

```
1 # Import libraries
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 # Data processing
5 import nltk
6 nltk.download('wordnets')
7 nltk.download('punkt')
8 from sklearn.feature_extraction.text import TfidfVectorizer
9 from sklearn.metrics.pairwise import cosine_similarity
10 # File saving
11 import pickle
12 # Used for text extraction
13 from io import StringIO
14 # Text extraction from a PDF
15 from pdfminer.converter import TextConverter
16 from pdfminer.layout import LAParams
17 from pdfminer.pdfdocument import PDFDocument
18 from pdfminer.pdfinterp import PDFPageInterpreter
19 from pdfminer.pdfpage import PDFPage
20 from pdfminer.pdfparser import PDFParser
21 # Open files
22 import urllib
23 # Preparing file for processing
24 # Read to import bytes
25 [nltk_data] Downloading package stopwords to
26 [nltk_data] C:\Users\Samuel\AppData\Local\Temp\nltk_data...
27 [nltk_data] Package stopwords is already up-to-date!
28 [nltk_data] Downloading package punkt to
29 [nltk_data] C:\Users\Samuel\AppData\Local\Temp\nltk_data...
30 [nltk_data] Package punkt is already up-to-date!
31 In [3]: # Extracts and parses text data from submitted document
32 def document_text(text):
33     job_descriptions = pd.read_csv('data/job_descriptions.pkl', index_col=0)
34     with open('pkl/job_descriptions_dec2020.pkl', 'rb') as f:
35         job_descriptions = pickle.load(f)
36     basic_documentdf = pd.DataFrame(data, columns = ['title', 'description'])
37     return basic_documentdf
38 In [4]: # Example of what this function produces
39 document = "https://rookielplay-uploads.s3.amazonaws.com/media/Bell-Resume_Apr2021_INJ3wOtc.pdf"
40 document_text = text(document)
41 document.to_csv('document.csv')
42 Out[4]: SAMUEL LOGAN BELLmedium.com@sambelloup | github.com/sambelloup | linkedin.com/in/sam-bell/sambelloup@gmail.com
43 1 751-652-7433 | Brooklyn, NY 11216Data scientist and machine learning engineer with a penchant for creat
44 ing software with a hunger to work and experience to learn TECHNICAL SKILLS (E3, E2), Rook, C#, Django,
45 b, Heroku, JavaScript, Linux, NumPy, Python, Pandas,scikit-learn, SQL, TableauRELEVANT EXPERIENCERookiePlaycha
46 nce and data analysis• Developing applications for users to upload documents stored in an S3 database• Eng
47 ineering job recommendation system backed by PostgreSQL• Utilizing APIs to funnel in data for live jobs ads
48 and computing growth in a PostgreSQL database• Volunteer Turnout Ratings - https://github.com/sambelloup/turnout_ratings
49 ing tree species growth in a PostgreSQL database• Utilize XGBoost, Decision Trees, and Random Forest
50 to create the best model for prediction• Manipulate features with Principal Component Analysis and then Lasso r
51 eature selection• OTHER EXPERIENCERokey Scouts of AmericaExploring ExecutivePeace CorpsEnglish Volunteer Manag
52 ed website updates for new programs and events• Taught English and computer skills to grades 4 through Middle
53 ONATION Schoolhouse York University, Speech-Language Pathology George Mason University, B.S. PsychologyNew Yo
54 rk, NYAug 2014 - May 2016, Fairfax, VAOct 2017 - Aug 2017Mar 2018 - Aug 2018, New York, NYSep 2017-Dec 2018, New York,
55 NYAug 2017 - May 2010, Fairfax, VAOct
56 In [5]: # Create dataframe of the text and label it resume
57 def compile_document_text(text):
58     job_descriptions = pd.read_csv('data/job_descriptions.pkl', index_col=0)
59     with open('pkl/job_descriptions_dec2020.pkl', 'rb') as f:
60         job_descriptions = pickle.load(f)
61     basic_documentdf = pd.DataFrame(data, columns = ['title', 'description'])
62     return basic_documentdf
63 In [6]: # Example
64 doc_df = compile_document_text(document_text)
65 doc_df.to_csv('document.csv')
66 Out[6]:
67 title description
68 0 resume SAMUEL LOGAN BELLmedium.com@sambelloup | git...
