# Web Scraper do Website JusBrasil pela biblioteca Puppeteer

Josué Teodoro Moreira

# Índice de conteúdo

Web Scraping e processamento de dados	3
Dos dados	3
Páginas de processos	3
Processos	4
Desenvolvimento	
Engenharia de Software	
Da State Machine	
Das Classes	
Spider	
ProcPagesSpider	5
ProcsSpider	
Middleware	
ProcPagesCleanerMiddleware	
ProcParserMiddleware	
Scraper	
Da execução	
Considerações finais	
C011014C14C0C0 1111410111111111111111111	

# Introdução

Assim que recebi o desafio, minha primeira ideia foi de utilizar uma API disponibilizada pela própria JusBrasil. Tal API existe, de fato, e se chama Digesto. Infelizmente, a Digesto é paga e requere diversos processos burocráticos que interferem com as regras definidas pelo desafio.

Dessa maneira, troquei minha atenção em seguida para o Web Scraping por meio da biblioteca Puppeteer. Na verdade, não pensei primeiramente na biblioteca Puppeteer, mas sim na Crawly (Elixir), depois na Selenium (Python), e por fim decidi continuar com a Puppeteer (Javascript/Typescript), já que durante minha entrevista foi citado que a LISS utilizava da tecnologia/linguagem Javascript.

# Web Scraping e processamento de dados

O web scraping é a prática de extrair informações por meio da leitura e processamento do conteúdo do híper texto (HTML) de um website. Nesse caso em específico, se trata da automatização do processo de entrar no website da JusBrasil, e utilizá-lo como um ser humano normal, porém controlado por um "robô", enquanto extrai as informações presentes.

# Dos dados

Existem diversas maneiras de processar o conteúdo de um website e filtrar somente o necessário. No caso específico do desafio, percebi que o site da JusBrasil tem o conteúdo de estilização (CSS) muito bem estruturado, então seria a melhor maneira de separar os dados presentes:

# Páginas de processos

Antes de extrair os dados dos processos em si, percebi que diversos processos dão nomes diferentes para a mesma empresa. Na prática, isso cria vários resultados de pesquisa diferentes para uma única empresa.

Ante isso, julguei necessário extrair e agregar todas as páginas que levavam a processos dispersos:

- Paginação: `.Pagination`
- Botão de próximo:
  - Desativado: `li.pagination-item > a[aria-label="Próximo"]`
  - Ativado: `li.pagination-item > a[aria-label="Próximo"]`
- Lista de resultados: `.EntitySnippet-item`
- Url da página de processos únicos: `.EntitySnippet-header > .EntitySnippet-header > .EntitySnippet-anchor-wrapper > a`

#### Processos

Como dito anteriormente, após extrair e agregar os links das páginas de processos dispersos, posso finalmente entrar nesses links e agregar os processos.

- Quantidade de processos: `.navbar-link.active > h2.navinline > span`
- Lista de processos: `a.LawsuitCardPersonPage`
- Link do processo: `a.LawsuitCardPersonPage-title—link`
- Número do processo: `span.LawsuitCardPersonPage-header-processNumber`
- Corte do processo: `p.LawsuitCardPersonPage-body-row-item-text[role="body-court"]`
- Procedimento/Classe do processo: `p.LawsuitCardPersonPage-body-row-item-text[role="body-kind"]`
- Partes do processo: `strong.LawsuitCardPersonPage-header-processInvolved`

# **Desenvolvimento**

Assim que separei os "selectors" que iria utilizar na extração dos dados, passei à programação em si. Criei o projeto NodeJS pel Yarn com versionamento Git e transpilação para Typescript, e fiz upload no Github: <a href="https://github.com/J0sueTM/jusbrasil\_scraper">https://github.com/J0sueTM/jusbrasil\_scraper</a>. A documentação do código a seguir possui hyperlinks que direcionam ao github, inclusive à linha exata que estou citando.

# Engenharia de Software

Como estilo de arquitetura, apesar de eu prefirir pessoalmente modelos hexagonais e/ou programação funcional, julguei propício para esse software um modelo orientado a objeto, por meio de abstrações, herança, encapsulamento, etc, com armazenamento de dados e estado por meio de uma state machine (máquina de estado).

