

UNIVERSIDADE DO MINHO

MESTRADO EM ENGENHARIA INFORMÁTICA

SCRIPTING NO PROCESSAMENTO DE LINGUAGEM NATURAL

---

# Website Processing

---

***Autores:***

Frederico Pinto  
Rui Vieira

A73639  
A74658

1 de Julho de 2019

# Conteúdo

<b>1</b>	<b>Introdução</b>	<b>2</b>
<b>2</b>	<b>Ficheiros TMX</b>	<b>3</b>
2.1	Estrutura . . . . .	3
<b>3</b>	<b>Objetivos</b>	<b>4</b>
<b>4</b>	<b>Implementação</b>	<b>5</b>
4.1	<i>linguee.pt</i> . . . . .	5
4.2	<i>lexico.com</i> . . . . .	9
<b>5</b>	<b>Apresentação do Resultado Final</b>	<b>11</b>
<b>6</b>	<b>Conclusão</b>	<b>13</b>
<b>7</b>	<b>Referências</b>	<b>14</b>

# 1 Introdução

Com o objetivo de melhorar os nossos conhecimentos sobre os temas estudados na unidade curricular de *Scripting no Processamento de Linguagem Natural* foi nos proposto uma série de enunciados, optamos então por escolher o enunciado 6, que se baseia na criação de ficheiros *TMX* a partir de técnicas de *web scraping* utilizando o *Beautiful Soup*.

Neste relatório, iremos numa fase inicial falar sobre a estrutura dos ficheiros *TMX* e posteriormente aprofundar os objetivos propostos. Por fim vamos falar sobre a nossa proposta de implementação para cumprir com sucesso esses objetivos.

## 2 Ficheiros TMX

*Translation Memory*, *TM*, é uma tecnologia de linguagem que permite a tradução de segmentos (frases, parágrafos ou frases) de documentos, pesquisando segmentos semelhantes numa base de dados e sugerindo correspondências de idiomas diferentes que se encontram presentes nessa base de dados.

É nesse contexto que surgem os ficheiros *TMX*, (*Translation Memory eXchange*), que são um padrão *XML* aberto e neutro de fornecedor. Têm como objetivo facilitar a troca de dados sobre *translation memory* entre ferramentas ou fornecedores de tradução apresentando pouca ou até nenhuma perda de dados importantes durante esse processo, o que é uma grande vantagem.

### 2.1 Estrutura

Como foi falado acima, o *TMX* apresenta uma sintaxe parecida ao *XML*, apresentando elementos e atributos únicos, contendo também várias regras para construir um ficheiro válido. Sendo uma quantidade grande de elementos, atributos e regras, colocamos nas referências uns *links* para se poder verificar todos os elementos, atributos e regras que o *TMX* contém.

Contudo, num sentido bastante abstrato, um ficheiro *TMX* possui um elemento base *<tmx>* que possui dois elementos, um *<header>* que contém atributos sobre o ficheiro, e um elemento *<body>* que contém a informação relevante. Esse elemento *<body>* contém uma lista de elementos *<tu>* que são *translation units*, possuindo cada um, elementos *<props>* e *<note>*, que representam informação sobre essa *translation unit*. Para além disso cada *<tu>* possui elementos *<tuv>* que são *translation unit variant*, este elemento possui um filho que é o elemento *<seg>* que contém o segmento de texto que o *<tuv>* possui.

Apresentamos de seguida um exemplo de um ficheiro *TMX*, contendo os elementos que falamos acima.

```

1 <?xml version="1.0" encoding="UTF-8" ?>
2 <tmx version="1.4">
3   <header adminlang="en"
4     creationtool="Heartsome Dictionary Editor"
5     creationtoolversion="1.0"
6     datatype="tbx"
7     o-tmf="unknown"
8     segtype="block"
9     srclang="en"/>
10  <body>
11    <tu origin="tbx" tuid="1108600011738">
12      <tuv xml:lang="en">
13        <prop type="administrativeStatus">admittedTerm-admn-sts </
14        prop>
15        <prop type="termType">entryTerm</prop>
16        <prop type="usageNote">Colloquial use term</prop>
17        <note>Informal salutation</note>
18        <seg>Hello</seg>
19      </tuv>
20      <tuv xml:lang="es">
21        <prop type="administrativeStatus">admittedTerm-admn-sts </
22        prop>
23        <prop type="termType">entryTerm</prop>
24        <prop type="usageNote">Termino de uso coloquial</prop>
25        <note>Saludo informal</note>
26        <seg>Hola</seg>
27      </tuv>
28    </tu>
  </body>
</tmx>

```

Listing 1: Exemplo de um ficheiro *TMX*.