69 In [7]: # RAKE algorithm to determine key phrases in a body of text
70 # by analyzing the frequency of word appearance
71 # and its co-occurrence with other words in the text.
72 # RAKE outputs a list, the words are joined into a single string
73 def text_to_bagofwords(df):
74     df['rake_key_words'] = ''
75     for index, row in df.iterrows():
76         r = Rake()
77         r.extract_keywords_from_text(row['description'])
78         key_words_dict_scores = r.get_word_scores()
79         rake_key_words = list(key_words_dict_scores.keys())
80     # Transform key words into bag of words
81     df['bag_of_words'] = ''
82     for index, row in df.iterrows():
83         words = ''
84         row['bag_of_words'] = words
85         verbose_documentdf = df
86         return verbose_documentdf
87 In [8]: # Example
88 bowdf = text_to_bagofwords(doc_df)
89 bowdf.to_csv('bagofwords.csv')
90 Out[8]:
91 title description rake_key_words bag_of_words
92 0 resume SAMUEL LOGAN BELLmedium.com@sambelloup | influential_preditcor_taloring_data... influential_preditcor_taloring_data... de...
93 In [9]: # Combine resume keywords with the pre-processed job description data
94 def join_and_condense(df):
95     with open('pkl/job_descriptions_dec2020.pkl', 'rb') as f:
96         job_descriptions = pickle.load(f)
97     job_descriptions = job_descriptions.append(df)
98     recommend_df = job_descriptions[['title', 'bag_of_words']]
99     recommend_df['recommend_df_rank_index'] = recommend_df['bag_of_words'].rank(ascending=False)
100     return recommend_df
101 In [10]: # Example
102 fullidf = join_and_condense(bowdf)
103 fullidf.to_csv('fullidf.csv')
104 Out[10]:
105 title bag_of_words
106 0 training_manager provide limited express interest additional tr...
107 1 training_manager workshops including schedules integrity respon...
108 2 training_manager oversee translate diverse advance excel...
109 3 training_manager public speaking - enjoy getting involved exper...
110 4 training_manager luxury business pr teams powerpoint paced emi...
111 - - -
112 Data Entry Specialist enter heavy volume northernmost us entry pre...
113 Pathway Advisor students choose guided academic pathway debt c...
114 Director of Food and Beverage 3 years jumped around forbes five star service...
115 Director of Operations improvements post meeting follow 1 administ...
116 resume influential_preditcor_taloring_data models de...
117 566 rows x 2 columns
118 In [11]: # Returns a matrix of values that represents how closely each document
119 # matches to each other on a scale of 0 to 1
120 def vectorize_text(df):
121     count = CountVecorizer()
122     count_matrix = count.fit_transform(df['bag_of_words'])
123     idf = TfidfTransformer()
124     tfidf_matrix = idf.fit_transform(count_matrix)
125     return cosine_sim
126 In [12]: # Example
127 sim = vectorize_text(fullidf)
128 sim.to_csv('sim.csv')
129 Out[12]: array([[0.05136845, 0.05346523, 0.08437563, 0.08924484, 0.04750557,
130 0.02939465, 0.05192616, 0.06359663, 0.07882312, 0.07639692, 0.06306039,
131 0.0514302, 0.06307053, 0.0649312, 0.07686972, 0.07193194,
132 0.06026831, 0.06847265, 0.06528984, 0.08248703, 0.06078019,
133 0.06349691, 0.06384693, 0.06268244, 0.06268244, 0.06268244,
134 0.0405917, 0.03338903, 0.0337621, 0.04549692, 0.042661,
135 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
136 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
137 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
138 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
139 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
140 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
141 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
142 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
143 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
144 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
145 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
146 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
147 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
148 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
149 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
150 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
151 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
152 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
153 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
154 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
155 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
156 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
157 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
158 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
159 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
160 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
161 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
162 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
163 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
164 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
165 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
166 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
167 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
168 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
169 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
170 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
171 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
172 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
173 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
174 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
175 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
176 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
177 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
178 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
179 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
180 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
181 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
182 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
183 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
184 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
185 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
186 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
187 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
188 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
189 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
190 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
191 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
192 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
193 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
194 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
195 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
196 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
197 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
198 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
199 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
200 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
201 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
202 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
203 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