#### Da State Machine

A máquina de estado é um conceito inicialmente matemático, mas bastante presente na engenharia de software. Uma máquina de estado é um agregado de informações que tem um valor dependente do momento; ou seja, no início da execução da aplicação, o estado pode ser 2, mas três operações a seguir, o estado pode ser 244.

Na minha aplicação, o estado se refere aos dados após cada execução de uma Spider (explicada a seguir). Está definida <u>neste arquivo</u>, baseados nas interfaces (DTOs) definidas <u>neste arquivo</u>.

Superficialmente, a state machine nesta aplicação é uma interface com os dados globalmente acessíveis, guardando informações geradas após os processamentos adjuntos das Spiders e dos Middlewares.

#### Das Classes

Classes, no âmbito da programação orientada a objeto, é uma representação abstrata de um conceito real, que possui propriedades reais, as quais são igualmente abstraídas na programação. A seguir segue as que foram inseridas por mim no Software.

# Spider

Na terminologia do Web Scraping, uma spider é a implementação de um "robô" que realiza ambas as operações de simulação de um ser humano, e a extração de dados. Baseado nisso, criei a classe base Spider.

A classe Spider é uma classe abstrata, ou seja, não existe com o intuito de realizar operações, mas para ser base na construção de classes mais complexas.

- Atributos
  - baseUrl: O link base. Ex: www.jusbrasil.com.br/busca-processual.
  - o urls: Uma lista com os links a serem processados.
  - ∘ name: Um nome de identificação.
  - o cacheFilePath: O arquivo onde os dados extraídos serão salvos.
  - o middlewares: Uma lista com os middlewares que processarão os dados.
- Métodos
  - o crawl(): Executa a extração dos dados.
  - o saveCompany(): Salva os dados apenas da última empresa que extraiu.
  - o save(): Salva todos os dados de todas as empresas extraídas.
  - <u>load()</u>: Caso a operação já tenha sido executada antes, e um cachê foi encontrado, carrega os dados já salvos na memória, salvando tempo e processamento.
  - shouldLoad(): Verifica se existe um cachê da Spider em questão.
  - run(): Executa os processos definidos pela Spider já implementada.

# ProcPagesSpider

Essa é a primeira Spider executada. <u>ProcPagesSpider</u> é responsável por extrair os links das páginas dispersas que referem à mesma empresa. O motivo, como explicado anteriormente, é que devido aos diferentes nomes utilizados nos múltiplos processos, acabaram espalhados e separados em páginas dispersas.

A classe ProcPagesSpider herda as características da classe abstrata Spider, porém implementa os métodos especificamente:

#### Implementação do método crawl()

```
for await (const [idx, url] of this.urls.entries()) {
   await appState.page.goto(url);

let curCompany: string = companies[idx];
```

Iteração de todos os urls de pesquisas para cada empresa.

Lê todos os urls de páginas de processos encontrados e os agrega na array curPages.

```
let disNextButton = await appState.page.$$eval(
    disNextBtnSlc,
    (elms) ⇒ elms
);
let pagination = await appState.page.$$eval(
    paginationSlc,
    (elms) ⇒ elms
);
if (disNextButton.length ≥ 1 || pagination.length ≤ 0) {
    finished = true;
    break;
}
```

Verifica se existe um botão de próximo. Se existe, continua o loop. Caso não, significa que terminamos de ler os urls de páginas de processos e podemos passar para a próxima empresa.

```
await appState.page.waitForSelector(actNextBtnSlc, { visible: true });
await appState.page.click(actNextBtnSlc);
```

Caso o botão de Próximo exista, ele é clicado, levando-nos para a próxima página.

```
appState.procPages.set(curCompany, curPages);
this.middlewares[0].run();
this.saveCompany(curCompany, curPages);
}
```

Assim que chega no fim da paginação, roda o middleware ProcPagesCleaner, o qual será explicado posteriormente, e salva os dados em cachê.