### 3 Objetivos

Após análise da estrutura do ficheiro *TMX*, podemos definir os objetivos do nosso trabalho, ou seja, que informação podemos encontrar em *websites* e colocar no ficheiro *TMX* criado.

Decidimos então, produzir um *script* capaz de traduzir uma ou mais palavras para idiomas escolhidos pelo utilizador, que para além disso tem que ser capaz de criar uma *translation unit* por cada significado semântico que essa palavra possa ter, incluindo a sua significação e palavras sinónimas.

## 4 Implementação

Como proposto no enunciado para resolver este problema temos que fazer *web scraping* em *Python* utilizando o *Beautiful Soup*. O *Beautiful Soup* é uma biblioteca para o *Python* que retira informação de ficheiros *HTML* e *XML*, utilizando um *parser* que permite obter maneiras de navegar, procurar e modificar a árvore analisada. Sem esta ferramenta, um *developer* perderia horas analisando toda a informação "à mão".

De acordo com os objetivos propostos na secção anterior e após realizarmos uma pesquisa, concluímos que necessitamos de fazer *scraping* em dois *websites*. Iremos utilizar o *Linguee* para extrair as traduções em todas os idiomas que este fornece e o *Lexico* para verificar a significação das palavras e os seus sinónimos.

De seguida, iremos falar sobre qual foi a nossa estratégia de implementação para cada um deles. Contudo, é de salientar que as palavras que o nosso *script* recebe têm que ser em inglês, pois é a língua mais suportada do *Linguee*, sendo possível traduzir para os seguintes idiomas.

- |             |               |            |
|-------------|---------------|------------|
| • Alemão    | • Holandês    | • Húngaro  |
| • Francês   | • Polaco      | • Eslovaco |
| • Espanhol  | • Sueco       | • Búlgaro  |
| • Chinês    | • Dinamarquês | • Esloveno |
| • Russo     | • Finlandês   | • Lituano  |
| • Japonês   | • Grego       | • Letão    |
| • Português | • Checo       | • Estónio  |
| • Italiano  | • Romeno      | • Maltês   |

A procura no *Lexico* é efetuada com a palavra em inglês, pois apresenta um excelente suporte.

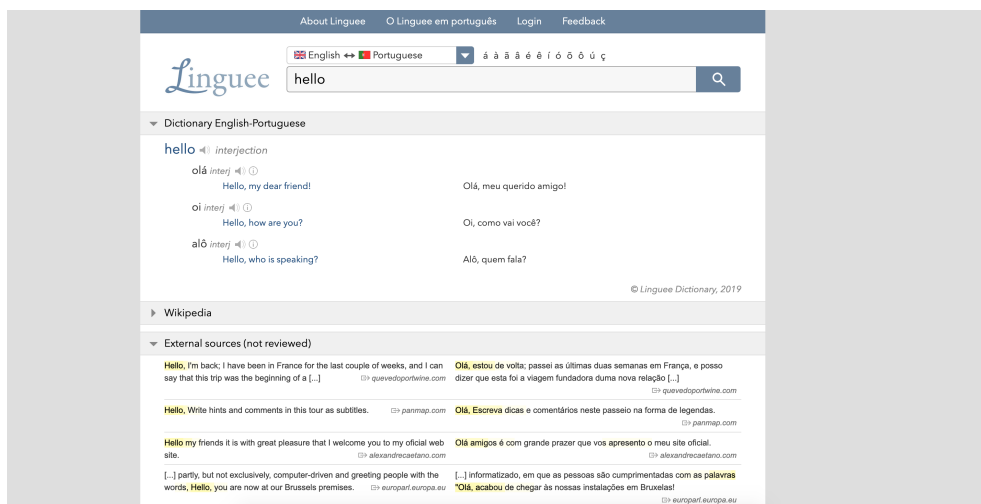
### 4.1 *linguee.pt*

Em primeiro lugar é importante analisar o *URL* que faz o pedido ao *endpoint* pela informação e quais são os campos a alterar de maneira a efetuarmos nós os pedidos com a informação que desejamos.

