

SCRIPTING NO PROCESSAMENTO DE LINGUAGEM NATURAL

UNIVERSIDADE DO MINHO

MESTRADO INTEGRADO EM ENGENHARIA INFORMÁTICA

Trabalho Prático 2

spaCy's POS Tagging

Grupo 7:

A73831 - João Barreira

A77364 - Mafalda Nunes

29 de Novembro de 2018



Resumo

O presente relatório tem com objetivo a aprendizagem das funcionalidades da ferramenta *spaCy*. Para tal, apresentar-se-á uma descrição da mesma, mais especificamente da funcionalidade de POS *tagging*, bem como um pequeno exemplo que demonstre como utilizar a ferramenta.

Este trabalho pretende dar resposta ao trabalho prático 2, proposto na unidade curricular SPLN, do Mestrado em Engenharia Informática, da Universidade do Minho.

Conteúdo

1 spaCy

O spaCy é uma biblioteca *open-source* para o processamento avançado de linguagem natural em *Python*, que foi desenvolvida especificamente para ser utilizada em ambientes de produção, ajudando a construir aplicações que processam grandes volumes de texto.

Dentro do contexto de processamento de linguagem natural, esta ferramenta pode ser utilizada em muitas vertentes, desde a extração de conhecimento, desenvolvimento de sistemas de compreensão de linguagem natural ou pré-processamento de texto para *deep-learning*.

1.1 Arquitetura

As estruturas de dados centrais do spaCy são o *Doc* e o *Vocab*.

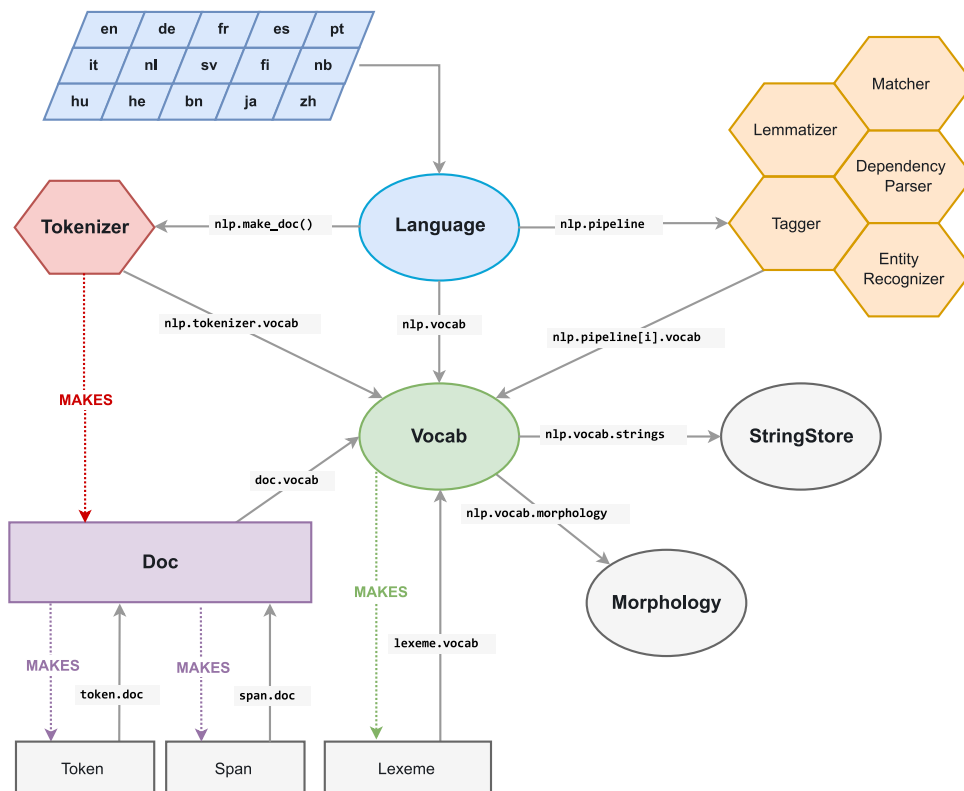


Figura 1: Esquema da arquitetura do spaCy

O objeto *Doc* possui uma sequência de *tokens* (palavra, pontuação, espaço em branco, etc.) e de *spans* (porções do objeto *Doc*), bem como todas as respectivas anotações. Os objetos *Span* e *Token* são vistas que apontam para o próprio *Doc*. O *Doc* é construído por um *Tokenizer*, sendo depois modificado pelos componentes de um *Pipeline* (figura ??). O objeto *Language* coordena estes componentes, ao pegar num texto limpo e enviá-lo para o *pipeline*, retornando um texto anotado. Para além disso, também permite organizar treinamento e serialização de dados.

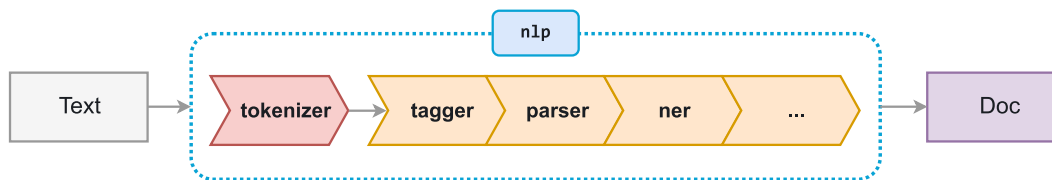


Figura 2: Esquema representativo da construção do *Doc*

Relativamente aos *pipelines*, é de se destacar os seguintes:

- *Lemmatizer* - permite determinar a forma base de palavras;
- *Tagger* - anota *part-of-speech* (POS) *tags* em objetos *Doc*;
- *Matcher* - faz *match* de sequências de *tokens*, tendo em conta regras padrão, similares a expressões regulares;
- *DependencyParser* - anota dependências sintáticas em objetos *Doc*;
- *EntityRecognizer* - anota nomes de entidades, como pessoas ou produtos, em objetos *Doc*.