#### Implementação dos métodos save() e load()

```
protected save() {
    // formatting would be nicer with:
    // const rawData = JSON.stringify(Object.fromEntries(this.dispProcUrls))
    // but I couldn't reverse it.

    const rawData = JSON.stringify(Array.from(appState.procPages.entries()));
    fs.writeFileSync(this.cacheFilePath, rawData);
}

protected load() {
    const rawData = fs.readFileSync(this.cacheFilePath, "utf8");
    appState.procPages = new Map(JSON.parse(rawData));
}
```

Os métodos em si são simples, e apenas parseiam os dados para JSON e salva em cachê, e vice-versa.

Após finalizado, o dado final é salvo desta maneira:

```
didas-do-brasil-ltda > 🕒 urls-de-paginas-de-processos.json > ...
 https://www.jusbrasil.com.br/processos/nome/608681872/adidas-do-brasil-ltda
 "https://www.jusbrasil.com.br/processos/nome/608681872/adidas-do-brasil-ltda",
"https://www.jusbrasil.com.br/processos/nome/499994541/adidas-do-brasil-ltda-solucao-extinto-processo-por-desistencia",
 "https://www.jusbrasil.com.br/processos/nome/134027405/adidas-do-brasil-ltda-1",
 https://www.jusbrasil.com.br/processos/nome/101711178/adidas-do-brasil-ltdabrasil-ltda",
 "https://www.jusbrasil.com.br/processos/nome/298950069/adidas-do-brasil-ltda-adidas",
 "https://www.jusbrasil.com.br/processos/nome/220524653/adidas-adidas-do-brasil-ltda",
  https://www.jusbrasil.com.br/processos/nome/94959649/adidas-do-brasil-ltda-filial'
 "https://www.jusbrasil.com.br/processos/nome/298936907/requerente-adidas-do-brasil-ltda",
 "https://www.jusbrasil.com.br/processos/nome/285912547/adidas-do-brasil-ltda-comercial-vulcabras-ltda",
 "https://www.jusbrasil.com.br/processos/nome/379337940/adidas-do-brasil-ltda-solucao-julgado-procedente-em-parte-o-pedido",
 "https://www.jusbrasil.com.br/processos/nome/186597337/adidas-com-adidas-do-brasil-ltda",
"https://www.jusbrasil.com.br/processos/nome/84436664/adidas-do-brasil-comercio-de-artigos-de-esporte-ltda",
 "https://www.jusbrasil.com.br/processos/nome/183269869/2013-adidas-do-brasil-ltda",
"https://www.jusbrasil.com.br/processos/nome/592932675/adidas-do-brasil-ltda-primeiro-a-s",
 "https://www.jusbrasil.com.br/processos/nome/319081550/adidas-do-brasil-ltda-loja",
"https://www.jusbrasil.com.br/processos/nome/383842440/adidas-do-brasil-ltda-autoridade",
 https://www.jusbrasil.com.br/processos/nome/360643433/adidas-do-brasil-ltda-3","
  https://www.jusbrasil.com.br/processos/nome/616464485/adidas-do-brasil-ltda-citado","
 "https://www.jusbrasil.com.br/processos/nome/586642628/adidas-d-brasil-ltda"
 "https://www.jusbrasil.com.br/processos/nome/282000011/representante-legal-da-adidas-do-brasil-ltda", "https://www.jusbrasil.com.br/processos/nome/302556013/adidas-do-brasil-ltda-e-filial-is",
 "https://www.jusbrasil.com.br/processos/nome/534706423/adidas-do-brasil-ltda-reu-re",
 https://www.jusbrasil.com.br/processos/nome/88529648/adidas-dop-brasil-ltda",
 https://www.jusbrasil.com.br/processos/nome/137314640/adidas-di-brasil-ltda"
 "https://www.jusbrasil.com.br/processos/nome/430817924/adidas-do-brasil-ltda-loja-online-reebok",
 https://www.jusbrasil.com.br/processos/nome/159332162/adidas-do-brasil-ltda-na-pessoa-de-seu-representante-legal",
 https://www.jusbrasil.com.br/processos/nome/299053205/adidas-do-brasil-ltda-atacadistas-e-fabricantes-de-calcados"
 https://www.jusbrasil.com.br/processos/nome/360071418/adidas-do-brasil-ltda-a-c-gustavo-stussi-neves",
  https://www.jusbrasil.com.br/processos/nome/130513025/adidas-do-brasil-com-de-artigos-de-esporte-ltda",
 "https://www.jusbrasil.com.br/processos/nome/594948313/adidas-do-brasil-ltda-solucao-julgado-improcedente-o-pedido",
 "https://www.jusbrasil.com.br/processos/nome/ɔy4y4xɔıɔ/auıuas up brasil-comercio-de-artigos-de-esportes-ltda",
"https://www.jusbrasil.com.br/processos/nome/142225597/adidas-do-brasil-ltda-soc-advoqados-ricardo-marfori-sampaio",
  https://www.jusbrasil.com.br/processos/nome/388250982/adidas-do-brasil-ltda-soc-advogados-ricardo-marfori-sa"
 "https://www.jusbrasil.com.br/processos/nome/155803318/adidas-do-brasil-comercio-de-artigos-esportivos-ltda",
 https://www.jusbrasil.com.br/processos/nome/167861546/adidas-franchise-brasil-servicos-ltda","
 https://www.jusbrasil.com.br/processos/nome/357146594/adidas-do-brasil-com-de-artigos-de-esportes-ltda",
 "https://www.jusbrasil.com.br/processos/nome/345849658/adidas-do-brasil-com-de-arts-de-esportes-ltda",
 "https://www.jusbrasil.com.br/processos/nome/595685772/adidas-do-brasil-comercio-de-artigo-esporte-ltda",
"https://www.jusbrasil.com.br/processos/nome/481136710/adidas-brasil-ltda-e-outro-s",
 "https://www.jusbrasil.com.br/processos/nome/155160582/adidas-brasil-ltda-e-outro-s",
"https://www.jusbrasil.com.br/processos/nome/155160582/adidas-brasil-ltda-wwwadidascombr",
  https://www.jusbrasil.com.br/processos/nome/166977664/adyen-do-brasil-ltda-adidas
```

