

```

1 from IPython.display import HTML, display
2
3 def set_css():
4     display(HTML('''
5         <style>
6             pre {
7                 white-space: pre-wrap;
8             }
9         </style>
10     '''))
11 get_ipython().events.register('pre_run_cell', set_css)

```



▼ Importing Libraries

```

1 !python -m spacy download en_core_web_sm
2 !pip install pdfminer.six

```

```

1 import numpy as np
2 import pandas as pd
3 from gensim.models import Word2Vec
4 import nltk
5 from nltk.corpus import stopwords
6 from nltk.stem import WordNetLemmatizer
7 from nltk.corpus import wordnet
8 import spacy
9 from spacy import displacy
10 import re
11 import string
12 from pdfminer.high_level import extract_text
13 import en_core_web_sm
14 from IPython.display import clear_output
15
16 nltk.download('stopwords')
17 nltk.download('punkt')
18 nltk.download('wordnet')
19 nltk.download('averaged_perceptron_tagger')
20 nlp = en_core_web_sm.load()

```

```

[ ] [nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /root/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]   date!

```

▼ Text extraction (PDF)

```
1 def pdf2text(n_files):
2     pdfs = []
3     for i in range(1, n_files + 1):
4         pdfs.append(extract_text('job{}.pdf'.format(i), maxpages = 2))
5     return pdfs
```



```
1 jobs = pdf2text(13)
```



▼ Data cleaning and processing

Regex

```
1 alphabets = '([A-Za-z])'
2 prefixes = '(Mr|Mrs|Ms|Dr|Pr)[.]'
3 suffixes = '(Inc|Ltd|Jr|Sr|Co)'
4 acronyms = '([a-zA-Z]{2,})'
5 websites = '([.](com|net|org|io|gov|fr|uk|usa|esp))'
6 starters = '(Mr\.|Mrs\.|Ms\.|Dr\.|Pr\.|However|But|The\s|This\s|That\s|Those\s|Their\s|
```



Lemmatization

```
1 lemmatizer = WordNetLemmatizer()
2
3 def nltk2wn_tag(nltk_tag):
4     if nltk_tag.startswith('J'):
5         return wordnet.ADJ
6     elif nltk_tag.startswith('V'):
7         return wordnet.VERB
8     elif nltk_tag.startswith('N'):
9         return wordnet.NOUN
10    elif nltk_tag.startswith('R'):
11        return wordnet.ADV
12    else:
13        return None
14
15 def lemmatize_sentence(sentence):
16     nltk_tagged = nltk.pos_tag(nltk.word_tokenize(sentence))
17     wn_tagged = map(lambda x: (x[0], nltk2wn_tag(x[1])), nltk_tagged)
18     res_words = []
```

```

18 res_words = []
19 for word, tag in wn_tagged:
20     if tag is None:
21         res_words.append(word)
22     else:
23         res_words.append(lemmatizer.lemmatize(word, tag))
24 return ' '.join(res_words)

```



Split function

```

1 def split_into_sentences(text):
2     text = text.replace('\n', ' ')
3     text = ' '.join([word for word in text.split(' ') if word != ''])
4     text = text.replace('\uf0b7', '').replace('\x0c', '').replace('\uf054', '').replace('
5     text = re.sub('\([a-zA-Z]{1}\)', '', text)
6     text = re.sub(prefixes, '\\1<prd>', text)
7     text = re.sub(suffixes + '[.]', '\\1<prd>', text)
8     text = re.sub('[.]' + alphabets, '<prd>\\1', text)
9     text = re.sub('([1-9a-zA-Z])[.]([1-9a-zA-Z])', '\\1<prd>\\2', text)
10    text = re.sub(alphabets + '[.]' + alphabets, '\\1<prd>\\2', text)
11    text = re.sub(alphabets + '[.]' + alphabets + '[.]' + alphabets, '\\1<prd>\\2<prd>\\3
12    text = re.sub('\.{2,}', '<prd>', text)
13    text = re.sub(websites, '<prd>\\1', text)
14    text = re.sub(starters, '<stop>\\1', text)
15    for acr in re.findall(acronyms, text): text = text.replace(acr, '<prd>'.join(acr.split('
16    text = text.replace('.', '.<stop>').replace('!', '!<stop>').replace('?', '?<stop>').r
17    text = text.replace('<prd>', '.')
18
19    sentences = text.split('<stop>')
20    sentences = [sent.strip() for sent in sentences]
21    sentences = [sent for sent in sentences if sent not in ['', ' ']]
22    sentences = [sent for sent in sentences if len(sent) > 1]
23    sentences = list(map(lambda sent: ' '.join([word for word in sent.split() if word.lower
24
25    sentences = [lemmatize_sentence(sent) for sent in sentences]
26    sentences = [sent.lower() for sent in sentences if sent != '']
27
28    return sentences

```



Extracting sentences to a dataframe

```

1 sentences = []
2 docs = []
3 for i in range(13):
4     sent_doc = split_into_sentences(jobs[i])
5     sentences += sent_doc
6     for k in range(len(sent_doc)): docs.append('Job {}'.format(i + 1))
7

```

```

8 df = pd.DataFrame()
9 df['Document'] = docs
10 df['Sentence'] = sentences
11 df.to_csv('sentences.csv')
12
13 df.head()

```



	Document	Sentence
0	Job 1	job vacancy notice
1	Job 1	software development engineer – time metrology
2	Job 1	international bureau weights measures (bimp)...
3	Job 1	bimp base sèvres , outskirts paris (france) ...
4	Job 1	information find www.bimp.org .

Cleaning sentences

```

1 def clean_sentence(sentence):
2     sentence = re.sub('\.([a-z0-9])', '<dot>\\1', sentence)
3     sentence = ' '.join([word for word in sentence.split() if word not in string.punctu
4     sentence = sentence.replace('<dot>', '.')
5     if '@' in sentence:
6         sentence = sentence.replace('@ ', '@')
7     if ' ' in sentence:
8         sentence = sentence.replace(' ', '\ ')
9
10    return sentence
11
12 df['Sentence'] = df['Sentence'].apply(clean_sentence)
13 df.head()

```



	Document	Sentence
0	Job 1	job vacancy notice
1	Job 1	software development engineer – time metrology
2	Job 1	international bureau weights measures bimp int...
3	Job 1	bimp base sèvres outskirts paris france intern...
4	Job 1	information find www.bimp.org

Tagging important chunks

```

1 skills = list(pd.read_table('skills.txt', sep = ',').columns)
2 job_titles = list(pd.read_csv('job_titles.txt').columns)
3 job_titles = [title.lower() for title in job_titles]

```

```

4 contract_types = ['full-time', 'part-time', 'fixed-term', 'temporary', 'internship']
5 degrees = ['associate degree', "bachelor's degree", "master's degree", 'doctoral degree']
6
7 def tag_words(sentence):
8     URL = '^[^@](((https?):\\/\\/)?(www.))[a-z0-9]+\\.[a-z]{2,})'
9     NUMBER = '\\s([1-9]{1,})'
10    EMAIL = '([a-zA-Z0-9_\\.+~]+@[a-zA-Z0-9-]+\\. [a-zA-Z0-9-\\.]+)'
11
12    #EMAIL
13    sentence = re.sub(EMAIL, '<email>\\1</email>', sentence)
14
15    # URL
16    sentence = re.sub(URL, ' <url>\\1</url>', sentence)
17
18    # JOB TITLES
19    sorted_list = sorted(job_titles, key = len)
20    sorted_list.reverse()
21    for title in sorted_list:
22        if title in sentence:
23            sentence = sentence.replace(title, ' <job_title>{}</job_title>'.format(title))
24            break
25
26    # SKILLS
27    sentence = sentence.split()
28    sentence = list(map(lambda word: ' <skill>{}</skill>'.format(word) if word in skill_list else word, sentence))
29    sentence = ' '.join(sentence)
30
31    # COUNTRIES, CITIES, STATES
32    doc = nlp(sentence)
33    for ent in doc.ents:
34        if ent.label_ == 'GPE':
35            sentence = sentence.replace(ent.text, ' <loc>{}</loc>'.format(ent.text))
36
37    # DATE
38    doc = nlp(sentence)
39    for ent in doc.ents:
40        if (ent.label_ == 'DATE') & ('<' not in ent.text) & ('>' not in ent.text):
41            sentence = sentence.replace(ent.text, ' <date>{}</date>'.format(ent.text))
42
43    # DEGREE
44    for degree in degrees:
45        if degree in sentence:
46            sentence = sentence.replace(degree, ' <degree>{}</degree>'.format(degree))
47            break
48
49    # COMPANY
50    doc = nlp(sentence)
51    for ent in doc.ents:
52        if ent.label_ == 'ORG':
53            sentence = sentence.replace(ent.text, ' <company>{}</company>'.format(ent.text))
54
55    # CONTRACT_TYPE
56    sentence = sentence.split()
57    sentence = list(map(lambda word: ' <contract_type>{}</contract_type>'.format(word) if word in contract_types else word, sentence))
58    sentence = ' '.join(sentence)
59

```

```

59
60     return sentence

```



```

1 df['Sentence'] = df['Sentence'].apply(tag_words)
2 df.head()

```



	Document	Sentence
0	Job 1	job vacancy notice
1	Job 1	<skill>software</skill> development engineer —...
2	Job 1	<company>international bureau weights measures...
3	Job 1	bipm base sèvres outskirts <company>paris fran...
4	Job 1	information find <url>www.bipm.org</url>

Extracting tags

```

1 def extract_tag(tag, sentence):
2     elements = []
3     for element in sentence.split('{>'.format(tag)):
4         if element.endswith('</'):
5             elements.append(element.replace('</', ''))
6     elements = [element for element in elements if element != '']
7     if len(elements) > 0: return set(elements)

```



```

1 df['Email'] = [extract_tag('email', sent) for sent in df['Sentence'].values]
2 df['URL'] = [extract_tag('url', sent) for sent in df['Sentence'].values]
3 df['Job Title'] = [extract_tag('job_title', sent) for sent in df['Sentence'].values]
4 df['Skills'] = [extract_tag('skill', sent) for sent in df['Sentence'].values]
5 df['Location'] = [extract_tag('loc', sent) for sent in df['Sentence'].values]
6 df['Company'] = [extract_tag('company', sent) for sent in df['Sentence'].values]
7 df['Date'] = [extract_tag('date', sent) for sent in df['Sentence'].values]
8 df['Degree'] = [extract_tag('degree', sent) for sent in df['Sentence'].values]
9 df['Contract Type'] = [extract_tag('contract_type', sent) for sent in df['Sentence'].va
10
11 df.head()

```



	Document	Sentence	Email	URL	Job Title	St
0	Job 1	job vacancy notice	None	None	None	

▼ Creating model dataframe

measures...

```
1 df_model = df.copy()
2 df_model.drop(columns = ['Email', 'URL', 'Date', 'Degree', 'Contract Type'], inplace =
3 df_model['Job Title'] = df_model['Job Title'].apply(lambda x: 0 if x is None else 1)
4 df_model['Skills'] = df_model['Skills'].apply(lambda x: 0 if x is None else 1)
5 df_model['Location'] = df_model['Location'].apply(lambda x: 0 if x is None else 1)
6 df_model['Company'] = df_model['Company'].apply(lambda x: 0 if x is None else 1)
7 df_model['Sentence'] = df_model['Sentence'].str.replace('<.*>', '', regex = True)
8
9 df_model.head()
```



	Document	Sentence	Job Title	Skills	Location
0	Job 1	job vacancy notice	0	0	0
1	Job 1	development engineer – time metrology	0	1	0
2	Job 1	whose mandate provide basis coherent system m...	0	0	0
3	Job 1	bipm base sèvres outskirts staff 70	0	0	0
4	Job 1	information find	0	0	0