



UNIVERSIDADE FEDERAL DE SANTA CATARINA  
CAMPUS REITOR JOÃO DAVID FERREIRA LIMA  
PROGRAMA DE GRADUAÇÃO EM SISTEMAS DE INFORMAÇÃO

Augusto Silva de Oliveira

**SACI - UMA FERRAMENTA PARA EXTRAÇÃO DE ALGORITMOS  
E SUAS COMPLEXIDADES DE TEMPO E ESPAÇO**

Florianópolis, Santa Catarina – Brasil  
2023



Augusto Silva de Oliveira

**SACI - UMA FERRAMENTA PARA EXTRAÇÃO DE ALGORITMOS  
E SUAS COMPLEXIDADES DE TEMPO E ESPAÇO**

Trabalho de Conclusão de Curso submetido ao Programa de Graduação em Sistemas de Informação da Universidade Federal de Santa Catarina para a obtenção do Grau de Bacharel em Sistemas de Informação.

**Orientador(a):** Carina Friedrich Dorneles, Dr.

**Coorientador(a):** Ricardo José Pfitscher, Dr.

Florianópolis, Santa Catarina – Brasil

2023

Catálogo na fonte pela Biblioteca Universitária da Universidade Federal de Santa Catarina.  
Arquivo compilado às 00:36h do dia 8 de julho de 2024.

Augusto Silva de Oliveira

SACI - Uma ferramenta para extração de algoritmos e suas complexidades de tempo e espaço / Augusto Silva de Oliveira; Orientador(a), Carina Friedrich Dorneles, Dr.; Coorientador(a), Ricardo José Pfitscher, Dr. - Florianópolis, Santa Catarina - Brasil, 05 de Julho de 2024.

148 p.

Trabalho de Conclusão de Curso - Universidade Federal de Santa Catarina, INE - Departamento de Informática e Estatística, CTC - Centro Tecnológico, Programa de Graduação em Sistemas de Informação.

Inclui referências

1. Scraper, 2. Dataset, 3. Complexidade, I. Carina Friedrich Dorneles, Dr. II. Ricardo José Pfitscher, Dr. III. Programa de Graduação em Sistemas de Informação IV. SACI - Uma ferramenta para extração de algoritmos e suas complexidades de tempo e espaço

CDU 02:141:005.7

Augusto Silva de Oliveira

**SACI - UMA FERRAMENTA PARA EXTRAÇÃO DE ALGORITMOS  
E SUAS COMPLEXIDADES DE TEMPO E ESPAÇO**

Este(a) Trabalho de Conclusão de Curso foi julgado adequado(a) para obtenção do Título de Bacharel em Sistemas de Informação, e foi aprovado em sua forma final pelo Programa de Graduação em Sistemas de Informação do INE – Departamento de Informática e Estatística, CTC – Centro Tecnológico da Universidade Federal de Santa Catarina.

Florianópolis, Santa Catarina – Brasil, 05 de Julho de 2024.

---

**Álvaro Junio Pereira Franco, Dr.**

Coordenador(a) do Programa de  
Graduação em Sistemas de Informação

**Banca Examinadora:**

---

**Carina Friedrich Dorneles, Dr.**

Orientador(a)  
Universidade Federal de Santa  
Catarina – UFSC

---

**Ricardo José Pfitscher, Dr.**

Coorientador(a)  
Universidade Federal de Santa  
Catarina – UFSC

---

**Roberto Willrich, Dr.**

Universidade Federal de Santa Catarina –  
UFSC

---

**Alexandre Gonçalves da Silva, Dr.**

Universidade Federal de Santa Catarina –  
UFSC

## **AGRADECIMENTOS**

Primeiramente, gostaria de expressar minha profunda gratidão à minha orientadora, Carina Friedrich Dorneles, por sua orientação, paciência e apoio contínuos ao longo deste trabalho. Sua expertise, conselhos valiosos e incentivo constante foram fundamentais para a realização deste projeto.

Agradeço também ao meu coorientador, Ricardo José Pfitscher, por seu suporte e contribuições significativas. Suas sugestões e feedbacks foram essenciais para o desenvolvimento deste trabalho.

À minha família, expresso minha mais sincera gratidão pelo amor incondicional, paciência e encorajamento durante toda a minha jornada acadêmica.

Aos meus amigos, que estiveram ao meu lado nos momentos de dificuldades e de conquistas, meu muito obrigado. Obrigado por todas as conversas, risadas. Vocês tornaram essa jornada mais leve e divertida.

## RESUMO

Determinar a complexidade de tempo de execução e de espaço em um código escrito é uma tarefa que pode ser complexa e que muitas vezes é essencial, considerando-se que deseja-se obter uma aplicação eficiente. Apesar de ser um aspecto importante na formação do desenvolvedor, a análise de algoritmos quanto a sua complexidade é uma tarefa difícil que por vezes é negligenciada. Algumas ferramentas são capazes de fazer sugestões para o preenchimento automático do código conforme o programador esteja escrevendo, mas, até o momento atual, ela infelizmente não possui capacidade de considerar as complexidades em suas sugestões. Por isso é necessário que o desenvolvedor tenha conhecimento sólido sobre análise de algoritmos para garantir que as soluções propostas sejam eficientes. Tendo isso em vista, esse trabalho propõe a criação de uma ferramenta de extração de dados, capaz de extrair códigos de algoritmos juntos de suas respectivas complexidades a fim de formar uma base de dados. Essa base de dados poderá ser usada para o treinamento de ferramentas, seja através de aprendizado de máquina ou outras técnicas, que visam determinar a complexidade de um trecho de código e auxiliar os desenvolvedores em suas rotinas. Ao fim do trabalho obteve-se uma base de dados com 501 algoritmos distintos acompanhados de suas complexidades, com boa precisão e qualidade dos dados, comprovado pelas métricas utilizadas.

**Palavras-chaves:** Scraper. Dataset. Complexidade.

## LISTA DE FIGURAS

Figura 1	–	Fluxo de atuação da ferramenta . . . . .	27
Figura 2	–	Fluxo de tarefas de extração . . . . .	28
Figura 3	–	Estrutura do SACI . . . . .	29
Figura 4	–	Estrutura das páginas implementadas . . . . .	30
Figura 5	–	Menu Geeks4Geeks . . . . .	39
Figura 6	–	Página de algoritmos . . . . .	39
Figura 7	–	Complexidades formato 1 . . . . .	40
Figura 8	–	Complexidades formato 2 . . . . .	40
Figura 9	–	Complexidades formato 3 . . . . .	41
Figura 10	–	Matriz confusão . . . . .	54
Figura 11	–	Matriz confusão, tempo . . . . .	56
Figura 12	–	Matriz confusão, espaço . . . . .	56
Figura 13	–	Matriz confusão OpenAI, tempo . . . . .	57
Figura 14	–	Matriz confusão OpenAI, espaço . . . . .	57
Figura 15	–	Caracterização dos Nodos . . . . .	59
Figura 16	–	Complexidades extraídas . . . . .	60
Figura 17	–	Distribuição de Algoritmos . . . . .	61
Figura 18	–	Linear Search . . . . .	62
Figura 19	–	Problema na extração . . . . .	63
Figura 20	–	Classe das complexidades . . . . .	64



## LISTA DE TABELAS

Tabela 1	–	Comparação de trabalhos . . . . .	26
Tabela 2	–	Métricas das complexidades de tempo e espaço . . . . .	58

## LISTA DE CÓDIGOS

Código 1	–	Exemplo de arquivo de análises . . . . .	31
Código 2	–	Exemplo de arquivo de dados . . . . .	32
Código 3	–	Esquema dos dados . . . . .	35
Código 4	–	Cabeçalhos utilizados . . . . .	41
Código 5	–	Regex de extração . . . . .	50
Código 6	–	Prompt . . . . .	51
Código 7	–	Regex de classificação de complexidades . . . . .	65

## LISTA DE ALGORITMOS

1	Código mostra a lista de páginas em alto nível . . . . .	31
2	Código mostra o logger em alto nível . . . . .	35
3	Código calcula métricas de desempenho em alto nível . . . . .	35
4	Código mostra a sessão http em alto nível . . . . .	37
5	Código mostra o objeto resposta em alto nível . . . . .	37
6	Código mostra a função principal do arquivo runner em alto nível . . . .	38
7	Código mostra a função principal do arquivo spider em alto nível . . . .	43
8	Código mostra as funções responsáveis por extrair URLs em alto nível .	43
9	Código mostra as funções responsáveis por extrair algoritmos em alto nível . . . . .	44
10	Código mostra a função de extração principal em alto nível . . . . .	46
11	Código mostra a função 'extrairCodigosReferencias' em alto nível . . .	47
12	Função 'extrairCodigo' . . . . .	48
13	Função 'extrairComentarios' . . . . .	48
14	Código mostra a função 'extrairComplexidadeDeReferencia' em alto nível	49
15	Código mostra a função 'fallback' em alto nível . . . . .	49
16	Código mostra a sessão http em alto nível . . . . .	51
17	Código mostra a função 'extrairComplexidadeDeReferencia' em alto nível	52

## SUMÁRIO

<b>1</b>	<b>INTRODUÇÃO</b>	<b>12</b>
1.1	MOTIVAÇÃO	13
1.2	OBJETIVOS GERAL E ESPECÍFICOS	14
1.3	ESTRUTURA DO TRABALHO	14
<b>2</b>	<b>CONCEITOS FUNDAMENTAIS</b>	<b>15</b>
2.1	WEB SCRAPING	15
2.2	COMPLEXIDADE DE ALGORITMOS	16
<b>2.2.1</b>	<b>Tempo de execução polinomial</b>	<b>17</b>
<b>2.2.2</b>	<b>Análise assintótica</b>	<b>18</b>
2.3	EXPRESSÕES REGULARES	19
2.4	MODELOS DE LINGUAGEM DE GRANDE ESCALA	20
2.5	CONSIDERAÇÕES FINAIS	21
<b>3</b>	<b>TRABALHOS RELACIONADOS</b>	<b>22</b>
3.1	LEARNING BASED METHODS FOR CODE RUNTIME COMPLEXITY PREDICTION	22
3.2	QFEX - UM CRAWLER PARA BUSCA E EXTRAÇÃO DE QUESTIONÁRIOS DE PESQUISA EM DOCUMENTOS HTML	22
3.3	FERRAMENTA PARA COLETA E COMPARAÇÃO DE DADOS DE PUBLICAÇÕES ACADÊMICAS DOS PROFESSORES COM O CURRÍCULO LATTES	23
3.4	CRAWLEX: UMA FERRAMENTA PARA EXTRAÇÃO DE DADOS NA WEB CONFIGURÁVEL ATRAVÉS DE EXEMPLOS	24
3.5	DESENVOLVIMENTO DE API PARA PREDIÇÃO DA COMPLEXIDADE DE TEMPO DE EXECUÇÃO DE CÓDIGOS POR MEIO DE AUTOML	24
3.6	ANÁLISE DOS TRABALHOS	25
<b>4</b>	<b>SACI - SCRAPER DE ALGORITMOS E COMPLEXIDADES</b>	<b>27</b>
4.1	VISÃO GERAL	27
<b>4.1.1</b>	<b>Fluxo do SACI</b>	<b>27</b>
<b>4.1.2</b>	<b>Fluxo de extração</b>	<b>28</b>
4.2	IMPLEMENTAÇÃO	28
4.3	ESTRUTURA GERAL	29
<b>4.3.1</b>	<b>Spiders individuais</b>	<b>30</b>
<b>4.3.2</b>	<b>Results</b>	<b>31</b>
<b>4.3.3</b>	<b>Observability</b>	<b>34</b>

---

4.3.4	<b>Schema</b> . . . . .	<b>35</b>
4.3.5	<b>Session</b> . . . . .	<b>36</b>
4.3.6	<b>Runner</b> . . . . .	<b>38</b>
4.4	GEEKS4GEEKS . . . . .	38
4.4.1	<b>Settings</b> . . . . .	<b>41</b>
4.4.2	<b>Scraper</b> . . . . .	<b>42</b>
4.5	EXTRACT . . . . .	44
4.5.1	<b>Extração com expressões regulares</b> . . . . .	<b>44</b>
4.5.2	<b>Extração com OpenAI</b> . . . . .	<b>50</b>
<b>5</b>	<b>EXPERIMENTOS</b> . . . . .	<b>53</b>
5.1	CONFIGURAÇÃO DO AMBIENTE . . . . .	53
5.2	MÉTRICAS . . . . .	53
5.2.1	<b>Histogramas</b> . . . . .	<b>53</b>
5.2.2	<b>Matriz confusão</b> . . . . .	<b>54</b>
5.2.3	<b>Acurácia</b> . . . . .	<b>54</b>
5.2.4	<b>Revocação</b> . . . . .	<b>55</b>
5.2.5	<b>Precisão</b> . . . . .	<b>55</b>
5.2.6	<b>F1-Score</b> . . . . .	<b>56</b>
5.3	RESULTADOS . . . . .	56
5.4	FONTE DE DADOS . . . . .	59
5.4.1	<b>Conjunto de dados extraídos</b> . . . . .	<b>59</b>
<b>6</b>	<b>CONCLUSÕES E TRABALHOS FUTUROS</b> . . . . .	<b>66</b>
	<b>REFERÊNCIAS</b> . . . . .	<b>69</b>
	<b>APÊNDICES</b> . . . . .	<b>71</b>

## 1 INTRODUÇÃO

A complexidade de algoritmos mede o consumo de recursos computacionais, como tempo de execução e o espaço ocupado em memória, necessários para execução de um programa ou algoritmo em relação ao tamanho da entrada; ou seja, é a medida de eficiência desse programa (CORMEN *et al.*, 2009). Trata-se de um tema muito importante durante o ciclo de desenvolvimento de uma aplicação porque está diretamente atrelado à forma como será seu desempenho. Uma vez conhecendo a complexidade de certo trecho de código, pode-se prever sua eficiência e otimizá-lo (SEDGEWICK, 2009); isso garante o melhor uso possível dos recursos computacionais, visto que os mesmos são finitos. Infelizmente, por falta de hábito, prática ou conhecimento, a tarefa de identificar a complexidade acaba sendo relegada ou pode tornar-se uma tarefa difícil para o programador, embora existam diversas variáveis que norteiam essa identificação, como tempo de execução, tamanho de entrada, espaço ocupado na memória, análises de melhor e pior caso, dentre outras (HALL; CHAPMAN, 2015). Considerar tudo isso pode ser demasiado complexo; então, por que não usar uma ferramenta que auxilie nessa tarefa?

A carência de dados nesse âmbito pode ser corroborada com uma breve pesquisa sobre o tema 'base de dados de algoritmos e complexidades', a qual retorna poucos resultados. Sikka *et al.* (SIKKA *et al.*, 2020a) trabalham com um banco de dados restrito, composto por 931 códigos em Java. O baixo número de amostras utilizado prejudica não somente o desempenho dos métodos de classificação, mas principalmente seu poder de generalização; além disso, o banco de dados não contempla algoritmos de ordem superior a  $O_n^2$ . A baixa amostragem de classes ineficientes dificulta a separação binária de eficiência do classificador, que é mais relevante do que a separação entre várias classes eficientes. Pfitscher *et al.* (PFITSCHER *et al.*, 2023) ampliam o total de códigos do banco de dados para 1325 amostras a partir de raspagem de dados por web scraping do site Geeks for Geeks <sup>1</sup>, que contém códigos em Java acompanhados de sua complexidade. O SACI busca ampliar ainda mais as bases de dados e amostras de algoritmos disponíveis para a comunidade científica e desenvolvedores, criando um repositório robusto e abrangente.

A ferramenta proposta neste trabalho é construída para extrair dados sobre fontes pré-selecionadas, das quais os algoritmos são extraídos junto de suas respectivas complexidades. O resultado da extração é um *dataset* composto de algoritmos e suas complexidades que é disponibilizado de forma aberta em plataforma online, como o Kaggle e o Github, onde esses dados podem vir a servir de base para análises e trabalhos futuros. Dentre os possíveis trabalhos pode-se citar uma ferramenta que, dado um determinado código, ela retorne sua respectiva complexidade. Hoje em dia

---

<sup>1</sup> <https://www.geeksforgeeks.org/>

já existem ferramentas que fazem tarefas similares a essa, como é o caso do GitHub Copilot<sup>2</sup> que sugere trechos de código em tempo real durante a escrita de código; o qual, além de seu atual funcionamento, também poderia sugerir novos trechos com melhoras em relação à complexidade e ao desempenho. Assim sendo, o trabalho aqui descrito visa construir uma ferramenta chamada SACI, capaz de extrair as informações necessárias e gerar uma base de dados de algoritmos e suas complexidades a partir do site Geeks for Geeks<sup>3</sup>. O SACI é desenvolvido na filosofia de *web scraping* (ZHAO, 2017), uma técnica para extrair informações da internet e salvá-las em um formato estruturado para análises futuras.

Um exemplo interessante de implementação prática é apresentado por Rodenbusch (RODENBUSCH, 2023), onde foi desenvolvida uma API especificamente para prever padrões em dados complexos, chamada RTCC (*Running Time Complexity Calculator*). Em seu trabalho, foram utilizados os resultados gerados pelo SACI, além de outras bases de dados, para criar uma nova base que, combinada com ferramentas de aprendizado de máquina para o treinamento de modelos, resultou em uma capacidade de predição com 79,05% de acurácia para classe de complexidade e 88,02% para classe de eficiência. Uma ferramenta capaz de inferir tais complexidades necessita ser provida de uma inteligência e, nesse caso, uma Inteligência Artificial (IA) poderia ser treinada com essa finalidade; mas, para tanto, são necessários dados, ou seja, algoritmos e complexidades. Os dados são parte crucial do processo de treinamento, pois fornecem as informações necessárias para que a IA possa aprender e melhorar continuamente (WANG *et al.*, 2022).

## 1.1 MOTIVAÇÃO

Atualmente, não existe no mercado uma ferramenta amplamente estabelecida a ponto de ser considerada um padrão da indústria para tal fim e que compartilhe todas estas variáveis de identificação; somente algumas que geram sugestões de trechos de código sem considerar a complexidade.

A criação de uma ferramenta como o SACI, que extrai algoritmos e suas complexidades, pode preencher essa lacuna. Além de facilitar a identificação e análise das complexidades, a base de dados resultante da ferramenta contém informações valiosas para a comunidade de pesquisadores e desenvolvedores. Tal base, pode ser usada para construção de modelos preditivos que estimem a complexidade de algoritmos, pode ajudar a otimizar algoritmos já existentes ou criar novos, também pode ser utilizada para realizar um análise comparativa dos algoritmos e suas complexidade e para futuros estudos teóricos.

---

<sup>2</sup> <https://github.com/features/copilot>

<sup>3</sup> <https://www.geeksforgeeks.org/>

## 1.2 OBJETIVOS GERAL E ESPECÍFICOS

Desenvolver uma ferramenta de *scraping* capaz de extrair algoritmos e suas complexidades, com o objetivo de gerar uma base de dados com o conteúdo extraído.

Objetivos específicos:

1. Realizar a extração com expressões regulares
2. Realizar a extração com OpenAI
3. Tratar os dados extraídos e padronizá-los
4. Caracterizar os dados extraídos

## 1.3 ESTRUTURA DO TRABALHO

Este trabalho está organizado como segue. Na Seção 2 são descritos alguns conceitos importantes utilizados no trabalho. Na Seção 3 são apresentados alguns trabalhos relacionados, envolvendo o conceito de *scraper*. Na Seção 4 é apresentado o Saci, o *web Scraper* desenvolvido para extrair algoritmos e suas complexidades. Na Seção 5 são descritos as métricas e os experimentos realizados. Finalmente, na Seção 6 são apresentadas algumas conclusões e os trabalhos futuros.



## 2 CONCEITOS FUNDAMENTAIS

Durante a elaboração, deste trabalho, foram utilizados alguns conceitos que são definidos neste capítulo: tais como os temas sobre *web scraping*, que é a base para a construção da ferramenta. Complexidade de algoritmos, um dos principais atributos que extraídos das páginas. Modelos de linguagem em grande escala que foram utilizados duran o processo de *scraping* e, por fim, expressões regulares que foram muito utilizadas para definir padrões e ajudar na extração.

### 2.1 WEB SCRAPING

A web é uma fonte valiosa de informação para diversos setores. Nela, encontramos muita informação útil e, ao mesmo tempo, muita coisa dispensável (DIOUF *et al.*, 2019a). Portanto, é preciso uma forma de coletar e extrair utilidade desses dados. Nesse sentido, a técnica de *web scraping*, largamente conhecida como uma maneira útil de coleta de *big data* (MOONEY; WESTREICH; EL-SAYED, 2015), e comumente descrita como o ato de extrair dados da internet e salvá-los em um arquivo ou base de dados para futuro uso e/ou análises, certamente ajudará nessa tarefa.

Há diversas maneiras de realizar essas extrações. Esse processo pode até ser feito manualmente por uma pessoa, mas isso costuma não ser eficiente porque demanda muito tempo. Hoje existem diversas ferramentas customizáveis para tornar o processo cem por cento automatizado, sendo capazes de simular a interação de um humano com o conteúdo (ZHAO, 2017). Tal simulação de comportamento humano costuma ser necessária em alguns casos, visto que diversas páginas costumam bloquear o acesso quando percebem que o usuário está agindo de forma estranha, identificando-o como um *scraper*. Técnicas como *fingerprinting* e análise de comportamento são comumente utilizadas para diferenciar pessoas reais de *scrapers*; algumas medidas também podem ser tomadas para dificultar a vida dos *scrapers*, como páginas com conteúdo dinâmico e o uso de *captcha*.

O processo de *web scraping* costuma ser dividido em duas etapas: a coleta dos dados puros e não estruturados e o tratamento e extração de informação interessante a partir desses dados; como, por exemplo, podemos fazer uma requisição HTTP a alguma página e depois navegar pelos elementos do DOM (*Document Object Model*, representação dos objetos que formam um documento HTML.) extraindo informações relevantes para depois tratar os dados e salvá-los em algum banco de dados.

É possível citar uma grande variedade de cenários em que o *web scraping* é utilizado, desde aplicações simples até uma análise de preço de ações automatizada até um *scraper* de mensagens em redes sociais. Muitas aplicações que fazem uso de algum modelo de IA contam com uma etapa de extração de dados. Nessa etapa, em alguns casos, temos o uso de *scraping*. Se analisarmos o contexto da web em alto nível,

é bem provável que boa parte dos metadados das páginas web sejam coletados por *scrapers*, por exemplo, quase todas as páginas web estão sofrendo ação de *scrapers* constantemente para construção dos motores de busca, como o Google (SNYDER, 2003).

Existem diversas abordagens que podem ser utilizadas para construir um scraper; abaixo estão algumas abordagens citadas na obra de (DIOUF *et al.*, 2019b):

1. **Mimetismo:** Funciona através de um conjunto de regras pré-definidas e customizáveis para o *scraper*. A localização dos dados a serem coletados é pré configurada no *scraper* e as ações acontecem normalmente através de seletores DOM (expressões usadas para manipular e/ou selecionar os objetos que formam um documento HTML.). É uma estratégia relativamente eficiente, mas não é recomendada para ser usada em diversos sites diferentes, não lida bem com heterogeneidade.
2. **Medidas de peso:** Baseada em um algoritmo genérico que analisa a árvore DOM da página e define um peso para cada palavra em cada galho. Através de dedução o algoritmo consegue extrair o texto de todos os nodos. A principal vantagem desse método é que não precisa de nenhuma especificidade, ou seja, se adapta a múltiplas páginas, mas os dados resultantes costumam ter muito ruído.
3. **Diferencial:** Baseia-se no argumento de que duas páginas do mesmo site só irão divergir em conteúdo do corpo ou em inglês *body*. Seguindo essa lógica, barras de menu, colunas e outros elementos seriam idênticos entre páginas, assim só o conteúdo que fosse diferente seria coletado. Funciona bem quando o site é estruturado em um padrões, por exemplo, utilizando *templates*, ou seja, modelos pré-definidos que definem a disposição dos elementos na página, como cabeçalhos, rodapés, barras de navegação e colunas. Em sites de notícias, por exemplo, esses *templates* garantem uma uniformidade visual e funcional, permitindo que apenas o conteúdo específico de cada página varie. Dessa forma, a coleta de dados pode focar diretamente nas áreas de conteúdo que mudam, otimizando a extração de informações e aumentando a precisão do processo.
4. **Aprendizado de máquina:** Tem como princípio treinar um algoritmo utilizando uma base com uma grande quantidade de páginas web analisadas à mão. Assim, após treinamento o algoritmo seria capaz de identificar onde está o conteúdo desejado e extraí-lo.

## 2.2 COMPLEXIDADE DE ALGORITMOS

Temas como algoritmos e estrutura de dados são frequentemente relegados por alguns programadores, seja por desinformação ou por acharem que não os afeta

diretamente. Enquanto isso, debates de qual linguagem é mais rápida se tornou uma pauta altamente discutida, mas totalmente sem importância se a construção do código for feita alheia à noção de algoritmos.

Mas o que são algoritmos? Algoritmos podem ser definidos como ferramentas ou um conjunto de regras que resolvem algum problema, seja esse problema calcular uma rota mais curta, organizar uma lista de nomes em ordem alfabética, filtrar algum histórico etc. (CORMEN *et al.*, 2001). Existem muitos problemas que os algoritmos nos ajudam a resolver. Para isso, é preciso saber qual técnica utilizar, qual categoria o problema se encaixa. Alguns exemplos de categoria de algoritmos: algoritmo de busca, algoritmo recursivo, algoritmo de ordenação, algoritmo guloso etc.

A complexidade de algoritmos busca analisar a eficiência de um algoritmo ignorando a linguagem ou hardware em que ele está implementado; comparar baseado em tempo não é uma boa maneira de medir isso. Por exemplo, é possível que um código em JavaScript execute mais rápido do que em Python, isso ocorre pelas diferenças estruturais entre as linguagens de programação. Tal fato perde importância se o algoritmo selecionado para uso for ruim em termos de eficiência. Portanto, a complexidade de algoritmos é um aspecto crucial no estudo e construção de códigos bem estruturados e computacionalmente eficientes; sendo uma medida do desempenho de uma solução para um determinado problema (SIKKA *et al.*, 2020a).

Em função disso é que foram definidos os conceitos de complexidade de tempo e espaço; um algoritmo é mais eficiente que outro se apresentar uma menor complexidade de tempo e espaço. Para determinar estes dois fatores precisamos entender alguns conceitos como análise do código baseado na velocidade que ele executa suas computações, contagem de instruções fundamentais que o código executa, análise de pior caso, análise de comportamento assintótico, entre outros.

### 2.2.1 Tempo de execução polinomial

A partir da análise de algoritmos discretos e discussões em meados de 1960 surgiu um conceito sobre como quantificar um tempo de execução 'razoável'. Para certos problemas um bom algoritmo é aquele que possui uma propriedade de escalonamento melhor, que significa que conforme o tamanho de entrada cresce com base em um fator constante, o algoritmo desacelera conforme outro fator constante. Pode-se formular essa propriedade de escalonamento da seguinte maneira: suponha que em um algoritmo, para cada entrada de tamanho  $N$ , seu tempo de execução seja limitado por  $cN^d$  etapas computacionais, onde  $c$  e  $d$  são constantes positivas. Em outras palavras, o tempo de execução é, no máximo, proporcional a  $N^d$ . Esse comportamento é conhecido como tempo de execução limitado a uma função polinomial. Assim, a definição proposta é que um algoritmo é eficiente se tiver um tempo de execução polinomial (KLEINBERG; TARDOS, 2006).

No entanto essa definição pode ser muito específica, nem todo algoritmo com tempo de execução proporcional a um polinômio é considerado eficiente. Além disso, um algoritmo com tempo de execução não polinomial, como  $O2^n$ , pode ser aceitável dependendo do caso em que está aplicado. A definição de tempo polinomial é comumente mais usada porque há uma grande diferença entre o crescimento de funções polinomiais e exponenciais. Isso ajuda a dizer se existem algoritmos eficientes para certos problemas, avaliar o desempenho dos algoritmos e a viabilidade dos problemas de forma mais objetiva e aceita por todos. Além disso, permite investigar cientificamente se existem ou não algoritmos eficientes (KLEINBERG; TARDOS, 2006).

### 2.2.2 Análise assintótica

A análise assintótica trata da análise de desempenho de algoritmos à medida que o tamanho da entrada aumenta; em vez de se concentrar em medidas exatas de tempo de execução, ela se preocupa com o comportamento geral do algoritmo à medida que o tamanho da entrada tende ao infinito. Ela é baseada na ideia de que, para tamanhos de entrada suficientemente grandes, o fator dominante que determina o desempenho do algoritmo é a taxa de crescimento da função que descreve o tempo de execução em relação ao tamanho da entrada; essa taxa de crescimento é chamada de ordem de crescimento assintótica.

A ordem de crescimento assintótica é geralmente expressa usando a notação Big O; por exemplo, se um algoritmo tem uma ordem de crescimento assintótica de  $O(n)$ , isso significa que o tempo de execução do algoritmo cresce proporcionalmente ao tamanho da entrada. Se a ordem de crescimento assintótica for  $O(n^2)$ , o tempo de execução crescerá quadraticamente com o tamanho da entrada, e assim por diante. Isso nos permite fazer comparações entre diferentes algoritmos; algoritmos com ordens de crescimento assintóticas menores são considerados mais eficientes, pois seu tempo de execução cresce mais lentamente à medida que o tamanho da entrada aumenta. Existem diversas classes de complexidades, abaixo estão alguns exemplos das complexidades mais utilizadas usando a notação Big O:

1. Constante:  $O(1)$  - significa que o tempo de execução do algoritmo não aumenta à medida que o tamanho da entrada aumenta. Algoritmos com complexidade de tempo constante são eficientes, pois executam em tempo constante, independentemente do tamanho da entrada.
2. Linear:  $O(n)$  - é quando, o tempo de execução do algoritmo aumenta de forma proporcional ao tamanho da entrada. Algoritmos com complexidade de tempo linear têm um desempenho que aumenta de maneira proporcional ao tamanho da entrada. À medida que o tamanho da entrada aumenta, o tempo de execução do algoritmo também aumenta na mesma proporção.

3. Exponencial:  $O(2^n)$  - neste caso, o tempo de execução do algoritmo cresce exponencialmente com o tamanho da entrada. Algoritmos com complexidade de tempo exponencial se tornam rapidamente inviáveis para entradas grandes, pois o tempo de execução aumenta drasticamente à medida que o tamanho da entrada aumenta exponencialmente.
4. Polinomial:  $O(n^k)$  - onde  $n$  é o tamanho da entrada e  $k$  é uma constante positiva. Algoritmos com complexidade de tempo polinomial têm um desempenho aceitável para entradas moderadamente grandes, pois o tempo de execução aumenta de acordo com um polinômio do tamanho da entrada. No entanto, é importante notar que a mesma ainda pode se tornar impraticável para entradas muito grandes, especialmente para valores grandes de  $k$ .
5. Fatorial:  $O(n!)$  - significa que o tempo de execução do algoritmo cresce de forma fatorial em relação ao tamanho da entrada. Algoritmos com complexidade fatorial são considerados inviáveis pois o tempo de execução aumenta extremamente rápido à medida que o tamanho da entrada aumenta.
6. Logarítmica:  $O(\log n)$  - aquela em, que o tempo de execução do algoritmo cresce de forma logarítmica em relação ao tamanho da entrada. Algoritmos com complexidade de tempo logarítmica são considerados muito eficientes, pois o tempo de execução aumenta de forma relativamente lenta à medida que o tamanho da entrada aumenta, pode-se afirmar que é uma das melhores complexidades em termos de eficiência e escalabilidade para problemas computacionais.

## 2.3 EXPRESSÕES REGULARES

A grande maioria dos programadores já se deparou com algo chamado expressões regulares em algum momento de sua vida, onde precisou filtrar padrões em um bloco de texto. Expressão regular, popularmente conhecida como *regex*, é um método conciso e flexível de se especificar um padrão e identificar uma cadeia de caracteres na qual se tem interesse (SERVIAN, 2020). Para definir uma expressão não é necessário listar todos os elementos do conjunto; existem algumas operações que ajudam nessa tarefa e dentre elas citam-se:

- Alternância: Representada por uma barra vertical '|' possibilita-nos separar duas alternativas; por exemplo, 'cara|casa' identifica cara ou casa.
- Agrupamento: Representado por parênteses, define a ordem de precedência e o escopo; utilizando o exemplo anterior, 'ca(r|s)a' identificaria o mesmo padrão.
- Quantificação: Quantificadores são caracteres que especificam quantas vezes o elemento anterior pode acontecer ou repetir. Representados por 3 símbolos, o '?'

faz com que o elemento ocorra zero ou uma vez, o ‘\*’ faz com que o elemento ocorra zero ou mais vezes, o ‘+’ faz com que o elemento ocorra uma ou mais vezes.

As expressões regulares constituem-se um poderoso meio de buscar, substituir e validar padrões de texto. Elas são amplamente suportadas em diversas linguagens de programação e aplicadas em uma variedade de tarefas, desde validação de formulários até análise de dados textuais complexos. Dominar o uso das expressões regulares pode aumentar a eficiência e a flexibilidade do processo de desenvolvimento de software.

## 2.4 MODELOS DE LINGUAGEM DE GRANDE ESCALA

Os Modelos de Linguagem de Grande Escala (LLMs, do inglês ‘Large Language Models’), como o GPT-4, representam o estado da arte no campo da Inteligência Artificial para o processamento de linguagem natural. Estes modelos são treinados em vastos conjuntos de dados textuais através de técnicas de *deep learning* (KERNER, 2023) e são capazes de realizar uma ampla gama de tarefas relacionadas à linguagem, desde responder perguntas e compor textos até tradução de idiomas e geração de código, entre outras tarefas como:

- Compreensão de leitura
- Geração de texto
- Tradução de idiomas
- Sumarização
- Resposta a perguntas
- Análise de sentimentos
- Autocompletar texto

Os LLMs (Modelos de Linguagem de Grande Escala) são baseados em redes neurais chamadas de transformadores, que contém bilhões de parâmetros treinados a partir de grandes bases de dados textuais. Essas redes conseguem analisar toda a sequência do texto, aprendendo padrões complexos e relações contextuais dentro dos dados, resultando em respostas mais precisas e relevantes. Além disso, eles são capazes de usar ‘transfer learning’, o que significa que podem aplicar conhecimento de uma tarefa para realizar outras tarefas de linguagem, mesmo com poucos ou nenhum dado adicional específico para essas novas tarefas. No entanto, por serem tão poderosos, esses modelos também levantam questões éticas e desafios, como a possibilidade

de gerar desinformação, a necessidade de minimizar viés e garantir a equidade, e o desafio de gerenciar direitos autorais e a originalidade do conteúdo gerado (ZHAO *et al.*, 2023).

Todas essas funcionalidades tornam os LLMs ferramentas muito úteis para integrar em aplicações. Em Python, uma biblioteca chamada Langchain facilita essa integração, atuando como um intermediário que simplifica a interação entre os modelos de linguagem e o software do desenvolvedor. Com a API do OpenAI integrada ao Langchain, os desenvolvedores podem criar aplicações que utilizam o poder do GPT-4 para entender e responder a consultas, realizar tarefas automatizadas e até gerar conteúdo dinâmico. Usar LLMs em um projeto não só amplia as funcionalidades de uma aplicação, mas também abre novas possibilidades para interfaces de usuário mais naturais e interativas, análises de dados mais aprofundadas e automação de processos relacionados à linguagem.

## 2.5 CONSIDERAÇÕES FINAIS

Neste capítulo, foram abordados conceitos fundamentais que sustentam a construção da ferramenta proposta neste trabalho. Abordaram-se as técnicas de *web scraping*, fundamentais para a coleta automatizada de dados da web. A importância da complexidade de algoritmos, apresentando como a análise assintótica e a compreensão dos tempos de execução polinomial são essenciais para a eficiência computacional. Discutiu-se a utilização de expressões regulares e seu papel na identificação de padrões textuais. Por fim, examinou-se os Modelos de Linguagem de Grande Escala e suas possíveis aplicações.

Esses conceitos fornecem uma base sólida para o desenvolvimento e implementação das técnicas e ferramentas descritas nos próximos capítulos. A compreensão profunda desses temas é crucial para a análise dos trabalhos relacionados, que será discutida no próximo capítulo. Nele, revisaremos estudos e pesquisas relevantes, contextualizando e comparando-os com a proposta deste trabalho.

### 3 TRABALHOS RELACIONADOS

Na elaboração da ferramenta de extração, objeto desse estudo, considerou-se as abordagens de quatro obras de desenvolvimento de ferramentas que, embora não realizem a extração das complexidades, possuem relação com a que está sendo elaborada, visto que as mesmas também têm um objetivo, características específicas e dificuldades inerentes. Assim sendo, detalhar resumidamente em que consiste cada uma dessas quatro ferramentas e traçar um comparativo de suas características com as da ferramenta apresentada neste estudo.

#### 3.1 LEARNING BASED METHODS FOR CODE RUNTIME COMPLEXITY PREDICTION

Este trabalho de (SIKKA *et al.*, 2020b) aborda a predição da complexidade de tempo de execução de um código de programação. Os autores destacam que prever a complexidade de tempo, até mesmo para humanos, é uma tarefa difícil, exigindo uma análise sutil e amplo conhecimento de algoritmos; por isso propõem abordar esse problema como uma tarefa de aprendizado de máquina e verificam sua viabilidade por meio de análises detalhadas. Como não há conjuntos de dados de código aberto disponíveis para essa tarefa, os autores propõem seu próprio conjunto de dados chamado 'CoRCoD: Code Runtime Complexity Dataset', que foi extraído de plataformas de programação online.

Os resultados apresentaram problemas de acurácia devido à base de dados criada ser pequena em relação ao tamanho considerado padrão atualmente; o modelo treinado teve dificuldades em prever as complexidades devido à falta de dados e variedade, um problema comum para algoritmos de *deep learning* que tendem a ter melhores resultados com mais dados.

Soluções como essa, potencialmente aplicam-se em avaliação automática de tarefas de programação, ferramentas integradas de análise de código estático, dentre outros. Este trabalho contribui com a disponibilização de um novo conjunto de dados anotado, bem como com a aplicação de técnicas de aprendizado de máquina para a previsão da complexidade de tempo de execução de código, proporcionando um avanço nesse campo.

#### 3.2 QFEX - UM CRAWLER PARA BUSCA E EXTRAÇÃO DE QUESTIONÁRIOS DE PESQUISA EM DOCUMENTOS HTML

De acordo com Mathias e Dorneles (2021), o (MATHIAS; DORNELES, 2021) é um *web crawler* construído para coletar da internet questionários de pesquisa e extrair as questões a fim de formar uma base de dados. Um dos desafios encontrados na



confeção desse *crawler* é um problema muito comum: a grande variedade de layouts e formas que as páginas web são desenhadas; logo, esse trabalho focou apenas na extração em páginas com conteúdo estático. Para lidar com a diversidade das páginas foi utilizada duas heurísticas; uma para a identificação dos questionários e outra para a extração, focando-se em procurar por palavras relacionadas como *survey* ou 'questionário', elementos HTML de formulário, entre outras coisas.

Sua implementação foi feita em Java com o auxílio de uma biblioteca chamada Crawler4j que possui métodos para navegar nos elementos HTML e trabalhar com JSON. Como resultado foram encontrados 2262 links de questionários, os quais passaram por uma limpeza e, ao fim, restaram 510 links totalizando 5765 perguntas. Todos os dados extraídos foram inicialmente postos em JSON e armazenados em um banco de dados relacional; a precisão foi boa, maior que 94 por cento; ou seja, 2137 dos links realmente possuíam questionários.

Os resultados desse trabalho podem ser usados em diversas áreas como enriquecer bases de conhecimento, avaliar a qualidade de questionários de pesquisa, gerar novos questionários etc.

### 3.3 FERRAMENTA PARA COLETA E COMPARAÇÃO DE DADOS DE PUBLICAÇÕES ACADÊMICAS DOS PROFESSORES COM O CURRÍCULO LATTES

(BRANCO, 2018) destaca que esta ferramenta visa criar um *crawler* capaz de extrair informações de outros sites como o ResearchGate para complementar o currículo Lattes com dados que não existem no mesmo. Assim, a aplicação iria integrar e unificar os dados do Lattes com as novas obras que forem encontradas. O currículo Lattes é uma poderosa ferramenta dentro do mundo acadêmico; através dela é possível avaliar o nível de um pesquisador e diversas instituições a utilizam quando da seleção de candidatos para ocupar suas vagas, daí a importância de se manter o currículo sempre atualizado.

O *crawler* desenvolvido possui uma arquitetura simples; ele se conecta a um servidor *proxy* para utilizar as ferramentas de pesquisa, acessa os repositórios e extrai os dados das publicações, faz a análise do currículo Lattes para comparar com os dados extraídos, e grava os dados em uma base de dados. Para acessar o Lattes era necessário resolver um *captcha*; então o autor optou por fazer isso manualmente e deixar os dados pré baixados em formato XML.

A implementação foi feita em JAVA, com auxílio da IDE Eclipse, e foi utilizado o PostgreSQL como banco de dados. As principais bibliotecas utilizadas foram 'jsoup', uma biblioteca desenvolvida para trabalhar com páginas HTML, com uma estrutura de manipulação e extração de dados; e 'Java-string-similarity', uma biblioteca que implementa algumas maneiras de medir a similaridade e distância entre strings.

A aplicação final foi desenvolvida especificamente para certos sites, portanto, não

funciona bem de forma genérica, sendo, um dos problemas, o fato de que pequenas mudanças no site podem implicar no mal funcionamento do código. Também enfrentou-se alguns problemas com relação a bloqueio no acesso ao site devido à detecção do *crawler*. Para trabalhos futuros, seria interessante o uso de uma IA na parte de comparação de dados a fim de obter melhores resultados e que demandem menos análise manual.

### 3.4 CRAWLEX: UMA FERRAMENTA PARA EXTRAÇÃO DE DADOS NA WEB CONFIGURÁVEL ATRAVÉS DE EXEMPLOS

Segundo (LESSA, 2022), o CrawlEX é uma ferramenta que busca facilitar a vida do usuário com a criação de um *crawler* genérico e configurável para a procura e extração de artigos da internet. Ele possui uma interface visual com diversas telas desenvolvidas com Python; na tela de início são exibidas algumas instruções quanto ao uso de cada opção, a tela de fornecer exemplos é onde o usuário entra com os artigos que serão a base utilizada pelo *crawler*, a tela de carregar fonte permite ao usuário selecionar um site já configurado e verificar os dados coletados, a tela de executar um *crawler* permite a seleção de páginas web já coletadas e o envio de argumentos que devem ser procurados nas mesmas, a tela de configurações apresenta algumas opções quanto ao envio de relatórios.

A implementação foi feita em Python, utilizou de bibliotecas como 'pyqt', 'beautifulsoup', 'sqlite', 'html5lib', entre outras. Após fornecer os exemplos, o programa gera uma expressão regular, identifica automaticamente uma possível URL que possua mais artigos e a extração é feita a partir de *tags* do 'beautifulsoup' pré-definidas.

Experimentos iniciais foram feitos passando 1, 2 e 3 exemplos para o programa. Com 1 exemplo o *crawler* teve dificuldade em encontrar os artigos; com 2 o desempenho foi bom, mais de 78 por cento de efetividade; já com 3, esse número abaixou um pouco e um dos motivos para isso é que, com mais exemplos ocorre uma maior generalização na pesquisa dos URLs.

### 3.5 DESENVOLVIMENTO DE API PARA PREDIÇÃO DA COMPLEXIDADE DE TEMPO DE EXECUÇÃO DE CÓDIGOS POR MEIO DE AUTOML

Em (RODENBUSCH, 2023), o autor realiza a combinação de duas bases de dados resultantes de *scraping* e gera uma base chamada de 'Crawleds', logo após ocorre outra combinação gerando uma terceira nomeada de 'Merged', que consta com 1668 amostras as quais são utilizadas como base principal para o trabalho. A partir do uso de métodos de aprendizado de máquina e a base foram gerados modelos de classificação, voltados para predição de complexidade em trechos de código, tais modelos atingiram 88,38% de acurácia para cálculo de eficiência e 77,58% para cálculo

de classe de complexidade. A partir destes foi desenvolvida então uma API a fim de disponibilizar uma interface para que o usuário final possa fazer fácil uso dos modelos treinados, na API é possível enviar trechos de código em Java e obter a classe de eficiência e complexidade predita.

### 3.6 ANÁLISE DOS TRABALHOS

Ao comparar quatro das obras apresentadas que implementaram alguma ferramenta de extração e contrastando-as com o trabalho atual, é perceptível que, embora haja distintos focos de extração, a maioria dos trabalhos é direcionada de forma especializada. A Tabela 1 ilustra essas diferenças.

O CrawlEX se destaca por tentar abordar uma perspectiva mais generalizada; no entanto, ainda necessita de configurações prévias, o que pode influenciar sua eficiência que não atinge a totalidade. Independentemente da abordagem, é comum que todos utilizem bibliotecas especializadas para navegar na estrutura DOM da página e nos elementos HTML. Essa dependência de configuração, seja ela feita pelo usuário ou já incorporada, é crucial, visto que cada página da web possui suas particularidades.

A presença ou ausência de uma interface gráfica parece estar diretamente ligada ao público-alvo de cada trabalho. Trabalhos como ECL e CrawlEX, que possuem interface gráfica, parecem ser orientados para usuários finais. Por outro lado, obras como o qFEx, CoRCoD e o SACI estão mais voltadas para a extração de dados para alimentar bases, não necessitando, portanto, de um elemento visual.

Em relação às linguagens de programação adotadas, duas predominam: Python e Java. O que é particularmente interessante é que, independentemente da linguagem escolhida, as bibliotecas auxiliares empregadas, como 'beautifulsoup' e 'jsoup', visam solucionar problemas semelhantes, demonstrando um padrão de desafios comuns na área de extração de dados web.

Comparação de trabalhos					
Trabalho	Tipo	Item Alvo	Técnica	Algoritmo de Extração	Interface Gráfica
qFEx	Focado	Questionários	Heurísticas	Do autor	Não
ECL	Focado	Publicações	Similaridade textual	Do autor	Sim
CrawlEX	Genérico	Artigos	Comparação textual	Do autor e Ratcliff/Obershelp	Sim
CoRCoD	Focado	Algoritmos e Complexidades	API e Scraping	Do autor	Não
SACI	Focado	Algoritmos e Complexidades	Expressões regulares e OpenAI	Do autor	Não

Tabela 1 – Comparação de trabalhos

## 4 SACI - SCRAPER DE ALGORITMOS E COMPLEXIDADES

Este capítulo apresenta o SACI (Scraper de algoritmos e complexidades), um *web scraper* desenvolvido para a extração de algoritmos e suas complexidades. Inicialmente, é apresentada uma visão geral de seu funcionamento com códigos em alto nível, em seguida são mostrados detalhes da implementação sobre a fonte de dados usada na extração, o Geeks4Geeks <sup>1</sup>, e, por fim, são apresentados alguns experimentos que mostram alguns números interessantes da extração realizada. Os arquivos com o código fonte podem ser encontrados no repositório<sup>2</sup> do GitHub.

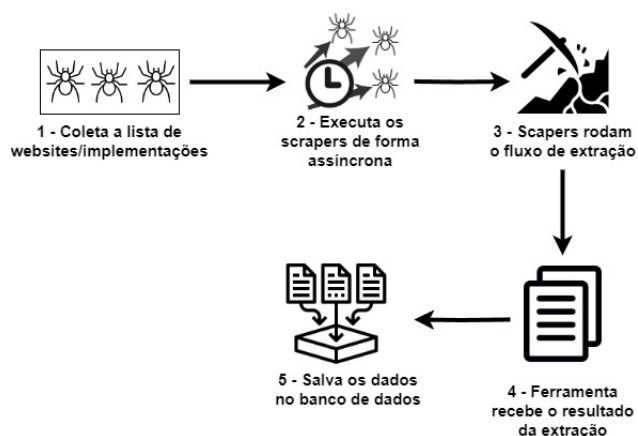
### 4.1 VISÃO GERAL

Esta seção apresenta o fluxo da ferramenta e de extração do *scrapers* desenvolvidos na ferramenta SACI.

#### 4.1.1 Fluxo do SACI

Para o desenvolvimento da base de dados foi construída uma ferramenta de extração utilizando conceitos de *web scraping*. Essa ferramenta é capaz de extrair os dados desejados, como os códigos e complexidades dos algoritmos das páginas escolhidas para serem extraídas e que tiveram suas implementações realizadas. A Figura 1 demonstra o fluxo no qual a ferramenta atua:

Figura 1 – Fluxo de atuação da ferramenta



Fonte: elaborada pelo autor

Como já mencionado, a ferramenta realiza a extração somente das páginas web

<sup>1</sup> <https://www.geeksforgeeks.org/>

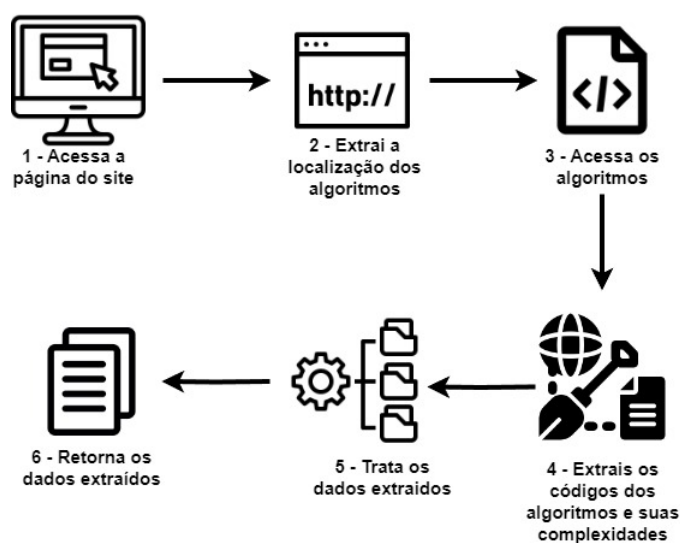
<sup>2</sup> <https://github.com/Augustives/SACI>

para as quais foi feita a implementação. Isso, devido ao fato que cada página contém suas peculiaridades exigindo, assim, certa especificidade para ser possível realizar a extração dos dados. Para tanto, existe uma lista de 'strings' representando websites, definida no arquivo '\_\_init\_\_.py', mais detalhes nas Sub-Seções 4.3 e 4.3.1. Conforme a Figura 1, a ferramenta começa lendo a lista de websites na etapa 1, executa cada implementação de forma assíncrona com o intuito de diminuir o tempo necessário, etapas 2 e 3, e, por fim, retorna todos os resultados em conjunto na etapa 4. Os dados retornados e já padronizados podem ser salvos em uma base de dados MongoDB ou em um arquivo JSON local durante a etapa 5.

### 4.1.2 Fluxo de extração

A Figura 2 apresenta o fluxo de extração executado, individualmente, em cada uma das páginas.

Figura 2 – Fluxo de tarefas de extração



Fonte: elaborada pelo autor

Conforme a Figura 2, o fluxo de extração dos dados para cada página/*website* se resume em acessar a página do site na etapa 1, localizar onde estão os alvos da extração e coletar os URLs que direcionam para esses dados, etapa 2 e 3, acessar cada URL e extrair os dados (códigos dos algoritmos, suas complexidades), etapa 4. Após extraídos, os dados são tratados e retornados nas etapas 5 e 6.

## 4.2 IMPLEMENTAÇÃO

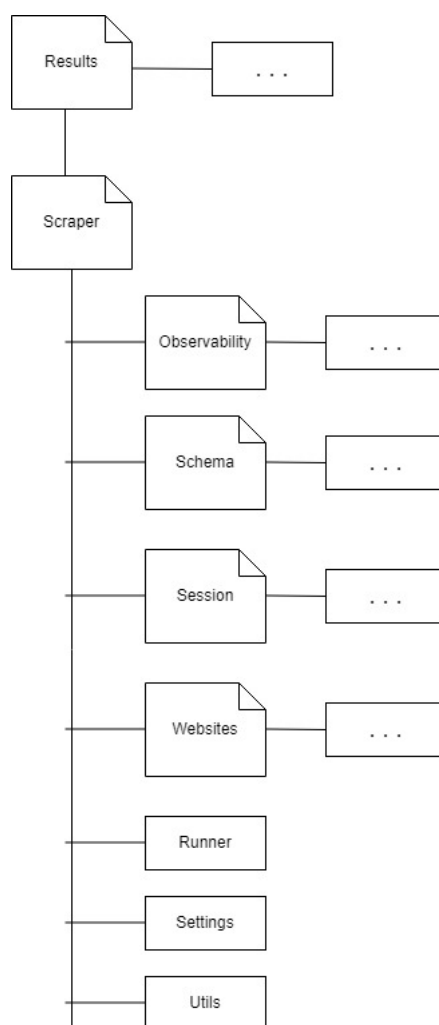
Para facilitar reprodutibilidade, SACI é explicado em detalhes conforme foi implementado. Primeiramente, é descrita a estrutura de pastas, e em seguida como cada

uma está organizada e o seu papel dentro da arquitetura do SACI.

O SACI foi desenhado de forma que é possível a inclusão de novas páginas. Dependendo do interesse do desenvolvedor, basta realizar a implementação seguindo os moldes definidos na Sub-Seção 4.3. A Figura 3 demonstra a estrutura da ferramenta em alto nível, onde cada nota, por exemplo 'Results', representa uma pasta e cada bloco representa um arquivo:

### 4.3 ESTRUTURA GERAL

Figura 3 – Estrutura do SACI



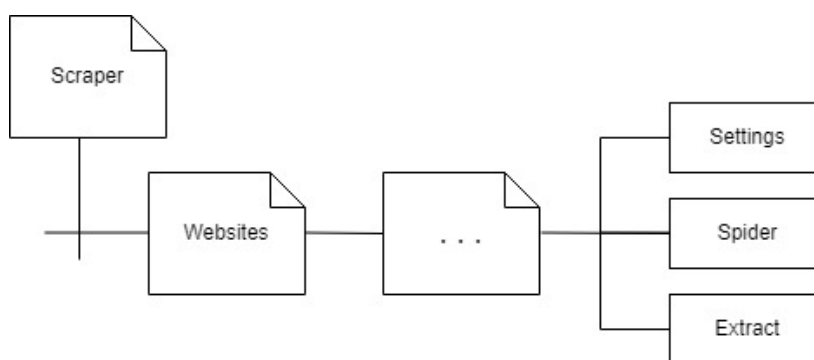
Fonte: elaborada pelo autor

A pasta 'Results', é destinada a armazenar os resultados e análises de cada página localmente. Dentro de 'Observability' estão os códigos responsáveis pelos 'logs' da ferramenta e o código responsável pela geração de métricas e análise do desempenho da extração; por exemplo, quantos algoritmos tiveram suas complexidades corretamen-

te extraídas. Em 'Schema' encontra-se o esquema JSON, desenvolvido com auxílio da biblioteca Pydantic na qual os dados resultantes são validados. A pasta 'Session' contém a implementação de uma sessão HTTP, responsável por todas as requisições feitas no código e pela construção do objeto response, que é o retorno para cada requisição feita. Na pasta 'Websites' está a lógica específica para cada página que foi implementada. Até o momento final deste trabalho apenas o site Geeks4Geeks foi implementado. No arquivo 'runner', está a lógica para executar todos os algoritmos de extração para cada página de forma assíncrona e armazenar os resultados de cada um na base de dados ou localmente. Em 'settings', estão descritas algumas configurações globais, como a variável que define se os resultados ficam localmente em JSON ou são enviados para o banco de dados. Por último, em 'utils', existem códigos que auxiliam no *scraping* de forma geral, por exemplo, a função de 'retry', que executa novamente funções caso alguma exceção seja levantada. Essa função é útil para refazer requisições que apresentem algum problema, seja por fatores de segurança da página ou outros.

### 4.3.1 Spiders individuais

Figura 4 – Estrutura das páginas implementadas



Fonte: elaborada pelo autor

Cada página implementada deve ter seu código em uma pasta inserida dentro de 'websites'. Apesar de não ser obrigatório, é recomendado que essas implementações sigam o padrão definido na Figura 4 e contemham os seguintes três arquivos: 'settings', 'spider' e 'extract'.

Em 'settings', estão alguns detalhes específicos para cada site, como URLs, cabeçalhos, expressões regulares etc. Em 'extract', estão os métodos de extração, funções para extrair os algoritmos, as complexidades, os títulos etc. No arquivo 'spider', fica a lógica geral da extração, onde são feitas as requisições e, em seguida, chamados os métodos de 'extract' para realizar a coleta dos dados.



Nada impede que mais arquivos sejam adicionados, caso o desenvolvedor deseje deixar o código mais modular e organizado. A única obrigatoriedade é realizar o registro da implementação no arquivo `'__init__.py'`, dentro da pasta `'websites'`, conforme o exemplo no Algoritmo 1. Apesar de ser chamada de lista, o objeto em código é um dicionário onde ficam associados o nome de cada página com o método que roda o fluxo de extração e retorna os dados.

---

**Algoritmo 1** Código mostra a lista de páginas em alto nível

---

```
Import executarSpider from página1 as executarPágina1  
Import executarSpider from página2 as executarPágina2
```

```
websites ← {'página1': executarPágina1, 'página2': executarPágina2}
```

---

### 4.3.2 Results

A pasta `'results'` é apenas um espaço para serem guardados os arquivos gerados pela ferramenta que são os resultados das extrações e os resultados das análises feitas sobre cada extração, todos em JSON. O Código 1 apresenta um exemplo, onde são mostrados alguns dos valores e características que são retornados com os resultados da execução do código, contém valores fictícios.

```
1 ##### RESULTS #####  
  
3 'Total Algorithms': 20,  
  'Time Complexity Extracted': '10 extracted - 50%',  
5 'Space Complexity Extracted': '10 extracted - 50%',  
  'URLs with problem': [  
7     'https://www.geeksforgeeks.org/interval-tree/',  
     ...  
9 ],  
  'Distinct Time Complexities': 10,  
11 'Distinct Space Complexities': 10,  
  'Non Trusted Time Complexities': 2,  
13 'Non Trusted Space Complexities': 2,  
  'Time Complexity': {  
15     'True Positive': 10,  
     'False Positive': 0,  
17     'True Negative': 0,  
     'False negative': 0,  
19     'Accuracy': 1.0,  
     'Precision': 1.0,  
21     'Recall': 1.0,  
     'F1 Score': 1.0  
23 },
```

```

25     'Space Complexity': {
26         'True Positive': 10,
27         'False Positive': 0,
28         'True Negative': 0,
29         'False negative': 0,
30         'Accuracy': 1.0,
31         'Precision': 1.0,
32         'Recall': 1.0,
33         'F1 Score': 1.0
34     },
35     'Time Complexity Classification': {
36         'Constant complexity': 2,
37         'Linear complexity': 2,
38         'Exponential complexity': 2,
39         'Polynomial complexity': 2,
40         'Factorial complexity': 1,
41         'Log complexity': 1
42     },
43     'Space Complexity Classification': {
44         'Constant complexity': 2,
45         'Linear complexity': 2,
46         'Exponential complexity': 2,
47         'Polynomial complexity': 2,
48         'Factorial complexity': 1,
49         'Log complexity': 1
50     }
51 ]

```

Código 1 – Exemplo de arquivo de análises

Os dados que formam a base de dados também ficam dentro da pasta 'results' e são retornados em formato JSON, formam uma lista de objetos que representam os dados extraídos, o Código 2 exemplifica isso.

```

1 [
2   {
3     "name": "Knapsack Problem",
4     "time_complexity": "O(2N)",
5     "trustable_time_complexity": true,
6     "space_complexity": "O(N)",
7     "trustable_space_complexity": true,
8     "url": "https://www.geeksforgeeks.org/0-1-knapsack-problem-dp-10/",
9     "codes": {
10      "C": {
11        "code": "...",
12        "comments": "...",
13      }
14    }
15  }
16 ]

```

```
15     "C#": {
16         "code": "...",
17         "comments": "... "
18     },
19     "C++": {
20         "code": "...",
21         "comments": "... "
22     },
23     "Python3": {
24         "code": "...",
25         "comments": "... "
26     },
27     "Python": {
28         "code": "...",
29         "comments": "... "
30     },
31     "PHP": {
32         "code": "...",
33         "comments": "... "
34     },
35     "Java": {
36         "code": "...",
37         "comments": "... "
38     },
39     "Javascript": {
40         "code": "...",
41         "comments": "... "
42     }
43 },
```

Código 2 – Exemplo de arquivo de dados

Descrição das variáveis do Código 2:

1. **name**: Nome/Título do algoritmo.
2. **time\_complexity**: Complexidade de tempo do algoritmo.
3. **trustable\_time\_complexity**: Booleano indicando se essa complexidade é confiável, foi corretamente extraída, houve ou não uso do método de 'fallback'.
4. **space\_complexity**: Complexidade de espaço do algoritmo.
5. **trustable\_space\_complexity**: Booleano indicando se essa complexidade é confiável, foi corretamente extraída, houve ou não uso do método de 'fallback'.
6. **url**: URL de onde os dados foram extraídos.

7. **codes**: Dicionário contendo o código do algoritmo em diferentes linguagens de programação.

### 4.3.3 Observability

Dentro de 'observability', existe um *logger* implementado, Algoritmo 2, responsável por produzir mensagens de erro ou avisos, de forma a ajudar no momento de depurar o código e obter mais informações em tempo de execução. Também há códigos para geração de métricas de avaliação, por exemplo o Algoritmo 3, onde se verifica quantos algoritmos foram coletados, a complexidade de tempo e espaço de cada um, o total que poderia ter sido extraído e quais deles tiveram erros, gera-se uma matriz de confusão, histogramas e outros gráficos.

**Algoritmo 2** Código mostra o logger em alto nível

- 1: **Function** constroiLogger(nomeLogger)
  - 2: Cria o logger
  - 3: Configura o logger
  - 4: **return** logger
- 
- 5: logger ← constroiLogger('scraperLogger')

**Algoritmo 3** Código calcula métricas de desempenho em alto nível

- 1: **Function** calculaMetricas(verdadeiroPositivo, falsoPositivo, verdadeiroNegativo, falsoNegativo)
- 2: Calcula a acurácia:  $(\text{verdadeiroPositivo} + \text{verdadeiroNegativo}) / (\text{verdadeiroPositivo} + \text{verdadeiroNegativo} + \text{falsoPositivo} + \text{falsoNegativo})$
- 3: Calcula a precisão:  $\text{verdadeiroPositivo} / (\text{verdadeiroPositivo} + \text{falsoPositivo})$
- 4: Calcula a sensibilidade (recall):  $\text{verdadeiroPositivo} / (\text{verdadeiroPositivo} + \text{falsoNegativo})$
- 5: Calcula a medida F1:  $(2 * \text{precisão} * \text{recall}) / (\text{precisão} + \text{recall})$
- 6: **return** métricas

**4.3.4 Schema**

Os dados retornados por cada página implementada devem seguir um esquema pré-definido. Na pasta 'schema', ele foi implementado com o auxílio da biblioteca Pydantic, uma biblioteca que fornece de forma prática maneiras de gerar esquemas de dados e validá-los. Antes do retorno do resultado definitivo, todos os dados passam por uma validação no formato apresentado no Código 3.

```

1 class ScrapedCode(BaseModel):
2     '''The scraped codes of a given algorithm'''
3
4     code: str = Field(
5         default='', description='The code of the algorithm in a certain
6             language'
7     )
8     comments: str = Field(default='', description='The first few comments
9         in the code')
10
11 class ScrapedAlgorithm(BaseModel):
12     '''The final result of the scraping operations, an extracted algorithm
13         with its attributes'''
14
15     name: str = Field(default='', description='The name of the scraped
16         algorithm')

```

```
15     time_complexity: Optional[str] = Field(  
        default='', description='The time complexity of the scraped  
            algorithm '  
    )  
17     trustable_time_complexity: bool = Field(  
        default=True,  
19     description='If the given time complexity used a trustable  
            extraction method',  
    )  
21     space_complexity: Optional[str] = Field(  
        default='', description='The space complexity of the algorithm '  
23     )  
25     trustable_space_complexity: bool = Field(  
        default=True,  
        description='If the given space complexity used a trustable  
            extraction method',  
27     )  
    url: str = Field(default='', description='The space complexity of the  
        algorithm')  
29     codes: dict[str, ScrapedCode] = Field(  
        default=dict(), description='The algorithm implementations '  
31     )
```

Código 3 – Esquema dos dados

Com a criação destes dois esquemas conseguimos atingir um nível de validação que evita problemas nos resultados. Isso é fundamental para garantir que os dados coletados durante o processo sejam consistentes e estruturados em um padrão.