# ProcsSpider

Esta é a segunda Spider a ser executada. <u>ProcsSpider</u> é responsável por utilizar dos links processados pela ProcPagesSpider, e finalmente extrair os dados dos processos em si.

#### Implementação do método crawl()

```
for await (const [company, page] of appState.procPages.entries()) {
  console.info(`Extraindo processos de ${company}...`);

  const progressBar = new cliProgress.MultiBar(
    {
      clearOnComplete: false,
      hideCursor: true,
      format: "{bar} >> {pagename} << {value}/{total}",
    },
    cliProgress.Presets.shades_classic
);</pre>
```

Inicia-se iterando as páginas de processos extraídas pela ProcPagesSpider. Desta vez, como algumas das extrações demoravam muito tempo, adicionei um barra de progresso, que indica quantas processos e páginas de processos foram extraídos, quantos faltam, e de qual empresa os processos estão sendo lidos.

```
await appState.page.goto(url);

let isPageUnavailable = await appState.page.evaluate(() ⇒ {
   const elms = Array.from(document.querySelectorAll("*"));
   return elms.some((elm) ⇒
      elm.textContent?.includes("Página não encontrada")
   );
});

if (isPageUnavailable) {
   continue;
}
```

Assim que o robô entra no website, é verificado se não ocorreu nenhum erro. Caso sim, a operação na página atual é encerrada, e a próxima é carregada.

```
const procCountSlc: string =
   ".navbar-link.active > h2.navinline > span";
let totalPageProcCount: number = await appState.page.$eval(
   procCountSlc,
   (elm) ⇒ {
    let rawNum: string = elm.textContent?.replace(/\D/g, '');

   const defaultThreshold: number = 1000;
   return Math.min(+rawNum, defaultThreshold);
}
);
const pageBar = progressBar.create(totalPageProcCount, 0);
```

