

Extração de informações da Plataforma Lattes

Alleff Dymytry Pereira de Deus

UERGS

Porto Alegre, Brasil

alleffdymytry@gmail.com

Resumo—Este documento reúne as informações necessárias para o entendimento da aplicação desenvolvida bem como as formas de modificação da mesma para obter outras informações. Dentro deste documento serão apresentadas as funções realizadas bem como as adversidades enfrentadas e os meios de contornar estas adversidades.

I. INTRODUÇÃO

As informações da vida acadêmica de cada professor possui um grande apreço para o desenvolvimento e propagação no incentivo a educação, estas informações normalmente contidas na plataforma Lattes. Todavia estas informações acabam sendo centralizadas e não possuindo uma separação em relação a cada instituição, sendo assim a extração das informações se faz necessária para a utilização destas informações em outras plataformas. As informações extraídas podem ser utilizadas em novas plataformas, podendo ser utilizadas e disponibilizadas em formas diferentes, alguns exemplos: tamanho diferente, filtros por áreas, classes de texto, etc.

O objetivo do presente trabalho é o desenvolvimento de uma aplicação para a extração das informações dos professores da Universidade Estadual do Rio Grande do Sul para a sua utilização em uma plataforma própria da universidade. As informações extraídas são armazenadas em arquivos para uso futuro na plataforma.

II. FUNDAMENTAÇÃO TEÓRICA

Para o desenvolvimento da aplicação proposta os seguintes conhecimentos e ferramentas foram necessários.

A. Python

A linguagem de programação Python foi utilizada para o desenvolvimento da aplicação devido as seguintes vantagens: alto desempenho, orientada a objetos, grande acervo de bibliotecas e integração com as mais diversas plataformas.

B. HTML

A linguagem de programação de páginas web HTML (HyperText Markup Language - Linguagem de Marcação de Hipertexto) é padrão na representação das mais diversas páginas na internet. A estrutura do HTML é dividida em blocos, sempre possuindo um bloco inicial e um bloco final, estes blocos representam os elementos que são possíveis de serem incluídos na páginas.

C. Selenium

O Selenium é um conjunto de bibliotecas e ferramentas que permitem a automação em navegadores web. É possível emular ações programadas dentro dos navegadores, assim é possível testar funções desenvolvidas e analisar possíveis falhas. O selenium possui biblioteca para o seu uso com a linguagem python,

III. MATERIAIS E MÉTODOS

As etapas da aplicação foram desenvolvidas utilizando a linguagem Python em conjunto com a biblioteca Selenium. A aplicação foi desenvolvida e testada em uma máquina com um processador AMD Ryzen 5 3400G de 3,7GHz, memória RAM de 16GB e placa de vídeo Radeon RX 590 de 8GB de VRAM.

O fluxo de funcionamento da aplicação proposta pode ser notada na Figura 1, onde para cada professor é realizado o fluxo do início ao fim.

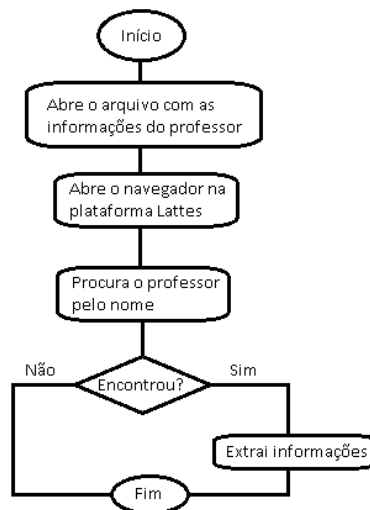


Figura 1. Fluxograma da aplicação.

A. Abrir Arquivos

A etapa de abertura dos arquivos se refere ao processo de abrir um arquivo previamente disposto para a aplicação contendo o nome completo de cada professor. O arquivo que contém o nome de cada professor deve ser salvo no formato de .csv e cada nome deve estar em uma nova linha, um exemplo

de como as informações deve ser organizadas esta disposta na Figura 2.

```
Name
Adriana Barni Truccolo
Daiana Bortoluzzi Baldoni
Edilma Machado de Lima
Fani Averbuh Tesseler
Laila Gabrielle Naymaer Gonçalves
Mariana Vargas Ferreira
Martha Giudice Narvaz
Pedro Luis Fagundes do Amaral
```

Figura 2. Exemplo do arquivo de informações do professor.

A primeira linha do arquivo deve conter a informação de cada coluna do arquivo, neste caso o nome do professor, caso mais informações sejam adicionadas deve-se separar cada informação por vírgulas e colocar uma descrição do que foi adicionado, como por exemplo se a data de nascimento for adicionada, a primeira linha da Figura 2 seria: "Name,dataNasc".