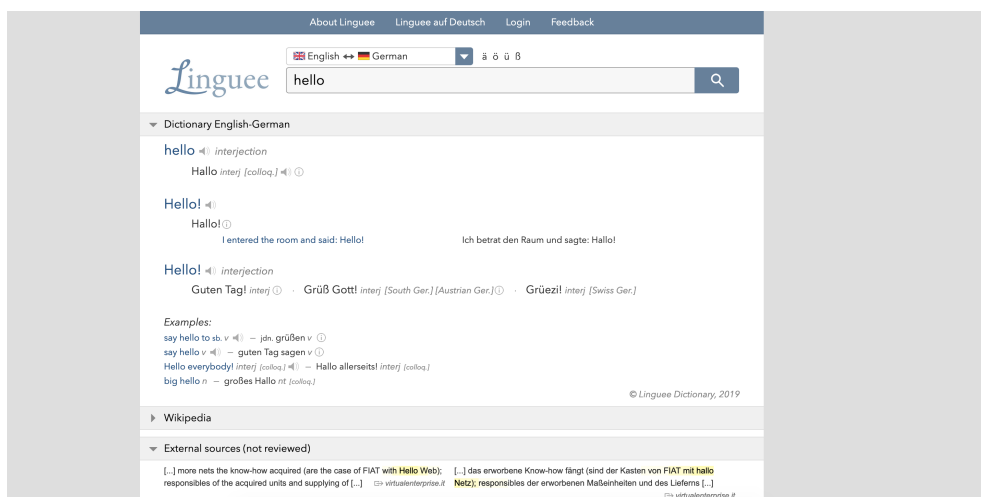
1 <https://www.linguee.com/english-portuguese/search?query=hello>

Listing 2: Exemplo de URL de pedidos ao Linguee.

Esse pedido devolve a seguinte página *HTML*:



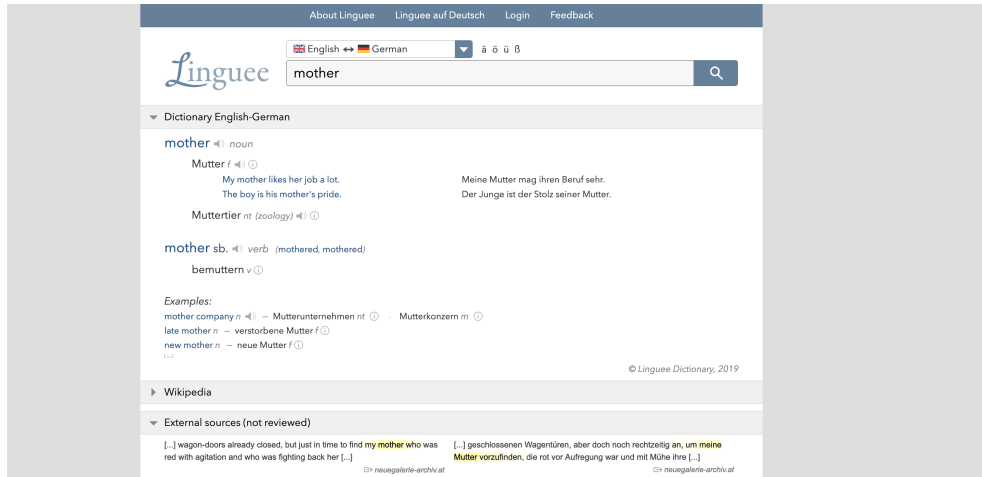
Após analisarmos o *URL* podemos concluir que existem componentes que sendo alterados, são pedidos válidos e contêm a informação necessária para obtermos aquilo que desejamos. O primeiro caso é o idioma para o qual queremos traduzir, o segundo componente do *URL* (*english-portuguese*), se substituirmos por *english-german*, obtemos a seguinte página, que apresenta a tradução de *hello* para alemão.



Para podermos variar o idioma da tradução, guardamos um *array* com todos os idiomas que o *Linguee* suporta e iteramos o mesmo, para fazer variar

o *URL* do pedido *HTTP*.

Por fim, e o mais óbvio, se alterarmos a variável *query* presente no *URL* para qualquer outra palavra em inglês, obtemos a tradução dessa palavra para o idioma indicado no *URL*. Alterando a variável *query* para *mother*, obtemos o seguinte resultado.