#### 4.3.5 Session

Na pasta 'session', está grande parte do funcionamento da ferramenta que é baseado em requisições HTTP. Para isso, foi implementado um tipo de sessão, apresentada em alto nível no Algoritmo 4, que utiliza como base a biblioteca Aiohttp que, por sua vez, é capaz de lidar com requisições assíncronas. Nela, ocorrem algumas validações de URL, um meio de definir os cabeçalhos que serão utilizados, um método de fazer requisições para páginas dinâmicas (que precisam de JavaScript) e um método para requisições normais. Também foi definido um objeto 'response' para cada resposta de requisição recebida onde são executados alguns métodos para se ter os dados da resposta em mãos. Com o objeto resposta é possível realizar buscas através de XPATH e através dos elementos HTML da resposta.

---

**Algoritmo 4** Código mostra a sessão http em alto nível

---

```
1: Class SessãoHttp
2: Init(sessão, jsSessão, cabeçalhosPadrão)
3: self.sessão ← sessão
4: self.jsSessão ← jsSessão
5: self.cabeçalhosPadrão ← cabeçalhosPadrão

6: Getter cabeçalhosPadrão
7: Setter cabeçalhosPadrão

8: StaticMethod validarUrl(url)
9: if not urlValida() then
10:     raise Erro
11: end if

12: Method validaArgumentosRequest(**kwargs)
13: Extrai metodo(GET, POST...) dos kwargs
14: if metodo is None or metodo not in ARGUMENTOS_NECESSARIOS then
15:     raise Erro
16: end if
17: validarUrl(url)

18: Async Method request(**kwargs)
19: Valida argumentos
20: Realiza HTTP request
21: if request falhou then
22:     Log erro
23: end if
24: Cria objeto resposta

25: Async Method jsRequest(**kwargs)
26: Valida argumentos
27: Realiza JS request
28: Cria objeto resposta
```

---

**Algoritmo 5** Código mostra o objeto resposta em alto nível

---

```
1: Class Resposta
2: Init(url, status, conteudo, cabeçalhos, resposta, soup)
3: self.url ← url
4: self.status ← status
5: self.conteudo ← conteudo
6: self.cabeçalhos ← cabeçalhos
7: self.resposta ← resposta
8: self.soup ← BeautifulSoup(content)
```

---

### 4.3.6 Runner

O arquivo 'runner', apresentado no Algoritmo 6, apesar de conter pouco código, é parte essencial do SACI. Nele fica a implementação responsável por gerenciar a extração de todas as páginas implementadas de forma assíncrona, tornando, assim, a tarefa mais rápida.

---

**Algoritmo 6** Código mostra a função principal do arquivo runner em alto nível

---

- 1: **Async Function** rodarScrapers(url)
  - 2: resultadosScrapers ← Executa scrapers e agrega resultados em uma lista
  - 3: Salva os resultados em arquivo local ou em uma base de dados
- 

## 4.4 GEEKS4GEEKS

O Geeks4Geeks<sup>3</sup> é um site com conteúdo voltado para programadores com uma variedade de materiais para aprendizado e, também, de entretenimento. repleto de artigos, tutoriais, aulas, seminários, é um ótimo lugar para quem procura aprender algo novo.

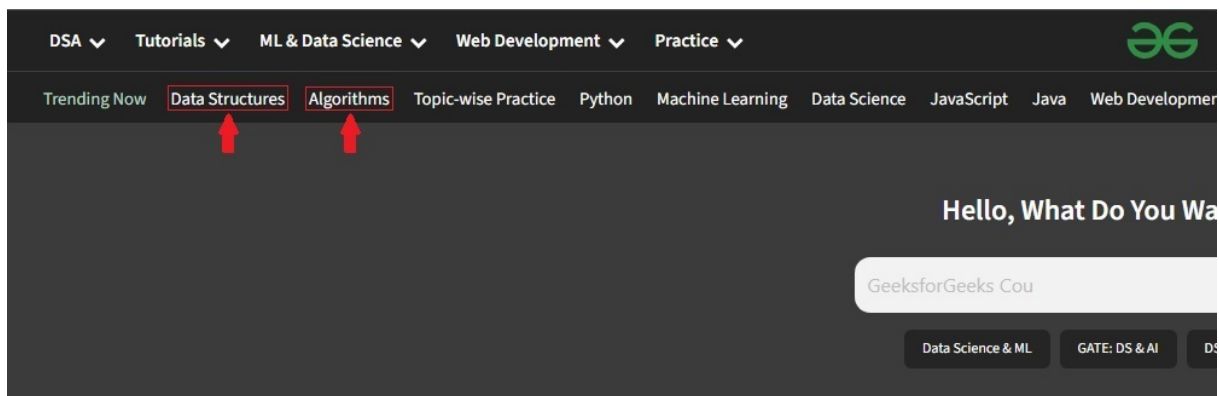
No topo do site, Figura ??, ficam evidenciados os principais tópicos: algoritmos de estrutura de dados, tutoriais, ciência de dados e aprendizado de máquina, desenvolvimento web, entre outros. Destes, o que interessa para o trabalho são os referentes a algoritmos e estrutura de dados. Dentro deles existem subtópicos e links, Figura ??, que direcionam para páginas com implementações de código, em múltiplas linguagens, acompanhadas das complexidades de tempo e espaço Figura ??, Figura 8 e Figura 9. Por conta dessas características, essas duas seções, apontadas na Figura ??, foram escolhidas como fontes de dados a serem extraídos pelo SACI. Nessas, existe uma variedade de links de algoritmos, organizados conforme seus tipos e temas.

---

<sup>3</sup> <https://www.geeksforgeeks.org/>



Figura 5 – Menu Geeks4Geeks



Fonte: elaborada pelo autor

Figura 6 – Página de algoritmos

**Topics:**

- [Analysis of Algorithms](#)
- [Searching and Sorting](#)
- [Greedy Algorithms](#)
- [Dynamic Programming](#)
- [Pattern Searching](#)
- [Backtracking](#)
- [Divide and Conquer](#)
- [Geometric Algorithms](#)
- [Mathematical Algorithms](#)
- [Bit Algorithms](#)
- [Graph Algorithms](#)
- [Randomized Algorithms](#)
- [Branch and Bound](#)
- [Quizzes](#)

**Analysis of Algorithms:**

1. [Asymptotic Analysis](#)
2. [Worst, Average and Best Cases](#)
3. [Asymptotic Notations](#)
4. [Lower and Upper Bound Theory](#)
5. [Introduction to Amortized Analysis](#)
6. [What does 'Space Complexity' mean?](#)
7. [Polynomial Time Approximation Scheme](#)
8. [Accounting Method | Amortized Analysis](#)
9. [Potential Method in Amortized Analysis](#)

**Searching and Sorting:**

1. [Introduction to Searching Algorithms](#)
2. [Introduction to Sorting Algorithm](#)
3. [Stable and Unstable Sorting Algorithms](#)
4. [Lower bound for comparison based sorting algorithms](#)
5. [Can Run Time Complexity of a comparison-based sorting algorithm be less than  \$N \log N\$ ?](#)
6. [Which sorting algorithm makes minimum number of memory writes?](#)

**Greedy Algorithms:**

1. [Introduction to Greedy Algorithms](#)

Fonte: elaborada pelo autor

Figura 7 – Complexidades formato 1

```
// Driver code
int main()
{
    int n = 5;
    printPascal(n);
    return 0;
}
```

**Output**

```
1
1 1
1 2 1
1 3 3 1
1 4 6 4 1
```

**Time Complexity:**  $O(n^2)$   
**Auxiliary Space:**  $O(1)$

So method 3 is the best method among all, but it may cause integer overflow for large values of  $n$  as it multiplies two integers to obtain values.

Fonte: elaborada pelo autor

Figura 8 – Complexidades formato 2

```
console.log(`The smallest distance is ${closest(P, n)}`);
}
main();
```

**Output**

```
The smallest distance is 1.41421
```

**Time Complexity:** Let Time complexity of above algorithm be  $T(n)$ . Let us assume that we use a  $O(n \log n)$  sorting algorithm. The above algorithm divides all points in two sets and recursively calls for two sets. After dividing, it finds the strip in  $O(n)$  time. Also, it takes  $O(n)$  time to divide the  $P_y$  array around the mid vertical line. Finally finds the closest points in strip in  $O(n)$  time. So  $T(n)$  can be expressed as follows

$$T(n) = 2T(n/2) + O(n) + O(n) + O(n)$$

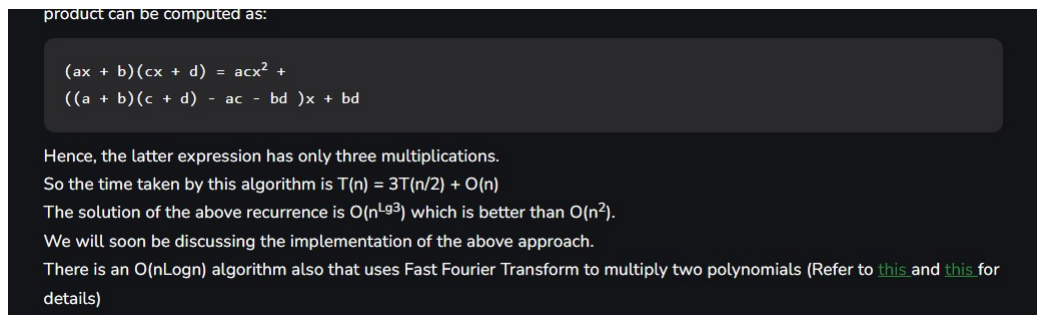
$$T(n) = 2T(n/2) + O(n)$$

$$T(n) = T(n \log n)$$

**Auxiliary Space:**  $O(\log n)$ , as implicit stack is created during recursive calls

Fonte: elaborada pelo autor

Figura 9 – Complexidades formato 3



Fonte: elaborada pelo autor

Ao acessar esses links, tem-se as implementações dos algoritmos e suas complexidades. Em geral as complexidades são apresentadas após os blocos de código com suas respectivas implementações e se apresentam de diversas formas, podendo ser de uma maneira mais direta como na Figura ?? ou dentro de textos como na Figura 8 e na Figura 9.

#### 4.4.1 Settings

Dentro do arquivo 'settings' tem-se, além das expressões regulares, os cabeçalhos que utilizados nas requisições HTTP, os URLs do site e a lista de linguagens de programação coletadas. Os cabeçalhos são muito importantes, pois são utilizados em todas as requisições e transmitem informações sobre a fonte da requisição e como ela deve ser tratada. Quando algum serviço de *scraping* faz uma requisição a um servidor, ele pode ser detectado e bloqueado com base em seus cabeçalhos. Vários sites têm mecanismos de defesa para identificar padrões e bloquear tais serviços que não se comportam como navegadores padrões. Por isso, os cabeçalhos são modificados para se parecerem com navegadores legítimos, 'disfarçando' sua verdadeira natureza e obtendo mais chances de acessar o conteúdo do site sem ser bloqueado.

```

1 HEADERS = {
    'Host': 'www.geeksforgeeks.org',
3    'Accept': 'text/html, application/xhtml+xml, application/xml;q=0.9,image/
      avif,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=
      b3;q=0.9',
    'Accept-Encoding': 'gzip, deflate',
5    'Upgrade-Insecure-Requests': '1',
    'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit
      /537.36 (KHTML, like Gecko) Chrome/104.0.0.0 Safari/537.36',
7    'Sec-Fetch-Site': 'none',
    'Sec-Fetch-Mode': 'navigate',
9    'Sec-Fetch-User': '?1',

```

```
11     'Sec-Fetch-Dest': 'document',  
13     'Sec-Ch-Ua': ' 'Chromium ';v='105 ', 'Not)A;Brand ';v='8 ',  
     'Sec-Ch-Ua-Mobile': '?0',  
     'Sec-Ch-Ua-Platform': ' 'Windows ',  
 }
```

Código 4 – Cabeçalhos utilizados

Os cabeçalhos, apresentados no Código 4, são modificados para garantir que a requisição HTTP se pareça com a de um navegador legítimo, usado no dia a dia por um ser humano, para isso seguiu-se estas estratégias:

1. User-Agent: Identifica o cliente que está fazendo a requisição. Utilizou-se uma string de User-Agent de um navegador popular, como o Google Chrome, para fazer a requisição parecer que vem de um usuário humano, em vez de um serviço de *scraping* automatizado.
2. Accept e Accept-Encoding: Indicam os tipos de conteúdo que o cliente pode processar. Configurou-se esses cabeçalhos para aceitarem uma variedade de tipos de conteúdo e codificações comuns para navegadores, como HTML, XML, e imagens.
3. Sec-Fetch-Site, Sec-Fetch-Mode, Sec-Fetch-User, Sec-Fetch-Dest: Fazem parte da política de segurança dos navegadores modernos e informam ao servidor o contexto da navegação. Ajustou-se estes para imitar os padrões de navegação de um usuário humano.
4. Host e Upgrade-Insecure-Requests: Especifica o domínio ao qual a requisição está sendo feita, enquanto Upgrade-Insecure-Requests indica ao servidor para aceitar a conexão através de HTTPS. Ajustou-se estes para se alinhar com as práticas de navegação segura.

Ao ajustar esses cabeçalhos para simular um navegador padrão, é possível evitar muitos bloqueios e defesas automáticas implementadas pelos sites, aumentando significativamente a eficácia do *scraper*.

#### 4.4.2 Scraper

Para o site do Geeks4Geeks, a extração seguiu em linhas gerais o fluxo já apresentado na Figura 2 e representado em alto nível no Algoritmo 7. A função principal do arquivo inicializa a sessão HTTP, faz o login na página a fim de evitar problemas como a detecção do *scraper* ou restrições a usuários não cadastrados, procura os URLs do algoritmo e, em seguida, faz o acesso aos algoritmos.

---

**Algoritmo 7** Código mostra a função principal do arquivo spider em alto nível

---

- 1: **Async Function** executarSpider(url)
  - 2: sessão ← SessãoHttp()
  - 3: fazLogin(sessão)
  - 4: urlsAlgoritmos ← encontrarUrlsAlgoritmos(sessão)
  - 5: **return** acessarTodosAlgoritmos(sessão, urlsAlgoritmos)
- 

A extração dos URLs se dá em três funções, apresentadas na Figura 8. Uma responsável por acessar a página do site que contém os URLs, uma para fazer a extração e uma para filtrar os URLs extraídos.

---

**Algoritmo 8** Código mostra as funções responsáveis por extrair URLs em alto nível

---

- 1: **Function** extrairUrlsAlgoritmos(response)
  - 2: Encontra e retorna todos os URLs que contém algoritmos
  
  - 3: **Function** filtraUrlsAlgoritmos(urlsAlgoritmos)
  - 4: **return** Lista de URLs filtradas, sem URLs desnecessárias
  
  - 5: **Async Function** encontrarUrlsAlgoritmos(sessão)
  - 6: urlsAlgoritmos ← Acessa página que contém as URLs
  - 7: urlsAlgoritmos ← extrairUrlsAlgoritmos(urlsAlgoritmos)
  - 8: urlsAlgoritmos ← filtraUrlsAlgoritmos(urlsAlgoritmos)
  - 9: **return** Lista de URLs de algoritmos
- 

A extração dos algoritmos é dividida em quatro etapas: primeiro, as URLs são percorridas e é chamada a função 'acessarAlgoritmo', linha 12 do Algoritmo 9, sobre cada uma, a qual faz a extração e adéqua os dados ao esquema, e por último os resultados são agregados em uma lista. A função 'extrairDados', linha 5 do Algoritmo 9, é encontrada no arquivo 'Extract', descrito na Seção 4.5.

---

**Algoritmo 9** Código mostra as funções responsáveis por extrair algoritmos em alto nível

---

```
1: Function extrairDadosAlgoritmo(response)
2: if not response then
3:   return [ ]
4: end if
5: dados ← extrairDados(response)
6: if not dados then
7:   raise Erro
8: end if
9: return Dados brutos dos algoritmos

10: Function algoritmosParaEsquema(dadosBrutosAlgoritmo)
11: return Lista de Algoritmos adequados ao esquema

12: Async Function acessarAlgoritmo(sessão, url)
13: response ← sessão.request(url)
14: dadosBrutosAlgoritmo ← extrairDadosAlgoritmo(response)
15: algoritmo ← algoritmosParaEsquema(dadosBrutosAlgoritmos)
16: return Algoritmo dentro do esquema definido

17: Async Function acessarTodosAlgoritmos(sessão, urlsAlgoritmos)
18: lista ← [ ]
19: for url in urlsAlgoritmos do
20:   algoritmo ← acessarAlgoritmo(sessão, url)
21:   Inseire algoritmo na lista
22: end for
23: return Lista com os dados do algoritmo
```

---

## 4.5 EXTRACT

Foram desenvolvidas duas implementações para a extração. Uma das versões faz uso de expressões regulares como o principal meio de identificação e extração das complexidades. Já a outra, no momento da extração, utiliza inteligência artificial através da API (*Application Programming Interface*) do OpenAI para a identificação e extração das complexidades. Todos os resultados das extrações podem ser encontrados no repositório<sup>4</sup> do GitHub e também no Kaggle<sup>5</sup>.

### 4.5.1 Extração com expressões regulares

Cada página do site pode conter mais de um código/implementação. Esses códigos ficam localizados em um elemento HTML chamado aqui de bloco. Cada bloco

---

<sup>4</sup> <https://github.com/Augustives/SACI>

<sup>5</sup> <https://www.kaggle.com/datasets/augustives/algorithms-and-their-complexities>

contém o algoritmo codificado em diversas linguagens de programação e uma complexidade associada ao algoritmo. Ao acessar à página o método procura pelas seções que contém esses blocos, linha 2 do Algoritmo 10, guarda as referências de posição desses blocos na árvore DOM junto dos códigos. Em seguida, procura em volta de cada referência pelas complexidades de tempo e espaço, cada uma através de uma seleção prévia de expressões regulares, linhas 7 a 18 do Algoritmo 10. Por fim, extraímos o título do algoritmo na linha 19 do Algoritmo 10 e adicionamos um objeto resultado a uma lista, linhas 20 a 32 do Algoritmo 10. Para a extração dos títulos dos códigos, linha 19, também utiliza-se a referência de localização dos blocos de código. É extraído o título principal da página que geralmente é o título do primeiro algoritmo e os restantes, para retornar uma lista de dicionários Python contendo o nome do algoritmo, a complexidade de tempo, a complexidade de espaço, o URL e o código em todas as linguagens, linhas 20 a 32. Essa lista é tratada pelo *scraper* na última etapa retornando os dados de acordo com o esquema JSON já definido na Sub-Seção 4.3.4. As chaves 'tempoConfiavel' e 'espaçoConfiavel' do objeto resultado tem relação com o uso ou não do método 'fallback' apresentado mais a frente. Caso o método seja utilizado essas chaves recebem como valor o booleano 'True', caso contrário 'False'.

---

**Algoritmo 10** Código mostra a função de extração principal em alto nível

---

```
1: Function extrairDados(response)

2: codigos, referenciaDom ← extrairCodigosReferencias(response)
3: if not codigos then
4:   Log aviso
5:   return []
6: end if

7: complexidades ← []
8: for ref in referenciaDom do
9:   try:
10:  complexidadeTempo ← extrairComplexidadeDeReferencia(ref, regexTempo)
11:  catch Erro:
12:  complexidadeTempo ← extrairComplexidadeComFallback(ref)

13:  try:
14:  complexidadeEspaço ← extrairComplexidadeDeReferencia(ref, regexEspaço)
15:  catch Erro:
16:  complexidadeEspaço ← extrairComplexidadeComFallback(ref)

17:   Insere complexidades na lista
18: end for

19: titulos ← extrairTitulos(ref)

20: resultados ← []
21: for (nome, complexidade, codigo) in zip(nomes, complexidades, codigos) do
22:   resultado ← {
23:     'nome': nome,
24:     'complexidadeTempo': complexidade['complexidadeTempo'],
25:     'tempoConfiavel': complexidade['tempoConfiavel'],
26:     'complexidadeEspaço': complexidade['complexidadeEspaço'],
27:     'espaçoConfiavel': complexidade['espaçoConfiavel'],
28:     'url': url,
29:     'codigo': codigo,
30:   }
31:   Insere resultado na lista
32: end for

33: return resultados
```

---

Para extrair os blocos de código e suas referências, existe a função 'extrairCodigosReferencias', utilizada na linha 2 do Algoritmo 10. Esses blocos podem ser localizados com a ajuda de uma biblioteca chamada 'BeautifulSoup' que recebe conteúdo HTML e o transforma em 'soups', objetos que podem ser usados para fazer pesquisas e navegar no DOM da página.



Das linhas 4 a 7 do Algoritmo 11, através do objeto `response` e sua `'soup'`, é feita uma pesquisa com a função `'find_all'` na qual retornam todas as `'divs'` que possuem a classe `'responsive-tabs'`; esses são os blocos de código. Em seguida, é percorrido cada bloco e, também, uma lista de linguagens (`['C', 'C#', 'C++', 'Python3'...]`).

Das linhas 8 a 15 do Algoritmo 11, procuramos em cada bloco o conteúdo com o algoritmo através da função `'find_next'`, estes localizados em elementos `'h2'` com a classe `'tabtitle'`. Caso exista o algoritmo, o código é localizado em elementos `'td'` com a classe `'code'` e extraído. Além do código, são coletados os comentários que ficam dentro do mesmo.

---

**Algoritmo 11** Código mostra a função `'extrairCodigosReferencias'` em alto nível

---

```
1: Function extrairCodigosReferencias(response)
2: codigos ← [ ]
3: referencias ← [ ]

4: blocosCodigo ← response.soup.find_all('div', {'class': 'responsive-tabs'})
5: for bloco in blocosCodigo do
6:     listaLinguagens ← [ ]
7:     Insere referencia do bloco na lista de referencias

8:     for linguagem in linguagensProcuradas do
9:         algoritmo ← bloco.find_next('h2', {'class': 'tabtitle'}, string=lang)
10:        if algoritmo then:
11:            codigo ← extrairCodigo(algoritmo.find_next('td', {'class': 'code'}))
12:            comentarios ← extrairComentarios(codigo)
13:        end if

14:        Insere codigo na lista de linguagens
15:    end for

16:    Insere código em suas linguagens na lista de codigos
17: end for

18: return codigos, referencias
```

---

As funções que extraem o código, Algoritmo 12, e os comentários, Algoritmo 13, são simples: apenas percorrem os elementos e suas linhas e coletam o desejado. Ambas se utilizam da função `'find_all'`, linhas 3 e 5, pertencente à biblioteca `'BeautifulSoup'`, que resumidamente encontra todos os elementos dentro do `'soup'`, os quais correspondem aos parâmetros de pesquisa usados.

---

Código mostra a função 'extrairCodigo' em alto nível

---

```
1: Function extrairCodigo(elementoCodigo)
2: codigo ← ""
3: linhasCodigo ← elementoCodigo.find_all('div', {'class': 'line'})

4: for linha in linhasCodigo do
5:     trechosCodigo ← linha.find_all('code')
6:     for trecho in trechosCodigo do
7:         codigo ← codigo + trecho
8:     end for

9: end for

10: return codigo
```

---

---

Código mostra a função 'extrairComentarios' em alto nível

---

```
1: Function extrairComentarios(codigo)
2: comentariosCodigo ← ""

3: for linha in codigo do
4:     comentariosCodigo ← comentariosCodigo + comentario
5: end for

6: return comentariosCodigo
```

---

No momento da extração das complexidades, são utilizados a referência da localização dos códigos, linhas 10 e 14 do Algoritmo 10, métodos do 'BeautifulSoup' e expressões regulares. A extração conta com dois caminhos: o padrão com regras mais específicas e um que chamamos de 'fallback', no qual são utilizadas regras mais generalizadas e inclusivas a fim de conseguir extrair a informação caso ocorra alguma falha durante o fluxo padrão.

O método padrão, representado no Algoritmo 14, procura por todas as palavras *Time Complexity* ou *Space Complexity* nas proximidades do bloco de códigos, linha 4. Para cada palavra existe uma lista de expressões regulares, para cada palavra encontrada elas são inseridas em uma lista de 'matches', linha 3 e 6. Então, procura o 'match' mais próximo da referência, linha 8, com um cálculo simples: para cada 'match' da lista subtrai-se a linha em que a referência está com a linha que o 'match' encontra-se, dessa forma, o menor valor dentre todos os calculados é o mais próximo. A seguir, o método percorre outra lista de expressões regulares, dessa vez focada em capturar os valores das complexidades. Para cada expressão regular de valor, itera os elementos HTML, para os quais é feita uma pesquisa, linha 11. Caso nenhuma complexidade consiga ser extraída, é levantado um erro de extração, linha 17.

Sendo o erro levantado, o fluxo de 'fallback', representado no Algoritmo 15, tenta extrair novamente a complexidade, e uma vez utilizado esse método, uma 'flag' é adicionada aos resultados indicando seu uso.

---

**Algoritmo 14** Código mostra a função 'extrairComplexidadeDeReferencia' em alto nível

---

```
1: Function extrairComplexidadeDeReferencia(referencia, regex)
2: for regex in regexes['palavra'] do
3:   matches ← [ ]
4:   palavra ← referencia.find_next(string=compile(regex))

5:   if palavra then:
6:     Insere palavra em matches
7:   end if

8:   match_mais_proximo ← Pega o match mais próximo da ref

9:   for regex in regexes['valor'] do
10:    for elemento in elementosHTML['valor'] do
11:      complexidade ← busca_regex(regex, match_mais_proximo)

12:      if complexidade then:
13:        return complexidade
14:      end if
15:    end for
16:  end for

17:  raise FalhaExtraçãoComplexidade
18: end for
```

---

---

**Algoritmo 15** Código mostra a função 'fallback' em alto nível

---

```
1: Function fallback(referencias, regex)
2: for referencia in referencias do
3:   if buscaComRegex(regex, referencia) then:
4:     return complexidade
5:   end if
```

---

A Figura 5 apresenta as expressões regulares criadas e utilizadas na extração dos dados. Aquelas que buscam pelas palavras, linhas 3, 16 e 17, tentam encontrar *time complexity*, *space complexity* e suas variações de escrita. As responsáveis pelos valores, linhas 5 a 10 e 19 a 28, são uma combinação da palavra com o valor da complexidade, buscando, por exemplo, algo como *Time Complexity: O(log N)*. Todas foram criadas com base nos diversos padrões com que as complexidades são descritas nos textos das páginas do site. Por fim, tem-se a que é usada no método

'fallback, linhas 11 e 29.

```

TIME_COMPLEXITY_REGEX = {
2   'word': [
4     r'[Tt]ime [Cc]omplexit\\w*\\s?:|[Tt]ime [Cc]omplexit\\w*\\s?',
6   ],
8   'value': [
10    r'[Tt]ime [Cc]omplexit\\w*\\s?:.*?(O\\s*(.*?\\))',
12    r'[Tt]ime [Cc]omplexit\\w*\\s?.*?(O\\s*(.*?\\))',
14    r'[Tt]ime [Cc]omplexit\\w*\\s.*?\\sis\\s.*?(O\\s*(.*?\\))',
16    r'[Tt]ime [Cc]omplexit\\w*\\s?:?.*?\\sis\\s(\\w*)',
18  ],
20  'fallback': r'O\\s*(.*?\\)',
22 }

AUXILIARY_SPACE_REGEX = {
24   'word': [
26    r'[Aa]uxiliary [Ss]pace\\s?:|[Aa]uxiliary [Ss]pace\\s?',
28    r'[Ss]pace [Cc]omplexit\\w*\\s?:|[Ss]pace [Cc]omplexit\\w*\\s?',
30  ],
32   'value': [
34    r'[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?.*?(O\\s*(.*?\\))',
36    r'[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?.*?(O\\s*(.*?\\))',
38    r'[Ss]pace [Cc]omplexit\\w*\\s?:.*?(O\\s*(.*?\\))',
40    r'[Ss]pace [Cc]omplexit\\w*\\s?.*?(O\\s*(.*?\\))',
42    r'[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?.*?\\sis\\s.*?(O\\s
44      *(.*?\\))',
46    r'[Ss]pace [Cc]omplexit\\w*\\s.*?\\sis\\s.*?(O\\s*(.*?\\))',
48    r'[Aa]uxiliary [Ss]pace\\w*\\s?:?.*?\\sis\\s(\\w*)',
50    r'[Ss]pace [Cc]omplexit\\w*\\s?:?.*?\\sis\\s(\\w*)',
52  ],
54  'fallback': r'O\\s*(.*?\\)',
56 }

```

Código 5 – Regex de extração

## 4.5.2 Extração com OpenAI

Para a extração com OpenAI utilizou-se da API oferecida <sup>6</sup>, e para sua implementação foi utilizada uma biblioteca do Python chamada de Langchain <sup>7</sup>. Ela oferece uma integração fácil com modelos de linguagem ou LLMs (*Large Language Models*), como o GPT-3 e outros similares, além de oferecer aos desenvolvedores um conjunto de

<sup>6</sup> <https://openai.com/index/openai-api/>

<sup>7</sup> <https://www.langchain.com/>

ferramentas e abstrações que facilitam a criação de aplicações interativas e complexas. A premissa central da biblioteca Langchain é permitir a construção de correntes de 'prompts' e respostas; ou seja, sequências de interações com modelos de linguagem que podem ser encadeadas para realizar tarefas sofisticadas. Dentro do arquivo *utils* do SACI está definida uma classe que faz o uso desta biblioteca e realiza essa comunicação com a OpenAI.

---

**Algoritmo 16** Código mostra a sessão http em alto nível
 

---

```

1: Class LlmSearcher
2: throttler ← Throttler(rate_limit=100, period=60)

3: Init()
4: self.LLM ← OpenAI(openai_api_key=...)
5: self.PROMPT ← PromptTemplate.from_template(template_do_prompt)
6: self.LLM_CHAIN ← LLMChain(prompt=self.PROMPT, llm=self.LLM)

7: Method search(complexidade, texto)
8: async with throttler:
9: resposta ← self.LLM_CHAIN.run(complexity=complexidade, text=texto)
10: return json.loads(resposta)
    =0
  
```

---

A classe, representada no Algoritmo ??, atua como um *singleton*, isto é, só existe uma instância sua durante a vida útil da aplicação. Dentro dela estão definidos o LLM que conecta com o modelo de linguagem da OpenAI, o 'PROMPT' que define como o modelo irá responder e a LLM\_CHAIN que une o 'PROMPT' com o modelo em si. Além disso, existe um 'throttler' que limita quantas 'requests' são feitas por minuto, pois, a API da OpenAI possui certas restrições.

```

1 I am using you during a scraping operation in which I need to extract from
   a piece of text the values of complexity/code complexity.
2 Consider that sometimes the space complexity is also called as auxiliary
   space.
   Consider that both complexities often come in this format: 'Some type of
   complexity: value of complexity'.
3
4 But they can also come inside the text like: 'The type of complexity of the
   given ... is ...'.
   Give the answer in JSON format with no line breaks, with a key called '
   complexity' and the value for the key is your answer.
5
6 If you cant determine the answer give the json with a null in the value.
   What is the {complexity} that is written in the following text: '{text}'?
  