B. Abrir o Navegador

A biblioteca Selenium possui funções para a abertura automática do navegador web escolhido na aplicação, neste caso o navegador Chromium foi o utilizado, contudo para tal abertura é necessário baixar o arquivo de driver do Selenium para o navegador escolhido. O arquivo deve ser colocado na mesma pasta que a aplicação.

O navegador tem sua página diretamente direcionada para a página de pesquisa de pesquisadores da plataforma Lattes, facilitando a busca dos professores. Com a janela do navegador na página correta, foi necessário realizar uma análise no código fonte da página em questão para obter o valor da caixa de inserção de texto, assim sendo possível colocar o nome do professor na pesquisa. Com o valor HTML do elemento, foi necessário realizar a busca utilizando a função de pesquisa de elementos do Selenium, esta pesquisa pode ser feita com diversos valores, como por exemplo: o identificador do elemento, a classe, o caminho xpath e entre outros. A forma como a análise foi realizada pode ser vista na Figura 3.

A extração das informações contidas em uma página web ocorre de forma manual, sendo necessário analisar previamente as informações que devem ser extraídas.

C. Pesquisa Professor

O elemento de inserção de texto deve ser encontrado na página web previamente, assim que o elemento é localizado é realizado a inserção do nome completo do professor e a inserção da tecla "ENTER" é realizada em seguida, esta inserção da tecla "ENTER" realiza a pesquisa, assim não é necessário realizar a procura do botão de pesquisa na página.

A pesquisa é realizada e a página é recarregada com as informações do site, as informações disponibilizadas são os perfis dos professores na base do Lattes. Estes perfis são clicáveis, abrindo uma janela para o redirecionamento para o currículo da pessoa em questão. O currículo aberto gera uma

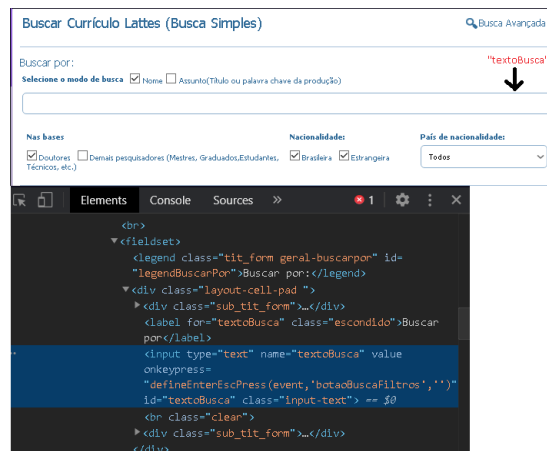


Figura 3. Análise do código HTML.

nova página no navegador, sendo assim é necessário fazer a troca de página para poder obter as informações do professor. Caso a pessoa não seja encontrada dentro de um tempo limite, a aplicação avisa que a pessoa em questão não foi encontrada e começa a pesquisar o próximo.

A troca de páginas segue o seguinte fluxo de operação: 1) deve-se obter todas as páginas abertas em uma lista, 2) salvar em uma variável o valor da página atual, para realizar o retorno posteriormente, 3) selecionar a página que não tem o valor da página original, 4) realizar a troca de página. A troca é realizada automaticamente pelas funções do Selenium. Com a página trocada é possível obter as informações da pessoa.

As informações extraídas da página são tratadas, para não conter espaços em branco e informações desnecessárias ou dispensáveis no atual momento. As informações tratadas são armazenadas em um arquivo no formato .json, o arquivo possui o nome da pessoa encontrada. Um exemplo do tratamento e organização das informações extraídas pode ser notada na Figura 4.

```
{
  "NOME": "Adriana Barni Truccolo",
  "PROFESOR": {
    {
      "ANO": "2020 - Atual",
      "TITULO": "Contação de Histórias on line como estratégia de comunicação afetiva com crianças em tempos de pandemia da COVID",
      "SITUACAO": "Situação: Em andamento; Natureza: Pesquisa."
    },
    {
      "ANO": "2020 - Atual",
      "TITULO": "Contação de histórias como estratégia de humanização à hospitalização infantil",
      "SITUACAO": "Situação: Desativado; Natureza: Pesquisa."
    }
  ]
}
```

Figura 4. Exemplo das informações extraídas.

As informações armazenadas em um arquivo no formato .json pode ser utilizada em diferentes aplicações, além de possuir um formato de descrição leve, onde caso necessário que uma nova informação seja adicionada manualmente, não seja um tarefa árdua, já que o padrão de inserção de informações pode ser copiado das informações anteriores.