A seguir, o indicador da quantidade de processos naquela página é parseado e lido.

```
const procUrl: string =
  elm.querySelector(procUrlSlc)?.getAttribute("href") ??
  "Indisponível";
const procNum: string =
  elm
    .querySelector(procNumSlc)
    ?.textContent.replace("Processo n° ", "") ?? "Indisponível";
const procCourtAndLocality: string[] = elm
  .querySelector(procCourtAndLocalitySlc)
  ?.textContent?.split(" · ") ?? [
  "Indisponível",
 "Indisponível",
1:
const procProcedure: string =
  elm.querySelector(procProcedureSlc)?.textContent ??
  "Indisponível";
const procParties: string[] = elm
  .querySelector(procPartiesSlc)
  ?.textContent?.split("x")
const procCourt: string = procCourtAndLocality[0];
var procLocality: string = procCourtAndLocality[1];
var procUF: string = procCourt?.slice(-2) ?? "Indisponível";
if (procCourtAndLocality[1] ≠ "Indisponivel") {
  procLocality = procCourtAndLocality[1] ?? "Indisponível";
 procUF = procLocality?.split(",")[1]?.slice(1) ?? "Indisponível";
// clean parties names' blank spaces
procParties[0] =
  procParties[0] == undefined
    ? "Indisponível"
    : procParties[0]?.slice(0, -1);
procParties[1] =
  procParties[1] == undefined
   ? "Indisponível"
    : procParties[1]?.slice(1);
```

Após tudo isso, os dados são finalmente extraídos.

```
curPageProcCount = curPageProcs.length;
   pageBar.update(curPageProcCount, { pagename: `Página ${idx + 1}` });
}
this.savePage(company, idx + 1, curPageProcs);
procs = procs.concat(curPageProcs);
progressBar.remove(pageBar);
compBar.update(idx + 1, { pagename: company.split(" ")[0] });
}
progressBar.stop();
console.info(`Total extraído: ${procs.length}\n`);
appState.procs.set(company, procs);
this.saveCompany(company, procs);
```

Em seguida, os processos lidos são adicionados à lista de processos da empresa atual, processados pelo middleware ProcParserMiddleware (que será explicado posteriormente), e por fim salvos em um arquivo JSON, como na imagem a seguir:

```
"url": "https://www.jusbrasil.com.br/processos/372539554/processo-n-000XXXX-4520208260048-do-tjsp",
"procNumber": "000XXXX-45.2020.8.26.0048", "court": "TJSP",
"locality": "Foro de Atibaia, SP",
"uf": "SP",
"procedure": "Cível · Obrigações",
  "Rafael Giuliano Santos Moura da Silveira",
"Under Armour Brasil Comércio e Distribuição de Artigos Esportivos LTDA"
"url": "https://www.jusbrasil.com.br/processos/333092448/processo-n-100XXXX-0620208260048-do-tjsp",
"procNumber": "100XXXX-06.2020.8.26.0048",
"court": "TJSP",
"locality": "Foro de Atibaia, SP",
"uf": "SP",
"procedure": "Cível · Procedimento do Juizado Especial Cível",
"parties":
  "Rafael Giuliano Santos Moura da Silveira",
"Under Armour Brasil Comércio e Distribuição de Artigos Esportivos LTDA"
"url": "https://www.jusbrasil.com.br/processos/423207093/processo-n-000XXXX-9420208190036-do-tjrj",
"procNumber": "000XXXX-94.2020.8.19.0036", "court": "TJRJ",
"uf": "RĴ",
"procedure": "Dano Moral - Outros/ Indenização por Dano Moral",
  "Romulo Nascimento dos Santos Lopes",
"Under Armour Brasil Comércio e Distribuição de Artigos Esportivos LTDA"
"url": "https://www.jusbrasil.com.br/processos/429051938/processo-n-002XXXX-8320208190001-do-tjrj",
"procNumber": "002XXXX-83.2020.8.19.0001",
```

### Middleware

Um middleware é um processo menor adicionado à uma sequência já existente de processos maiores; No caso específico desse Software, os processos maiores são as Spiders, que executam a extração dos dados, enquanto os processos menores, os middlewares, executam a limpeza, parseamento dos dados, etc.