O objeto *Vocab* possui um conjunto de tabelas de *look-up*, que tornam a informação comum disponível entre documentos. O *Vocab* é constituído por entradas do tipo *Lexeme*, que são um tipo de palavras sem contexto, ao contrário de uma palavra *token*, que possui uma *POS tag*, um parser de dependências, entre outros. O *Vocab* permite também aceder ao objeto *Morphology*, que possibilita a associação de características linguísticas, como o *lemma*, o tipo de nome, o tempo verbal, entre outros, aos *tokens*, tendo em conta a palavra e a respetiva *POS tag*. Ao centralizar vetores de palavras, atributos léxicos e *strings*, o spaCy evita armazenar múltiplas cópias desses dados, pelo que poupa memória e assegura que existe uma única fonte de verdade.

Também para poupar memória, o spaCy codifica todas as *strings* para valores de *hash*, que são acessíveis a partir do *StringStore*, que corresponde a um *map* de *strings* de e para valores de *hash*. Internamente, o spaCy apenas trabalha com valores de *hash* (figura ??).

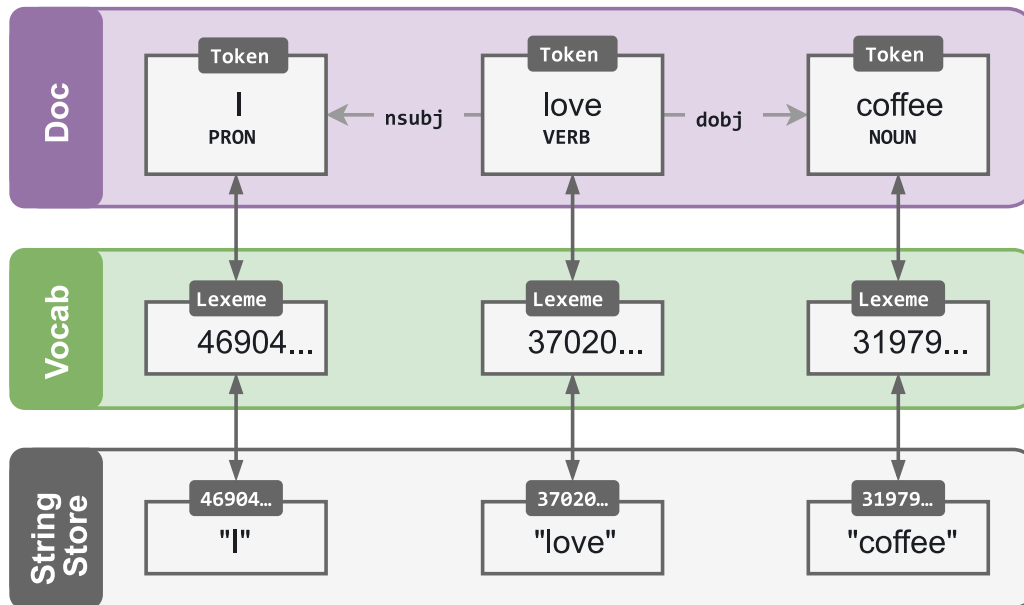


Figura 3: Representação esquemática dos objetos *Doc*, *Vocab* e *String Store*

1.2 Funcionalidades

A biblioteca spaCy disponibiliza diversas funcionalidades, sendo que algumas estão associadas a conceitos linguísticos e outras a capacidades mais gerais de *machine learning*. Apesar de algumas destas funcionalidades trabalharem de forma independente, outras requerem o carregamento de modelos estatísticos, para que o spaCy seja capaz de prever anotações linguísticas. Atualmente, o spaCy providencia modelos estatísticos para 8 linguagens, que podem diferir em tamanho, velocidade, utilização de memória, correção e os dados que incluem.

As principais funcionalidades desta biblioteca são sumarizadas de seguida:

- ***Tokenization*** – Separação de textos em segmentos significativos, como palavras, pontuação, entre outros. Esta divisão é efetuada tendo em conta regras específicas de cada linguagem. O *input* corresponde a um texto limpo e o *output* um objeto *Doc*.
- ***POS Tagging*** – atribuir tipos de palavras ou propriedades gramaticais (como verbo ou nome) a *tokens*, que constituem um *Doc*. O spaCy recorre a modelos estatísticos para prever quais as *tags* ou etiquetas que mais provavelmente se aplicam ao contexto corrente. As anotações linguísticas do tipo POS são obtidas a partir do atributo *pos_* do *Token*.
- ***Dependency Parsing*** – Processo de obtenção de relações de dependência entre os elementos de uma frase através de um *parser*. Este *parser* possibilita a iteração do texto *tokenizado* (*Doc*) pelas suas frases nominais, sendo que as dependências são dispostas numa árvore de fácil manipulação. É ainda possível visualizar mais facilmente as dependências de um ou mais *Docs* através de um grafo, cujos arcos possuem uma legenda que caracteriza a ligação sintática entre os *tokens*.
- ***Lemmatization*** – redução de palavras à sua formas base, sendo frequentemente utilizado para standardizar palavras com significados similares. O *lemma* de 'era' e 'foi', e.g., é 'ser'. Pode-se aceder ao *lemma* através do atributo *lemma_* do *Token*.
- ***Sentence Boundary Detection (SBD)*** – Processo de separação de frases de um texto separado pelos seus *tokens* (*Doc*). É utilizado o *parsing* de dependências para uma divisão mais correta. Possibilita a definição manual dos limites frásicos (i.e. *boundaries*) em detrimento dos pré-definidos '.', '!' e '?'.
- ***Named Entity Recognition (NER)*** - atribuição de etiquetas ou categorias pré-definidas a nomes correspondentes a objetos do mundo real, como pessoas, companhias, localizações, países, produtos, entre outros. O spaCy reconhece vários tipos de nomes de entidades num documento, através da solicitação de uma previsão ao modelo. Como estes modelos são estatísticos e muito dependentes dos textos com que foram treinados, nem sempre esta previsão é correta, podendo ser necessário posterior processamento específico para o caso em estudo. Os nomes das entidades estão disponíveis na propriedade *ents* de um *Doc* ou na propriedade *ent_type_* de um *Token*. Quando se pretende

alterar anotações de entidades, tem de se o fazer a nível do documento, de forma a assegurar que a sequência de anotações de *tokens* permanece consistente.