```

Código 6 – Prompt

O 'PROMPT' utilizado é descrito no Algoritmo 6, com ele é possível explicar o

contexto em que o SACI está, demonstrando alguns exemplos de como os dados irão chegar, seus formatos e diferenças, sendo ainda explicado como a resposta deverá ser retornada e qual pergunta é feita. Com isso o modelo de linguagem consegue entender que está sendo feito a pesquisa pelos valores de complexidade do texto que ele está recebendo.

O método de extração das complexidades sobre as referências, representado pelo Algoritmo 17, também sofreu uma pequena alteração. Não é mais iterado a lista de 'regexes' responsáveis pela extração dos valores de complexidade, como era feito no Algoritmo 14, pois esse trabalho agora é feito pelo LLM na linha 10.

---

**Algoritmo 17** Código mostra a função 'extrairComplexidadeDeReferencia' em alto nível

---

```
1: Function extrairComplexidadeDeReferencia(referencia, regex)
2: for regex in regexes['palavra'] do
3:     matches ← [ ]
4:     palavra ← referencia.find_next(string=compile(regex))

5:     if palavra then:
6:         Insere palavra em matches
7:     end if

8:     match_mais_proximo ← Pega o match mais próximo da ref

9:     for elemento in elementosHTML['valor'] do
10:         complexidade ← llm_searcher.search(complexidade, match_mais_proximo)

11:         if complexidade then:
12:             return complexidade
13:         end if
14:     end for

15:     raise FalhaExtraçãoComplexidade
```

---

## 5 EXPERIMENTOS

Para cada página em que foi implementada um *scraper* foram realizadas avaliações a fim de avaliar se a extração das informações teve um bom resultado ou não. Tentou-se garantir que a quantidade de extrações fosse boa e que as informações extraídas estivessem corretas. A seguir, é apresentada a configuração do ambiente que executou o *scraper*, descrevem as métricas utilizadas, caracterizam os dados extraído e detalham os experimentos e seus respectivos resultados a partir do uso do SACI.

### 5.1 CONFIGURAÇÃO DO AMBIENTE

Os experimentos foram realizados em uma máquina local, com o uso de Python e Visual Studio Code para rodar o *scraper* e com o auxílio do Libre Office para cálculos estatísticos e geração de gráficos. Configurações da máquina: Sistema operacional Microsoft Windows 10, processador AMD Ryzen 5600X, placa de Vídeo: GTX 1060 3GB, 16GB RAM, 1TB SSD.

### 5.2 MÉTRICAS

Para caracterizar a fonte de dados foram criados histogramas e gráficos para exibir o comportamento e a distribuição dos dados extraídos. Esses gráficos ajudam a visualizar padrões e tendências que podem não ser imediatamente aparentes através de uma simples observação.

Para melhor entender os resultados de cada extração foram selecionadas métricas de avaliação. A partir dos resultados é criada uma matriz de confusão e calculada a acurácia, revocação, precisão e F1-Score. Utilizam melhor essas métricas em conjunto do que individualmente, pois cada uma oferece uma perspectiva única sobre a performance do processo. Estas serão usadas para comparar a extração realizada com expressões regulares e a feita com OpenAi.

#### 5.2.1 Histogramas

Um histograma é uma representação gráfica da distribuição de um conjunto de dados. Ele é utilizado para visualizar a frequência com que os diferentes valores ou intervalos de valores ocorrem em um conjunto de dados. Um histograma é composto por barras retangulares (também chamadas de *bins*) onde a altura de cada barra representa a frequência ou a contagem de dados que se enquadram em um determinado intervalo.

### 5.2.2 Matriz confusão

A matriz confusão é uma tabela que resume o desempenho de um modelo em dados de teste; com ela é possível avaliar sua qualidade, pois fornece uma visão detalhada dos acertos e erros do modelo. Ela é organizada em quatro quadrantes quando usada a classificação binária, que representam as quatro possíveis combinações de previsão e realidade:

Verdadeiro Positivo (VP): a previsão do modelo é positiva e a classe real também é positiva. Falso Positivo (FP): a previsão do modelo é positiva, mas a classe real é negativa. Falso Negativo (FN): a previsão do modelo é negativa, mas a classe real é positiva. Verdadeiro Negativo (VN): a previsão do modelo é negativa e a classe real também é negativa.

A partir da matriz de confusão na Figura 10, podem ser calculadas diversas métricas de avaliação, tais como acurácia, precisão, revocação e F1-Score.

Figura 10 – Matriz confusão

		Valor Predito	
		Sim	Não
Real	Sim	Verdadeiro Positivo (TP)	Falso Negativo (FN)
	Não	Falso Positivo (FP)	Verdadeiro Negativo (TN)

Fonte: elaborada pelo autor

### 5.2.3 Acurácia

A acurácia mede a proporção de classificações corretas em relação ao total de classificações feitas pelo modelo. Em outras palavras, a acurácia mede a taxa de sucesso geral do modelo em classificar as amostras corretamente. Por exemplo, se um *scraper* é configurado para extrair informações de um conjunto de 100 páginas da web e consegue obter corretamente os dados de 80 delas, então a acurácia desse seria de 80%.

Vale ressaltar que a acurácia é apenas uma métrica superficial e deve ser avaliada em conjunto com outras métricas. Por exemplo, suponha que o *scraper* esteja coletando dados de duas categorias diferentes, A e B, onde a categoria A representa 90% das páginas e a categoria B representa apenas 10%. Nesse caso, um *scraper*



que sempre extrai informações da categoria A terá uma acurácia de 90%, mas não será útil na prática, já que não está realmente obtendo os dados corretos da categoria B. Nesse cenário, a acurácia alta pode ser enganosa, pois o *scraper* parece estar funcionando bem ao considerar todas as páginas, mas na verdade está falhando em capturar informações importantes da categoria B.

$$\text{acurácia} = (\text{número de classificações corretas}) / (\text{total de amostras})$$

#### 5.2.4 Revocação

A revocação ou 'recall', é uma métrica usada para medir a proporção de verdadeiros positivos que foram identificados corretamente em relação ao total de amostras positivas. A revocação é especialmente relevante em cenários em que os falsos negativos têm um alto custo; ou seja, em casos em que classificar erroneamente uma amostra negativa como positiva pode ter consequências graves.

Consideremos um cenário no qual temos um *scraper* que faz a extração de preços de ações, em que uma amostra pode ser classificada como apta a ser comprada (positivo) ou não (negativo). O custo de um falso negativo é alto, pois pode resultar em uma compra indesejada e causar prejuízos significativos ao usuário. Nesse caso, é crucial que o sistema tenha uma alta revocação; ou seja, que identifique corretamente a maioria dos preços.

$$\text{revocação} = (\text{amostras positivas classificadas corretamente}) / (\text{total de amostras})$$

#### 5.2.5 Precisão

A precisão é uma métrica utilizada para avaliar a proporção de amostras classificadas como positivas que são realmente positivas em relação ao total de amostras classificadas como positivas. A precisão é também importante em situações em que o custo de um falso positivo é alto.

Vamos considerar um cenário em que temos um sistema de *scraper* de mensagens com o objetivo de detectar spam, no qual uma amostra pode ser classificada como spam (positivo) ou não (negativo). O custo de um falso positivo é alto porque pode resultar em mensagens importantes sendo filtradas erroneamente. Nesse caso, é crucial que o sistema tenha uma alta precisão; isto é, que classifique

corretamente a maioria das mensagens positivas como spam.

$$\text{precisão} = (\text{amostras classificadas como positivas corretamente}) / (\text{amostras classificadas como positivas})$$

### 5.2.6 F1-Score

O F1-Score é uma métrica que combina a precisão e a revocação (recall) em uma única medida. Essa métrica é especialmente útil em problemas de classificação desbalanceados, nos quais a acurácia pode ser enganosa quando uma classe é muito mais frequente do que a outra.

Por exemplo, em um conjunto de dados com 99 amostras corretamente extraídas por um *scraper* e apenas 1 amostra errônea, poderíamos dizer que o *scraper* em questão teria uma acurácia de 99%; mas, seu desempenho real na extração seria muito ruim. Nesse caso, o F1-Score seria uma medida mais apropriada para ser utilizada, pois leva em consideração tanto a precisão quanto a revocação do modelo, proporcionando uma visão mais equilibrada da sua eficácia na classificação das amostras positivas corretamente.

$$\text{F1-Score} = 2 * (\text{precisão} * \text{revocação}) / (\text{precisão} + \text{revocação})$$

## 5.3 RESULTADOS

A partir da análise do SACI na versão com expressões regulares, considerando as métricas citadas anteriormente, pode-se observar que os resultados obtidos mostraram um desempenho promissor na extração de ambas as complexidades, conforme demonstram as matrizes de confusão para tempo e espaço das Figuras 11 e 12.

		Valor Predito	
		Sim	Não
Real	Sim	391	11
	Não	88	11

Figura 11 – Matriz confusão, tempo

Fonte: elaborada pelo autor

		Valor Predito	
		Sim	Não
Real	Sim	419	1
	Não	47	34

Figura 12 – Matriz confusão, espaço

Fonte: elaborada pelo autor

A validação dos resultados foi feita de forma manual onde cada URL, onde ocorreu alguma extração de algoritmo, foi acessada e comparados os valores que foram extraídos com os valores presentes na página web.

Nota-se um bom desempenho na identificação de verdadeiros positivos para ambos os tipos de complexidade, com uma quantidade significativa de registros corretamente identificados pelo SACI. Porém, na complexidade de tempo observa-se um número considerável de falso-positivos, totalizando 88 casos. Isso indica que alguns dados errôneos foram incluídos na base de dados. A explicação para existência desse número um pouco mais elevado de falso-positivos é o uso do método de 'fallback', já que o mesmo parte do princípio de ser genérico e coletar a qualquer custo alguma complexidade.

Para falso-negativos, nota-se também sua existência significativa apenas para a complexidade de tempo, com 11 registros, o que sugere que alguns dados foram ignorados durante a extração. A baixa quantidade de verdadeiro-negativos na complexidade de tempo pode significar uma menor capacidade de filtrar dados irrelevantes.

Por outro lado, a complexidade de espaço apresentou um desempenho superior. A quantidade de verdadeiros negativos foi maior, com 34 registros, e os falsos negativos foram mínimos, com apenas 1 registro, mostrando uma melhor performance em comparação com a complexidade de tempo. Isso se deve ao fato de que a complexidade de espaço é apresentada de forma mais padronizada em todas as páginas, resultando em uma melhor identificação dos dados.

Na versão OpenAI notamos uma melhora geral em todos os parâmetros se comparados com a versão padrão, mas alguns apresentaram diferenças significativas.

		Valor Predito	
		Sim	Não
Real	Sim	451	6
	Não	32	12

Figura 13 – Matriz confusão OpenAI, tempo

Fonte: elaborada pelo autor

		Valor Predito	
		Sim	Não
Real	Sim	424	1
	Não	37	39

Figura 14 – Matriz confusão OpenAI, espaço

Fonte: elaborada pelo autor

Para a complexidade de tempo observou-se um grande incremento nos verdadeiro-positivos, de 391 para 451, além de uma diminuição significativa nos falso-positivos que reduziram de 88 para 32, sugerindo a diminuição da inclusão de dados incorretos na base. Para a complexidade de espaço notou-se também um aumento dos verdadeiro-positivos, de 419 para 424, além da redução dos

falso-positivos, de 47 para 37, e o aumento dos verdadeiro-negativos, de 34 para 39, evidenciando, assim, uma significativa melhora na capacidade da nova versão em evitar a inclusão de dados incorretos e discernir quando uma complexidade de espaço não está presente.

Tabela 2 – Métricas das complexidades de tempo e espaço

<b>Métrica</b>	<b>Tempo Re-gex</b>	<b>Espaço Re-gex</b>	<b>Tempo OpenAI</b>	<b>Espaço OpenAI</b>
Acurácia	0,8023	0,9041	0,9241	0,9241
Precisão	0,8162	0,8991	0,9337	0,9197
Revocação	0,9726	0,9976	0,9868	0,9976
F1 Score	0,8876	0,9458	0,9595	0,9571

Fonte: elaborada pelo autor

Os resultados alcançados pelo SACI são motivo de encorajamento. Eles evidenciaram alta precisão, revocação e acurácia em todas as complexidades analisadas. Notadamente, a implementação das técnicas da OpenAI trouxe melhorias substanciais sobre a versão padrão em vários aspectos críticos. Contudo, é importante reconhecer a dependência desses resultados à estrutura da fonte de dados. Mudanças futuras nos sites podem influenciar o desempenho da ferramenta.

Ao examinar as métricas de desempenho, confirma-se que o SACI se destaca como uma solução robusta e eficiente na coleta de complexidades. As altas taxas de precisão e revocação asseguram a habilidade em extrair informações com o mínimo de erros, um ponto crucial para a integridade dos dados coletados, garantindo a validade das informações e evitando dados imprecisos ou irrelevantes.

Além disso, a elevada acurácia e os índices de F1-Score reforçam a capacidade de realizar classificações precisas, otimizando o equilíbrio entre a detecção de dados relevantes e a exclusão de falso-positivos. A adoção de tecnologias avançadas em processamento de linguagem natural e aprendizado de máquina da OpenAI foi decisiva nesse avanço, proporcionando uma compreensão contextual mais refinada e uma extração mais acertada das complexidades.

A versão aprimorada do SACI, equipada com os modelos da OpenAI, mostrou uma redução expressiva tanto em falso-positivos quanto em falso-negativos, um diferencial que potencializa a confiabilidade dos dados extraídos. Tal performance sugere que, mesmo diante de estruturas de dados variáveis nos sites-alvo, a ferramenta possui a versatilidade para se adaptar e manter alto padrão de qualidade.

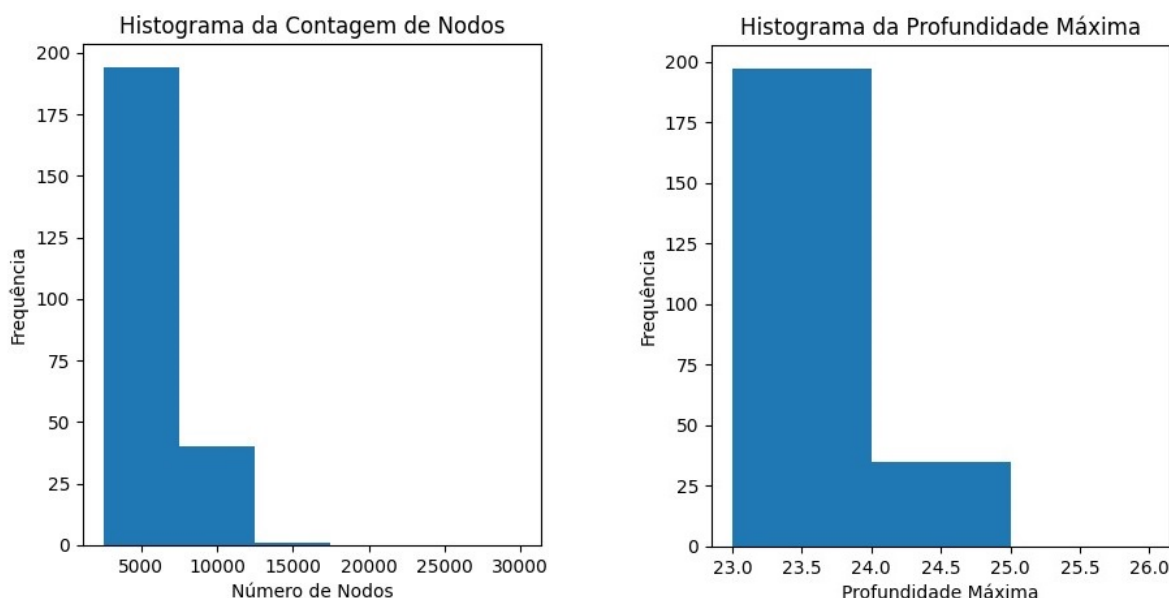
As métricas obtidas solidificam a confiança na ferramenta e pavimentam o caminho para seu contínuo desenvolvimento, com foco na adaptação para diversos domínios e na capacidade de enfrentar mudanças estruturais nos sites de origem.

## 5.4 FONTE DE DADOS

Todos os dados usados neste trabalho tiveram sua extração no dia '24/08/2023', conforme descrito previamente na Seção 4.4. Foi realizada uma análise sobre cada página onde houve a extração de uma ou mais complexidades. De um total de 235 URLs, com 501 algoritmos disponíveis através delas, para cada página, executou-se uma função que percorria todos os nodos da árvore HTML e retornava qual a profundidade máxima encontrada e quantidade de nodos existentes.

Tal análise foi feita porque o número de nodos HTML em uma página tem impacto direto na velocidade das operações de extração; quanto maior o número de nós, maior será o tempo necessário para percorrer a árvore de elementos e extrair as informações desejadas. A maioria das profundidades máximas fica em 24 nodos e, no caso da quantidade de nodos, entre 5 e 15 mil nodos é o mais frequente.

Figura 15 – Caracterização dos Nodos



Fonte: elaborada pelo autor

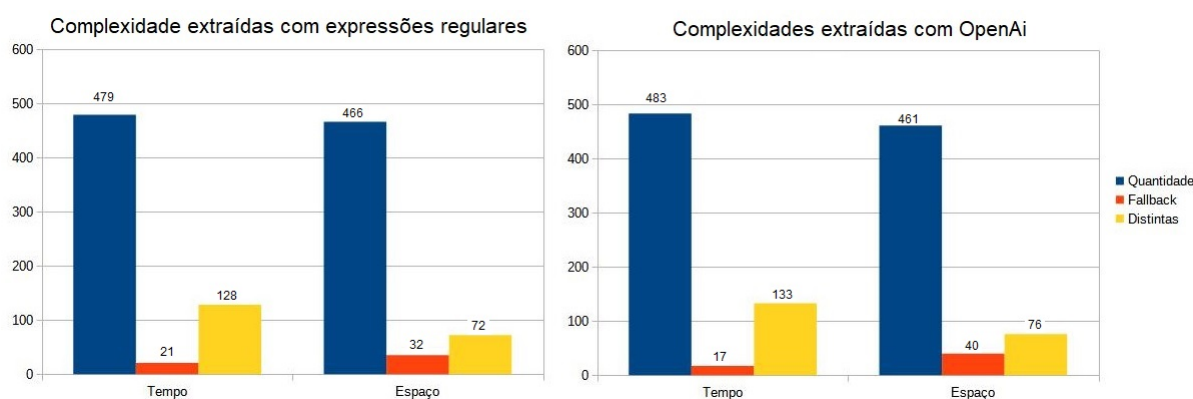
### 5.4.1 Conjunto de dados extraídos

A caracterização do conjunto de dados extraídos foi feita, primeiramente, avaliando o quanto a ferramenta de coleta utilizou da técnica de 'fallback' para extrair as complexidades. Os títulos dos algoritmos e os outros atributos coletados não fazem

parte da caracterização. Conforme a Figura ??, 479 algoritmos tiveram suas complexidades de tempo extraídas e 466 suas complexidades de espaço; dentre estes, 21 fizeram uso do 'fallback' para extrair a complexidade de tempo e 32 para a de espaço. Foram obtidas 128 complexidades de tempo distintas e 72 complexidades de espaço distintas.

Um dos motivos para não se ter todas as complexidades de tempo extraídas com sucesso é a falta de um padrão na página. Em algumas páginas as complexidades estão escritas de maneira diferente, em locais diferentes, fazendo, assim, com que a expressão regular não as identifique. No caso da complexidade de espaço, alguns algoritmos realmente não a possuem e, em outros casos, ocorreu o mesmo problema que com a extração da complexidade de tempo.

Figura 16 – Complexidades extraídas



Fonte: elaborada pelo autor

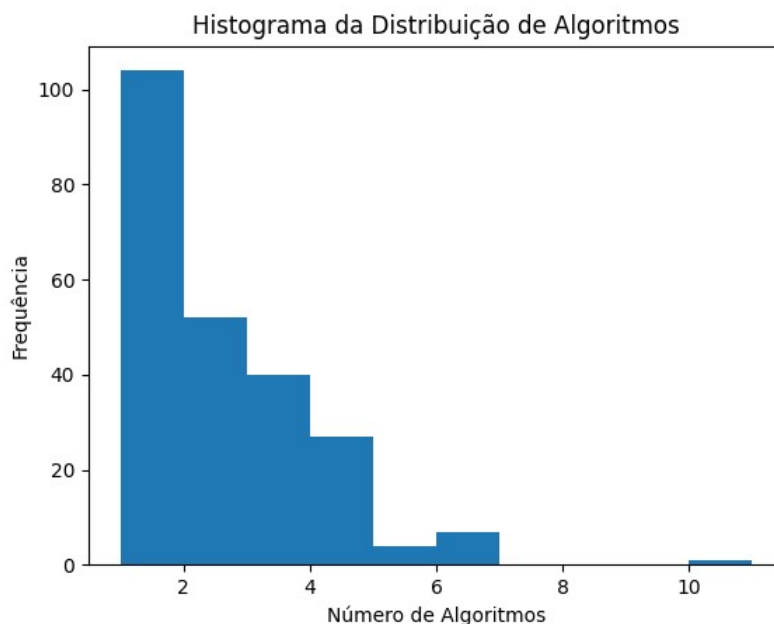
Na nova versão do método de extração, que incorpora as capacidades da OpenAI, observou-se um incremento na eficiência da extração de dados. Como exibido na Figura ??, esta versão avançada conseguiu extrair com sucesso as complexidades de tempo de 483 algoritmos, e as complexidades de espaço de 461 algoritmos, superando um pouco a marca anterior.

Além disso, a dependência da técnica de 'fallback' foi reduzida para apenas 17 casos na extração das complexidades de tempo indicando, com isso, uma melhoria. Nas complexidades de espaço houve um pequeno aumento de 32 para 40 casos. Em termos de diversidade, o SACI identificou 133 complexidades de tempo distintas e 76 de espaço, um aumento em relação às 128 e 72 variantes respectivamente encontradas na versão que utilizava expressões regulares.

O número de URLs se manteve o mesmo em ambas as versões e é menor que o número de algoritmos extraídos, pois, cada página pode conter mais de um algoritmo; para exemplificar, foi gerado um histograma na Figura 17. Existem em média 2.27

algoritmos em cada página, sendo o máximo 10 e o mínimo 1.

Figura 17 – Distribuição de Algoritmos



Fonte: elaborada pelo autor

Ao realizar um teste de uma das URLs do site<sup>1</sup> com ambas as versões, notamos que o *scraper* retorna como valores de complexidade ' $O(N)$ ' para a de tempo e ' $O(1)$ ' para a de espaço. Acessando o site através da própria URL exposta no retorno, pode-se confirmar que os dados foram extraídos com sucesso, conforme os valores da página exibida na Figura 18.

<sup>1</sup> <https://www.geeksforgeeks.org/linear-search/>

Figura 18 – Linear Search

Below is the implementation of the above approach:

```
C    C++   Java   Python3   C#    PHP    Javascript
# Python3 code to linearly search x in arr[].
# If x is present then return its location,
# otherwise return -1

def search(arr, N, x):
    for i in range(0, N):
        if (arr[i] == x):
            return i
    return -1

# Driver Code
if __name__ == "__main__":
    arr = [2, 3, 4, 10, 40]
    x = 10
    N = len(arr)

    # Function call
    result = search(arr, N, x)
    if(result == -1):
        print("Element is not present in array")
    else:
        print("Element is present at index", result)
```

### Output

```
Element is present at index 3
```

**Time complexity:**  $O(N)$

**Auxiliary Space:**  $O(1)$

Fonte: elaborada pelo autor

Um exemplo de caso com problemas seria o do algoritmo encontrado no site <sup>2</sup>, onde as complexidades de tempo aparecem em uma lista, cada uma associada a uma determinada operação, o que causa problemas quando o código tenta realizar a extração, apresentado na Figura 19.

<sup>2</sup> <https://www.geeksforgeeks.org/introduction-and-array-implementation-of-queue/>



Figura 19 – Problema na extração

**Complexity Analysis:**

- **Time Complexity**

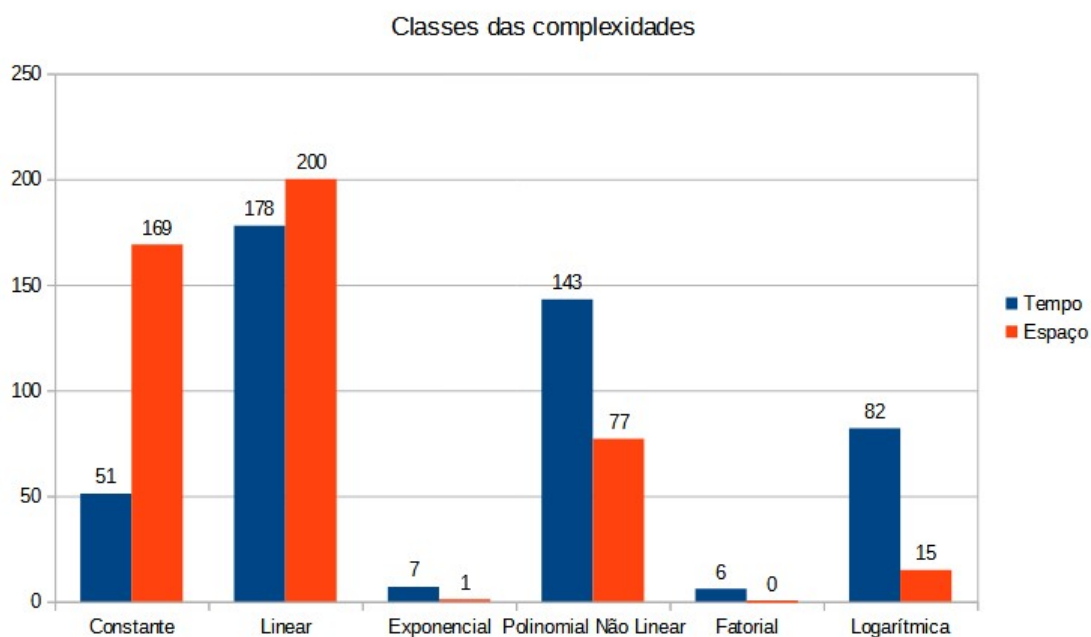
Operations	Complexity
Enqueue(insertion)	O(1)
Deque(deletion)	O(1)
Front(Get front)	O(1)
Rear(Get Rear)	O(1)
IsFull(Check queue is full or not)	O(1)
IsEmpty(Check queue is empty or not)	O(1)

- **Auxiliary Space:**  
O(N) where N is the size of the array for storing elements.

Fonte: elaborada pelo autor

Para se ter um panorama mais detalhado, também foi analisado a distribuição das classes de complexidade extraídas e chegou-se ao resultado apresentado na Tabela 20.

Figura 20 – Classe das complexidades



Fonte: elaborada pelo autor

Conforme Tabela 20, ao analisar as complexidades de tempo, observa-se que a classe *linear* é a predominante com 178 entradas. Ela é seguida pelas classes *polinomial* com 143 entradas; *logarítmica* com 82; e *constante* com 51. As classes *fatorial* e *exponencial* apresentam números mais baixos, com 6 e 7 entradas, respectivamente. Isso sugere que, embora haja um esforço para representar todas as classes na base de dados, algumas têm uma presença reduzida no site analisado.

Um fato notável é que as complexidades *fatorial* e *exponencial* são, geralmente, menos eficientes, resultando em sua menor utilização em aplicações práticas. Esta tendência pode ser interpretada como uma inclinação da comunidade em favor de abordagens que evitem tais complexidades ineficientes. Ao explorar as classes de complexidade de espaço, conforme indicado na mesma tabela, percebe-se que a classe *linear* se sobressai, contabilizando 200 entradas. A classe *constante* também apresenta uma boa representatividade, com 169 entradas, seguida pela *polinomial*, com 77. As classes *exponencial* e *fatorial* são menos frequentes, ambas com apenas 1 e 0 entrada, respectivamente, enquanto a *logarítmica* registra 15.

Esta distribuição reforça a ideia de que há uma preferência por complexidades mais práticas e gerenciáveis em termos de uso de espaço. A baixa representatividade das complexidades *exponencial* e *fatorial* reflete a natureza ineficiente destas quando se trata de consumo de memória. A presença mais discreta da complexidade *logarítmica*, por sua vez, pode indicar sua aplicação em contextos mais específicos.

Em resumo, os dados mostram uma clara inclinação da base para soluções

que equilibram eficiência e praticidade, priorizando algoritmos que otimizem o uso de recursos. A versão com OpenAI não apresentou diferenças significativas. Algumas entradas diminuíram e outras aumentaram, mas a ordem de representatividade se manteve.

A classificação das complexidades em ambas versões foi feita através de expressões regulares, exibidas no Código 7 utilizando-se de um bloco de código de condicionais. Caso a expressão capturasse algo na complexidade, significa que tal complexidade é pertencente a tal categoria.

```
1  Constante = r'O(1)'
   Linear = r'O(w)'
3  Exponencial = r'O(d^?w)'
   Fatorial = r'O(*!.*)'
5  Logarítmica = r'O(* log.*)'
   Polinomial = r'O(w[s*[\^*\|+]?s*[\d\w\s]+.*))' e r'O(\d+(\w*\w))'
```

Código 7 – Regex de classificação de complexidades

## 6 CONCLUSÕES E TRABALHOS FUTUROS

A criação de ferramentas capazes de identificar a complexidade de um algoritmo e, através disso, gerar *insights* para o usuário final sobre o código que está sendo analisado, requer uma base de dados com algoritmos e suas complexidades, de forma que tais ferramentas possam aproveitar-se e aprender com eles a partir de técnicas de inteligência artificial.

O SACI apresentou resultados bastantes satisfatórios. A quantidade de códigos e complexidades coletadas das páginas do site formam uma base de dados adequada e as métricas calculadas comprovam a eficácia da ferramenta.

Durante o desenvolvimento, foi perceptível a dificuldade existente na construção de um *scraper* genérico, principalmente devido à grande diversidade de layouts de páginas. No entanto, a forma como se desenhou a ferramenta permitiu a reutilização de certos trechos de códigos, facilitando, com isso, a implementação futura para outras páginas. O maior obstáculo, porém, consistia em detectar e extrair os elementos desejados. Neste caso, a estratégia escolhida foi o uso de expressões regulares juntamente com o método 'fallback', utilizado quando as expressões regulares não obtinham êxito na extração.

Quanto a questões de impedimentos por 'anti-bots' ou sistemas do gênero, estes não implicaram em problemas para a página coletada, bastando o simples envio dos cabeçalhos adequados e a criação de login para evitar futuros transtornos. Todavia, é preciso atentar-se para a quantidade de requisições enviadas simultaneamente, pois a ferramenta faz isso de forma assíncrona e pode ser facilmente notado caso exceda um número normal. Assim sendo, criou-se um argumento para limitar as requisições e, também, para adicionar um intervalo entre as mesmas, mitigando, assim, possíveis eventos de alerta.