Após analisarmos como poderíamos navegar pelo *Linguee* a nível de pedidos *HTTP*, temos agora que analisar como chegar à informação necessária presente na página *HTML* com o *Beautiful Soup*. De seguida mostramos um excerto da página que contém informação necessária.

```
1 <html>
2 <head>
3 ...
4 </head>
5 <body>
6 ...
7 <div id="dictionary">
8 <h1>...</h1>
9 <div class="isMainTerm">
10 <div class="exact">
11 <div class="lemma">
12 <div>
13 <h2 class="line lemma_desc" lid="EN: mother25752">
14 <span class="tag_lemma">
15 <a class="dictLink" rel="nofollow" href="/english-german
/translation/mother.html">mother</a>
16 <span class="tag_wordtype">noun</span>
17 </span>
18 </h2>
19 <div class="lemma_content">
20 <div class="meaninggroup sortablemg" gid="0">
21 <div class="translation_lines">
22 <div class="translation sortablemg featured">
23 <h3 class="translation_desc">
24 <span class="tag_trans" bid="10000250407"
lid="DE: Mutter24716"><a id="dictEntry10000250407"
href="/german-english/translation/Mutter.
html" class="dictLink featured">Mutter</a>
25 </span>
26 </h3>
27 </div>
28 <div class="translation sortablemg featured">
29 <h3 class="translation_desc">
30 <span class="tag_trans" bid="10000247790"
lid="DE: Muttertier60356">
```



```

32 english/translation/Muttertier.html" <a id="dictEntry10000247790" href="/german-
33 class="dictLink featured">Muttertier </a>
34 </span>
35 </h3>
36 </div>
37 </div>
38 </div>
39 <div class="lemma">
40 (que contem outras traduções para diferentes significados
41 semanticos da palavra)
42 </div>
43 </div>
44 </div>

```

Listing 3: Página HTML fornecida pelo Linguee.

Após analisar todo o *HTML* recebido, para irmos buscar a informação importante, precisamos de ter o seguinte *workflow*.

```

1
2 # Fazer o pedido com o URL previamente construido
3 response = requests.get(url).content
4
5 # Analise da pagina html pelo BeautifulSoup
6 soup = BS(response, 'html.parser')
7
8 # Encontrar a DIV que contem o id dictionary
9 dictionary = soup.find('div', id="dictionary")
10
11 # Encontrar dentro da DIV dictionary a DIV com classe exact
12 exact = dictionary.find('div', 'exact')
13
14 # Encontrar dentro da DIV exact todas as DIVS com classe
15 lemma
16 lemma = exact.findAll('div', 'lemma')
17
18 # Por cada DIV encontrada acima
19 for lem in lemma:
20     # Encontrar o span que contem a classe tag_lemma para ir
21     buscar a tag_wordtype
22     word = lem.find('span', 'tag_lemma')
23
24     # Descobrir o span que contem o tipo semantico da palavra
25     que esta presente neste DIV lem
26     tag_wordtype = word.find('span', 'tag_wordtype')
27
28     word_type = tag_wordtype.text
29
30     # Ir buscar a primeira tradução que estiver contida na
31     DIV lem
32     translation = lem.find('span', 'tag_trans')

```

Listing 4: Implementação com o BeautifulSoup para o *Linguee*.