- ***Similarity*** – Comparação de dois textos e cálculo de valor que indique quão similares são. A função de similitude compara palavras de um dado texto com base em vetores de palavras específicos para cada linguagem que podem ser customizados para um melhor desempenho em contextos de aplicação específicos.

- ***Text Classification*** – atribuição de categorias ou etiquetas a todo o documento ou partes de um documento.

- ***Rule-based matching*** – Processo de busca de correspondências entre *tokens* e determinados padrões. Estes padrões podem ser bastante complexos, tendo algumas parecenças com o funcionamento expressões regulares. No entanto, podem também tirar partido do conhecimento léxico desta ferramenta, filtrando, por exemplo, por bases morfológicas (i.e. *Lemma*) ou classes gramaticais (i.e. verbos, substantivos, etc.). Possibilita ainda a formulação de padrões mais complexos através de quantificadores, *wildcards*, expressões regulares. Podem ainda ser definidas, de um modo simples, funções que serão chamadas caso se dê uma ou mais correspondências.

- ***Training*** – atualização e melhoria de modelos estatísticos de previsões. Para treinar um modelo, são necessários dados de treino, i.e., exemplos de texto e etiquetas (POS *tag*, *named entity*, etc.) que se pretende que o modelo seja capaz de prever. O modelo pode ser melhorado através do processamento de novos textos limpos, em que o utilizador, que conhece a correção da resposta, dá *feedback* ao modelo sobre a sua previsão na forma de um gradiente de erro da função de perda que calculam a diferença entre o exemplo de treino e o *output* esperado (figura ??). Quanto maior a diferença, mais significativo o gradiente e as atualizações ao modelo.

- ***Serialization*** – Funcionalidade responsável por guardar objetos do *spaCy* para ficheiros em disco, bem como por carregar objetos previamente guardados. Este processo torna-se especialmente importante quando usado para guardar objetos do tipo *Doc*, por exemplo, que possuem todas as informações retiradas do texto original, bem como outras características inseridas manualmente como entidades ou vetores de palavras customizados. Para isso, é utilizado a biblioteca *Pickle*,

built-in da linguagem *Python*.

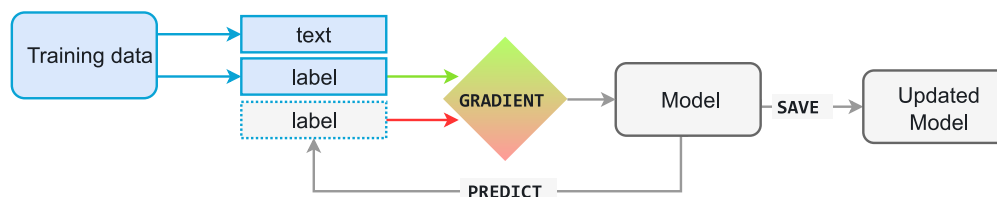


Figura 4: Esquema representativo do funcionamento do processo de treino

1.3 Vantagens

Existem diversas ferramentas do género do spaCy que podem ser utilizadas para efetuar processamento de linguagem natural. Assim, há que apresentar as principais distinções entre estas, de forma a perceber qual será a que trará mais vantagens na sua utilização.

Em primeiro lugar, pode-se efetuar uma breve comparação entre as funcionalidades oferecidas pelo spaCy e outras ferramentas semelhantes, como SyntaxNet, NLTK e CoreNLP.

| | spaCy | SyntaxNet | NLTK | CoreNLP |
|-------------------------|--------|-----------|--------|---------|
| Programming language | Python | C++ | Python | Java |
| Neural network models | ✓ | ✓ | ✗ | ✓ |
| Integrated word vectors | ✓ | ✗ | ✗ | ✗ |
| Multi-language support | ✓ | ✓ | ✓ | ✓ |
| Tokenization | ✓ | ✓ | ✓ | ✓ |
| Part-of-speech tagging | ✓ | ✓ | ✓ | ✓ |
| Sentence segmentation | ✓ | ✓ | ✓ | ✓ |
| Dependency parsing | ✓ | ✓ | ✗ | ✓ |
| Entity recognition | ✓ | ✗ | ✓ | ✓ |
| Coreference resolution | ✗ | ✗ | ✗ | ✓ |

Tabela 1: Funcionalidades de diversas ferramentas NLP

Como se pode perceber através da análise da tabela ??, o spaCy disponibiliza todas as principais funcionalidades associadas o processamento de linguagem natural, exceto a *Conference Resolution*, que corresponde à tarefa de encontrar todas as expressões que se referem à mesma entidade num texto. Apenas o CoreNLP possui essa capacidade, sendo que, em contrapartida, não possui *Integrated word vectors*, que possibilitam a previsão de quão similares são dois objetos. Assim, percebe-se que o spaCy será a melhor alternativa em muitos dos casos.

Outro fator fundamental na comparação entre o spaCy e outras ferramentas semelhantes, é o desempenho, i.e., o tempo que o spaCy demora a gerar respostas. Dois artigos revistos em 2015 (REF) confirmam que o spaCy oferece o *parser* sintático mais rápido do mundo e que a sua precisão está dentro de 1% do melhor disponível. Os poucos sistemas mais precisos são pelo menos 20 vezes mais lentos. Alguns valores de precisão e velocidade são apresentados na tabela ?. Informações mais detalhadas relativamente à velocidade são apresentadas na tabela ?.

| System | Year | Language | Accuracy | Speed (wps) |
|------------|------|---------------|----------|-------------|
| spaCy v2.x | 2017 | Python/Cython | 92.6 | n/a |
| spaCy v1.x | 2015 | Python/Cython | 91.8 | 13,963 |
| ClearNLP | 2015 | Java | 91.7 | 10,271 |
| CoreNLP | 2015 | Java | 89.6 | 8,602 |
| MATE | 2015 | Java | 92.5 | 550 |
| Turbo | 2015 | C++ | 92.4 | 349 |

Tabela 2: Precisão e velocidade de várias ferramentas NLP

A partir das informações das tabelas ? e ?, percebe-se que o spaCy deverá ser uma boa escolha, também em termos de precisão e desempenho.

Há, contudo, que se ter em atenção o tipo de texto que se irá utilizar e o tipo de processamento que se pretende efetuar, uma vez que os valores de precisão, o tamanho do modelo e a velocidade de processamento dependem de vários fatores, como, por exemplo, a linguagem ou modelo utilizado.

| System | Absolute (ms per doc) | | | Relative (to spaCy) | | |
|---------|-----------------------|-------|-------|---------------------|------|-------|
| | Tokenize | Tag | Parse | Tokenize | Tag | Parse |
| spaCy | 0.2ms | 1ms | 19ms | 1x | 1x | 1x |
| CoreNLP | 0.18ms | 10ms | 49ms | 0.9x | 10x | 2.6x |
| ZPar | 1ms | 8ms | 850ms | 5x | 8x | 44.7x |
| NLTK | 4ms | 443ms | n/a | 20x | 443x | n/a |

Tabela 3: Informações detalhadas sobre velocidade de diversas ferramentas NLP

1.4 Utilização

O processo de instalação da biblioteca é bastante simples, podendo ser realizado através do *pip* (*package manager* da linguagem *Python*), utilizado, neste exemplo, na sua versão 3.X.

```
$ pip3 install spacy
```

Os modelos correspondentes às linguagens podem ser também facilmente obtidos através do seguinte comando abaixo, exemplificado para um dos modelos da língua inglesa.

```
$ python3 -m spacy download en_core_web_sm
```

Após o processo de instalação da biblioteca e de *download* das linguagens, o primeiro passo será importar a biblioteca. Depois, poder-se-á carregar a linguagem já descarregada e produzir um *Doc* a partir de um qualquer texto. Será então possível começar imediatamente a tirar proveito das funcionalidades providenciadas pela biblioteca, utilizando o objeto do tipo *Doc* resultante do processamento do texto inserido.

```
>>> import spacy
>>> nlp = spacy.load('en_core_web_sm')
>>> doc = nlp('This text can now be processed using spaCy!')
```

1.5 Aplicações

Com todas as funcionalidades do spaCy, consegue-se concretizar diversos tipos de processamento de texto de forma geralmente simples. Algumas possibilidades são as que se indicam de seguida:

- desenhar os gráficos de dependências num texto;
- reconhecer entidades num texto;
- extrair relações entre frases e entidades utilizando as funcionalidades de *entity recognizer* e de *dependency parse*;
- obter marcações de um conjunto personalizado de entidades, alterando-se para isso anotações de entidades com base, e.g., numa lista de nomes desse conjunto e juntando-se entidades num único *token*;
- corrigir ou personalizar atributos nos objetos *Doc*, *Span* e/ou *Token*;
- treinar o *Named Entity Recognizer*, o *Dependency Parser*, o *Part-of-speech Tagger* ou o *Text Classifier* do spaCy, partindo de um modelo já treinado ou de um em branco;

Estes exemplos apresentados são de implementação relativamente fácil, tendo-se concretizado alguns dos mais simples.

1.5.1 Grafos de Dependências

Relativamente ao desenho dos grafos de dependências, é possível desenvolver uma função que os imprima para ficheiros com formato *svg* ou para *html*, sendo esta apresentada no excerto de código do anexo ??) com o nome de *generate_dependencies_graph*.

Para obter as dependências dentro de cada frase, basta separar as mesmas do objeto *Doc* e invocar sobre estas a função *serve* (mostra no *browser*) ou *render* (devolve conteúdo a guardar num documento *svg* ou *html*) com o estilo *dep*. Existem outras formas de realizar esta tarefa, mas esta é a mais simples e ainda permite personalizar o grafo gerado, como mudar a cor do fundo, tipo de letra, entre outros.

Esta função pode ser utilizada da seguinte forma:

```
>>> import spacy
>>> from spacy import displacy
>>> nlp = spacy.load('pt')
>>> doc = nlp('Pretende-se gerar o gráfico de dependências do
spaCy!')
>>> generate_dependencies_graph(doc)
```

Serving on port 5000...
Using the 'dep' visualizer

Acedendo-se ao endereço `127.0.0.1:5000/` no *browser*, obtém-se o resultado apresentado na figura ???. Também se poderia solicitar o output em ficheiro(s) *html* ou *svg*, aterando o argumento *type*. Nesse caso, o conteúdo dos ficheiros seria retornado pela função. Os restantes argumentos, também opcionais, conferem diferentes características estéticas ao *output*.

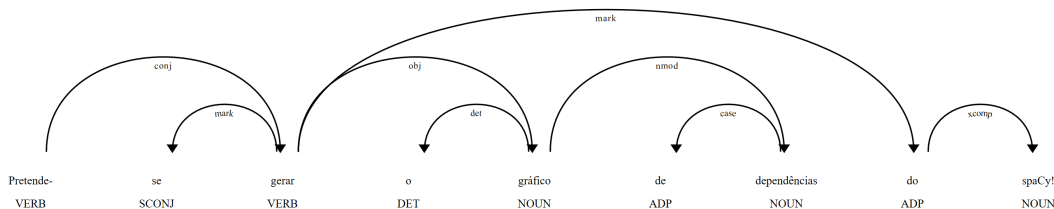


Figura 5: *Output* de invocação da função `generate_dependencies_graph`

1.5.2 Reconhecimento de Entidades

Para se proceder ao reconhecimento de entidades, desenvolveu-se a função `generate_tagged_text` do excerto de código do anexo ???). Essa função permite obter resultados no *browser*, em ficheiros *html* ou na *shell*.

Para desenvolver esta funcionalidade, basta invocar sobre o *Doc* a função `displacy.serve` (mostra no *browser*) ou `displacy.render` (devolve conteúdo a guardar num documento *html*) com o estilo *ent*. Para se obter o *output* em formato de texto, percorreu-se todos os *tokens* do *Doc* e, caso este tivesse uma entidade associada (atributo *ent_type*), a mesma seria escrita junto ao *token*.

Esta função pode ser utilizada da seguinte forma, para obter o *output* no *browser*:

```
>>> import spacy
>>> from spacy import displacy
>>> nlp = spacy.load('pt')
>>> doc = nlp('Hoje o Manuel gostava de ir ao cinema do Braga
Parque!')
>>> generate_tagged_text(doc)
```

Serving on port 5000...
Using the 'ent' visualizer

Acedendo-se ao endereço 127.0.0.1:5000/ no *browser*, obtém-se o resultado apresentado na figura ???. Também se poderia solicitar o output em ficheiro *html*, aterando o argumento *type*. Nesse caso, o conteúdo do ficheiro seria retornado pela função. O atributo *entities* permite indicar quais as entidades que se pretendem etiquetar. Caso nenhuma seja indicada, todas são etiquetadas. O argumento *color*, também opcional, confere diferentes cores aos *labels* das entidades.

Hoje o **Manuel** **PER** gostava de ir ao cinema do **Braga Parque** **ORG** !

Figura 6: *Output* de invocação da função *generate_tagged_text*

Caso se pretenda obter o resultado em texto, pode-se realizar o seguinte comando:

```
>>> generate_tagged_text(doc, 'text')

'Hoje o Manuel{PER} gostava de ir a o cinema do Braga{ORG}
Parque{ORG} !'
```

1.5.3 Especificação de *Tokens*

Com o spaCy consegue-se, de forma bastante simples, alterar informação associada a determinado *Token*, ou gerar uma nova divisão para determinada palavra. A palavra inglesa *don't*, e.g., é geralmente dividida em *do* e *n't*

(*not*), mas tal não se verifica no Reino Unido, onde esta deve permanecer sempre um único *Token*.

A função *add_tokenizer_exceptions*, apresentada no excerto de código do anexo ??, demonstra como personalizar um *Token*, ao adicionar ao *Tokenizer* do modelo casos especiais. O argumento *tokens* dessa função deverá ser um dicionário, em que a chave é o *Token* original e o valor é uma lista de dicionários. Esta lista de dicionários corresponde aos atributos do *Token*, devendo conter o nome do atributo e o respetivo valor como chave e valor, respetivamente. Note-se que o atributo *ORTH* deve ser sempre especificado e, todos concatenados, devem formar o *Token* original.

A referida função pode ser utilizada da seguinte forma:

```
>>> import spacy
>>> nlp = spacy.load('en_core_web_lg')
>>> doc = nlp("We don't need to go now!")
>>> [print((token.orth_, token.lemma_, token.tag_,
token.pos_)) for token in doc]
('We', '-PRON-', 'PRP', 'PRON')
('do', 'do', 'VBP', 'VERB')
('n't', 'not', 'RB', 'ADV')
('need', 'need', 'VB', 'VERB')
('to', 'to', 'TO', 'PART')
('go', 'go', 'VB', 'VERB')
('now', 'now', 'RB', 'ADV')
('!', '!', '.', 'PUNCT')

>>> tokens = {"don't": [{'ORTH': 'do', 'LEMMA': 'do',
'POS': 'VERB'}, {'ORTH': "n't", 'LEMMA': 'not'}]}
>>> doc = nlp("We don't need to go now!")
>>> [print((token.orth_, token.lemma_, token.tag_,
token.pos_)) for token in doc]
('We', '-PRON-', 'PRP', 'PRON')
("don't", "don't", '', 'VERB')
('need', 'need', 'VB', 'VERB')
('to', 'to', 'TO', 'PART')
('go', 'go', 'VB', 'VERB')
('now', 'now', 'RB', 'ADV')
('!', '!', '.', 'PUNCT')
```

1.5.4 Treino do *dependency parser*

Como foi dito anteriormente, o *spaCy* possibilita o treino de novos modelos estatísticos utilizados nas suas mais diversas funcionalidades. Uma dessas funcionalidades, demonstrada neste exemplo, é o *dependency parser*.

Em primeiro lugar, é criada uma variável *TRAIN_DATA* que possui duas frases nominais, bem como o nome das dependências (*deps*) e quais os elementos principais dessas mesmas dependências (*heads*). Esta variável servirá como base para o processo de treino. De seguida, é criada a linguagem que poderá ser uma das pré-existentes, ou criada de raiz (i.e. *blank*). Depois, é adicionado o *parser* à *pipeline* e adicionadas as etiquetas com os nomes das dependências ao *parser* que é finalmente utilizado para treinar a linguagem utilizando a variável *TRAIN_DATA* criada anteriormente. O resultado é, a seguir, testado, através da impressão do resultado das dependências para um novo texto.

Adicionalmente, tira-se ainda partido da capacidade de serialização do *spaCy* (abordada anteriormente neste relatório) para guardar o modelo da linguagem que se encontra modificada após o processo de treino.

O código deste exemplo está presente no Anexo [??](#), encontrando-se também disponível no repositório do *GitHub* REF.

2 *spaCy*'s POS Tagging

A funcionalidade de POS *Tagging* do *spaCy* foi explicitada anteriormente, mas merece especial destaque por ser o principal tema do trabalho.

Um modelo consiste em dados binários e é produzido ao apresentar ao sistema exemplos suficientes para que ele faça previsões que são generalizadas a toda a linguagem. Há a necessidade de prever, e.g., a classe (POS) das várias palavras, podendo a mesma palavra ter diferentes classes em contextos distintos ('gosto' pode ser um verbo ou um nome, dependendo da frase em que está incluída). Este tipo de previsão permite efetuar novas conjecturas, como que tipo de palavra se espera encontrar após outra (e.g., a palavra que se segue ao determinante artigo definido 'o' será muito provavelmente um nome).

Para além de palavras com a mesma classe tenderem a seguir uma estrutura sintática semelhante, também são úteis em processamento baseado

em regras. Assim, consegue-se perceber a importância de determinar corretamente a classe (*part-of-speech*) de cada palavra num determinado texto de *input*.

O spaCy utiliza as POS *tags* do projeto *Penn Treebank* REF.

2.1 Utilização

Após ter sido instalada a biblioteca e feito o *download* das linguagens (processo descrito anteriormente), poder-se-á testar o processo de obtenção das *POS tags* para um determinado texto, de forma simples. Para tal, bastará introduzir um texto de input e – como apresenta o código abaixo –, imprimir, por exemplo, a informação sobre as *POS tags* de cada uma das suas palavras.

```
>>> import spacy
>>> nlp = spacy.load('en_core_web_sm')
>>> doc = nlp('This is a demonstration of a very cool
Python library!')
>>> for token in doc:
...     print(token.text, token.pos_)
This DET
is VERB
a DET
demonstration NOUN
of ADP
a DET
very ADV
cool ADJ
Python PROPN
library NOUN
! PUNCT
```

3 Exemplo Demonstrativo

Para construir um exemplo demonstrativo de algumas das funcionalidades do spaCy mais interativo, utilizou-se a biblioteca *flask* para correr a

aplicação no *browser* com uma interface gráfica baseada em *templates html*. Esta ferramenta foi utilizada apenas como auxiliar, não estando relacionada com o tema do trabalho, pelo que não se aprofundará o funcionamento da mesma. Contudo, se o utilizador preferir, pode realizar todas as operações a partir da linha de comandos, através da utilização de opções.

O código associado a esta parte da aplicação, à interação com o utilizador, é apresentado no anexo ??. A utilização da mesma é indicada pela própria aplicação quando se solicita ajuda (-h) ou se insere um comando inválido:

```
USAGE: python3 app.py [(-i <outputfile> | -p <outputfile>
| -t <outputfile> | -g <outputfile>) [-l <language>]
<inputfile>]
```

OPTIONS:

```
    -l    language - 'pt' (default) or 'en'
-i    returns table with information about tokens to
outputfile (or 'shell')
-p    generates bar chart with POS tag frequency to
outputfile
    -t    returns tagged text to outputfile (or 'shell')
    -g    generates dependencies graphs to outputfile.svg
* if no options are provided, a web server will be
initialized, where the input will be inserted and the
output presented.
```

Nos anexos REF a REF apresentam-se os resultados obtidos na aplicação *web*

Para se implementar as funcionalidades disponibilizadas pelo spaCy, utilizou-se as funções anteriormente referidas, que se encontram no anexo ??. No ambiente *web*, o reconhecimento de entidades permite fazer uma filtragem de quais serão etiquetadas. Nesse contexto, todos os *outputs* obtidos são do tipo *html*, sendo depois integrados nos *templates* existentes.

Para além das funções anteriormente referidas, podem-se encontrar no código apresentado no anexo ?? outras mais direccionadas ao tema principal do projeto: POS *tagging*.

3.1 Gráfico de frequência POS

A função *generate_pos_chart* é responsável pela geração de um gráfico que indica a frequência de cada POS *tag* num determinado *Doc*. Para se obter as frequências, basta utilizar a função *count_by* do *Doc* e passar-lhe como argumento o identificador de POS (*spacy.attrs.POS*). O resultado dessa invocação é um dicionário com o valor de *hash* do POS como chave e o número de vezes que se repete como valor. De seguida, basta obter o nome associado aos valores de *hash*. Caso seja solicitado o *output* do tipo *html*, é retornada uma lista de listas, sendo que cada lista interna armazena o par (POS, frequência). Essa lista poderá então ser utilizada para construir um *Google Chart*. Caso seja solicitado o *output* em formato *pict* (imagem), é guardado um gráfico de barras como imagem.

Esta função pode ser utilizada da seguinte forma:

```
>>> import spacy
>>> import matplotlib as mpl
>>> mpl.use('Agg')
>>> from matplotlib import pyplot as plt
>>> nlp = spacy.load('pt')
>>> doc = nlp('Hoje o João e a Maria tiveram de ir ao hospital
visitar a Teresa.')
>>> generate_pos_chart(doc, type='pict')
```

3.2 Tabela com informações sobre *Tokens*

A função *generate_information* permite obter várias informações sobre os *Tokens* de um documento. Para isso, basta percorrer os *tokens* do *Doc* e aceder aos atributos do mesmo, conforme explicitado previamente. Se o *output* for solicitado no tipo *html*, retorna-se uma tabela nesse formato. Caso contrário, retorna-se uma *PrettyTable* que poderá ser impressa no terminal.

```
>>> import spacy
>>> from prettytable import PrettyTable
>>> nlp = spacy.load('en_core_web_lg')
>>> doc = nlp('We have to go there now!')
>>> print(generate_information(doc, nlp.vocab, type='text'))
```