Por questões de tempo ou por escolhas quando da implementação, algumas funções da ferramenta não foram desenvolvidas completamente, estando, portanto, passíveis de melhorias em futuros projetos e desenvolvimentos. São elas:

1. Melhorar as métricas: Atualmente o SACI só está levando em questão as complexidades de tempo e espaço em suas métricas, é possível adicionar mais atributos a serem acompanhados.
2. Melhorar as técnicas de extração: As expressões regulares e lógicas utilizadas nas extrações dos dados desejados podem ser redesenhadas e melhoradas a fim de obter mais sucesso na tarefa. Também é possível usar inteligência artificial além do Chat GPT para as extrações, treinando um modelo especificamente para o caso de uso que identifique os dados relevantes da página e usar isso para extraí-los.

3. Fonte de dados: Realizar a implementação em mais fontes de dados a fim de coletar mais códigos e complexidades, gerando uma base de dados mais rica.

Além destas melhorias, a base de dados proveniente dos resultados já existentes na ferramenta contém informações valiosas para pesquisadores da área, que podem utilizar os 501 algoritmos extraídos em trabalhos que dependam de uma base de dados com tais características e, desta forma, melhorá-los, otimizá-los ou desenvolver algo novo a partir destes. Dentre os possíveis usos para essa base destacam-se:

1. Análise Comparativa de Algoritmos: Uma das aplicações mais óbvias que utiliza uma base de dados de algoritmos e suas complexidades, é a análise comparativa de algoritmos. Os dados fornecidos podem ser usados para comparar o desempenho de diferentes algoritmos em termos de tempo de execução e uso de espaço. Isso permite identificar algoritmos mais eficientes para determinados problemas e fornece *insights* sobre as características de desempenho de cada algoritmo.
2. Otimização de Algoritmos: É possível realizar trabalhos de otimização de algoritmos com as informações obtidas. Os dados podem revelar oportunidades para melhorar algoritmos existentes, seja reduzindo a complexidade de tempo e espaço ou refinando suas implementações. Por exemplo, é possível identificar partes do algoritmo que podem ser substituídas por abordagens mais eficientes ou adaptar o algoritmo para melhor lidar com entradas específicas para casos específicos.
3. Previsão de Desempenho de Algoritmos: A base de dados pode ser utilizada para construir modelos preditivos que estimem o desempenho de algoritmos com base em suas características e nos tamanhos dos conjuntos de dados de entrada. Isso pode ser útil para tomar decisões sobre o uso de um determinado algoritmo em uma situação específica. Os modelos preditivos podem ser usados para prever o tempo de execução ou o consumo de espaço de um algoritmo com base em dados históricos e características conhecidas do problema.
4. Estudo Teórico de Algoritmos: Também pode-se usar a base para fins teóricos. É possível realizar estudos para investigar as características dos algoritmos em relação a suas complexidades de tempo e espaço. Esses estudos podem envolver a análise de classes de algoritmos, a identificação de limites teóricos para determinados problemas computacionais ou a análise de propriedades gerais dos algoritmos em relação às suas complexidades.
5. Desenvolvimento de Novos Algoritmos: Os dados disponíveis na base de dados podem servir como inspiração e orientação para o desenvolvimento de novos

algoritmos. Ao estudar as complexidades de tempo e espaço de algoritmos existentes, é possível identificar lacunas ou problemas não resolvidos e propor novas soluções. Por exemplo, gerar um novo algoritmo de *sorting* baseado nas melhores características dos já existentes.

A base de dados construída oferece uma riqueza de informações que podem ser exploradas em diferentes áreas. A análise comparativa de algoritmos, a otimização, a previsão de desempenho, o estudo teórico e o desenvolvimento de novos algoritmos são apenas alguns exemplos dos trabalhos possíveis que podem ser realizados a partir dessa base de dados.

## REFERÊNCIAS

BRANCO, Arthur Machado. Ferramenta para coleta e comparação de dados de publicações acadêmicas dos professores com o currículo lattes. *In: UFSC*. Citado na p. 23.

CORMEN, Thomas H. *et al.* **Algoritmos: Estruturas de Dados e Algoritmos para Resolução de Problemas**. Rio de Janeiro: LTC, 2009. Citado na p. 12.

CORMEN, Thomas H. *et al.* **Introduction to Algorithms**. 2nd. [S.l.]: The MIT Press, 2001. ISBN 0262032937. Disponível em:

<<http://www.amazon.com/Introduction-Algorithms-Thomas-H-Cormen/dp/0262032937%3FSubscriptionId%3D13CT5CVB80YFWJEPWS02%26tag%3Dws%26linkCode%3Dxm%26camp%3D2025%26creative%3D165953%26creativeASIN%3D0262032937>>. Citado na p. 17.

DIOUF, Rabiyatou *et al.* Web Scraping: State-of-the-Art and Areas of Application. *In: 2019 IEEE International Conference on Big Data (Big Data)*. [S.l.: s.n.], 2019. P. 6040–6042. DOI: 10.1109/BigData47090.2019.9005594. Citado na p. 15.

DIOUF, Rabiyatou *et al.* Web scraping: state-of-the-art and areas of application. *In: IEEE. 2019 IEEE International Conference on Big Data (Big Data)*. [S.l.: s.n.], 2019. P. 6040–6042. Citado na p. 16.

HALL, Mark; CHAPMAN, Barry. Measuring algorithm complexity. **ACM SIGPLAN Notices**, ACM Press, v. 48, n. 1, p. 16–23, 2015. Citado na p. 12.

KERNER, Sean Michael. **What are large language models (LLMs)?** 2023.

Disponível em:

<<https://www.techtarget.com/whatis/definition/large-language-model-LLM>>. Acesso em: 5 nov. 2022. Citado na p. 20.

KLEINBERG, J.; TARDOS, É. **Algorithm Design**. [S.l.]: Pearson/Addison-Wesley, 2006. (Alternative Etext Formats). ISBN 9780321295354. Disponível em:

<<https://books.google.com.br/books?id=0iGhQgAACAAJ>>. Citado nas pp. 17, 18.

LESSA, Marcos Aurélio. CrawlEX: uma ferramenta para extração de dados na web configurável através de exemplos. *In: UFSC*. Citado na p. 24.

MATHIAS, Gilney; DORNELES, Carina. qFEx - um crawler para busca e extração de questionários de pesquisa em documentos HTML. *In: ANAIS do III Dataset Showcase Workshop*. Rio de Janeiro: SBC, 2021. P. 1–10. DOI: 10.5753/dsw.2021.17409.

Disponível em: <<https://sol.sbc.org.br/index.php/dsw/article/view/17409>>.

Citado na p. 22.

MOONEY, Stephen J; WESTREICH, Daniel J; EL-SAYED, Abdulrahman M. Epidemiology in the era of big data. **Epidemiology (Cambridge, Mass.)**, NIH Public Access, v. 26, n. 3, p. 390, 2015. Citado na p. 15.

PFITSCHER, Ricardo J. *et al.* Estimating Code Running Time Complexity with Machine Learning. *In: NALDI, Murilo C.; BIANCHI, Reinaldo A. C. (Ed.). Intelligent Systems*. Cham: Springer Nature Switzerland, 2023. P. 400–414. Citado na p. 12.

RODENBUSCH, Gabriel Braun. **Desenvolvimento de API para predição da complexidade de tempo de execução de códigos por meio de AutoML**. [S.l.: s.n.], 2023. TCC (graduação) - Universidade Federal de Santa Catarina, Campus Joinville, Engenharia Mecatrônica. Disponível em: <<https://repositorio.ufsc.br/handle/123456789/252875>>. Acesso em: 7 dez. 2023. Citado nas pp. 13, 24.

SEDGEWICK, Robert. Analyzing algorithms. **Communications of the ACM**, ACM Press, v. 52, n. 1, p. 139–148, 2009. Citado na p. 12.

SERVIAN, Alexandre. **Regex: Um guia prático para expressões regulares**. 2020. Disponível em: <<https://medium.com/xp-inc/regex-um-guia-pratico-para-express%C3%B5es-regulares-1ac5fa4dd39f>>. Acesso em: 10 dez. 2022. Citado na p. 19.

SIKKA, Jagriti *et al.* Learning Based Methods for Code Runtime Complexity Prediction. *In: SPRINGER. EUROPEAN Conference on Information Retrieval*. [S.l.: s.n.], 2020. P. 313–325. Citado nas pp. 12, 17.

SIKKA, Jagriti *et al.* Learning Based Methods for Code Runtime Complexity Prediction. *In: JOSE, Joemon M. et al. (Ed.). Advances in Information Retrieval*. Cham: Springer International Publishing, 2020. P. 313–325. Citado na p. 22.

SNYDER, Russell. **Web search engine with graphic snapshots**. [S.l.]: Google Patents, nov. 2003. US Patent 6,643,641. Citado na p. 16.



WANG, Meng *et al.* A Survey on Large-Scale Machine Learning. **IEEE Transactions on Knowledge and Data Engineering**, v. 34, n. 6, p. 2574–2594, 2022. DOI: 10.1109/TKDE.2020.3015777. Citado na p. 13.

ZHAO, Bo. Web scraping. **Encyclopedia of big data**, Springer Living ed. Cham, p. 1–3, 2017. Citado nas pp. 13, 15.

ZHAO, Wayne Xin *et al.* **A Survey of Large Language Models**. [S.l.: s.n.], 2023. arXiv: 2303.18223 [cs.CL]. Citado na p. 21.

# Apêndice A - Código Fonte Versão

## Expressões Regulares

- saci/requirements.txt

```
aiohttp
2 beautifulsoup4
black
4 pandas
pydantic
6 pymongo
pyppeteer
8 requests
scikit-learn
10 scipy
seaborn
12 mlxtend
requests-html
14 pip-tools
openai
16 langchain
asyncio-throttle
```

- saci/makefile

```
1 PYTHON :=
  ifeq ($(OS),Windows_NT)
3     PYTHON = venv/Scripts/python
     PIP = venv/Scripts/pip
5 else
     PYTHON = venv/bin/python3
7     PIP = venv/bin/pip
  endif
9
requirements: requirements.txt
11 $(PIP) install -r requirements.txt
```

```
13 run_scraper :
    $(PYTHON) -m main
15
17 scrape_url :
    $(PYTHON) -m scraper.scripts --script=scrape_url --url=$(url) --scraper=$(scraper)
19
19 results_analysis :
    $(PYTHON) -m scraper.scripts --script=results_analysis --scraper=$(scraper)
21
23 manual_results_boilerplate :
    $(PYTHON) -m scraper.scripts --script=manual_results_boilerplate --scraper=$(scraper)
```

- saci/main.py

```
1 from asyncio import get_event_loop
3 from scraper.runner import run_scrapers
5 if __name__ == "__main__":
    loop = get_event_loop()
7    loop.run_until_complete(run_scrapers())
```

- saci/scraper/\_\_init\_\_.py
- saci/scraper/database.py

```
1 import os
from typing import Dict, List
3
4 from pymongo import MongoClient
5 from pymongo.collection import Collection
6 from pymongo.operations import UpdateOne
7
8 from scraper.settings import DATABASE_NAME
9
10 FIELDS_TO_UPDATE = [
11     "name",
12     "time_complexity",
```

```
13     "trustable_time_complexity",
14     "space_complexity",
15     "trustable_space_complexity",
16 ]
17
18
19 class ScraperDatabase:
20     def __init__(self):
21         self.client: MongoClient = MongoClient(
22             os.environ.get("MONGO_CONNECTION_STRING"),
23             serverSelectionTimeoutMS=1
24         )
25         self.database = self.client[DATABASE_NAME]
26
27     def _get_collection(self, collection_name: str) -> Collection:
28         return self.database[collection_name]
29
30     def update_database(self, collection_name: str, data: List[Dict[str,
31         any]]):
32         collection = self._get_collection(collection_name)
33         operations = [
34             UpdateOne(
35                 filter={"url": item["url"], "codes": item["codes"]},
36                 update={"$set": {field: item[field] for field in
37                     FIELDS_TO_UPDATE}},
38                 upsert=True,
39             )
40             for item in data
41         ]
42         collection.bulk_write(operations)
```

- saci/scrapper/exceptions.py

```
1 class InvalidUrlException(Exception):
2     def __init__(self, url):
3         super().__init__(f"The provided URL '{url}' is not valid.")
4
5
6
7 class TooManyRetrysException(Exception):
8     pass
9
10
11 class FailedExtraction(Exception):
12     """Failed the extraction process"""
```

```
13     pass
15
16 class FailedComplexityExtraction(Exception):
17     """Failed to extract the complexity"""
18
19     pass
```

- saci/scrapper/runner.py

```
1 from asyncio import gather
2
3 from scrapper.database import ScrapperDatabase
4 from scrapper.settings import USE_MONGO_DATABASE
5 from scrapper.utils import write_results_to_json
6 from scrapper.websites import websites
7
8
9 async def run_scrapers():
10     scrapers_results = await gather(*[scrapper() for scrapper in websites.
11                                     values()])
12     for name, data in zip(websites.keys(), scrapers_results):
13         if USE_MONGO_DATABASE:
14             database = ScrapperDatabase()
15             database.update_database(name, [result.dict() for result in
16                                         data])
17         else:
18             write_results_to_json(
19                 name,
20                 sorted([result.dict() for result in data], key=lambda alg:
21                       alg["url"]),
22             )
```

- saci/scrapper/scripts.py

```
1 import argparse
2 import pprint
3 from asyncio import get_event_loop
4 from importlib import import_module
5 from scrapper.observability.metrics import (
6     make_results_analysis,
7     make_manual_results_boilerplate,
```

```
)
9
11 def extract_single_url(url: str, scraper: str):
    module = import_module(f"scraper.websites.{scraper}")
13     spider = getattr(module, "spider")
    loop = get_event_loop()
15     result = loop.run_until_complete(spider.run(url))
    pprint.PrettyPrinter(indent=4).pprint(result)
17
19 def main():
    parser = argparse.ArgumentParser()
21     parser.add_argument("--script")
    parser.add_argument("--url", help="URL to be extracted.")
23     parser.add_argument("--scraper", help="Scraper to be used.")
    args = parser.parse_args()
25
    actions = {
27         "scrape_url": lambda: extract_single_url(args.url, args.scraper),
        "results_analysis": lambda: make_results_analysis(args.scraper),
29         "manual_results_boilerplate": lambda:
            make_manual_results_boilerplate(
31             args.scraper
            ),
    }
33
    action = actions.get(args.script)
35     if action:
        action()
37     else:
        print(f"Unknown script: {args.script}")
39
41 if __name__ == "__main__":
    main()
```

- saci/scraper/settings.py

```
2 DATABASE_NAME = "COMPLEXITY_SCRAPER"
USE_MONGO_DATABASE = False
```

- saci/scrapper/utils.py

```
1 import json
3 from scrapper.exceptions import TooManyRetrysException
  from scrapper.observability.log import scrapper_log as log
5
7 def retry(
  times: int,
  raise_exception=True,
  return_value=None,
11 ):
  def func_wrapper(f):
13     async def wrapper(*args, **kwargs):
15         for _ in range(times):
17             try:
19                 return await f(*args, **kwargs)
21             except Exception:
23                 pass
25
27             if not raise_exception:
29                 return return_value
31             else:
33                 raise TooManyRetrysException
35
37         return wrapper
39
41     return func_wrapper
43
45 def remove_duplicates(data: list) -> list:
47     return list(set(data))
49
51 def open_results_from_json(file_path: str) -> dict:
53     with open(file_path, "r") as file:
55         file_contents = file.read()
57         results = json.loads(file_contents)
59     return results
61
63 def write_results_to_json(file_path: str, data: list):
65     with open(f"{file_path}.json", "w") as file:
67         json.dump(data, file)
```

- saci/observability/\_\_init\_\_.py
- 

- saci/observability/log.py

```
import platform
2
from logging import DEBUG, Filter, Formatter, StreamHandler, getLogger
4
from scraper.observability.settings import SCRAPER_LOG_FORMAT,
    SESSION_LOG_FORMAT
6
8 class AddScraperName(Filter):
    def filter(self, record):
10         if not hasattr(record, "scraper"):
                split_char = "\\\" if "Windows" in platform.platform() else "/"
12             record.scraper = record.pathname.split(split_char)[-2]
            return True
14
16 def build_scraper_log(logger_name):
    log_handler = StreamHandler()
18     log_handler.setLevel(DEBUG)
    log_handler.setFormatter(Formatter(SCRAPER_LOG_FORMAT))
20
    log = getLogger(logger_name)
22     log.addFilter(AddScraperName())
    log.addHandler(log_handler)
24     return log
26
28 def build_session_log(logger_name):
    log_handler = StreamHandler()
    log_handler.setLevel(DEBUG)
30     log_handler.setFormatter(Formatter(SESSION_LOG_FORMAT))
32
    log = getLogger(logger_name)
    log.addHandler(log_handler)
34     return log
36
38 scraper_log = build_scraper_log("scraper_log")
session_log = build_session_log("session_log")
```



- saci/observability/metrics.py

```
import math
2 import re
from collections import Counter
4 from typing import Any, Dict, List, Tuple, Union

6 import matplotlib.pyplot as plt
import numpy
8 import requests
from bs4 import BeautifulSoup, NavigableString, Tag
10 from mlxtend.plotting import plot_confusion_matrix

12 from scraper.utils import (
    open_results_from_json,
14     remove_duplicates,
    write_results_to_json,
16 )

18
def make_results_characterization(scraper: str) -> Dict[str, int]:
20     results = open_results_from_json(f"./results/{scraper}.json")

22     return {
        "Distinct Time Complexities": len(
24         {algorithm["time_complexity"] for algorithm in results}
        ),
        "Distinct Space Complexities": len(
26         {algorithm["space_complexity"] for algorithm in results}
28     ),
        "Non Trusted Time Complexities": sum(
30         [1 for algorithm in results if not algorithm["
            trustable_time_complexity"]]
        ),
        "Non Trusted Space Complexities": sum(
32         [1 for algorithm in results if not algorithm["
            trustable_space_complexity"]]
34     ),
    }

36
38 def make_completion_rate(scraper: str) -> Dict[str, int]:
    results = open_results_from_json(f"./results/{scraper}.json")

40
    total = len(results)
42     time_complexity, space_complexity = sum(
        1 for algorithm in results if not algorithm["time_complexity"]
44     ), sum(1 for algorithm in results if not algorithm["space_complexity"])
```

```
46     urls_with_problem = remove_duplicates(  
47         [  
48             algorithm["url"]  
49             for algorithm in results  
50             if not algorithm["time_complexity"] or not algorithm["  
51                 space_complexity"]  
52         ]  
53     )  
54     return {  
55         "Total Algorithms": total ,  
56         "Time Complexity Extracted": f"{total - time_complexity} extracted  
57             - {100 - (time_complexity / total) * 100}%",  
58         "Space Complexity Extracted": f"{total - space_complexity}  
59             extracted - {100 - (space_complexity / total) * 100}%",  
60         "URLs with problem": urls_with_problem ,  
61     }  
62  
63 def calculate_metrics(  
64     true_positive: int , false_positive: int , true_negative: int ,  
65     false_negative: int  
66 ) -> Dict[str , float]:  
67     accuracy = (true_positive + true_negative) / (  
68         true_positive + true_negative + false_positive + false_negative  
69     )  
70     precision = true_positive / (true_positive + false_positive)  
71     recall = true_positive / (true_positive + false_negative)  
72     f1_score = (2 * precision * recall) / (precision + recall)  
73  
74     return {  
75         "True Positive": true_positive ,  
76         "False Positive": false_positive ,  
77         "True Negative": true_negative ,  
78         "False negative": false_negative ,  
79         "Accuracy": accuracy ,  
80         "Precision": precision ,  
81         "Recall": recall ,  
82         "F1 Score": f1_score ,  
83     }  
84  
85 def plot_confusion(  
86     title: str ,  
87     true_positive: int ,  
88     false_negative: int ,  
89     false_positive: int ,  
90     true_negative: int ,
```

```
) :
90     array = numpy.array(
91         [[true_positive , false_negative], [false_positive , true_negative]]
92     )
93
94     fig , ax = plot_confusion_matrix(conf_mat=array)
95     plt.title("Space Complexity Confusion Matrix")
96     plt.show()
97
98
99 def make_confusion_matrix(scraper: str) -> Dict[str , Dict[str , Any]]:
100     manual_results = open_results_from_json(f"./results/manual_{scraper}.
101         json")
102     results = open_results_from_json(f"./results/{scraper}.json")
103
104     time_metrics = {
105         "true_positive": 0,
106         "false_positive": 0,
107         "true_negative": 0,
108         "false_negative": 0,
109     }
110
111     space_metrics = {
112         "true_positive": 0,
113         "false_positive": 0,
114         "true_negative": 0,
115         "false_negative": 0,
116     }
117
118     for result , manual_result in zip(results , manual_results):
119         # Time complexity
120         if time_complexity := result["time_complexity"]:
121             if time_complexity == manual_result["time_complexity"]:
122                 time_metrics["true_positive"] += 1
123             else:
124                 time_metrics["false_positive"] += 1
125         else:
126             if not manual_result["time_complexity"]:
127                 time_metrics["true_negative"] += 1
128             else:
129                 time_metrics["false_negative"] += 1
130
131         # Space complexity
132         if space_complexity := result["space_complexity"]:
133             if space_complexity == manual_result["space_complexity"]:
134                 space_metrics["true_positive"] += 1
135             else:
136                 space_metrics["false_positive"] += 1
```

```
136     else:
137         if not manual_result["space_complexity"]:
138             space_metrics["true_negative"] += 1
139         else:
140             space_metrics["false_negative"] += 1
141
142     # Calculating metrics
143     time_complexity_metrics = calculate_metrics(**time_metrics)
144     space_complexity_metrics = calculate_metrics(**space_metrics)
145
146     # Plotting the confusion matrices
147     plot_confusion("Time Complexity Confusion Matrix", **time_metrics)
148     plot_confusion("Space Complexity Confusion Matrix", **space_metrics)
149
150     # Return result
151     return {
152         "Time Complexity": time_complexity_metrics,
153         "Space Complexity": space_complexity_metrics,
154     }
155
156 def make_manual_results_boilerplate(scraper: str) -> None:
157     results = open_results_from_json(f"./results/{scraper}.json")
158
159     manual_results = [
160         {"url": alg["url"], "time_complexity": " ", "space_complexity": " "
161         }
162         for alg in results
163     ]
164
165     write_results_to_json(f"./results/manual_{scraper}", manual_results)
166
167 def make_complexitys_classification(scraper: str) -> Dict[str, Dict[str,
168 int]]:
169     results = open_results_from_json(f"./results/{scraper}.json")
170
171     time_classification = classify_complexity(
172         [result["time_complexity"] for result in results]
173     )
174     space_classification = classify_complexity(
175         [result["space_complexity"] for result in results]
176     )
177
178     return {
179         "Time Complexity Classification": time_classification,
180         "Space Complexity Classification": space_classification,
181     }
```

```
182
184 def classify_complexity(complexities: List[str]) -> Dict[str, int]:
    constant, linear, exponential, polynomial, factorial, logarithmic = 0,
    0, 0, 0, 0, 0
186 for complexity in complexities:
    if complexity:
188         if re.match(r"O(1)", complexity):
            constant += 1
190         elif re.match(r"O(\w)", complexity):
            linear += 1
192         elif re.match(r"O(.*!.*)", complexity):
            factorial += 1
194         elif re.match(r"O(.*log.*)", complexity):
            logarithmic += 1
196         elif re.match(r"O(\d(?:^?\w)", complexity):
            exponential += 1
198         elif re.match(r"O(\w\s*(\^*\+)?\s*(\d\w\s)+.*)", complexity)
            or re.match(
200             r"O(\d+(\w*\w)", complexity
            ):
                polynomial += 1
202
    return {
204         "Constant complexity": constant,
206         "Linear complexity": linear,
208         "Exponential complexity": exponential,
210         "Polynomial complexity": polynomial,
212         "Factorial complexity": factorial,
214         "Log complexity": logarithmic,
    }
216
218 def make_algorithms_histogram(scraper: str):
    results = open_results_from_json(f"./results/{scraper}.json")
220
    count = Counter([result["url"] for result in results]).values()
    plt.hist(count)
222     plt.xlabel("Número de Algoritmos")
    plt.ylabel("Frequência")
224     plt.title("Histograma da Distribuição de Algoritmos")
    plt.show()
226
228 def traverse_tree(
    node: Union[Tag, NavigableString], depth: int, depths: Dict[str, List[
        int]]
    ) -> None:
```

```
228 """
    Traverse the BeautifulSoup tree and populate the depths dictionary with
        tag names and their depths.
230 Args:
    node: The current BeautifulSoup node.
232 depth: Current depth of traversal.
    depths: Dictionary to populate with depths.
234 """
    if isinstance(node, Tag): # Using isinstance to check if node is a Tag
        type
236         if node.name not in depths:
            depths[node.name] = []
238         depths[node.name].append(depth)

240         for child in node.children:
            traverse_tree(child, depth + 1, depths)
242
244 def calculate_html_nodes_depth(scraper: str) -> None:
    results = open_results_from_json(f"./results/{scraper}.json")
246 urls = list(set(result["url"] for result in results))

248 all_max_depths: List[int] = []
    all_num_nodes: List[int] = []
250
252 for url in urls:
    response = requests.get(url)
    soup = BeautifulSoup(response.content, "html.parser")
254
    depths: Dict[str, List[int]] = {}
256    traverse_tree(soup, 0, depths)

258    max_depth = max(max(depths[tag]) for tag in depths)
    num_nodes = sum(len(depths[tag]) for tag in depths)
260
    all_max_depths.append(max_depth)
262    all_num_nodes.append(num_nodes)

264 _, axs = plt.subplots(1, 2, figsize=(10, 5))
    axs[0].hist(all_max_depths, bins=[23, 24, 25, 26])
266    axs[0].set_xlabel("Profundidade Máxima")
    axs[0].set_ylabel("Frequência")
268    axs[0].set_title("Histograma da Profundidade Máxima")

270    axs[1].hist(all_num_nodes, bins=[2500, 7500, 12500, 17500, 22500,
        30000])
    axs[1].set_xlabel("Número de Nodos")
```

```
272     axs[1].set_ylabel("Frequência")
273     axs[1].set_title("Histograma da Contagem de Nodos")
274
275     plt.suptitle("Combined Histogram for All URLs")
276     plt.show()
277
278 def make_results_analysis(scraper: str):
279     write_results_to_json(
280         f"./results/{scraper}_analysis",
281         [
282             {
283                 **make_completion_rate(scraper),
284                 **make_results_characterization(scraper),
285                 **make_confusion_matrix(scraper),
286                 **make_complexitys_classification(scraper),
287             }
288         ],
289     )
290     # calculate_html_nodes_depth(scraper)
291     # make_algorithms_histogram(scraper)
```

- saci/observability/settings.py

```
SCRAPER_LOG_FORMAT = (
2     "%(asctime)s - %(levelname)s - Scraper: %(scraper)s - Message: %(
3     message)s" # noqa
4 )
5 SESSION_LOG_FORMAT = (
6     "%(asctime)s - %(levelname)s - Session - Message: %(message)s" # noqa
7 )
```

- saci/schema/\_\_init\_\_.py

- saci/schema/data\_schemas.py

```
from typing import Optional
2
```

```
from pydantic import Field
4 from pydantic import BaseModel

6
class ScrapedCode(BaseModel):
8     """The scraped codes of a given algorithm"""

10     code: str = Field(
        default="", description="The code of the algorithm in a certain
            language"
12     )
    comments: str = Field(default="", description="The first few comments
        in the code")
14

16 class ScrapedAlgorithm(BaseModel):
    """The final result of the scraping operations, an extracted algorithm
        with its attributes"""

18
    name: str = Field(default="", description="The name of the scraped
        algorithm")
20    time_complexity: Optional[str] = Field(
        default="", description="The time complexity of the scraped
            algorithm"
22    )
    trustable_time_complexity: bool = Field(
24        default=True,
        description="If the given time complexity used a trustable
            extraction method",
26    )
    space_complexity: Optional[str] = Field(
28        default="", description="The space complexity of the algorithm"
    )
    trustable_space_complexity: bool = Field(
30        default=True,
        description="If the given space complexity used a trustable
            extraction method",
32    )
    url: str = Field(default="", description="The space complexity of the
        algorithm")
34    codes: dict[str, ScrapedCode] = Field(
        default=dict(), description="The algorithm implementations"
36    )
    )
```

- `saci/session/__init__.py`



- saci/session/exceptions.py

```
class InvalidUrlException(Exception):
2     def __init__(self):
        self.message = "Please provide a valid URL string!"
4         super().__init__(self.message)

6
class MissingMethodException(Exception):
8     def __init__(self):
        self.message = "Please provide a request method!"
10        super().__init__(self.message)

12
class UnsupportedMethodException(Exception):
14    def __init__(self):
        self.message = "Please provide a supported method!"
16        super().__init__(self.message)

18
class MissingArgumentException(Exception):
20    def __init__(self, argument):
        self.message = f"Please provide the missing argument: {argument}!"
22        super().__init__(self.message)

24
class InvalidArgumentType(Exception):
26    def __init__(self, type):
        self.message = f"Please provide the correct argument type: {type}!"
28        super().__init__(self.message)
```

- saci/session/http\_session.py

```
import asyncio
2 import inspect

4 from aiohttp import ClientSession
from requests_html import AsyncHTMLSession
6 from yarl import URL
```

```
8 from scraper.exceptions import TooManyRetrysException
9 from scraper.observability.log import session_log as log
10 from scraper.session.exceptions import (
11     InvalidArgumentType,
12     InvalidUrlException,
13     MissingArgumentException,
14     MissingMethodException,
15     UnsupportedMethodException,
16 )
17 from scraper.session.response import Response
18 from scraper.session.settings import MAX_REDIRECTS, REQUIRED_REQUEST_ARGS
19 from scraper.utils import retry
20
21
22 class HttpSession:
23     def __init__(self):
24         self._session = ClientSession()
25         self._js_session = AsyncHTMLSession()
26         self._default_headers = {}
27
28     @property
29     def default_headers(self):
30         return self._default_headers
31
32     @default_headers.setter
33     def default_headers(self, headers: dict):
34         if not isinstance(headers, dict):
35             raise InvalidArgumentType(headers)
36         self._default_headers = headers
37
38     @staticmethod
39     def validate_url(url: str):
40         if not URL(url).is_absolute():
41             raise InvalidUrlException(url)
42
43     def _validate_request_args(self, **kwargs):
44         method = kwargs.get("method")
45         if not method:
46             raise MissingMethodException
47         if method not in REQUIRED_REQUEST_ARGS:
48             raise UnsupportedMethodException
49         for arg in REQUIRED_REQUEST_ARGS[method]:
50             if arg not in kwargs:
51                 raise MissingArgumentException(arg)
52         self.validate_url(kwargs["url"])
53
54     def _make_headers(self, headers: dict) -> dict:
55         return {**self._default_headers, **(headers or {})}
```

```
56
@retry(times=3)
58 async def _http_request(self, **kwargs):
    response = await self._session.request(
60         max_redirects=MAX_REDIRECTS,
        **kwargs,
62         headers=self._make_headers(kwargs.get("headers")),
    )
64     await asyncio.sleep(2)
    return response

66
68 async def request(self, **kwargs) -> Response:
    self._validate_request_args(**kwargs)
    callbacks = kwargs.pop("callbacks", None)
70
    try:
72         aiohttp_response = await self._http_request(**kwargs)
    except TooManyRetrysException:
74         log.error(f'Request failed. URL: {kwargs["url"]}')
        return Response()
76
    response = await Response.create_response_object(aiohttp_response)
78
    if callbacks:
80         for callback in callbacks:
            if inspect.iscoroutinefunction(callback):
82                 response = await callback(response)
            else:
84                 response = callback(response)
86
    return response

88
89 async def js_script_request(self, **kwargs):
    self._validate_request_args(**kwargs)
90     script = kwargs.get("script")
92
    response = await self._js_session.request(kwargs)
    if script:
94         return await response.html.arender(script=script)
    else:
96         await response.html.arender()
        return response.html.full_text
```

- `saci/session/response.py`

```

1 from aiohttp import ClientResponse
  from bs4 import BeautifulSoup
3
5 class Response:
  def __init__(
7     self,
  url: str,
9     status: int,
  content: str,
11    headers: dict,
  original_response: ClientResponse,
13    soup: BeautifulSoup = None,
  ):
15    self.url = url
  self.status = status
17    self.content = content
  self.headers = headers
19    self.original_response = original_response
  self.soup = soup or BeautifulSoup(content, "html.parser")
21
  @classmethod
23  async def create_response_object(cls, aiohttp_response: ClientResponse)
  :
  content = await aiohttp_response.text()
25
  return cls(
27    url=aiohttp_response.url,
  status=aiohttp_response.status,
29    headers=aiohttp_response.headers,
  content=content,
31    original_response=aiohttp_response,
  )

```

- saci/session/settings.py

```

AVAILABLE_METHODS = ["get", "post"]
2
REGEX = {
4  "url": r"(https ?://(?:www\.|(?!www)) [a-zA-Z0-9][a-zA-Z0-9-]+[a-zA-Z0-9]
  \. [\^s]{2,} | www\. [a-zA-Z0-9][a-zA-Z0-9-]+[a-zA-Z0-9] \. [\^s]{2,} |
  https ?://(?:www\.|(?!www)) [a-zA-Z0-9]+\.[\^s]{2,} | www\. [a-zA-Z0-9]
  -9]+\.[\^s]{2,})" # noqa
}

```

```
6 REQUIRED_REQUEST_ARGS = {"get": ["url"], "post": ["url", "data"]}
8 MAX_REDIRECTS = 5
```

- saci/session/utils.py

```
1 from enum import Enum
3
4 class Methods(Enum):
5     GET = "get"
6     POST = "post"
```

- saci/websites/\_\_init\_\_.py

- saci/websites/geeks\_for\_geeks/\_\_init\_\_.py

- saci/websites/geeks\_for\_geeks/extract.py

```
1 from re import IGNORECASE, compile, search
2 from typing import Dict, List, Optional, Tuple, Union
3
4 from bs4.element import Tag
5
6 from scraper.exceptions import (
7     FailedComplexityExtraction,
8 )
9 from scraper.observability.log import scraper_log as log
10 from scraper.session.response import Response
11 from scraper.websites.geeks_for_geeks.settings import (
12     AUXILIARY_SPACE_REGEX,
13     COMMENTS_STARTING_STRINGS,
14     HTML_ELEMENTS_NAMES,
```

```
LANGUAGES,
16 TIME_COMPLEXITY_REGEX,
)
18

20 def get_sourceline(element):
    return element.sourceline if isinstance(element, Tag) else element.
        parent.sourceline
22

24 def search_regex(pattern: str, text: str, group_num: int = 0) -> Optional[
    str]:
    match = search(pattern, text)
26     return match.group(group_num) if match else None

28
def extract_code(code_table: Tag) -> str:
30     code_text = ""
    code_lines = code_table.find_all("div", {"class": "line"})
32     for line in code_lines:
        code_pieces = line.find_all("code")
34         for code in code_pieces:
            code_text += code.text
36         code_text += "\n"
    return code_text
38

40 def extract_code_comments(algorithm: str) -> str:
    algorithm_comments = ""
42     for line in algorithm.splitlines():
        if line[:2] in COMMENTS_STARTING_STRINGS:
44         algorithm_comments += line
            algorithm_comments += "\n"
46     else:
        return algorithm_comments
48

50 def extract_complexity_from_reference(
    reference: Tag, regex_map: Dict[str, List[str]]
52 ) -> str:
    matches = [
54         reference.find_next(string=compile(regex)) for regex in regex_map["
            word"]
    ]

56     matches = [match for match in matches if match]

58     closest_match = min(
```

```
60     matches ,
        key=lambda match: abs(getsourceline(reference) - getsourceline(
            match)) ,
62     )

64     for regex in regex_map["value"]:
        for element in HTML_ELEMENTS_NAMES:
66         complexity = search_regex(
            regex, closest_match.find_previous(element).text , 1
68         )
            if complexity:
70                 return complexity

72     raise FailedComplexityExtraction

74
def extract_name(dom_reference: Tag) -> str:
76     name = dom_reference.find_previous("h2")
    if name and "tabtitle" not in name.get("class", ""):
78         return name.text

80
def fallback_search(reference_list: List[str], pattern: str) -> Optional[
    str]:
82     for result in reference_list:
        if search(pattern, result, flags=IGNORECASE):
84         return search_regex(TIME_COMPLEXITY_REGEX["fallback"], result)
    return None

86
88 def extract_complexity_with_fallback(
    dom_reference: Tag,
90 ) -> Tuple[Optional[str], Optional[str]]:
    search_func = lambda x: search_regex(TIME_COMPLEXITY_REGEX["fallback"],
        x.text)
92     next_results = [p.text for p in dom_reference.find_all_next("p") if
        search_func(p)]
    prev_results = [
94         p.text for p in dom_reference.find_all_previous("p") if search_func
            (p)
        ]

96
    time_complexity = fallback_search(next_results, r"time") or
        fallback_search(
98         prev_results, r"time"
        )
    space_complexity = fallback_search(
100     next_results, r"auxiliary|space"
```

```
102     ) or fallback_search(prev_results, r"auxiliary|space")
104     return time_complexity, space_complexity
106
107 def extract_codes_and_references(
108     response: Response,
109 ) -> Tuple[List[Dict[str, str]], List[Tag]]:
110     codes, references = [], []
111     code_tabs = response.soup.find_all("div", {"class": "responsive-tabs"})
112
113     for tab in code_tabs:
114         code_map = {}
115         references.append(tab)
116
117         for lang in LANGUAGES:
118             algorithm = tab.find_next("h2", {"class": "tabtitle"}, string=
119                 lang)
120             if algorithm:
121                 code = extract_code(algorithm.find_next("td", {"class": "
122                     code"}))
123                 comments = extract_code_comments(code)
124
125                 code_map[lang] = {
126                     "code": code,
127                     "comments": comments,
128                 }
129
130             codes.append(code_map)
131
132     return codes, references
133
134 def extract_data(
135     response: Response,
136 ) -> List[Dict[str, Union[str, bool, Dict[str, str]]]]:
137     codes, dom_references = extract_codes_and_references(response)
138     if not codes:
139         log.warning(f"Failed to find codes. URL: {response.url}")
140         return []
141
142     complexities = []
143     for ref in dom_references:
144         try:
145             time_complexity = extract_complexity_from_reference(
146                 ref, TIME_COMPLEXITY_REGEX
147             )
148         except FailedComplexityExtraction:
```



```
148         time_complexity, _ = extract_complexity_with_fallback(ref)
150     try:
151         space_complexity = extract_complexity_from_reference(
152             ref, AUXILIARY_SPACE_REGEX
153         )
154     except FailedComplexityExtraction:
155         _, space_complexity = extract_complexity_with_fallback(ref)
156
157     complexities.append(
158         {
159             "time_complexity": time_complexity,
160             "trustable_time_complexity": bool(time_complexity),
161             "space_complexity": space_complexity,
162             "trustable_space_complexity": bool(space_complexity),
163         }
164     )
165
166     main_name = response.soup.find("div", {"class": "article-title"}).text
167     names = [main_name] + [
168         extract_name(ref) if extract_name(ref) else main_name
169         for ref in dom_references[1:]
170     ]
171
172     return [
173         {
174             "name": name,
175             "time_complexity": complexity["time_complexity"],
176             "trustable_time_complexity": complexity["
177                 trustable_time_complexity"],
178             "space_complexity": complexity["space_complexity"],
179             "trustable_space_complexity": complexity["
180                 trustable_space_complexity"],
181             "url": str(response.url),
182             "codes": code,
183         }
184         for name, complexity, code in zip(names, complexities, codes)
185     ]
```

- [saci/websites/geeks\\_for\\_geeks/login.py](https://saci/websites/geeks_for_geeks/login.py)

```
1 import os
3 from scraper.session.http_session import HttpSession
  from scraper.session.response import Response
```

```
5 from scraper.session.utils import Methods
7 from scraper.websites.geeks_for_geeks.settings import LOGIN_URL
9 def get_login_payload():
11     return {
13         "reqType": "Login",
15         "user": os.getenv("USERNAME"),
17         "pass": os.getenv("PASSWORD"),
19         "rem": False,
21         "to": "https://auth.geeksforgeeks.org/?to=https://www.geeksforgeeks
23         .org/",
25         "rem": "on",
27         "g-recaptcha-response": None,
29         "browserInfo": {},
31     }
33
35 async def follow_login(session: HttpSession):
37     await session.request(
39         method=Methods.GET.value,
41         data=get_login_payload(),
43         url=LOGIN_URL,
45     )
```

- `saci/websites/geeks_for_geeks/settings.py`

```
1 BASE_URL = "https://www.geeksforgeeks.org"
3 LOGIN_URL = "https://auth.geeksforgeeks.org/auth.php"
5 ALGORITHMS_LOCATION_URLS = [
7     f"{BASE_URL}/fundamentals-of-algorithms",
9     f"{BASE_URL}/data-structures",
11 ]
13 HEADERS = {
15     "Host": "www.geeksforgeeks.org",
17     "Accept": "text/html, application/xhtml+xml, application/xml;q=0.9,image/
19     avif,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=
21     b3;q=0.9",
23     "Accept-Encoding": "gzip, deflate",
25     "Upgrade-Insecure-Requests": "1",
27     "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit
29     /537.36 (KHTML, like Gecko) Chrome/104.0.0.0 Safari/537.36",
```

```

17     "Sec-Fetch-Site": "none",
18     "Sec-Fetch-Mode": "navigate",
19     "Sec-Fetch-User": "?1",
20     "Sec-Fetch-Dest": "document",
21     "Sec-Ch-Ua": '"Chromium";v="105", "Not)A;Brand";v="8" ',
22     "Sec-Ch-Ua-Mobile": "?0",
23     "Sec-Ch-Ua-Platform": '"Windows" ',
24 }
25
26 ALGORITHM_LINKS_JS_SCRIPT = """
27     () => {
28         const algorithm_links = []
29
30         const ordered_lists = document.getElementsByTagName('ol');
31         for (ol of ordered_lists) {
32             const hrefs = []
33             const a_tags = ol.getElementsByTagName('a')
34             for (a of a_tags) {
35                 hrefs.push(a.getAttribute('href'))
36             }
37             algorithm_links.push(hrefs)
38         }
39
40         return algorithm_links
41     }
42 """
43
44 TIME_COMPLEXITY_REGEX = {
45     "word": [
46         r"[Tt]ime [Cc]omplexit\\w*[\\s]?:[Tt]ime [Cc]omplexit\\w*[\\s]?",
47     ],
48     "value": [
49         r"[Tt]ime [Cc]omplexit\\w*\\s?:.*?(O\\s*\\(.*?\\))",
50         r"[Tt]ime [Cc]omplexit\\w*\\s?.*?(O\\s*\\(.*?\\))",
51         r"[Tt]ime [Cc]omplexit\\w*\\s.*?\\sis\\s.*?(O\\s*\\(.*?\\))",
52         r"[Tt]ime [Cc]omplexit\\w*\\s?:.*?\\sis\\s(\\w*)",
53     ],
54     "fallback": r"O\\s*\\(.*?\\)",
55 }
56
57 AUXILIARY_SPACE_REGEX = {
58     "word": [
59         r"[Aa]uxiliary [Ss]pace\\s?:|[Aa]uxiliary [Ss]pace\\s?",
60         r"[Ss]pace [Cc]omplexit\\w*\\s?:|[Ss]pace [Cc]omplexit\\w*\\s?",
61     ],
62     "value": [
63         r"[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?:.*?(O\\s*\\(.*?\\))",
64         r"[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?.*?(O\\s*\\(.*?\\))",

```

```

65     r"[Ss]pace [Cc]omplexit\\w*\\s?:.*?(O\\s*(.*?\\))",
66     r"[Ss]pace [Cc]omplexit\\w*\\s?:.*?(O\\s*(.*?\\))",
67     r"[Aa]uxiliary [Ss]pace\\s?[\\Cc]omplexit\\w*?.*?\\sis\\s.*?(O\\s
68         *\\(.*?\\))",
69     r"[Ss]pace [Cc]omplexit\\w*\\s.*?\\sis\\s.*?(O\\s*(.*?\\))",
70     r"[Aa]uxiliary [Ss]pace\\w*\\s?:?.*?\\sis\\s(\\w*)",
71     r"[Ss]pace [Cc]omplexit\\w*\\s?:?.*?\\sis\\s(\\w*)",
72 ],
73     "fallback": r"O\\s*(.*?\\)",
74 }
75 HTML_ELEMENTS_NAMES = ["p", "li", "ul", "i"]
76 COMMENTS_STARTING_STRINGS = ["//", "# ", "/*"]
77 LANGUAGES = ["C", "C#", "C++", "Python3", "Python", "PHP", "Java", "
    Javascript"]

```

- [saci/websites/geeks\\_for\\_geeks/spider.py](#)

```

1 from asyncio import gather
2
3 from scraper.exceptions import FailedExtraction
4 from scraper.observability.log import scraper_log as log
5 from scraper.schema.data_schemas import ScrapedAlgorithm
6 from scraper.session.http_session import HttpSession
7 from scraper.session.response import Response
8 from scraper.session.utils import Methods
9 from scraper.utils import remove_duplicates, retry
10 from scraper.websites.geeks_for_geeks.extract import extract_data
11 from scraper.websites.geeks_for_geeks.login import follow_login
12 from scraper.websites.geeks_for_geeks.settings import
13     ALGORITHMS_LOCATION_URLS, HEADERS
14
15 def parse_algorithm_schema(algorithms_data: list) -> list:
16     return [ScrapedAlgorithm(**algorithm) for algorithm in algorithms_data]
17
18 def filter_algorithms_urls(algorithms_urls: list) -> list:
19     return [href for href in algorithms_urls if "geeksquiz" not in href]
20
21
22 async def extract_algorithm_data(response: Response) -> list:
23     if not response:
24

```

```
        return []
26
    data = await extract_data(response)
28    if not data:
        raise FailedExtraction
30
    return data
32
34 def extract_algorithms_urls(response: Response) -> list:
    page_content = response.soup.find("div", {"class": "page_content"})
36    return [
        a.attrs.get("href")
38        for ol in page_content.find_all("ol")[1:]
        for a in ol.find_all("a")
40        if a.attrs.get("href")
    ]
42
44 @retry(times=3, raise_exception=False, return_value=[])
async def follow_algorithm(session: HttpSession, url: str):
46     return await session.request(
        method=Methods.GET.value,
48         url=url,
        callbacks=[extract_algorithm_data, parse_algorithm_schema],
50     )
52
54 async def follow_algorithms(session: HttpSession, algorithms_urls: list) ->
    list:
    algorithms = await gather(
        *[follow_algorithm(session, url) for url in algorithms_urls]
56     )
    return sum(algorithms, [])
58
60 async def follow_algorithms_urls(session: HttpSession) -> list:
    algorithm_location_urls = await gather(
62         *[
            session.request(
64                 method=Methods.GET.value,
                url=url,
66                 callbacks=[
                    extract_algorithms_urls,
68                 filter_algorithms_urls,
                    remove_duplicates,
70                 ],
            )
        ]
    )
```

```
72         for url in ALGORITHMS_LOCATION_URLS
73     ]
74 )
75 return sum(algorithm_location_urls , [])
76
77
78 async def run(url: str = None) -> list:
79     session = HttpSession()
80     session.default_headers = HEADERS
81
82     log.info('Starting "Geeks for Geeks" algorithms extraction')
83     await follow_login(session)
84
85     algorithms_urls = [url] if url else await follow_algorithms_urls(
86         session)
87     return await follow_algorithms(session, algorithms_urls)
```

## Apêndice B - Código Fonte Versão OpenAI

- saci/requirements.txt

```
1 aiohttp
2 beautifulsoup4
3 black
4 pandas
5 pydantic
6 pymongo
7 pyppeteer
8 requests
9 scikit-learn
10 scipy
11 seaborn
12 mlxtend
13 requests-html
14 pip-tools
15 openai
16 langchain
17 asyncio-throttle
```

- saci/makefile

```
1 PYTHON :=
2 ifeq ($(OS),Windows_NT)
3     PYTHON = venv/Scripts/python
4     PIP = venv/Scripts/pip
5 else
6     PYTHON = venv/bin/python3
7     PIP = venv/bin/pip
8 endif
9
10 requirements: requirements.txt
11 $(PIP) install -r requirements.txt
12
13 run_scraper:
```

```
$(PYTHON) -m main
15
scrape_url:
17 $(PYTHON) -m scraper.scripts --script=scrape_url --url=$(url) --scraper=$(
    scraper)

19 results_analysis:
    $(PYTHON) -m scraper.scripts --script=results_analysis --scraper=$(
        scraper)

21
manual_results_boilerplate:
23 $(PYTHON) -m scraper.scripts --script=manual_results_boilerplate --
    scraper=$(scraper)
```

- saci/main.py

```
1 from asyncio import get_event_loop
3 from scraper.runner import run_scrapers
5 if __name__ == "__main__":
    loop = get_event_loop()
7    loop.run_until_complete(run_scrapers())
```

- saci/scraper/\_\_init\_\_.py
- saci/scraper/database.py

```
1 import os
from typing import Dict, List
3
from pymongo import MongoClient
5 from pymongo.collection import Collection
from pymongo.operations import UpdateOne
7
from scraper.settings import DATABASE_NAME
9
FIELDS_TO_UPDATE = [
11     "name",
    "time_complexity",
13     "trustable_time_complexity",
    "space_complexity",
```



```
15     "trustable_space_complexity",
16 ]
17
18
19 class ScraperDatabase:
20     def __init__(self):
21         self.client: MongoClient = MongoClient(
22             os.environ.get("MONGO_CONNECTION_STRING"),
23             serverSelectionTimeoutMS=1
24         )
25         self.database = self.client[DATABASE_NAME]
26
27     def _get_collection(self, collection_name: str) -> Collection:
28         return self.database[collection_name]
29
30     def update_database(self, collection_name: str, data: List[Dict[str,
31         any]]):
32         collection = self._get_collection(collection_name)
33         operations = [
34             UpdateOne(
35                 filter={"url": item["url"], "codes": item["codes"]},
36                 update={"$set": {field: item[field] for field in
37                     FIELDS_TO_UPDATE}},
38                 upsert=True,
39             )
40             for item in data
41         ]
42         collection.bulk_write(operations)
```

- saci/scrapper/exceptions.py

```
1 class InvalidURLException(Exception):
2     def __init__(self, url):
3         super().__init__(f"The provided URL '{url}' is not valid.")
4
5
6
7 class TooManyRetrysException(Exception):
8     pass
9
10
11 class FailedExtraction(Exception):
12     """Failed the extraction process"""
13
14     pass
```

```
15 class FailedComplexityExtraction(Exception):
17     """Failed to extract the complexitys"""
19     pass
```

- saci/scrapper/runner.py

```
1 from asyncio import gather
3 from scraper.database import ScraperDatabase
4 from scraper.settings import USE_MONGO_DATABASE
5 from scraper.utils import write_results_to_json
6 from scraper.websites import websites
7
9 async def run_scrapers():
10     scrapers_results = await gather(*[scraper() for scraper in websites.
11     values()])
12     for name, data in zip(websites.keys(), scrapers_results):
13         if USE_MONGO_DATABASE:
14             database = ScraperDatabase()
15             database.update_database(name, [result.dict() for result in
16             data])
17         else:
18             write_results_to_json(
19                 name,
20                 sorted([result.dict() for result in data], key=lambda alg:
21                 alg["url"]),
22             )
```

- saci/scrapper/scripts.py

```
1 import argparse
2 import pprint
3 from asyncio import get_event_loop
4 from importlib import import_module
5 from scraper.observability.metrics import (
6     make_results_analysis,
7     make_manual_results_boilerplate,
8 )
9
```

```
11 def extract_single_url(url: str, scraper: str):
    module = import_module(f"scraper.websites.{scraper}")
13     spider = getattr(module, "spider")
    loop = get_event_loop()
15     result = loop.run_until_complete(spider.run(url))
    pprint.PrettyPrinter(indent=4).pprint(result)
17
19 def main():
    parser = argparse.ArgumentParser()
21     parser.add_argument("--script")
    parser.add_argument("--url", help="URL to be extracted.")
23     parser.add_argument("--scraper", help="Scraper to be used.")
    args = parser.parse_args()
25
    actions = {
27         "scrape_url": lambda: extract_single_url(args.url, args.scraper),
        "results_analysis": lambda: make_results_analysis(args.scraper),
29         "manual_results_boilerplate": lambda:
            make_manual_results_boilerplate(
31             args.scraper
            ),
    }
33
    action = actions.get(args.script)
35     if action:
        action()
37     else:
        print(f"Unknown script: {args.script}")
39
41 if __name__ == "__main__":
    main()
```

- saci/scraper/settings.py

```
1 DATABASE_NAME = "COMPLEXITY_SCRAPER"
2
USE_MONGO_DATABASE = False
```

- saci/scraper/utils.py

```
1 import json
import os

3
4
5 from asyncio_throttle import Throttler
6 from langchain.chains import LLMChain
7 from langchain.llms import OpenAI
8 from langchain.prompts import PromptTemplate

9 from scraper.exceptions import TooManyRetrysException
10 from scraper.observability.log import scraper_log as log
11

12
13 def retry(
14     times: int,
15     raise_exception=True,
16     return_value=None,
17 ):
18     def func_wrapper(f):
19         async def wrapper(*args, **kwargs):
20             for _ in range(times):
21                 try:
22                     return await f(*args, **kwargs)
23                 except Exception:
24                     pass
25
26                 if not raise_exception:
27                     return return_value
28             else:
29                 raise TooManyRetrysException
30
31         return wrapper
32
33     return func_wrapper
34
35
36 def remove_duplicates(data: list) -> list:
37     return list(set(data))
38
39
40 def open_results_from_json(file_path: str) -> dict:
41     with open(file_path, "r") as file:
42         file_contents = file.read()
43         results = json.loads(file_contents)
44     return results
45
46
47 def write_results_to_json(file_path: str, data: list):
```

```
with open(f"{file_path}.json", "w") as file:
    json.dump(data, file)

class LlmComplexitySearcher:
    _instance = None
    throttler = Throttler(rate_limit=100, period=60)

    def __init__(self):
        self.LLM = OpenAI(openai_api_key=os.environ.get("OPENAI_KEY", ""))
        self.PROMPT = PromptTemplate.from_template(
            """I am using you during a scraping operation in which I need
            to extract from a piece of text the values of complexity/
            code complexity.
            Consider that sometimes the space complexity is also called as
            auxiliary space.
            Consider that both complexities often come in this format: 'Some
            type of complexity: value of complexity'.
            But they can also come inside the text like: 'The type of
            complexity of the given ... is ...'.
            Give the answer in JSON format with no line breaks, with a key
            called "complexity" and the value for the key is your answer
            .
            If you cant determine the answer give the json with a null in
            the value.
            What is the {complexity} that is written in the following text:
            "{text}"?"""
        )
        self.LLM_CHAIN = LLMChain(prompt=self.PROMPT, llm=self.LLM)

    def __new__(cls, *args, **kwargs):
        if not cls._instance:
            cls._instance = super().__new__(cls)
        return cls._instance

    @retry(times=3, raise_exception=False, return_value={"complexity": None})
    async def search(self, complexity: str, text: str) -> dict:
        async with self.throttler:
            answer = self.LLM_CHAIN.run(complexity=complexity, text=text)
            return json.loads(answer)
```

- `saci/observability/__init__.py`

- saci/observability/log.py

```
import platform
2
from logging import DEBUG, Filter, Formatter, StreamHandler, getLogger
4
from scraper.observability.settings import SCRAPER_LOG_FORMAT,
    SESSION_LOG_FORMAT
6
8 class AddScraperName(Filter):
    def filter(self, record):
10         if not hasattr(record, "scraper"):
                split_char = "\\\" if "Windows" in platform.platform() else "/"
12                 record.scraper = record.pathname.split(split_char)[-2]
                return True
14
16 def build_scraper_log(logger_name):
    log_handler = StreamHandler()
18     log_handler.setLevel(DEBUG)
    log_handler.setFormatter(Formatter(SCRAPER_LOG_FORMAT))
20
    log = getLogger(logger_name)
22     log.addFilter(AddScraperName())
    log.addHandler(log_handler)
24     return log
26
28 def build_session_log(logger_name):
    log_handler = StreamHandler()
    log_handler.setLevel(DEBUG)
30     log_handler.setFormatter(Formatter(SESSION_LOG_FORMAT))
32
    log = getLogger(logger_name)
    log.addHandler(log_handler)
34     return log
36
38 scraper_log = build_scraper_log("scraper_log")
session_log = build_session_log("session_log")
```

- saci/observability/metrics.py

```
import math
```

```
2 import re
3 from collections import Counter
4 from typing import Any, Dict, List, Tuple, Union
5
6 import matplotlib.pyplot as plt
7 import numpy
8 import requests
9 from bs4 import BeautifulSoup, NavigableString, Tag
10 from mlxtend.plotting import plot_confusion_matrix
11
12 from scraper.utils import (
13     open_results_from_json,
14     remove_duplicates,
15     write_results_to_json,
16 )
17
18
19 def make_results_characterization(scraper: str) -> Dict[str, int]:
20     results = open_results_from_json(f"./results/{scraper}.json")
21
22     return {
23         "Distinct Time Complexities": len(
24             {algorithm["time_complexity"] for algorithm in results}
25         ),
26         "Distinct Space Complexities": len(
27             {algorithm["space_complexity"] for algorithm in results}
28         ),
29         "Non Trusted Time Complexities": sum(
30             [1 for algorithm in results if not algorithm["
31                 trustable_time_complexity"]]
32         ),
33         "Non Trusted Space Complexities": sum(
34             [1 for algorithm in results if not algorithm["
35                 trustable_space_complexity"]]
36     )
37
38 def make_completion_rate(scraper: str) -> Dict[str, int]:
39     results = open_results_from_json(f"./results/{scraper}.json")
40
41     total = len(results)
42     time_complexity, space_complexity = sum(
43         1 for algorithm in results if not algorithm["time_complexity"]
44     ), sum(1 for algorithm in results if not algorithm["space_complexity"])
45     urls_with_problem = remove_duplicates(
46         [
47             algorithm["url"]
```

```
48         for algorithm in results
49             if not algorithm["time_complexity"] or not algorithm["
50                 space_complexity"]
51         ]
52     )
53
54     return {
55         "Total Algorithms": total,
56         "Time Complexity Extracted": f"{total - time_complexity} extracted
57             - {100 - (time_complexity / total) * 100}%",
58         "Space Complexity Extracted": f"{total - space_complexity}
59             extracted - {100 - (space_complexity / total) * 100}%",
60         "URLs with problem": urls_with_problem,
61     }
62
63 def calculate_metrics(
64     true_positive: int, false_positive: int, true_negative: int,
65     false_negative: int
66 ) -> Dict[str, float]:
67     accuracy = (true_positive + true_negative) / (
68         true_positive + true_negative + false_positive + false_negative
69     )
70     precision = true_positive / (true_positive + false_positive)
71     recall = true_positive / (true_positive + false_negative)
72     f1_score = (2 * precision * recall) / (precision + recall)
73
74     return {
75         "True Positive": true_positive,
76         "False Positive": false_positive,
77         "True Negative": true_negative,
78         "False negative": false_negative,
79         "Accuracy": accuracy,
80         "Precision": precision,
81         "Recall": recall,
82         "F1 Score": f1_score,
83     }
84
85 def plot_confusion(
86     title: str,
87     true_positive: int,
88     false_negative: int,
89     false_positive: int,
90     true_negative: int,
91 ):
92     array = numpy.array(
93         [[true_positive, false_negative], [false_positive, true_negative]]
```



```
92     )
93
94     fig , ax = plot_confusion_matrix(conf_mat=array)
95     plt.title("Space Complexity Confusion Matrix")
96     plt.show()
97
98
99 def make_confusion_matrix(scraper: str) -> Dict[str, Dict[str, Any]]:
100     manual_results = open_results_from_json(f"./results/manual_{scraper}.
101         json")
102     results = open_results_from_json(f"./results/{scraper}.json")
103
104     time_metrics = {
105         "true_positive": 0,
106         "false_positive": 0,
107         "true_negative": 0,
108         "false_negative": 0,
109     }
110
111     space_metrics = {
112         "true_positive": 0,
113         "false_positive": 0,
114         "true_negative": 0,
115         "false_negative": 0,
116     }
117
118     for result, manual_result in zip(results, manual_results):
119         # Time complexity
120         if time_complexity := result["time_complexity"]:
121             if time_complexity == manual_result["time_complexity"]:
122                 time_metrics["true_positive"] += 1
123             else:
124                 time_metrics["false_positive"] += 1
125         else:
126             if not manual_result["time_complexity"]:
127                 time_metrics["true_negative"] += 1
128             else:
129                 time_metrics["false_negative"] += 1
130
131         # Space complexity
132         if space_complexity := result["space_complexity"]:
133             if space_complexity == manual_result["space_complexity"]:
134                 space_metrics["true_positive"] += 1
135             else:
136                 space_metrics["false_positive"] += 1
137         else:
138             if not manual_result["space_complexity"]:
139                 space_metrics["true_negative"] += 1
```

```

    else:
140         space_metrics["false_negative"] += 1

142     # Calculating metrics
    time_complexity_metrics = calculate_metrics(**time_metrics)
144     space_complexity_metrics = calculate_metrics(**space_metrics)

146     # Plotting the confusion matrices
    plot_confusion("Time Complexity Confusion Matrix", **time_metrics)
148     plot_confusion("Space Complexity Confusion Matrix", **space_metrics)

150     # Return result
    return {
152         "Time Complexity": time_complexity_metrics,
        "Space Complexity": space_complexity_metrics,
154     }

156 def make_manual_results_boilerplate(scraper: str) -> None:
158     results = open_results_from_json(f"./results/{scraper}.json")

160     manual_results = [
        {"url": alg["url"], "time_complexity": " ", "space_complexity": " "
162     }
        for alg in results
164     ]

    write_results_to_json(f"./results/manual_{scraper}", manual_results)
166

168 def make_complexitys_classification(scraper: str) -> Dict[str, Dict[str,
int]]:
    results = open_results_from_json(f"./results/{scraper}.json")
170

    time_classification = classify_complexity(
172         [result["time_complexity"] for result in results]
    )
174     space_classification = classify_complexity(
        [result["space_complexity"] for result in results]
176     )

178     return {
        "Time Complexity Classification": time_classification,
180         "Space Complexity Classification": space_classification,
    }
182

184 def classify_complexity(complexities: List[str]) -> Dict[str, int]:
```

```
constant, linear, exponential, polynomial, factorial, logarithmic = 0,
0, 0, 0, 0, 0
186 for complexity in complexities:
    if complexity:
188         if re.match(r"O\(\d\)", complexity):
            constant += 1
190         elif re.match(r"O\(\w\)", complexity):
            linear += 1
192         elif re.match(r"O\(.!\.!\.!\)", complexity):
            factorial += 1
194         elif re.match(r"O\(.! log .!\)", complexity):
            logarithmic += 1
196         elif re.match(r"O\(\d\^?\w\)", complexity):
            exponential += 1
198         elif re.match(r"O\(\w\s*[\^!*\+]? \s*[\d\w\s]+.\!\)", complexity)
            or re.match(
200             r"O\(\d+\(\w*\w\)", complexity
            ):
                polynomial += 1
202
203 return {
204     "Constant complexity": constant,
205     "Linear complexity": linear,
206     "Exponential complexity": exponential,
207     "Polynomial complexity": polynomial,
208     "Factorial complexity": factorial,
209     "Log complexity": logarithmic,
210 }
211
212
213 def make_algorithms_histogram(scraper: str):
214     results = open_results_from_json(f"./results/{scraper}.json")
215
216     count = Counter([result["url"] for result in results]).values()
217     plt.hist(count)
218     plt.xlabel("Número de Algoritmos")
219     plt.ylabel("Frequência")
220     plt.title("Histograma da Distribuição de Algoritmos")
221     plt.show()
222
223
224 def traverse_tree(
225     node: Union[Tag, NavigableString], depth: int, depths: Dict[str, List[
226         int]]
227 ) -> None:
228     """
229     Traverse the BeautifulSoup tree and populate the depths dictionary with
230     tag names and their depths.
```

```
230     Args:
231         node: The current BeautifulSoup node.
232         depth: Current depth of traversal.
233         depths: Dictionary to populate with depths.
234     """
235     if isinstance(node, Tag): # Using isinstance to check if node is a Tag
236         type
237         if node.name not in depths:
238             depths[node.name] = []
239         depths[node.name].append(depth)
240
241         for child in node.children:
242             traverse_tree(child, depth + 1, depths)
243
244 def calculate_html_nodes_depth(scraper: str) -> None:
245     results = open_results_from_json(f"./results/{scraper}.json")
246     urls = list(set(result["url"] for result in results))
247
248     all_max_depths: List[int] = []
249     all_num_nodes: List[int] = []
250
251     for url in urls:
252         response = requests.get(url)
253         soup = BeautifulSoup(response.content, "html.parser")
254
255         depths: Dict[str, List[int]] = {}
256         traverse_tree(soup, 0, depths)
257
258         max_depth = max(max(depths[tag]) for tag in depths)
259         num_nodes = sum(len(depths[tag]) for tag in depths)
260
261         all_max_depths.append(max_depth)
262         all_num_nodes.append(num_nodes)
263
264     _, axs = plt.subplots(1, 2, figsize=(10, 5))
265     axs[0].hist(all_max_depths, bins=[23, 24, 25, 26])
266     axs[0].set_xlabel("Profundidade Máxima")
267     axs[0].set_ylabel("Frequência")
268     axs[0].set_title("Histograma da Profundidade Máxima")
269
270     axs[1].hist(all_num_nodes, bins=[2500, 7500, 12500, 17500, 22500,
271         30000])
272     axs[1].set_xlabel("Número de Nodos")
273     axs[1].set_ylabel("Frequência")
274     axs[1].set_title("Histograma da Contagem de Nodos")
```

```
plt.suptitle("Combined Histogram for All URLs")
plt.show()

def make_results_analysis(scraper: str):
    write_results_to_json(
        f"./results/{scraper}_analysis",
        [
            {
                **make_competition_rate(scraper),
                **make_results_characterization(scraper),
                **make_confusion_matrix(scraper),
                **make_complexity_classification(scraper),
            }
        ],
    )
    # calculate_html_nodes_depth(scraper)
    # make_algorithms_histogram(scraper)
```

- saci/observability/settings.py

```
SCRAPER_LOG_FORMAT = (
    "%(asctime)s - %(levelname)s - Scraper: %(scraper)s - Message: %(
        message)s" # noqa
)
SESSION_LOG_FORMAT = (
    "%(asctime)s - %(levelname)s - Session - Message: %(message)s" # noqa
)
```

- saci/schema/\_\_init\_\_.py
- 

- saci/schema/data\_schemas.py

```
from typing import Optional
from pydantic import Field
from pydantic import BaseModel
```

```
6
class ScrapedCode(BaseModel):
8     """The scraped codes of a given algorithm"""

10     code: str = Field(
        default="", description="The code of the algorithm in a certain
            language"
12     )
    comments: str = Field(default="", description="The first few comments
        in the code")
14

16 class ScrapedAlgorithm(BaseModel):
    """The final result of the scraping operations, an extracted algorithm
        with its attributes"""

18     name: str = Field(default="", description="The name of the scraped
        algorithm")
20     time_complexity: Optional[str] = Field(
        default="", description="The time complexity of the scraped
            algorithm"
22     )
    trustable_time_complexity: bool = Field(
24         default=True,
        description="If the given time complexity used a trustable
            extraction method",
26     )
    space_complexity: Optional[str] = Field(
28         default="", description="The space complexity of the algorithm"
    )
30     trustable_space_complexity: bool = Field(
        default=True,
32         description="If the given space complexity used a trustable
            extraction method",
    )
34     url: str = Field(default="", description="The space complexity of the
        algorithm")
    codes: dict[str, ScrapedCode] = Field(
36         default=dict(), description="The algorithm implementations"
    )
```

- `saci/session/___init___py`

- saci/session/exceptions.py

```
class InvalidUrlException(Exception):
2     def __init__(self):
        self.message = "Please provide a valid URL string!"
4         super().__init__(self.message)

6
class MissingMethodException(Exception):
8     def __init__(self):
        self.message = "Please provide a request method!"
10        super().__init__(self.message)

12
class UnsupportedMethodException(Exception):
14    def __init__(self):
        self.message = "Please provide a supported method!"
16        super().__init__(self.message)

18
class MissingArgumentException(Exception):
20    def __init__(self, argument):
        self.message = f"Please provide the missing argument: {argument}!"
22        super().__init__(self.message)

24
class InvalidArgumentType(Exception):
26    def __init__(self, type):
        self.message = f"Please provide the correct argument type: {type}!"
28        super().__init__(self.message)
```

- saci/session/http\_session.py

```
import asyncio
2 import inspect

4 from aiohttp import ClientSession
from requests_html import AsyncHTMLSession
6 from yarl import URL

8 from scraper.exceptions import TooManyRetrysException
from scraper.observability.log import session_log as log
10 from scraper.session.exceptions import (
    InvalidArgumentType,
12     InvalidUrlException,
```

```
MissingArgumentException ,
14 MissingMethodException ,
    UnsupportedMethodException ,
16 )
from scraper.session.response import Response
18 from scraper.session.settings import MAX_REDIRECTS, REQUIRED_REQUEST_ARGS
from scraper.utils import retry
20
22 class HttpSession:
    def __init__(self):
24         self._session = ClientSession()
        self._js_session = AsyncHTMLSession()
26         self._default_headers = {}

    @property
28     def default_headers(self):
30         return self._default_headers

    @default_headers.setter
32     def default_headers(self, headers: dict):
34         if not isinstance(headers, dict):
            raise InvalidArgumentType(dict)
36         self._default_headers = headers

    @staticmethod
38     def validate_url(url: str):
40         if not URL(url).is_absolute():
            raise InvalidUrlException(url)
42
    def _validate_request_args(self, **kwargs):
44         method = kwargs.get("method")
        if not method:
46             raise MissingMethodException
        if method not in REQUIRED_REQUEST_ARGS:
48             raise UnsupportedMethodException
        for arg in REQUIRED_REQUEST_ARGS[method]:
50             if arg not in kwargs:
                raise MissingArgumentException(arg)
52         self.validate_url(kwargs["url"])

    def _make_headers(self, headers: dict) -> dict:
54         return {**self._default_headers, **(headers or {})}

56
    @retry(times=3)
58     async def _http_request(self, **kwargs):
        response = await self._session.request(
60             max_redirects=MAX_REDIRECTS,
```



```
        **kwargs ,
        headers=self._make_headers(kwargs.get("headers")),
    )
    await asyncio.sleep(2)
    return response

    66
    async def request(self, **kwargs) -> Response:
    68        self._validate_request_args(**kwargs)
        callbacks = kwargs.pop("callbacks", None)
    70
        try:
    72            aiohttp_response = await self._http_request(**kwargs)
        except TooManyRetrysException:
    74            log.error(f'Request failed. URL: {kwargs["url"]}')
            return Response()
    76
        response = await Response.create_response_object(aiohttp_response)
    78
        if callbacks:
    80            for callback in callbacks:
                if inspect.iscoroutinefunction(callback):
    82                    response = await callback(response)
                else:
    84                    response = callback(response)
    86
        return response

    88    async def js_script_request(self, **kwargs):
        self._validate_request_args(**kwargs)
    90        script = kwargs.get("script")

    92        response = await self._js_session.request(kwargs)
        if script:
    94            return await response.html.arender(script=script)
        else:
    96            await response.html.arender()
            return response.html.full_text
```

- `saci/session/response.py`

```
1 from aiohttp import ClientResponse
  from bs4 import BeautifulSoup
3
5 class Response:
```

```

7     def __init__(
9         self,
10        url: str,
11        status: int,
12        content: str,
13        headers: dict,
14        original_response: ClientResponse,
15        soup: BeautifulSoup = None,
16    ):
17        self.url = url
18        self.status = status
19        self.content = content
20        self.headers = headers
21        self.original_response = original_response
22        self.soup = soup or BeautifulSoup(content, "html.parser")
23
24    @classmethod
25    async def create_response_object(cls, aiohttp_response: ClientResponse)
26    :
27        content = await aiohttp_response.text()
28
29        return cls(
30            url=aiohttp_response.url,
31            status=aiohttp_response.status,
32            headers=aiohttp_response.headers,
33            content=content,
34            original_response=aiohttp_response,
35        )

```

- saci/session/settings.py

```

1 AVAILABLE_METHODS = ["get", "post"]
2
3 REGEX = {
4     "url": r"(https ?://(?:www\.|(?!www)) [a-zA-Z0-9][a-zA-Z0-9-]+[a-zA-Z0-9-9]\.[^\s]{2,}|www\.[a-zA-Z0-9][a-zA-Z0-9-]+[a-zA-Z0-9]\.[^\s]{2,}|https ?://(?:www\.|(?!www)) [a-zA-Z0-9]+\.[^\s]{2,}|www\.[a-zA-Z0-9]+\.[^\s]{2,})" # noqa
5 }
6
7 REQUIRED_REQUEST_ARGS = {"get": ["url"], "post": ["url", "data"]}
8
9 MAX_REDIRECTS = 5

```

- saci/session/utils.py

```
1 from enum import Enum
3
4 class Methods(Enum):
5     GET = "get"
6     POST = "post"
```

- saci/websites/\_\_init\_\_.py
- 

- saci/websites/geeks\_for\_geeks/\_\_init\_\_.py
- 

- saci/websites/geeks\_for\_geeks/extract.py

```
1 from re import IGNORECASE, compile, search
2 from typing import Dict, List, Optional, Tuple, Union
3
4 from bs4.element import Tag
5
6 from scraper.exceptions import FailedComplexityExtraction
7 from scraper.observability.log import scraper_log as log
8 from scraper.session.response import Response
9 from scraper.utils import LlmComplexitySearcher
10 from scraper.websites.geeks_for_geeks.settings import (
11     AUXILIARY_SPACE_REGEX,
12     COMMENTS_STARTING_STRINGS,
13     HTML_ELEMENTS_NAMES,
14     LANGUAGES,
15     TIME_COMPLEXITY_REGEX,
16 )
17
18 llm_searcher = LlmComplexitySearcher()
19
20 def get_sourceline(element):
```

```
22     return element.sourceline if isinstance(element, Tag) else element.  
        parent.sourceline  
  
24  
25 def search_regex(pattern: str, text: str, group_num: int = 0) -> Optional[  
26     str]:  
27     match = search(pattern, text)  
28     return match.group(group_num) if match else None  
  
29  
30 def extract_code(code_table: Tag) -> str:  
31     code_text = ""  
32     code_lines = code_table.find_all("div", {"class": "line"})  
33     for line in code_lines:  
34         code_pieces = line.find_all("code")  
35         for code in code_pieces:  
36             code_text += code.text  
37             code_text += "\n"  
38     return code_text  
  
39  
40  
41 def extract_code_comments(algorithm: str) -> str:  
42     algorithm_comments = ""  
43     for line in algorithm.splitlines():  
44         if line[:2] in COMMENTS_STARTING_STRINGS:  
45             algorithm_comments += line  
46             algorithm_comments += "\n"  
47         else:  
48             return algorithm_comments  
  
49  
50  
51 async def extract_complexity_from_reference(  
52     reference: Tag, regex_map: Dict[str, List[str]]  
53 ) -> str:  
54     matches = [  
55         reference.find_next(string=compile(regex)) for regex in regex_map["  
56         word"]  
57     ]  
  
58     matches = [match for match in matches if match]  
  
59  
60     closest_match = min(  
61         matches,  
62         key=lambda match: abs(get_sourceline(reference) - get_sourceline(  
63             match)),  
64     )  
  
65     for element in HTML_ELEMENTS_NAMES:
```

```
66     search_result = await llm_searcher.search(
67         regex_map["value"], closest_match.find_previous(element).text
68     )
69
70     complexity = search_result.get("complexity")
71     if complexity:
72         return complexity
73
74     raise FailedComplexityExtraction
75
76
77 def extract_name(dom_reference: Tag) -> str:
78     name = dom_reference.find_previous("h2")
79     if name and "tabtitle" not in name.get("class", ""):
80         return name.text
81
82
83 def fallback_search(reference_list: List[str], pattern: str) -> Optional[
84     str]:
85     for result in reference_list:
86         if search(pattern, result, flags=IGNORECASE):
87             return search_regex(TIME_COMPLEXITY_REGEX["fallback"], result)
88     return None
89
90 def extract_complexity_with_fallback(
91     dom_reference: Tag,
92 ) -> Tuple[Optional[str], Optional[str]]:
93     search_func = lambda x: search_regex(TIME_COMPLEXITY_REGEX["fallback"],
94         x.text)
95     next_results = [p.text for p in dom_reference.find_all_next("p") if
96         search_func(p)]
97     prev_results = [
98         p.text for p in dom_reference.find_all_previous("p") if search_func
99         (p)
100     ]
101
102     time_complexity = fallback_search(next_results, r"time") or
103         fallback_search(
104         prev_results, r"time"
105     )
106     space_complexity = fallback_search(
107         next_results, r"auxiliary|space"
108     ) or fallback_search(prev_results, r"auxiliary|space")
109
110     return time_complexity, space_complexity
```

```
def extract_codes_and_references(  
110     response: Response,  
) -> Tuple[List[Dict[str, str]], List[Tag]]:  
112     codes, references = [], []  
     code_tabs = response.soup.find_all("div", {"class": "responsive-tabs"})  
114  
     for tab in code_tabs:  
116         code_map = {}  
         references.append(tab)  
118  
         for lang in LANGUAGES:  
120             algorithm = tab.find_next("h2", {"class": "tabtitle"}, string=  
                 lang)  
             if algorithm:  
122                 code = extract_code(algorithm.find_next("td", {"class": "  
                     code"}))  
                 comments = extract_code_comments(code)  
124  
                 code_map[lang] = {  
126                     "code": code,  
                     "comments": comments,  
128                 }  
130             codes.append(code_map)  
132     return codes, references  
134  
async def extract_data(  
136     response: Response,  
) -> List[Dict[str, Union[str, bool, Dict[str, str]]]]:  
138     codes, dom_references = extract_codes_and_references(response)  
     if not codes:  
140         log.warning(f"Failed to find codes. URL: {response.url}")  
         return []  
142  
     complexities = []  
144     for ref in dom_references:  
         try:  
146             time_complexity = await extract_complexity_from_reference(  
                 ref, TIME_COMPLEXITY_REGEX  
148             )  
         except FailedComplexityExtraction:  
150             time_complexity, _ = extract_complexity_with_fallback(ref)  
152  
         try:  
             space_complexity = await extract_complexity_from_reference(  
154                 ref, AUXILIARY_SPACE_REGEX
```

```
    )
156     except FailedComplexityExtraction:
157         _, space_complexity = extract_complexity_with_fallback(ref)
158
159     complexities.append(
160         {
161             "time_complexity": time_complexity,
162             "trustable_time_complexity": bool(time_complexity),
163             "space_complexity": space_complexity,
164             "trustable_space_complexity": bool(space_complexity),
165         }
166     )
167
168     main_name = response.soup.find("div", {"class": "article-title"}).text
169     names = [main_name] + [
170         extract_name(ref) if extract_name(ref) else main_name
171         for ref in dom_references[1:]
172     ]
173
174     return [
175         {
176             "name": name,
177             "time_complexity": complexity["time_complexity"],
178             "trustable_time_complexity": complexity["
179                 trustable_time_complexity"],
180             "space_complexity": complexity["space_complexity"],
181             "trustable_space_complexity": complexity["
182                 trustable_space_complexity"],
183             "url": str(response.url),
184             "codes": code,
185         }
186         for name, complexity, code in zip(names, complexities, codes)
187     ]
```

- [saci/websites/geeks\\_for\\_geeks/login.py](#)

```
1 import os
2
3 from scraper.session.http_session import HttpSession
4 from scraper.session.response import Response
5 from scraper.session.utils import Methods
6 from scraper.websites.geeks_for_geeks.settings import LOGIN_URL
7
8
9 def get_login_payload():
```

```
11     return {
12         "reqType": "Login",
13         "user": os.getenv("USERNAME"),
14         "pass": os.getenv("PASSWORD"),
15         "rem": False,
16         "to": "https://auth.geeksforgeeks.org/?to=https://www.geeksforgeeks.org/",
17         "rem": "on",
18         "g-recaptcha-response": None,
19         "browserInfo": {},
20     }
21
22     async def follow_login(session: HttpSession):
23         await session.request(
24             method=Methods.GET.value,
25             data=get_login_payload(),
26             url=LOGIN_URL,
27         )
```

- `saci/websites/geeks_for_geeks/settings.py`

```
1 BASE_URL = "https://www.geeksforgeeks.org"
3 LOGIN_URL = "https://auth.geeksforgeeks.org/auth.php"
5 ALGORITHMS_LOCATION_URLS = [
6     f"{BASE_URL}/fundamentals-of-algorithms",
7     f"{BASE_URL}/data-structures",
8 ]
9
10 HEADERS = {
11     "Host": "www.geeksforgeeks.org",
12     "Accept": "text/html,application/xhtml+xml,application/xml;q=0.9,image/avif,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3;q=0.9",
13     "Accept-Encoding": "gzip, deflate",
14     "Upgrade-Insecure-Requests": "1",
15     "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/104.0.0.0 Safari/537.36",
16     "Sec-Fetch-Site": "none",
17     "Sec-Fetch-Mode": "navigate",
18     "Sec-Fetch-User": "?1",
19     "Sec-Fetch-Dest": "document",
20     "Sec-Ch-Ua": '"Chromium";v="105", "Not)A;Brand";v="8" ',
```



```
21     "Sec-Ch-Ua-Mobile": "?0",
22     "Sec-Ch-Ua-Platform": '"Windows"',
23 }
24
25 ALGORITHM_LINKS_JS_SCRIPT = """
26     () => {
27         const algorithm_links = []
28
29         const ordered_lists = document.getElementsByTagName('ol');
30         for (ol of ordered_lists) {
31             const hrefs = []
32             const a_tags = ol.getElementsByTagName('a')
33             for (a of a_tags) {
34                 hrefs.push(a.getAttribute('href'))
35             }
36             algorithm_links.push(hrefs)
37         }
38
39         return algorithm_links
40     }
41 """
42
43 TIME_COMPLEXITY_REGEX = {
44     "word": [
45         r"[Tt]ime [Cc]omplexit\\w+\\s?:|[Tt]ime [Cc]omplexit\\w+\\s?",
46     ],
47     "value": "time complexity",
48     "fallback": r"O\\s*\\(.*?\\)",
49 }
50
51 AUXILIARY_SPACE_REGEX = {
52     "word": [
53         r"[Aa]uxiliary [Ss]pace\\s?:|[Aa]uxiliary [Ss]pace\\s?",
54         r"[Ss]pace [Cc]omplexit\\w+\\s?:|[Ss]pace [Cc]omplexit\\w+\\s?",
55     ],
56     "value": "auxiliary/space complexity",
57     "fallback": r"O\\s*\\(.*?\\)",
58 }
59
60 HTML_ELEMENTS_NAMES = ["p", "li", "ul", "i"]
61
62 COMMENTS_STARTING_STRINGS = ["//", "# ", "/*"]
63
64 LANGUAGES = ["C", "C#", "C++", "Python3", "Python", "PHP", "Java", "
65     Javascript"]
```

- saci/websites/geeks\_for\_geeks/spider.py

```
from asyncio import gather
2
from scraper.exceptions import FailedExtraction
4 from scraper.observability.log import scraper_log as log
from scraper.schema.data_schemas import ScrapedAlgorithm
6 from scraper.session.http_session import HttpSession
from scraper.session.response import Response
8 from scraper.session.utils import Methods
from scraper.utils import remove_duplicates, retry
10 from scraper.websites.geeks_for_geeks.extract import extract_data
from scraper.websites.geeks_for_geeks.login import follow_login
12 from scraper.websites.geeks_for_geeks.settings import
    ALGORITHMS_LOCATION_URLS, HEADERS
14
def parse_algorithm_schema(algorithms_data: list) -> list:
16     return [ScrapedAlgorithm(**algorithm) for algorithm in algorithms_data]
18
def filter_algorithms_urls(algorithms_urls: list) -> list:
20     return [href for href in algorithms_urls if "geeksquiz" not in href]
22
async def extract_algorithm_data(response: Response) -> list:
24     if not response:
26         return []
28
29     data = await extract_data(response)
30     if not data:
31         raise FailedExtraction
32
33     return data
34
def extract_algorithms_urls(response: Response) -> list:
35     page_content = response.soup.find("div", {"class": "page_content"})
36     return [
37         a.attrs.get("href")
38         for ol in page_content.find_all("ol")[1:]
39         for a in ol.find_all("a")
40         if a.attrs.get("href")
41     ]
42
43
44 @retry(times=3, raise_exception=False, return_value=[])
async def follow_algorithm(session: HttpSession, url: str):
```

```
46     return await session.request(  
48         method=Methods.GET.value ,  
         url=url ,  
         callbacks=[extract_algorithm_data , parse_algorithm_schema] ,  
50     )  
  
52  
53 async def follow_algorithms(session: HttpSession , algorithms_urls: list) ->  
54     list:  
55     algorithms = await gather(  
56         *[follow_algorithm(session , url) for url in algorithms_urls]  
57     )  
58     return sum(algorithms , [])  
  
60 async def follow_algorithms_urls(session: HttpSession) -> list:  
61     algorithm_location_urls = await gather(  
62         * [  
63             session.request(  
64                 method=Methods.GET.value ,  
65                 url=url ,  
66                 callbacks=[  
67                     extract_algorithms_urls ,  
68                     filter_algorithms_urls ,  
69                     remove_duplicates ,  
70                 ] ,  
71             )  
72             for url in ALGORITHMS_LOCATION_URLS  
73         ]  
74     )  
75     return sum(algorithm_location_urls , [])  
  
76  
77  
78 async def run(url: str = None) -> list:  
79     session = HttpSession()  
80     session.default_headers = HEADERS  
  
81     log.info('Starting "Geeks for Geeks" algorithms extraction')  
82     await follow_login(session)  
  
83  
84     algorithms_urls = [url] if url else await follow_algorithms_urls(  
85         session)  
86     return await follow_algorithms(session , algorithms_urls)
```

## Apêndice C - Artigo SBC

### SACI - UMA FERRAMENTA PARA EXTRAÇÃO DE ALGORITMOS E SUAS COMPLEXIDADES DE TEMPO E ESPAÇO

Augusto Silva de Oliveira<sup>1</sup>, Carina Friedrich Dorneles<sup>2</sup>, Ricardo José Pfitscher<sup>3</sup>

<sup>1</sup>Universidade Federal de Santa Catarina (UFSC)

Departamento de Informática e Estatística (INE) – Florianópolis – SC – Brazil  
augusto.so@grad.ufsc.br, carina.dorneles@ufsc.br, ricardo.pfitscher@ufsc.br

**Abstract.** Determining the time and space complexity of a written code is a task that can be complex and is often essential, especially when aiming to achieve an efficient application. Despite being an important aspect in the education of developers, the analysis of algorithms regarding their complexity is a difficult task that is sometimes neglected. Some tools can make suggestions for code auto-completion as the programmer writes, but, unfortunately, they currently do not have the capability to consider complexities in their suggestions. Therefore, it is necessary for developers to have a solid understanding of algorithm analysis to ensure that the proposed solutions are efficient. With this in mind, this work proposes the creation of a data extraction tool capable of extracting algorithm codes along with their respective complexities in order to form a database. This database can be used for training tools, whether through machine learning or other techniques, aimed at determining the complexity of a piece of code and assisting developers in their routines. By the end of the work, a database with 501 distinct algorithms accompanied by their complexities was obtained, with good precision and data quality, as evidenced by the metrics used.

**Resumo.** Determinar a complexidade de tempo de execução e de espaço em um código escrito é uma tarefa que pode ser complexa e que muitas vezes é essencial, considerando-se que deseja-se obter uma aplicação eficiente. Apesar de ser um aspecto importante na formação do desenvolvedor, a análise de algoritmos quanto a sua complexidade é uma tarefa difícil que por vezes é negligenciada. Algumas ferramentas são capazes de fazer sugestões para o preenchimento automático do código conforme

o programador esteja escrevendo, mas, até o momento atual, ela infelizmente não possui capacidade de considerar as complexidades em suas sugestões. Por isso é necessário que o desenvolvedor tenha conhecimento sólido sobre análise de algoritmos para garantir que as soluções propostas sejam eficientes. Tendo isso em vista, esse trabalho propõe a criação de uma ferramenta de extração de dados, capaz de extrair códigos de algoritmos juntos de suas respectivas complexidades a fim de formar uma base de dados. Essa base de dados poderá ser usada para o treinamento de ferramentas, seja através de aprendizado de máquina ou outras técnicas, que visam determinar a complexidade de um trecho de código e auxiliar os desenvolvedores em suas rotinas. Ao fim do trabalho obteve-se uma base de dados com 501 algoritmos distintos acompanhados de suas complexidades, com boa precisão e qualidade dos dados, comprovado pelas métricas utilizadas.

## 1 INTRODUÇÃO

A complexidade de algoritmos mede o consumo de recursos computacionais, como tempo de execução e espaço em memória, em relação ao tamanho da entrada, sendo uma medida de eficiência essencial para o desempenho de uma aplicação (CORMEN, et al., 2009). Conhecer a complexidade de um trecho de código permite prever sua eficiência e otimizá-lo (SEDEGWICK, 2009). No entanto, identificar a complexidade pode ser uma tarefa difícil para muitos programadores devido à falta de hábito, prática ou conhecimento (HALL; CHAPMAN, 2015).

Há uma carência de dados sobre algoritmos e suas complexidades, como evidenciado por Sikka et al. (SIKKA et al., 2020a), que trabalharam com um banco de dados restrito de 931 códigos em Java. Pfitscher et al. (PFITSCHER et al., 2023) ampliaram esse banco de dados para 1325 amostras com dados extraídos do site Geeks for Geeks. O SACI visa expandir ainda mais esses dados, criando um repositório robusto e abrangente de algoritmos e suas complexidades, disponibilizado online.

A ferramenta proposta extrai algoritmos e suas complexidades de fontes pré-selecionadas, resultando em um dataset aberto para a comunidade científica e desenvolvedores. Exemplos de aplicação incluem ferramentas que retornam a complexidade de um código, como o GitHub Copilot, que poderia sugerir melhorias de desempenho. O SACI utiliza técnicas de *web scraping* para construir essa base de dados (ZHAO, 2017).

Um exemplo prático é apresentado por Rodenbusch (RODENBUSCH, 2023), com a API RTCC (*Running Time Complexity Calculator*), que utiliza dados do SACI para prever padrões de complexidade com alta precisão. Ferramentas de IA treinadas com esses dados são cruciais para melhorar a capacidade de predição e eficiência (WANG et al., 2022).

## 2 OBJETIVOS GERAL E ESPECÍFICOS

Desenvolver uma ferramenta de *scraping* capaz de extrair algoritmos e suas complexidades, com o objetivo de gerar uma base de dados com o conteúdo extraído.

Objetivos específicos:

1. Realizar a extração com expressões regulares
2. Realizar a extração com OpenAI
3. Tratar os dados extraídos e padronizá-los
4. Caracterizar os dados extraídos

## 3 TRABALHOS RELACIONADOS

Este estudo considerou quatro abordagens relacionadas ao desenvolvimento de ferramentas, apesar de não focarem na extração de complexidades, pois compartilham objetivos e características específicas.

### 3.1 Learning Based Methods for Code Runtime Complexity Prediction

Este trabalho de (SIKKA et al., 2020b) aborda a previsão da complexidade de tempo de execução de códigos utilizando aprendizado de máquina. Os autores criaram um conjunto de dados próprio, 'CoRCoD', devido à falta de dados abertos, mas enfrentaram problemas de acurácia devido ao tamanho reduzido do conjunto de dados.

### 3.2 qFEx - um crawler para busca e extração de questionários de pesquisa em documentos HTML

O (MATHIAS; DORNELES, 2021) é um web crawler desenvolvido para coletar questionários de pesquisa da internet. Utiliza heurísticas para identificar e extrair questionários, enfrentando desafios devido à diversidade de layouts de páginas web. Foi implementado em Java usando a biblioteca Crawler4j.

### 3.3 Ferramenta para coleta e comparação de dados de publicações acadêmicas dos professores com o currículo Lattes

Desenvolvido por (BRANCO, 2018), este *scraper* complementa o currículo Lattes com dados de sites como ResearchGate. Implementado em Java, utiliza bibliotecas como 'jsoup' e 'Java-string-similarity', mas enfrenta problemas de bloqueio e manutenção devido a mudanças nos sites.

### 3.4 CrawlEX: uma ferramenta para extração de dados na web configurável através de exemplos

O CrawlEX, descrito por (LESSA, 2022), é uma ferramenta configurável para extração de artigos da web, com uma interface visual em Python. Usa bibliotecas como 'beautifulsoup' e 'sqlite' para a extração baseada em exemplos fornecidos pelo usuário.

### 3.5 Desenvolvimento de API para predição da complexidade de tempo de execução de códigos por meio de AutoML

Em (RODENBUSCH, 2023), o autor combina bases de dados de *scraping* para criar modelos de classificação de complexidade de código, resultando em uma API que permite a predição de classes de eficiência e complexidade para códigos em Java, com acurácia de até 88,38%.

### 3.6 Análise dos trabalhos

Ao comparar quatro das obras apresentadas que implementaram alguma ferramenta de extração e contrastando-as com o trabalho atual, é perceptível que, embora haja distintos focos de extração, a maioria dos trabalhos é direcionada de forma especializada. A Tabela 1 ilustra essas diferenças.

O CrawlEX se destaca por tentar abordar uma perspectiva mais generalizada; no entanto, ainda necessita de configurações prévias, o que pode influenciar sua eficiência que não atinge a totalidade. Independentemente da abordagem, é comum que todos utilizem bibliotecas especializadas para navegar na estrutura DOM da página e nos elementos HTML. Essa dependência de configuração, seja ela feita pelo usuário ou já incorporada, é crucial, visto que cada página da web possui suas particularidades.

A presença ou ausência de uma interface gráfica parece estar diretamente ligada ao público-alvo de cada trabalho. Trabalhos como ECL e CrawlEX, que possuem interface gráfica, parecem ser orientados para usuários finais. Por outro lado, obras como o qFEx, CoRCoD e o SACI estão mais voltadas para a extração de dados para alimentar bases, não necessitando, portanto, de um elemento visual.

Em relação às linguagens de programação adotadas, duas predominam: Python e Java. O que é particularmente interessante é que, independentemente da linguagem escolhida, as bibliotecas auxiliares empregadas, como 'beautifulsoup' e 'jsoup', visam solucionar problemas semelhantes, demonstrando um padrão de desafios comuns na área de extração de dados web.