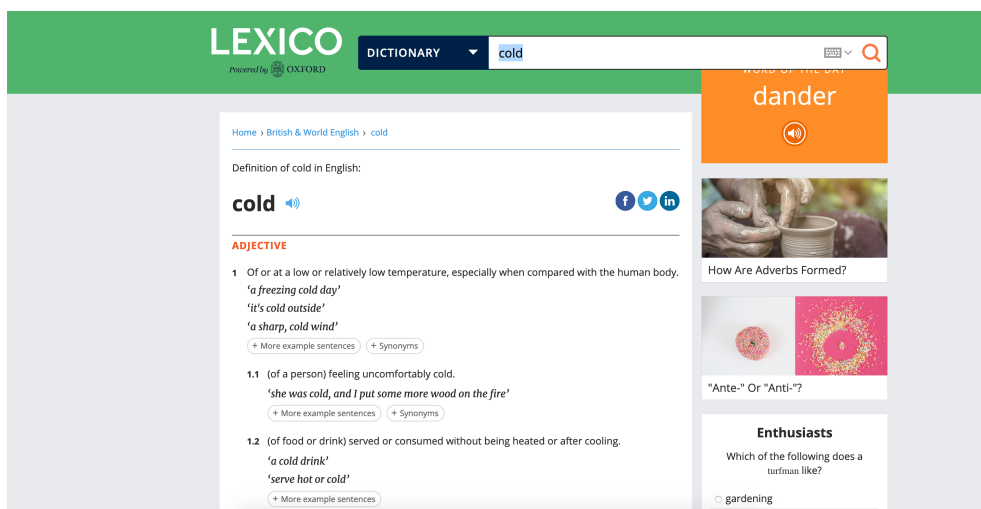
## 4.2 *lexico.com*

Tal como o *Linguee* numa primeira fase temos que analisar os *URLs* que fazem as conexões aos *endpoints* que queremos atingir. Sendo assim apresentamos de seguida, o *URL* que permite ver a definição de *hello* no *Lexico*.

```
1 https://www.lexico.com/en/definition/cold
```

Listing 5: Exemplo de URL de pedidos ao Lexico.

Esse pedido devolve a seguinte página:



Como podemos verificar, a única parte variável deste *URL* é a ultima parte, *cold*. Substituindo pela palavra desejada, obtemos a página que contém a definição dessa palavra.

Sendo assim, mostramos de seguida o *HTML* que compõe a página para sabermos como podemos obter a informação que necessitamos, que é a significação e as palavras sinónimos.

```
1 <html>
2 <head>
3   ...
4 </head>
5 <body>
6   ...
7   <section class="gramb">
8     <h3 class="ps pos"><span class="pos">adjective </span></h3><span class="
9       transitivity"></span>
10     <ul class="semb">
11       <li>
12         <div class="trg">
13           <p>
14             <span class="iteration">1</span>
15             <span class="ind">Of or at a low or relatively low temperature,
16               especially when compared with the human body.</span>
17           </p>
18           <div class="exg"> ... </div>
19           <div class="examples"> ... </div>
20           <div class="synonyms">
21             <div class="moreInfo"><button data-behaviour="ga-event-synonyms"
22               data-value="expand/collapse">Synonyms</button></div>
23             <div class="exg">
```

```

23         <div>
24             <strong class="syn">chilly </strong>
25             <span class="syn">, cool, freezing, icy, snowy, icy-
cold,
26                 glacial, wintry, crisp, frosty, frigid, bitter,
bitterly cold, biting, piercing, numbing,
27                 sharp, raw, polar, arctic, Siberian
28             </span>
29         </div>
30         <a data-behaviour="ga-event-synonyms"
31           data-value="view synonyms" href="/en/synonym/cold">
View synonyms
32     </a>
33 </div>
34 </div>
35 <ol class="subSenses">
36     ....
37     ....
38 </ol>
39 </div>
40 </li>
41 <li>
42     (Com a mesma estrutura do li acima, variando a informacao sobre o
significado e os sinonimos)
43 </li>
44 </section>
45 (contem varias seccoes com class gramb, cada uma para diferentes significados
semanticos da palavra)

```

Listing 6: Página HTML fornecida pelo Lexico.

Após efetuada uma análise ao *HTML* implementamos o seguinte *workflow* utilizando o *Beautiful Soup* para retirar a informação que necessitamos.

```

1
2 # Fazer o pedido com o URL previamente construido
response = requests.get(url).content
3
4
5 # Analise da pagina html pelo BeautifulSoup
soup = BS(response, 'html.parser')
6
7
8 # Encontrar todas as seccoes gramb
word_types = soup.findAll('section', 'gramb')
9
10
11 # Por cada seccao encontrada
for word_type in word_types:
12
13     # Encontrar o significado semantico da palavra na section
instance = word_type.find('h3', 'ps pos')
14
15
16     # Encontrar a lista na section que contem todos os
significados e sinonimos
17     ul = word_type.find('ul', 'semb')
18
19
20     # Encontrar todos os lis que essa lista contem
lis = ul.findAll('li')
21
22
23     # Por cada li encontrado
for li in lis:
24
25
26         # Encontrar a DIV com classe trg

```

```

27         div = li.find('div', 'trg')
28
29         # Se a DIV existir
30         if(div):
31             # Encontrar o paragrafo onde o significado da
palavra esta
32             p = div.find('p')
33
34             # Se existir
35             if(p):
36                 # Retirar o significado que se encontra
dentro do span
37                 span = p.find('span', 'ind')
38
39                 # Encontrar a DIV que contem os sinonimos que
esta dentro da DIV trg
40                 sysn = div.find('div', 'synonyms')
41
42                 # Se existir
43                 if(sysn):
44
45                     # Retirar o primeiro sinonimo
46                     strong = sysn.find('strong', 'syn')
47
48                     # Retirar os restantes sinonimos
49                     exg = sysn.find('span', 'syn')

```

Listing 7: Implementação com o BeautifulSoup para o *Lexico*.

É de salientar que para o mesmo significado semântico, possuímos vários significados e como consequência várias palavras sinónimas. Então para resolver isso, fazemos uma concatenação de significados, separando-os para o mesmo significado semântico pelos caracteres (*///*). O mesmo é aplicado para as palavras sinónimas.

## 5 Apresentação do Resultado Final

Para testar o nosso *script* precisamos de executar o seguinte comando.

```

1 ./scriptname --to pt,it,fr cold

```

Este comando, irá gerar um ficheiro *TMX* com toda a informação descrita acima e com as traduções da palavra *cold* para português, italiano e francês.

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <tmx version="1.4">
3   <header adminlang="en"
4     datatype="tbx"

```

```

5      o-tmf="unknown"
6      segtype="block"
7      srclang="en"/>
8  <body>
9      <tu tuid="1">
10         <prop type="word_type">adjective</prop>
11         <prop type="meaning">Of or at a low or relatively low
temperature, especially when compared with the human body. ||
Lacking affection or warmth of feeling; unemotional. || (of
the scent or trail of a hunted person or animal) no longer
fresh and easy to follow. || Without preparation or rehearsal
.</prop>
12         <prop type="synonyms">chilly, cool, freezing, icy, snowy,
icy-cold, glacial, wintry, crisp, frosty, frigid, bitter,
bitterly cold, biting, piercing, numbing, sharp, raw, polar,
arctic, Siberian || unfriendly, cool, inhospitable,
unwelcoming, unsympathetic, forbidding, stony, frigid, frosty
, glacial, lukewarm, haughty, supercilious, disdainful, aloof
, distant, remote, indifferent, reserved, withdrawn,
uncommunicative, unresponsive, unfeeling, unemotional,
dispassionate, passionless, wooden, impersonal, formal, stiff
, austere || Doesn't have synonyms || unprepared, unready,
inattentive, unwary, unwatchful, with one's defences down, by
surprise, cold, unsuspecting</prop>
13         <tuv xml:lang="en">
14             <seg>cold</seg>
15         </tuv>
16         <tuv xml:lang="fr">
17             <seg>froide</seg>
18         </tuv>
19         <tuv xml:lang="pt">
20             <seg>frio</seg>
21         </tuv>
22         <tuv xml:lang="it">
23             <seg>fredde</seg>
24         </tuv>
25     </tu>
26
27
28     <tu tuid="2">
29         <prop type="word_type">noun</prop>
30         <prop type="meaning">A low temperature; cold weather; a
cold environment. || A common infection in which the mucous
membrane of the nose and throat becomes inflamed, typically
causing running at the nose, sneezing, and a sore throat.</
prop>
31         <prop type="synonyms">Doesn't have synonyms || cold, dose
of flu, dose of influenza, respiratory infection, viral
infection, virus</prop>

```

```

32     <tuv xml:lang="en">
33         <seg>cold</seg>
34     </tuv>
35     <tuv xml:lang="fr">
36         <seg>rhume</seg>
37     </tuv>
38     <tuv xml:lang="pt">
39         <seg>constipação</seg>
40     </tuv>
41     <tuv xml:lang="it">
42         <seg>freddo</seg>
43     </tuv>
44 </tu>
45 </body>
46 </tmx>

```

Listing 8: Ficheiro resultando do comando.

## 6 Conclusão

A realização deste trabalho prático foi importante para o nosso desenvolvimento como futuros engenheiros informáticos, pois permitiu-nos aprofundar os conhecimentos sobre *web scraping* e sobre ferramentas que o auxiliam como o *Beautiful Soup*. Para além disso, apresentou-nos um formato de ficheiro que nunca tínhamos ouvido falar e que viemos a descobrir ser muito importante no âmbito das traduções.

Por fim, podemos concluir que o *Beautiful Soup* é uma ferramenta poderosa, capaz de nos ajudar na busca de informação necessária para a construção de um projeto de maneira eficaz e direta.

## 7 Referências

- Lista de Elementos do *TMX*
- Lista de Atributos do *TMX*
- *Linguee*
- *Lexico*