Assim como a classe Spider, a classe <u>Middleware</u> é uma classe abstrata que não possui nenhum atributo, mas apenas um método: <u>run()</u>, que executa o processo do middleware em si.

# ProcPagesCleanerMiddleware

A classe ProgPagesCleanerMiddleware herda as propriedados da classe Middleware, e implementa o método run(), realizando duas operações:

Remove quaisquer urls que não sejam estritamente das empresas em questão. Isso acontece devido a conflitos de nomes de empresas, como é o caso entre a Puma do Brasil (calçados) e a Puma do Brasil (petróleo).

```
let uniqueNewCompProcUrls = [...new Set(newCompProcUrls)];
appState.procPages.set(company, uniqueNewCompProcUrls);
```

Por fim, remove duplicatas e troca a lista antiga pela nova lista de urls de páginas de processos.

#### **ProcParserMiddleware**

A classe <u>ProcParserMiddleware</u> também herda as propriedades da classe Middleware, e implementa o método run() parseando os processos adquiridos de cada empresa, e salvando-os em suas pastas finais.

```
public run() {
    try {
      const parser = new Parser(this.opts);

    for (const [company, procs] of appState.procs) {
      let newProcs = procs.map((proc) ⇒ {
         if (proc.procNumber.includes("Processo") || proc.procNumber.includes("Indisponivel")) {
            proc.procNumber = proc.url?.split("processo-n-")[1]?.slice(0, 21)?.replaceAll('-', '.');
      }

      return proc;
    })
    appState.procs.set(company, newProcs)
    const parsedCsv = parser.parse(newProcs);
    fs.writeFileSync(`cache/${company.replaceAll(" ", "-")}/processos.csv`, parsedCsv);
    }
} catch (err) {
      console.error(`Não foi possível gerar o CSV :: ${err}`);
}
```

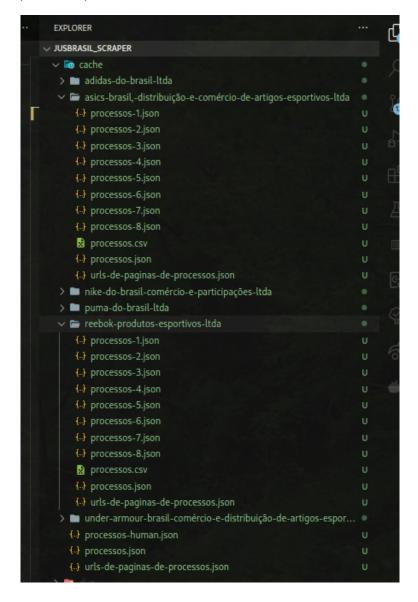
## Scraper

A class Scraper é a que controla todas as operações, incluindo as Spiders, os Middlewares e o driver do Puppeteer.

- Atributos
  - baseUrl: 0 url base.
  - o spiders: Uma lista com instâncias das spiders que serão executadas.
- Métodos
  - o createCache(): Cria as pastas necessárias para salvar os dados extraídos.
  - openBrowser(): Abre o navegador para realizar o Web Scraping.
  - o setSpiders(): Setta as spiders a serem executadas.
  - o scrape(): Inicia o processo de scraping.

# Da execução

Com os blocos básicos do Software detalhados, a inicialização e execução é realizada <u>neste arquivo</u>. Após executado, os arquivos são colocados em cachê, e prontos para análise.



# Considerações finais

Apesar de tudo, o desenvolvimento do Software foi um processo divertido e de altíssima aprendizagem. Eu nunca havia utilizado o puppeteer, e consegui produzir uma aplicação aparentemente viável, e que cumpriu seus objetivos de maneira muito melhor do que eu pessoalmente esperava.

Acerca de críticas, existem dois pontos que eu acredito que, se tivesse tempo, melhoraria e adicionaria ao Software: Threading e maior organização de arquitetura de software.

De qualquer forma, foi um prazer participar do desafio de estágio da LISS e espero que a minha solução seja de agrado para a vaga.