Tabela 1 - Comparação de trabalhos

Comparação de trabalhos					
Trabalho	Tipo	Item Alvo	Técnica	Algoritmo de Extração	Interface Gráfica
qFEx	Focado	Questionários	Heurísticas	Do autor	Não
ECL	Focado	Publicações	Similaridade textual	Do autor	Sim
CrawlEX	Genérico	Artigos	Comparação textual	Do autor e Ratcliff/Obershelp	Sim
CoRCoD	Focado	Algoritmos e Complexidades	API e Scra-ping	Do autor	Não
SACI	Focado	Algoritmos e Complexidades	Expressões regulares e OpenAI	Do autor	Não

## 4 SACI - SCRAPER DE ALGORITMOS E COMPLEXIDADES

### 4.1 Visão Geral

A ferramenta SACI extrai dados de algoritmos e suas complexidades de páginas web específicas, utilizando conceitos de *web scraping*. O fluxo da ferramenta é executado de forma assíncrona para aumentar a eficiência, e os dados extraídos são armazenados em uma base de dados MongoDB ou em um arquivo JSON local.

### 4.2 Implementação

A estrutura da ferramenta é modular, permitindo a inclusão de novas páginas. As principais pastas são:

1. **Results:** Armazena os resultados das extrações.
2. **Observability:** Contém os códigos responsáveis pelos logs e métricas de desempenho.
3. **Schema:** Define os esquemas JSON para validação dos dados.
4. **Session:** Implementa a sessão HTTP para requisições assíncronas.
5. **Websites:** Contém a lógica específica para cada página implementada.

### 4.3 Geeks4Geeks

O Geeks4Geeks <sup>1</sup> é um site com conteúdo voltado para programadores com uma variedade de materiais para aprendizado e, também, de entretenimento. repleto de

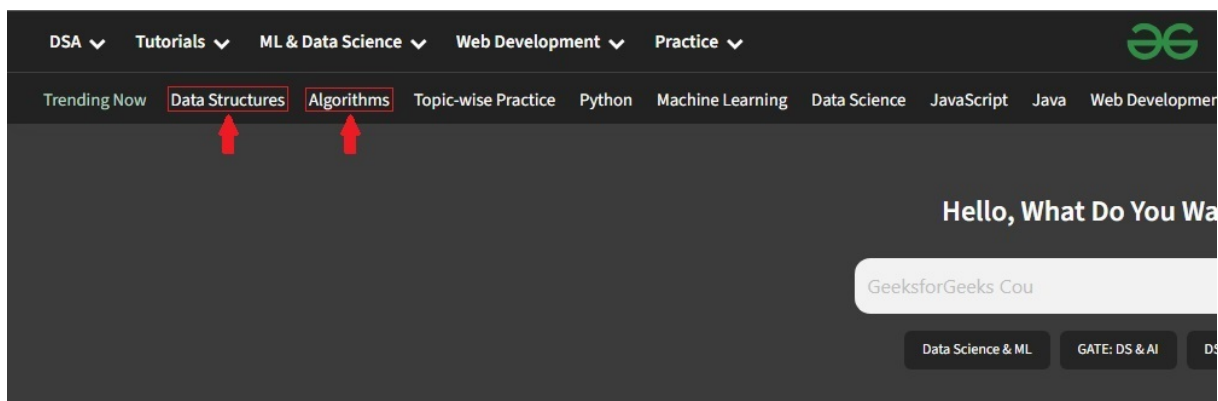
<sup>1</sup> <https://www.geeksforgeeks.org/>



artigos, tutoriais, aulas, seminários, é um ótimo lugar para quem procura aprender algo novo.

No topo do site, Figura 1, ficam evidenciados os principais tópicos: algoritmos de estrutura de dados, tutoriais, ciência de dados e aprendizado de máquina, desenvolvimento web, entre outros. Destes, o que interessa para o trabalho são os referentes a algoritmos e estrutura de dados. Dentro deles existem subtópicos e links, Figura 2, que direcionam para páginas com implementações de código, em múltiplas linguagens, acompanhadas das complexidades de tempo e espaço Figura . Por conta dessas características, essas duas seções, apontadas na Figura 3, foram escolhidas como fontes de dados a serem extraídos pelo SACI. Nessas, existe uma variedade de links de algoritmos, organizados conforme seus tipos e temas.

Figura 1 - Menu Geeks4Geeks



Fonte: elaborada pelo autor

Figura 2 - Página de algoritmos

**Topics:**

- Analysis of Algorithms
- Searching and Sorting
- Greedy Algorithms
- Dynamic Programming
- Pattern Searching
- Backtracking
- Divide and Conquer
- Geometric Algorithms
- Mathematical Algorithms
- Bit Algorithms
- Graph Algorithms
- Randomized Algorithms
- Branch and Bound
- Quizzes

**Analysis of Algorithms:**

1. Asymptotic Analysis
2. Worst, Average and Best Cases
3. Asymptotic Notations
4. Lower and Upper Bound Theory
5. Introduction to Amortized Analysis
6. What does 'Space Complexity' mean?
7. Polynomial Time Approximation Scheme
8. Accounting Method | Amortized Analysis
9. Potential Method in Amortized Analysis

**Searching and Sorting:**

1. Introduction to Searching Algorithms
2. Introduction to Sorting Algorithm
3. Stable and Unstable Sorting Algorithms
4. Lower bound for comparison based sorting algorithms
5. Can Run Time Complexity of a comparison-based sorting algorithm be less than  $N \log N$ ?
6. Which sorting algorithm makes minimum number of memory writes?

**Greedy Algorithms:**

1. Introduction to Greedy Algorithms

Fonte: elaborada pelo autor

Figura 3 - Complexidades

```
// Driver code
int main()
{
    int n = 5;
    printPascal(n);
    return 0;
}
```

**Output**

```
1
1 1
1 2 1
1 3 3 1
1 4 6 4 1
```

**Time Complexity:**  $O(n^2)$   
**Auxiliary Space:**  $O(1)$

So method 3 is the best method among all, but it may cause integer overflow for large values of n as it multiplies two integers to obtain values.

Fonte: elaborada pelo autor

#### 4.4 Scraper

O scraper segue um fluxo para localizar e extrair os códigos dos algoritmos e suas complexidades das páginas do Geeks4Geeks, utilizando expressões regulares e métodos 'fallback' para garantir a extração eficiente.

#### 4.5 Extract

Foram desenvolvidas duas implementações para a extração. Uma das versões faz uso de expressões regulares como o principal meio de identificação e extração das complexidades. Já a outra, no momento da extração, utiliza inteligência artificial através da API (*Application Programming Interface*) do OpenAI para a identificação e extração das complexidades. Todos os resultados das extrações podem ser encontrados no repositório<sup>2</sup> do GitHub e também no Kaggle<sup>3</sup>.

#### 4.6 Extração com expressões regulares

Cada página do site pode conter mais de um código/implementação. Esses códigos ficam localizados em um elemento HTML chamado aqui de bloco. Cada bloco contém o algoritmo codificado em diversas linguagens de programação e uma complexidade associada ao algoritmo. Ao acessar à página o método procura pelas seções que contém esses blocos, linha 2 do Algoritmo 1, guarda as referências de posição desses blocos na árvore DOM junto dos códigos. Em seguida, procura em volta de cada referência pelas complexidades de tempo e espaço, cada uma através de uma seleção prévia de expressões regulares, linhas 7 a 18 do Algoritmo 1. Por fim, extraímos o título do algoritmo na linha 19 do Algoritmo 1 e adicionamos um objeto resultado a uma lista, linhas 20 a 32 do Algoritmo 1. Para a extração dos títulos dos códigos, linha 19, também utiliza-se a referência de localização dos blocos de código. É extraído o título principal da página que geralmente é o título do primeiro algoritmo e os restantes, para retornar uma lista de dicionários Python contendo o nome do algoritmo, a complexidade de tempo, a complexidade de espaço, o URL e o código em todas as linguagens, linhas 20 a 32. Essa lista é tratada pelo *scraper* na última etapa retornando os dados de acordo com o esquema JSON. As chaves 'tempoConfiavel' e 'espaçoConfiavel' do objeto resultado tem relação com o uso ou não do método 'fallback' apresentado mais a frente. Caso o método seja utilizado essas chaves recebem como valor o booleano 'True', caso contrário 'False'.

<sup>2</sup> <https://github.com/Augustives/SACI>

<sup>3</sup> <https://www.kaggle.com/datasets/augustives/algorithms-and-their-complexities>

---

**Algoritmo 1 - Código mostra a função de extração principal em alto nível**

---

```
1: Function extrairDados(response)

2: codigos, referencisaDom ← extrairCodigosReferencias(response)
3: if not codigos then
4:   Log aviso
5:   return []
6: end if

7: complexidades ← []
8: for ref in referencisaDom do
9:   try:
10:  complexidadeTempo ← extrairComplexidadeDeReferencia(ref, regexTempo)
11:  catch Erro:
12:  complexidadeTempo ← extrairComplexidadeComFallback(ref)

13:  try:
14:  complexidadeEspaço ← extrairComplexidadeDeReferencia(ref, regexEspaço)
15:  catch Erro:
16:  complexidadeEspaço ← extrairComplexidadeComFallback(ref)

17:   Insere complexidades na lista
18: end for

19: titulos ← extrairTitulos(ref)

20: resultados ← []
21: for (nome, complexidade, codigo) in zip(nomes, complexidades, codigos) do
22:   resultado ← {
23:     'nome': nome,
24:     'complexidadeTempo': complexidade['complexidadeTempo'],
25:     'tempoConfiavel': complexidade['tempoConfiavel'],
26:     'complexidadeEspaço': complexidade['complexidadeEspaço'],
27:     'espaçoConfiavel': complexidade['espaçoConfiavel'],
28:     'url': url,
29:     'codigo': codigo,
30:   }
31:   Insere resultado na lista
32: end for

33: return resultados =0
```

---

No momento da extração das complexidades, são utilizados a referência da localização dos códigos, linhas 10 e 14 do Algoritmo 1, métodos do 'BeautifulSoup' e expressões regulares. A extração conta com dois caminhos: o padrão com regras mais específicas e um que chamamos de 'fallback', no qual são utilizadas regras mais generalizadas e inclusivas a fim de conseguir extrair a informação caso ocorra alguma falha

durante o fluxo padrão.

O método padrão, representado no Algoritmo 2, procura por todas as palavras *Time Complexity* ou *Space Complexity* nas proximidades do bloco de códigos, linha 4. Para cada palavra existe uma lista de expressões regulares, para cada palavra encontrada elas são inseridas em uma lista de 'matches', linha 3 e 6. Então, procura o 'match' mais próximo da referência, linha 8, com um cálculo simples: para cada 'match' da lista subtrai-se a linha em que a referência está com a linha que o 'match' encontra-se, dessa forma, o menor valor dentre todos os calculados é o mais próximo. A seguir, o método percorre outra lista de expressões regulares, dessa vez focada em capturar os valores das complexidades. Para cada expressão regular de valor, itera os elementos HTML, para os quais é feita uma pesquisa, linha 11. Caso nenhuma complexidade consiga ser extraída, é levantado um erro de extração, linha 17.

Sendo o erro levantado, o fluxo de 'fallback', tenta extrair novamente a complexidade, e uma vez utilizado esse método, uma 'flag' é adicionada aos resultados indicando seu uso.

---

Algoritmo 2 - Código mostra a função 'extrairComplexidadeDeReferencia' em alto nível

---

```

1: Function extrairComplexidadeDeReferencia(referencia, regex)
2: for regex in regexes['palavra'] do
3:   matches ← [ ]
4:   palavra ← referencia.find_next(string=compile(regex))

5:   if palavra then:
6:     Insere palavra em matches
7:   end if

8:   match_mais_proximo ← Pega o match mais próximo da ref

9:   for regex in regexes['valor'] do
10:    for elemento in elementosHTML['valor'] do
11:      complexidade ← busca_regex(regex, match_mais_proximo)

12:      if complexidade then:
13:        return complexidade
14:      end if
15:    end for
16:  end for

17:  raise FalhaExtraçãoComplexidade
18: end for

```

---

A Figura 4 apresenta as expressões regulares criadas e utilizadas na extração dos dados. Aquelas que buscam pelas palavras, linhas 3, 16 e 17, tentam encontrar *time complexity*, *space complexity* e suas variações de escrita. As responsáveis

pelos valores, linhas 5 a 10 e 19 a 28, são uma combinação da palavra com o valor da complexidade, buscando, por exemplo, algo como *Time Complexity: O(log N)*. Todas foram criadas com base nos diversos padrões com que as complexidades são descritas nos textos das páginas do site. Por fim, tem-se a que é usada no método 'fallback', linhas 11 e 29.

Figura 4 - Regex de extração

```

TIME_COMPLEXITY_REGEX = {
2   'word': [
4     r'[Tt]ime [Cc]omplexit\\w*\\s?:|[Tt]ime [Cc]omplexit\\w*\\s?',
6   ],
8   'value': [
10    r'[Tt]ime [Cc]omplexit\\w*\\s?:.*?(O\\s*\\(.*?\\))',
12    r'[Tt]ime [Cc]omplexit\\w*\\s?.*?(O\\s*\\(.*?\\))',
14    r'[Tt]ime [Cc]omplexit\\w*\\s.*?\\sis\\s.*?(O\\s*\\(.*?\\))',
16    r'[Tt]ime [Cc]omplexit\\w*\\s?:?.*?\\sis\\s(\\w*)',
18  ],
20  'fallback': r'O\\s*\\(.*?\\)',
22 }

AUXILIARY_SPACE_REGEX = {
24   'word': [
26     r'[Aa]uxiliary [Ss]pace\\s?:|[Aa]uxiliary [Ss]pace\\s?',
28     r'[Ss]pace [Cc]omplexit\\w*\\s?:|[Ss]pace [Cc]omplexit\\w*\\s?',
30  ],
32   'value': [
34     r'[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?.*?(O\\s*\\(.*?\\))',
36     r'[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?.*?(O\\s*\\(.*?\\))',
38     r'[Ss]pace [Cc]omplexit\\w*\\s?:.*?(O\\s*\\(.*?\\))',
40     r'[Ss]pace [Cc]omplexit\\w*\\s?.*?(O\\s*\\(.*?\\))',
42     r'[Aa]uxiliary [Ss]pace\\s?[\\[Cc\\]omplexit\\w*]?.*?\\sis\\s.*?(O\\s*\\(.*?\\))',
44     r'[Ss]pace [Cc]omplexit\\w*\\s.*?\\sis\\s.*?(O\\s*\\(.*?\\))',
46     r'[Aa]uxiliary [Ss]pace\\w*\\s?:?.*?\\sis\\s(\\w*)',
48     r'[Ss]pace [Cc]omplexit\\w*\\s?:?.*?\\sis\\s(\\w*)',
50  ],
52  'fallback': r'O\\s*\\(.*?\\)',
54 }

```

## 4.7 Extração com OpenAI

Para a extração com OpenAI utilizou-se da API oferecida <sup>4</sup>, e para sua implementação foi utilizada uma biblioteca do Python chamada de Langchain <sup>5</sup>. Ela oferece uma integração fácil com modelos de linguagem ou LLMs (*Large Language Models*), como o GPT-3 e outros similares, além de oferecer aos desenvolvedores um conjunto de ferramentas e abstrações que facilitam a criação de aplicações interativas e complexas. A premissa central da biblioteca Langchain é permitir a construção de correntes de 'prompts' e respostas; ou seja, sequências de interações com modelos de linguagem que podem ser encadeadas para realizar tarefas sofisticadas. Dentro do arquivo *utils* do SACI está definida uma classe que faz o uso desta biblioteca e realiza essa comunicação com a OpenAI.

---

Algoritmo 3 - Código mostra a sessão http em alto nível

---

```

1: Class LlmSearcher
2: throttler ← Throttler(rate_limit=100, period=60)

3: Init()
4: self.LLM ← OpenAI(openai_api_key=...)
5: self.PROMPT ← PromptTemplate.from_template(template_do_prompt)
6: self.LLM_CHAIN ← LLMChain(prompt=self.PROMPT, llm=self.LLM)

7: Method search(complexidade, texto)
8: async with throttler:
9: resposta ← self.LLM_CHAIN.run(complexity=complexidade, text=texto)
10: return json.loads(resposta)

```

---

A classe, representada no Algoritmo 3, atua como um *singleton*, isto é, só existe uma instância sua durante a vida útil da aplicação. Dentro dela estão definidos o LLM que conecta com o modelo de linguagem da OpenAI, o 'PROMPT' que define como o modelo irá responder e a LLM\_CHAIN que une o 'PROMPT' com o modelo em si. Além disso, existe um 'throttler' que limita quantas 'requests' são feitas por minuto, pois, a API da OpenAI possui certas restrições.

Figura 5 - Prompt do LLM

I am using you during a scraping operation in which I need to extract from a piece of text the values of complexity/code complexity.

2 Consider that sometimes the space complexity is also called as auxiliary space.

<sup>4</sup> <https://openai.com/index/openai-api/>

<sup>5</sup> <https://www.langchain.com/>

```
Consider that both complexitys often come in this format: 'Some type of
  complexity: value of complexity'.
4 But they can also come inside the text like: 'The type of complexity of the
  given ... is ...'.
Give the answer in JSON format with no line breaks, with a key called '
  complexity' and the value for the key is your answer.
6 If you cant determine the answer give the json with a null in the value.
What is the {complexity} that is written in the following text: '{text}'?
```

O 'PROMPT' utilizado é descrito na Figura 5, com ele é possível explicar o contexto em que o SACI está, demonstrando alguns exemplos de como os dados irão chegar, seus formatos e diferenças, sendo ainda explicado como a resposta deverá ser retornada e qual pergunta é feita. Com isso o modelo de linguagem consegue entender que está sendo feito a pesquisa pelos valores de complexidade do texto que ele está recebendo.

O método de extração das complexidades sobre as referências também sofreu uma pequena alteração. Não é mais iterado a lista de 'regexes' responsáveis pela extração dos valores de complexidade, como era feito no antes, pois esse trabalho agora é feito pelo LLM.

## 5 EXPERIMENTOS

Para cada página em que foi implementada um *scraper* foram realizadas avaliações a fim de avaliar se a extração das informações teve um bom resultado ou não. Tentou-se garantir que a quantidade de extrações fosse boa e que as informações extraídas estivessem corretas. A seguir, é apresentada a configuração do ambiente que executou o *scraper*, descrevem as métricas utilizadas, caracterizam os dados extraído e detalham os experimentos e seus respectivos resultados a partir do uso do SACI.

### 5.1 Métricas

Para caracterizar a fonte de dados foram criados histogramas e gráficos para exibir o comportamento e a distribuição dos dados extraídos. Esses gráficos ajudam a visualizar padrões e tendências que podem não ser imediatamente aparentes através de uma simples observação.

Para melhor entender os resultados de cada extração foram selecionadas métricas de avaliação. A partir dos resultados é criado uma matriz de confusão e calculada a acurácia, revocação, precisão e F1-Score. Utilizam melhor essas métricas em conjunto do que individualmente, pois cada uma oferece uma perspectiva única sobre a performance do processo. Estas serão usadas para comparar a extração realizada com expressões regulares e a feita com OpenAi.



## 5.2 Resultados

A partir da análise do SACI na versão com expressões regulares, considerando as métricas citadas anteriormente, pode-se observar que os resultados obtidos mostraram um desempenho promissor na extração de ambas as complexidades, conforme demonstram as matrizes de confusão para tempo e espaço da Figura 6.

Figura 6 - Resultados Regex

		Valor Predito	
		Sim	Não
Real	Sim	391	11
	Não	88	11

Fonte: elaborada pelo autor

		Valor Predito	
		Sim	Não
Real	Sim	419	1
	Não	47	34

Fonte: elaborada pelo autor

A validação dos resultados foi feita de forma manual onde cada URL, onde ocorreu alguma extração de algoritmo, foi acessada e comparados os valores que foram extraídos com os valores presentes na página web.

Nota-se um bom desempenho na identificação de verdadeiros positivos para ambos os tipos de complexidade, com uma quantidade significativa de registros corretamente identificados pelo SACI. Porém, na complexidade de tempo observa-se um número considerável de falso-positivos, totalizando 88 casos. Isso indica que alguns dados errôneos foram incluídos na base de dados. A explicação para existência desse número um pouco mais elevado de falso-positivos é o uso do método de 'fallback', já que o mesmo parte do princípio de ser genérico e coletar a qualquer custo alguma complexidade.

Para falso-negativos, nota-se também sua existência significativa apenas para a complexidade de tempo, com 11 registros, o que sugere que alguns dados foram ignorados durante a extração. A baixa quantidade de verdadeiro-negativos na complexidade de tempo pode significar uma menor capacidade de filtrar dados irrelevantes.

Por outro lado, a complexidade de espaço apresentou um desempenho superior. A quantidade de verdadeiros negativos foi maior, com 34 registros, e os falsos negativos foram mínimos, com apenas 1 registro, mostrando uma melhor performance em comparação com a complexidade de tempo. Isso se deve ao fato de que a complexidade de espaço é apresentada de forma mais padronizada em todas as páginas, resultando em uma melhor identificação dos dados.

Na versão OpenAI notamos uma melhora geral em todos os parâmetros se comparados com a versão padrão, mas alguns apresentaram diferenças significativas.

Figura 7 - Resultados OpenAI

		Valor Predito	
		Sim	Não
Real	Sim	451	6
	Não	32	12

Fonte: elaborada pelo autor

		Valor Predito	
		Sim	Não
Real	Sim	424	1
	Não	37	39

Fonte: elaborada pelo autor

Tabela 2 - Métricas das complexidades de tempo e espaço

Métrica	Tempo Re-gex	Espaço Re-gex	Tempo OpenAI	Espaço OpenAI
Acurácia	0,8023	0,9041	0,9241	0,9241
Precisão	0,8162	0,8991	0,9337	0,9197
Revocação	0,9726	0,9976	0,9868	0,9976
F1 Score	0,8876	0,9458	0,9595	0,9571

Fonte: elaborada pelo autor

Os resultados alcançados pelo SACI são motivo de encorajamento. Eles evidenciaram alta precisão, revocação e acurácia em todas as complexidades analisadas. Notadamente, a implementação das técnicas da OpenAI trouxe melhorias substanciais sobre a versão padrão em vários aspectos críticos. Contudo, é importante reconhecer a dependência desses resultados à estrutura da fonte de dados. Mudanças futuras nos sites podem influenciar o desempenho da ferramenta.

A versão aprimorada do SACI, equipada com os modelos da OpenAI, mostrou uma redução expressiva tanto em falso-positivos quanto em falso-negativos, um diferencial que potencializa a confiabilidade dos dados extraídos. Tal performance sugere que, mesmo diante de estruturas de dados variáveis nos sites-alvo, a ferramenta possui a versatilidade para se adaptar e manter alto padrão de qualidade.

### 5.3 Fonte de dados

Todos os dados usados neste trabalho tiveram sua extração no dia '24/08/2023', conforme descrito previamente na Seção 4.4. Foi realizada uma análise sobre cada página onde houve a extração de uma ou mais complexidades. De um total de 235

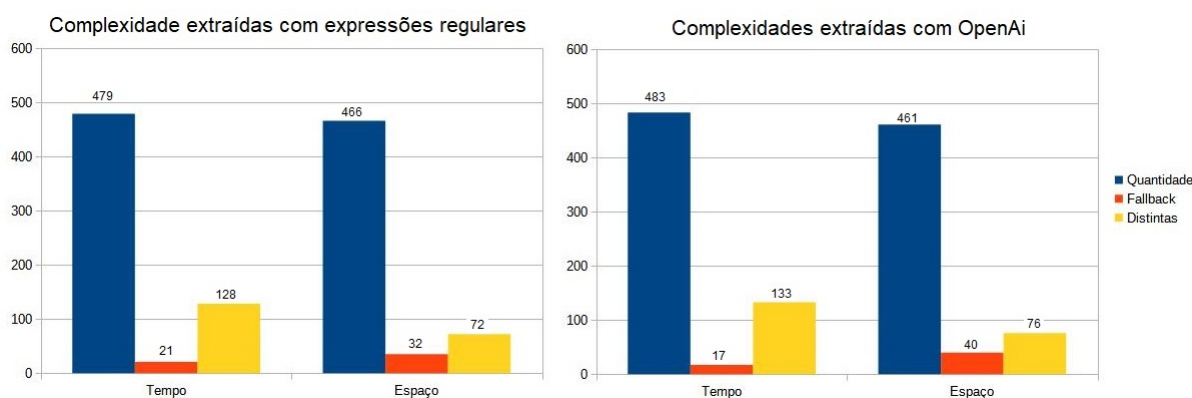
URLs, com 501 algoritmos disponíveis através delas, para cada página, executou-se uma função que percorria todos os nodos da árvore HTML e retornava qual a profundidade máxima encontrada e quantidade de nodos existentes.

#### 5.4 Conjunto de dados extraídos

A caracterização do conjunto de dados extraídos foi feita, primeiramente, avaliando o quanto a ferramenta de coleta utilizou da técnica de 'fallback' para extrair as complexidades. Os títulos dos algoritmos e os outros atributos coletados não fazem parte da caracterização. Conforme a Figura 8, 479 algoritmos tiveram suas complexidades de tempo extraídas e 466 suas complexidades de espaço; dentre estes, 21 fizeram uso do 'fallback' para extrair a complexidade de tempo e 32 para a de espaço. Foram obtidas 128 complexidades de tempo distintas e 72 complexidades de espaço distintas.

Um dos motivos para não se ter todas as complexidades de tempo extraídas com sucesso é a falta de um padrão na página. Em algumas páginas as complexidades estão escritas de maneira diferente, em locais diferentes, fazendo, assim, com que a expressão regular não as identifique. No caso da complexidade de espaço, alguns algoritmos realmente não a possuem e, em outros casos, ocorreu o mesmo problema que com a extração da complexidade de tempo.

Figura 8 - Complexidades extraídas



Fonte: elaborada pelo autor

Na nova versão do método de extração, que incorpora as capacidades da OpenAI, observou-se um incremento na eficiência da extração de dados. Como exibido na Figura 8, esta versão avançada conseguiu extrair com sucesso as complexidades de tempo de 483 algoritmos, e as complexidades de espaço de 461 algoritmos, superando um pouco a marca anterior.

Além disso, a dependência da técnica de 'fallback' foi reduzida para apenas 17 casos na extração das complexidades de tempo indicando, com isso, uma melhoria. Nas complexidades de espaço houve um pequeno aumento de 32 para 40 casos. Em termos de diversidade, o SACI identificou 133 complexidades de tempo distintas e 76 de espaço, um aumento em relação às 128 e 72 variantes respectivamente encontradas na versão que utilizava expressões regulares.

A classificação das complexidades em ambas versões foi feita através de expressões regulares, utilizando-se de um bloco de código de condicionais. Caso a expressão capturasse algo na complexidade, significa que tal complexidade é pertencente a tal categoria.

```

1  Constante = r'O\ (1\)'
   Linear = r'O\ (w\)'
3  Exponencial = r'O\ (d\^?\w\)'
   Fatorial = r'O\ (.!\.!\)'
5  Logarítmica = r'O\ (.!log.!\)'
   Polinomial = r'O\ (w\s*[\^*\+]?[s*[\d\w\s]+.!\)' e r'O\ (d+\ (w\*w\)'

```

## 6 CONSIDERAÇÕES FINAIS E TRABALHOS FUTUROS

A criação de ferramentas para identificar a complexidade de algoritmos e gerar *insights* para o usuário requer uma base de dados de algoritmos e suas complexidades. O SACI apresentou resultados satisfatórios, formando uma base de dados adequada com eficácia comprovada pelas métricas calculadas.

Durante o desenvolvimento, foram enfrentadas dificuldades na construção de um *scraper* genérico devido à diversidade de layouts de páginas. A reutilização de trechos de código facilitou futuras implementações. A detecção e extração dos elementos desejados foram realizadas com expressões regulares e um método *fallback*. Problemas com *anti-bots* foram mitigados com o envio de cabeçalhos adequados e limitações de requisições.

Para trabalhos futuros, melhorias podem incluir:

1. **Melhorar as métricas:** Adicionar mais atributos além das complexidades de tempo e espaço.
2. **Melhorar as técnicas de extração:** Refinar expressões regulares e usar inteligência artificial para identificar e extrair dados relevantes.
3. **Expandir fontes de dados:** Implementar a coleta de mais códigos e complexidades de diversas fontes.

A base de dados resultante possui 501 algoritmos que podem ser utilizados por pesquisadores para diversos fins, como:

1. **Análise Comparativa de Algoritmos:** Comparar o desempenho de diferentes algoritmos.
2. **Otimização de Algoritmos:** Identificar oportunidades para melhorar algoritmos existentes.
3. **Previsão de Desempenho de Algoritmos:** Construir modelos preditivos para estimar o desempenho dos algoritmos.
4. **Estudo Teórico de Algoritmos:** Investigar características teóricas dos algoritmos.
5. **Desenvolvimento de Novos Algoritmos:** Inspirar e orientar o desenvolvimento de novos algoritmos.

A base de dados construída oferece uma riqueza de informações que podem ser exploradas em diversas áreas, proporcionando avanços significativos em análise, otimização e desenvolvimento de algoritmos.

## 7 REFERÊNCIAS

CORMEN, Thomas H. et al. Algoritmos: Estruturas de Dados e Algoritmos para Resolução de Problemas. Rio de Janeiro: LTC, 2009.

SEDGEWICK, Robert. Analyzing algorithms. *Communications of the ACM*, ACM Press, v. 52, n. 1, p. 139–148, 2009.

HALL, Mark; CHAPMAN, Barry. Measuring algorithm complexity. *ACM SIGPLAN Notices*, ACM Press, v. 48, n. 1, p. 16–23, 2015.

SIKKA, Jagriti et al. Learning Based Methods for Code Runtime Complexity Prediction. In: SPRINGER. EUROPEAN Conference on Information Retrieval. [S.l.: s.n.], 2020. P. 313–325.

PFITSCHER, Ricardo J. et al. Estimating Code Running Time Complexity with Machine Learning. In: NALDI, Murilo C.; BIANCHI, Reinaldo A. C. (Ed.). *Intelligent Systems*. Cham: Springer Nature Switzerland, 2023. P. 400–414.

ZHAO, Bo. Web scraping. *Encyclopedia of big data*, Springer Living ed. Cham, p. 1–3, 2017.

RODENBUSCH, Gabriel Braun. Desenvolvimento de API para predição da complexidade de tempo de execução de códigos por meio de AutoML. [S.l.: s.n.], 2023. TCC (graduação) - Universidade Federal de Santa Catarina, Campus Joinville, Engenharia Mecatrônica.

WANG, Meng et al. A Survey on Large-Scale Machine Learning. *IEEE Transactions on Knowledge and Data Engineering*, v. 34, n. 6, p. 2574–2594, 2022. DOI: 10.1109/TKDE.2020.3015777.

SIKKA, Jagriti et al. Learning Based Methods for Code Runtime Complexity Prediction. In: JOSE, Joemon M. et al. (Ed.). *Advances in Information Retrieval*. Cham: Springer International Publishing, 2020. P. 313–325.

MATHIAS, Gilney; DORNELES, Carina. qFEx - um crawler para busca e extração de questionários de pesquisa em documentos HTML. In: *ANAIS do III Dataset Showcase Workshop*. Rio de Janeiro: SBC, 2021. P. 1–10. DOI: 10.5753/dsw.2021.17409.

BRANCO, Arthur Machado. Ferramenta para coleta e comparação de dados de publicações acadêmicas dos professores com o currículo lattes.

LESSA, Marcos Aurélio. *CrawLEX: uma ferramenta para extração de dados na web configurável através de exemplos*.