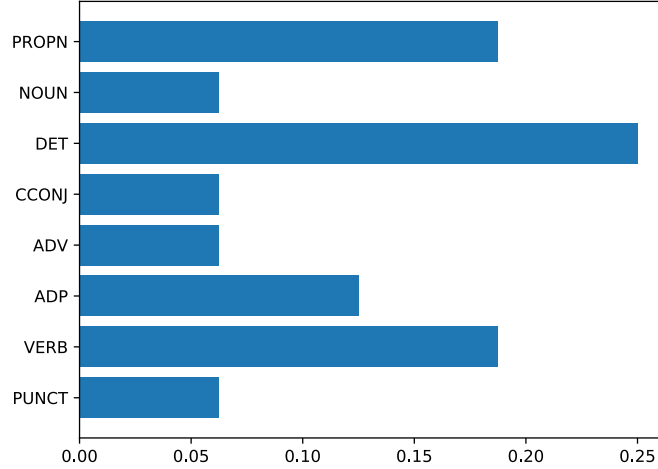


Figura 7: *Output* de invocação da função *generate_pos_chart*

| Text | Lemma | POS | TAG | DEP | SHAPE | MORPHOLOGIAL INFO | IS_ALPHA | IS_STOP |
|-------|--------|-------|-----|--------|-------|--|----------|---------|
| We | -PRON- | PRON | PRP | nsubj | Xx | {'PronType': 'prs'} | True | False |
| have | have | VERB | VP | ROOT | xxxx | {'VerbForm': 'fin', 'Tense': 'pres'} | True | False |
| to | to | PART | TO | aux | xx | {'PartType': 'inf', 'VerbForm': 'inf'} | True | False |
| go | go | VERB | VB | xcomp | xx | {'VerbForm': 'inf'} | True | False |
| there | there | ADV | RB | advmod | xxxx | {'Degree': 'pos'} | True | False |
| now | now | ADV | RB | advmod | xxx | {'Degree': 'pos'} | True | False |
| ! | ! | PUNCT | . | punct | ! | {'PunctType': 'peri'} | False | False |

Figura 8: *Output* de invocação da função *generate_info*

3.3 Resultados

Os resultados obtidos através da execução via linha de comandos são iguais aos previamente recolhidos para a execução das funções individualmente. No contexto *wed*, apresentam-se nas imagens do anexo REF exemplos de funcionamento de cada uma das funcionalidades.

Anexos

A *Working Example*: código associado à utilização de spaCy

B *Working Example:* código associado à interação
com o utilizador

C Exemplo de aplicação: Treino do *dependency parser* (e serialização)

D *Working Example: interface web*

Scripting no Processamento de Linguagem Natural

spaCy's POS Tagging

Input Form

Language: English

Entities: ☒ ORG ☒ MONEY ☒ QUANTITY ☒ ORDINAL ☒ LAW ☒ WORK_OF_ART ☒ EVENT ☒ PRODUCT ☒ LANGUAGE ☒ DATE ☒ PERCENT ☒ TIME ☒ GPE ☒ NORP ☒ FAC ☒ CARDINAL ☒ PERSON ☒ LOC

Add/Modify Token

Select File: Escolher ficheiro en.txt

Input Text

Article 13 approved: What is the EU copyright law amendment?

Critics of Article 13, the proposed EU directive on copyright, warned that it could censor internet users

Politics

In September the European Parliament voted on – and approved – an amended version of the EU copyright law that was rejected by the same body in July of this year.

The divisive legislation was put to the ultimate democratic test in Strasbourg, with 438 voting in favour of the measures, 226 voting against and 39 abstaining.

Figura 9: Página inicial onde se indica a linguagem, as entidades a etiquetar e o *input*

Token Exception

| | |
|----------|------------|
| Token | don't |
| ORTH | do,n't |
| IS_ALPHA | True,False |
| IS_ASCII | |
| IS_DIGIT | |
| IS_LOWER | |
| IS_PUNCT | |
| IS_SPACE | |
| IS_TITLE | |
| IS_UPPER | |
| LIKE_URL | |

* Concatenated ORTH elements must be equal to the inserted Token

** All fields, except Token, must be separated with ','

Confirm Cancel

Figura 10: Modal da página inicial, onde se especifica um caso especial de *Token*

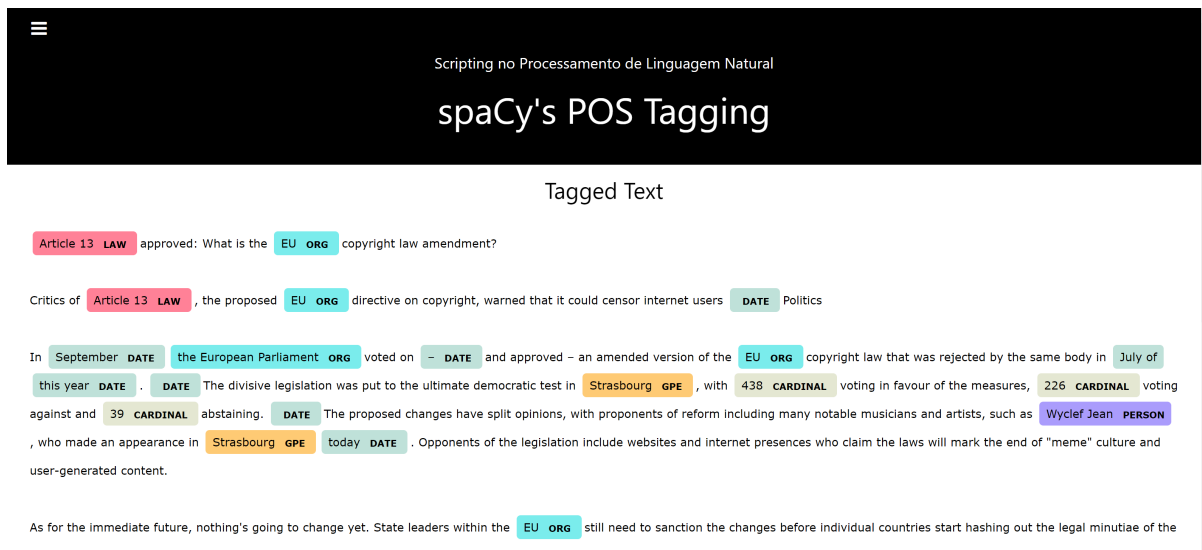


Figura 11: Página que apresenta o texto de *input* com as entidades etiquetadas

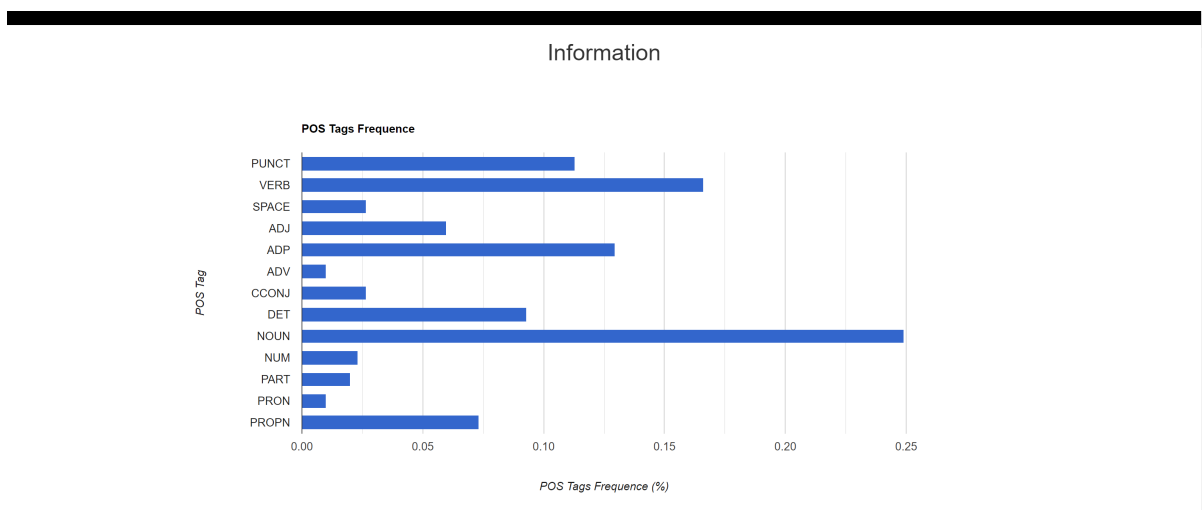


Figura 12: Excerto de página que apresenta a taxa de ocorrência de cada POS

| Text | Lemma | POS | TAG | DEP | SHAPE | MORPHOLOGICAL INFO | IS_ALPHA | IS_STOP |
|-----------|-----------|-------|-----|----------|-------|---|----------|---------|
| Article | article | PROPN | NNP | nsubj | Xxxxx | {'NounType': 'prop', 'Number': 'sing'} | True | False |
| 13 | 13 | NUM | CD | nummod | dd | {'NumType': 'card'} | False | False |
| approved | approve | VERB | VRN | ROOT | xxxx | {'VerbForm': 'part', 'Tense': 'past', 'Aspect': 'perf'} | True | False |
| : | : | PUNCT | : | punct | : | | False | False |
| What | what | NOUN | WP | attr | Xxxx | {'PronType': 'int rel'} | True | False |
| is | be | VERB | VBZ | ccomp | xx | {'VerbForm': 'fin', 'Tense': 'pres', 'Number': 'sing', 'Person': 3} | True | False |
| the | the | DET | DT | det | xxx | | True | False |
| EU | eu | PROPN | NNP | compound | XX | {'NounType': 'prop', 'Number': 'sing'} | True | False |
| copyright | copyright | NOUN | NN | compound | xxxx | {'Number': 'sing'} | True | False |
| law | law | NOUN | NN | compound | xxx | {'Number': 'sing'} | True | False |
| amendment | amendment | NOUN | NN | nsubj | xxxx | {'Number': 'sing'} | True | False |
| ? | ? | PUNCT | . | punct | ? | {'PunctType': 'per'} | False | False |

Figura 13: Excerto de página que apresenta informações sobre cada token do texto inicial

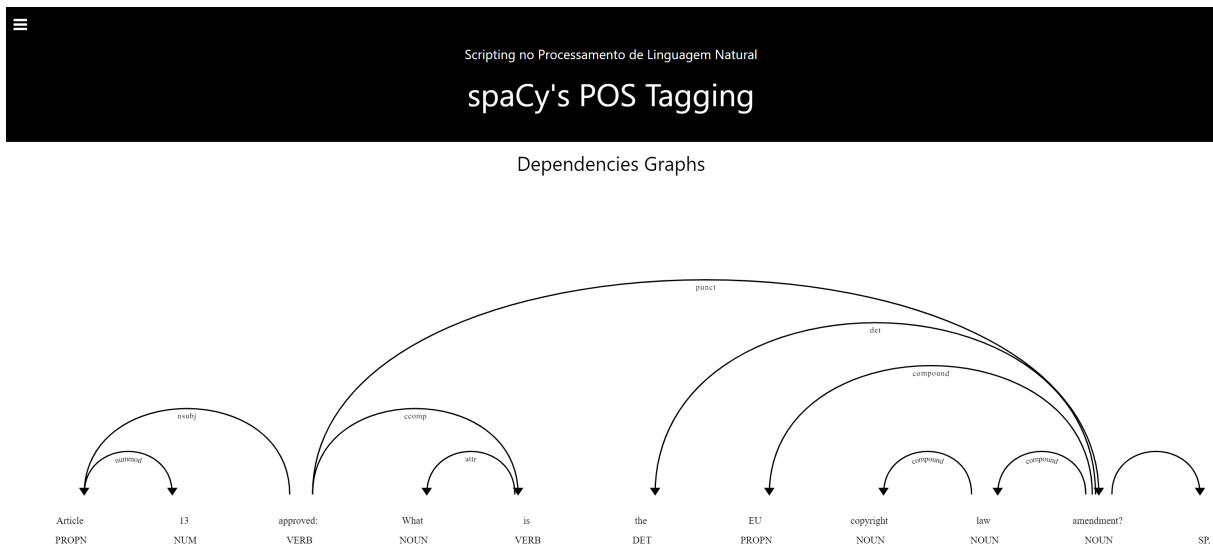


Figura 14: Página que apresenta grafos de dependências de cada frase do texto inicial

Referências

- [1] Autor, “Título”, *website name*, [websitelink](#).