IV. RESULTADOS E DISCUSSÕES

Com o desenvolvimento da aplicação proposta, foi possível obter as informações de projetos de pesquisa e projetos de extensão de cada professor da Universidade Estadual do Rio

Grande do Sul. Contudo algumas dificuldades foram enfrentadas, sendo estas dificuldades discutidas nesta seção.

A. Descrição dos Projetos

As informações descritas nos perfis de cada professor possuem algumas informações em comum, como por exemplo: 1) ano do projeto, 2) nome do projeto e 3) situação do projeto. Contudo a descrição é uma informação que nem todos os projetos possuem, fazendo com que não se possa utilizar uma única técnica de extração e tratamento das informações, um exemplo desta situação pode ser vista na Figura 5. Para contornar esta adversidade foi necessário realizar uma análise nas informações extraídas, onde mesmo que um projeto não contivesse uma descrição, o título de "Descrição:" seria inserido junto de cada projeto extraído. A solução para contorno da falta de descrição nos projetos, faz com que seja possível utilizar um padrão de extração para todos os projetos de cada professor.

2015 - 2015	Educação Integral e as Avaliações em Larga Escala: Desafios Contemporâneos em Tempos de Inclusão no Ciclo da Alfabetização
	Projeto certificado pelo(a) coordenador(a) Rochelle da Silva Santana em 21/07/2017.
	Descrição: 20 horas semanais de atividades durante todo o período (março/dezembro de 2015), desempenhando funções de forma satisfatória, momentos de leitura e discussão dos temas propostos, de construção e organização dos questionários e visitas às escolas.
	Situação: Concluído; Natureza: Pesquisa.
	Integrantes: Laila Gabrielle Naysmaer Gonçalves - Integrante / Rochelle da Silva Santana - Coordenador.
Projetos de extensão	
2014 - 2015	Pedagogias da Igualdade
	Projeto certificado pelo(a) coordenador(a) Martha Gudice Narvaz em 23/06/2019.
	Situação: Concluído; Natureza: Extensão.
	Integrantes: Laila Gabrielle Naysmaer Gonçalves - Integrante / Martha Gudice Narvaz - Coordenador.
2013 - 2014	A Boniteza de um Sorriso no Alegrete: Urges e Comunidade no Enfrentamento da Violência Contra as Mulheres e as Meninas
	Projeto certificado pelo(a) coordenador(a) Martha Gudice Narvaz em 21/07/2017.
	Descrição: Desenvolvimento de atividades no período de 15 de abril de 2013 a 15 de agosto de 2014, perfazendo um total de 768 horas.
	Situação: Concluído; Natureza: Extensão.
	Integrantes: Laila Gabrielle Naysmaer Gonçalves - Integrante / Martha Gudice Narvaz - Coordenador.

Figura 5. Exemplo das informações extraídas.

B. Nomes Repetidos

As pessoas cadastradas na plataforma Lattes podem possuir nomes parecidos ou até mesmo iguais, com isso existe a possibilidade de acesso de um perfil que não pertence ao professor correto, ver Figura 6. Esta situação pode se apresentar de duas excessões diferentes, sendo elas: 1) existe pessoas com partes do nome iguais, contudo o nome completo pode distinguir quem é o correto e 2) quando o nome completo é igual a de outra pessoa ou pessoas.

A primeira excessão foi solucionada realizando a pesquisa do elemento que possuía o nome completo inserido no campo de pesquisa, onde o currículo somente era aberto se o nome estivesse exatamente igual. Já a segunda excessão não possui uma solução, devido as múltiplas possibilidades, assim se faz necessário a utilização de outra informação do que o nome, poderia se utilizar o identificador único do perfil do professor.

A solução para a segunda excessão necessita de um pré-processamento no arquivo de informações dos professores, onde seria necessário os identificadores de todos os perfis de cada professor na plataforma Lattes, assim caso o professor não fosse encontrado pelo seu nome, seria procurado pelo seu identificador único.

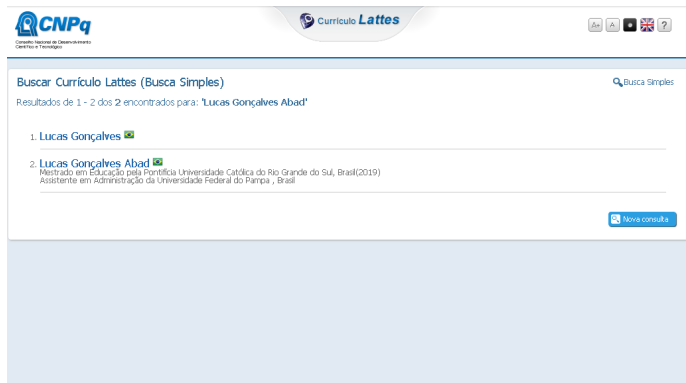


Figura 6. Exemplo de perfis com nomes iguais.

C. Nome Incorreto

As informações dos professores são carregadas do arquivo que foi processado previamente, com isso é necessário que todas as informações que estão neste arquivo estejam corretas. Todavia existe a possibilidade de que o nome armazenado no arquivo não esteja correto ou não seja o nome cadastrado na plataforma Lattes. Sendo assim é necessário o tratamento do arquivo previamente. Caso o nome esteja correto mas o seu perfil não seja encontrado, foi necessário deixar de pesquisar o perfil do professor em questão, já que não há uma forma de comprovar se o nome está correto ou se o nome do arquivo e o nome do perfil na plataforma são os mesmos, um exemplo deste problema pode ser visto na Figura 7.

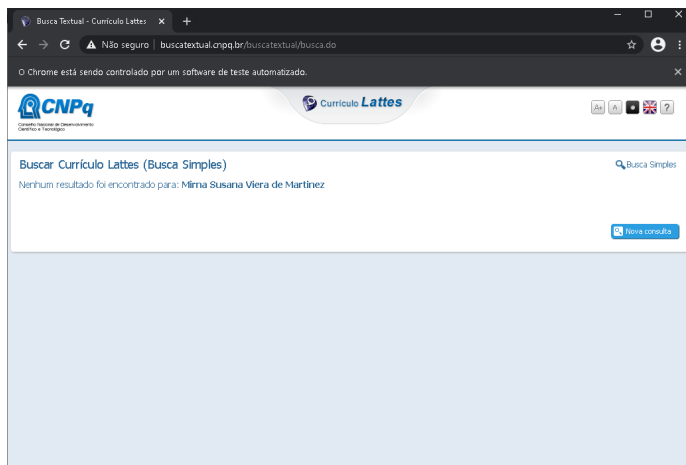


Figura 7. Exemplo de nome incorreto.

D. Pesquisa não carregada

As pesquisas realizadas em sua grande maioria sempre retornava uma página com os perfis encontrados ou uma página sem perfis. Contudo durante os testes ocorreram erros no carregamento dos perfis, onde a busca de um professor retornava positiva porém o seu perfil não estava acessível para a aplicação, um exemplo deste erro pode ser notado na Figura 8.

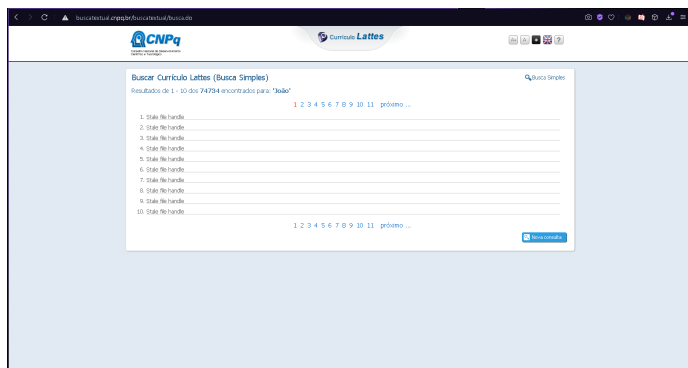


Figura 8. Exemplo de nome incorreto.

O erro dos perfis não serem disponibilizados não foram resolvidos devido a possibilidade de um mal carregamento das informações da plataforma, sendo necessário um melhor estudo na ocorrência deste erro, para assim estipular uma forma de contornar este erro. Uma possibilidade seria a de caso este erro ocorra, realizar a pesquisa do professor novamente.

V. CONCLUSÃO

As informações extraídas da plataforma Lattes devem ser tratadas, pois podem possuir espaços e informações desnecessárias devido a construção da página web.

A aplicação desenvolvida pode ser utilizada para obter as informações de um conjunto de professores, mesmo que seja necessário somente a informação de um professor, se estiver dentro de um arquivo seguindo o padrão apresentado neste trabalho, não haverá problemas.

A aplicação desenvolvida não pode ser utilizada para a extração de qualquer informação do perfil de um professor, já que ele foi pensado até o presente momento para somente extrair as informações de projetos de pesquisa e extensão de forma automática. Para a extração de novas informações é necessária a refatoração do código fonte da aplicação para que seja extraída a informação desejada.

APÊNDICE

Os códigos da aplicação podem ser encontrados na plataforma de versionamento de arquivos Github no link <https://github.com/DarkDym/GDL>.

REFERÊNCIAS

- [1] Selenium, <https://www.selenium.dev/documentation/en/>.
- [2] Python, <https://www.python.org>.
- [3] HTML, <https://pt.wikipedia.org/wiki/HTML>.
- [4] Plataforma Lattes, <http://buscatextual.cnpq.br/buscatextual/busca.do>.