204 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
205 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
206 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
207 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
208 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
209 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
210 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
211 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
212 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
213 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
214 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
215 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
216 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
217 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
218 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
219 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
220 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
221 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
222 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
223 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
224 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
225 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
226 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
227 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
228 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
229 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
230 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
231 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
232 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
233 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
234 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
235 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
236 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
237 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
238 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
239 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
240 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
241 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
242 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
243 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
244 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
245 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
246 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
247 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
248 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
249 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
250 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
251 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
252 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
253 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
254 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
255 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
256 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
257 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
258 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
259 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
260 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
261 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
262 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
263 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
264 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
265 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
266 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
267 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
268 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
269 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
270 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
271 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
272 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
273 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
274 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
275 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
276 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
277 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
278 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
279 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
280 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
281 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
282 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
283 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
284 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
285 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
286 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
287 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
288 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
289 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
290 0.07281639, 0.06716127, 0.05646931, 0.05148525, 0.05237976,
291 0.07281639, 0.06716127, 0.05646931, 0.05148525
```



Out[39]: ['Research Data Scientist',  
'Ecommerce Category Manager',  
'Actuary Analyst',  
'Program Coordinator',  
'Data Science Internship',  
'Research Associate/Associate Scientist- Immunology',  
'Controls, Modeling, And Simulation Engineer',  
'Remote Mortgage Underwriter',  
'Training Manager',  
'Motion Graphics Designer',  
'Profit Growth Management Business Partner',  
'Manager, User Experience',  
'Fugusol Manager',  
'Scientist Ii, Discovery Science',  
'Pac-3 Systems Engineer - Reliability',  
'Associate I Corporate Finance',  
'Quality Engineer',  
'Foreign Language Instructor',  
'HR Manager',  
'Data Analyst/Report Writer',  
'Project Analyst',  
'Project Manager',  
'Supply Chain Business Analyst',  
'Area Manager',  
'Lead Software Engineer',  
'Accounting Clerk',  
'Supply Chain Analyst',  
'Sales Rep',  
'Admissions Coordinator',  
'Accounts Payable Specialist',  
'Data Scientist',  
'Marketing Operations Manager',  
'Principal Software Engineer',  
'Consulting Analyst',  
'Document Management Analyst',  
'Molding Process Engineer',  
'Data Analyst',  
'Health Program Manager',  
'Business Analyst',  
'Account Manager',  
'Reag Project Analyst - Real Estate Services Group',  
'Product Designer',  
'Business Analyst - Financial Aid Experience',  
'Mechanical Engineering Visiting Scientist',  
'Customer Onboarding Specialist',  
'Mechanical Engineer',  
'Leadership Development Coach',  
'Quality Control Engineer',  
'Budget Analyst',  
'Scientist',  
'Healthcare Access And Value Policy Analyst',  
'Financial Planning Analyst',  
'Executive Director',  
'Director Of Nursing',  
'Customer Service Representative',  
'Seasonal Sales & Service Professional ',  
'Analyst Finance I Us',  
'Associate Portfolio Manager',  
'Missile Assembly / Test Engineer Associate',  
'Electrical Engineer',  
'Population Health Data Analyst',  
'Communications Specialist',  
'Project Management Intern',  
'Manufacturing Planner',  
'Sr Population Health Analyst',  
'Gameplay Engineer - C++',  
'Mechanical Engineering Assistant Professor',  
'Presentation Designer',  
'Application Analyst',  
'Gallery Client Services Associate',  
'Remote Product Design Co-op',  
'Aerospace Engineer']

In [ ]: