$ — соседние темы t_i

6.2. Построение графа

Сначала необходимо получить первое приближение G_0 искомого графа G . Вершинами обоих графов являются исходные темы онлайн-курса. Для построения дуг используются данные из Wikidata: с помощью запросов на языке SPARQL [17] проверяется, как темы соотносятся между собой. Тема t_1 считается предком темы t_2 , если для них имеет место отношение "instance of", "subclass of" или "part of". Однако полученный граф не отражает действительную структуру курса, поскольку никак не учитывает его уроки: в таком графе нельзя сопоставить уроку какой-либо путь. Вдобавок, этот граф не удовлетворяет условиям ацикличности, поэтому необходимо удалить некоторые дуги для приведения его к ациклическому виду. Таким образом, G_0 не отражает последовательность изучения тем в рамках курса. Например, в контексте Wikidata темы "I/O" и "String" никак не связаны, однако в процессе прохождения курса тема "I/O" должна изучаться строго после темы "String".

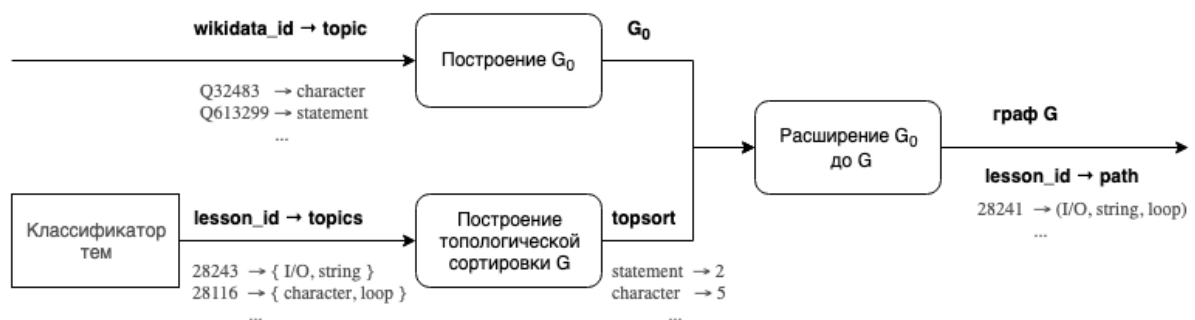
В рамках расширения графа тем новыми дугами предполагается, что каждый урок соответствует некоторому пути в искомом графе, то есть

$$\forall l \in \mathcal{L} \exists (t_1, t_2, \dots, t_n) \in G : l = \{t_1, t_2, \dots, t_n\}$$

Основная идея построения графа состоит в том, чтобы сначала найти топологическую сортировку расширенного графа, а затем по ней построить сам расширенный граф. Использование топологической сортировки гарантирует, что полученный граф будет соответствовать последовательности изучения тем курса.

6.3. Расширение графа

Рис. 5: Схема работы алгоритма



1. Составить словарь $\forall t_i \in \mathcal{T} \{t_i : count(t_i)\}$
2. В каждом уроке отсортировать темы t_i по убыванию $count(t_i)$
3. Составить словарь $\forall t_i \in \mathcal{T} \{t_i : neighbors(t_i)\}$; в дальнейшем считаем, что $neighbors(t_i)$ являются предками t_i
4. Отсортировать уроки по убыванию количества тем в них
5. Пронумеровать уроки уровнями в соответствии с топологической сортировкой: $\forall t_i \in \mathcal{T} : level[t_i] \in [0..N]$

После шага 5 мы знаем уровни вершин (т.е. длины путей от одной из корневых вершин). Осталось добавить в граф G_0 ребра в соот-

ветствии с уровнями и уроками (как упоминалось выше, каждый урок соответствует некоторому пути)

6. Проходим по всем уровням от 0 до последнего: для каждой темы добавляем ребра из его предков, если тема ещё недостижима из предка.

Достижимость вершины проверяется с помощью поиска в глубину или алгоритма Флойда-Уоршалла.

В результате выполнения алгоритма заданный граф G_0 расширяется новыми дугами.

7. Апробация

7.1. Привязка заданий к темам

В ходе апробации было проведено сравнение нескольких перечисленных методов классификации по темам. Для этого использовалась композиция метрик Precision, Recall и F1-score. Поскольку наилучшим образом справился с задачей метод Decision Tree, ниже на Рис. 6 приведена таблица значений метрик по каждому из классов, определенных при помощи данного классификатора.

После того, как каждому заданию курса при помощи классификации был сопоставлен набор тем, данному набору при помощи приведенного выше алгоритма был сопоставлен некоторый путь в графе тем.

7.2. Построение графа тем

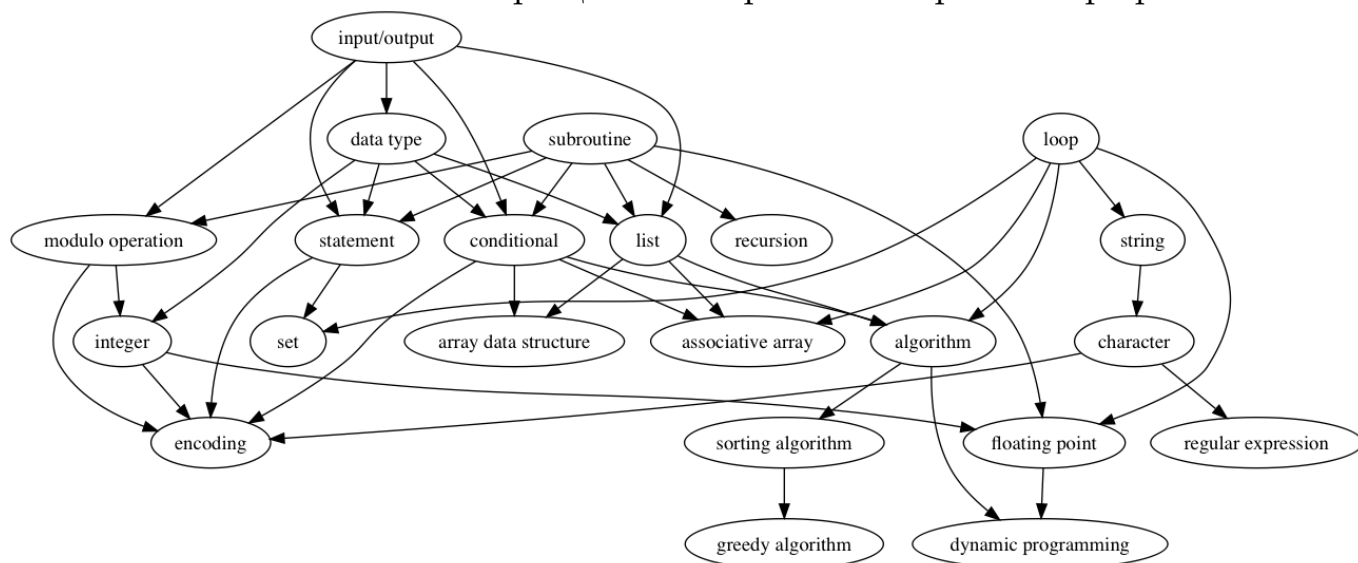
Был построен граф тем, каждой из вершин которого соответствует задание курса «Adaptive Python». Во избежание утечки коммерческих данных и для наглядности исходный граф тем был значительно сокращён. Более полноразмерный вариант графа, а также соответствие заданий и путей в нём, приведены вместе с исходным кодом проекта на ресурсе GitHub ⁶.

⁶<https://github.com/shabatar/knowledge-spaces>

Рис. 6: Результат Decision Tree Classifier для данных.



Рис. 7: Сокращённый вариант построения графа тем



Заключение

В ходе данной работы были получены следующие результаты:

1. Проведён анализ предметной области. Сделан обзор аналогичных существующих решений, используемых в них методов анализа данных и машинного обучения, а также метрик для оценки результатов их работы.
2. Реализован алгоритм многоклассовой классификации извлечённых данных. Он затем используется для привязки задач онлайн-курса к той или иной вершине в графе тем. Были протестированы варианты реализации алгоритма с различными классификаторами. В результате наилучшим образом себя показал метод решающего дерева.
3. Реализован метод автоматического построения графа тем по учебному контенту:
 - Используя результаты работы классификатора на предыдущем шаге, задачи курса привязываются к темам в графе
 - Используя идею топологической сортировки, реализован алгоритм расширения полученного графа новыми вершинами

Список литературы

- [1] Adaptive Python. — <https://stepik.org/course/568/>. — Дата обращения: 07.04.2019.
- [2] Blackboard Learn. — <https://www.blackboard.com/about-us>. — Дата обращения: 07.04.2019.
- [3] Cormen Thomas H. Algorithms Unlocked. — MIT Press, 2013. — ISBN: 978-0-262-51880-2. — Access mode: <http://mitpress.mit.edu/books/algorithms-unlocked>.
- [4] Coursera, learning platform. — <https://coursera.org/>. — Дата обращения: 07.04.2019.
- [5] DOT language. — <http://www.graphviz.org/about/>. — Дата обращения: 07.04.2019.
- [6] Design Patterns: Elements of Reusable Object-oriented Software / Erich Gamma, Richard Helm, Ralph Johnson, John Vlissides. — Boston, MA, USA : Addison-Wesley Longman Publishing Co., Inc., 1995. — ISBN: 0-201-63361-2.
- [7] Doignon Jean-Paul, Falmagne Jean-Claude. Knowledge Spaces and Learning Spaces. — 2015. — 11.
- [8] Handbook of educational data mining / C. Romero, S. Ventura, M. Pechenizkiy, R.S.J. Baker, de. Chapman and Hall/CRC data mining and knowledge discovery series. — CRC Press, 2011. — ISBN: 978-1-4398-045-7-5.
- [9] Hyperskill, learning platform. — <https://hyperskill.org/>. — Дата обращения: 07.04.2019.
- [10] IntelliJ Labs. — <https://www.jetbrains.com/company/>. — Дата обращения: 07.04.2019.

- [11] KERTKEIDKACHORN Natthawut, ICHISE Ryutaro. An Automatic Knowledge Graph Creation Framework from Natural Language Text // IEICE Transactions on Information and Systems. — 2018. — 01. — Vol. E101.D. — P. 90–98.
- [12] Knowledge Graphs: In Theory and Practice / Nitish Aggarwal, Saeedeh Shekarpour, Sumit Bhatia, Amit Sheth. — 2017. — 11.
- [13] Kotlin, programming language. — <https://kotlinlang.org/docs/>. — Дата обращения: 07.04.2019.
- [14] Müller A.C., Guido S. Introduction to Machine Learning with Python: A Guide for Data Scientists. — O'Reilly Media, Incorporated, 2016. — ISBN: 9781449369415. — Access mode: <https://books.google.ru/books?id=qjUVogEACAAJ>.
- [15] Python. — <https://www.python.org/about/>. — Дата обращения: 07.04.2019.
- [16] Rubin L. F. Syntax-Directed Pretty Printing—A First Step Towards a Syntax-Directed Editor // IEEE Transactions on Software Engineering. — 1983. — March. — Vol. SE-9, no. 2. — P. 119–127.
- [17] SPARQL Query Language. — <https://www.w3.org/TR/rdf-sparql-query/>. — Дата обращения: 07.04.2019.
- [18] Scikit-Learn. — <http://scikit-learn.org/>. — Дата обращения: 07.04.2019.
- [19] Shearer C. The CRISP-DM model: the new blueprint for data mining // J Data Warehouse. — 2000. — 01. — Vol. 5. — P. 13–22.
- [20] Stepik, learning platform. — <https://stepik.org/>. — Дата обращения: 07.04.2019.
- [21] Taxonomy of educational objectives. The classification of educational goals. Handbook 1: Cognitive domain / B. S. Bloom, M. B. Engelhart, E. J. Furst et al. — New York : Longmans Green, 1956.

- [22] Supraja S., Hartman Kevin, Tatinati Sivanagaraja, Khong Andy W. H. Toward the Automatic Labeling of Course Questions for Ensuring their Alignment with Learning Outcomes. — 2017. — Access mode: http://educationaldatamining.org/EDM2017/proc_files/papers/paper_44.pdf.
- [23] Wikidata Knowledge Base. — <https://www.wikidata.org/wiki/Wikidata:Introduction>. — Дата обращения: 07.04.2019.
- [24] ast Python library. — <https://docs.python.org/3/library/ast.html>. — Дата обращения: 07.04.2019.
- [25] openpyxl Python library. — <https://openpyxl.readthedocs.io/en/stable/>. — Дата обращения: 07.04.2019.
- [26] pandas Python library. — <https://pandas.pydata.org/about.html>. — Дата обращения: 07.04.2019.