

Uncovering Key Drivers in Selecting Professionals for Global Software Development and the Gig Economy

Erivelton Antonio dos Santos , Dalton Garcia Borges de Souza , and Carlos Eduardo Sanches da Silva 

This document contains essential **Supporting Information** concerning the **Algorithms** presents in the paper.

APPENDIX A CCSS ALGORITHM

This Appendix A presents the criteria clusters' semantic similarity (CCSS) algorithm1 used in Python framework.

APPENDIX B ALGORITHM FOR THE 3D BUBBLE CHART

This Appendix B presents the algorithm 2, an interactive 3D hierarchical clustering graph illustrating the hierarchical structure of CC.

Algorithm 1 CCSS between the clusters

```

1  # Data Structures
2  import numpy as np
3  import pandas as pd
4
5  # Keyword extraction
6  from sklearn.feature_extraction.text import CountVectorizer
7  from sentence_transformers import SentenceTransformer
8  from sklearn.metrics.pairwise import cosine_similarity
9
10 # Loading the data
11 df = pd.read_excel('Clusters definitions_vf01.xlsx')
12 df.head()
13
14 # Importing the clustering algorithms - initialize our model and tokenizer
15 model_name = 'sentence-transformers/all-mpnet-base-v2'
16 from transformers import AutoTokenizer, AutoModel
17 import torch
18 tokenizer = AutoTokenizer.from_pretrained(model_name)
19 model = AutoModel.from_pretrained(model_name)
20
21 # Declaring the variables - tokenize the sentences
22 sentences = df['Cluster definition']
23 sentences
24 sentences_list = list(df['Cluster definition'])
25 sentences_list
26 tokens = {'input_ids': [], 'attention_mask': []}
27 for sentence in sentences_list:
28     new_tokens = tokenizer.encode_plus(sentence, max_length=384,
29                                       truncation=True, padding='max_length', return_tensors='pt',
30                                       return_attention_mask=True)
31
32     tokens['input_ids'].append(new_tokens['input_ids'][0])
33     tokens['attention_mask'].append(new_tokens['attention_mask'][0])
34
35 tokens['input_ids']
36
37 # reformat list of tensors into single tensor
38 tokens['input_ids'] = torch.stack(tokens['input_ids'])
39 tokens['attention_mask'] = torch.stack(tokens['attention_mask'])
40
41 # Checking the variables
42 tokens['input_ids']
43 type(tokens['input_ids'])
44 tokens['input_ids'].shape
45
46 # Making the operations - Processing these tokens through our model
47 outputs = model(**tokens)
48 outputs.keys()
49
50 # The dense vector declarations of our text are contained within the outputs 'last_hidden_state' tensor
51 embeddings = outputs.last_hidden_state
52 embeddings
53 embeddings.shape
54
55 # Resize our attention_mask tensor
56 attention = tokens['attention_mask']
57 attention.shape
58 mask = attention.unsqueeze(-1).expand(embeddings.shape).float()
59 mask
60
61 # Multiply the two tensors to apply the attention masks
62 mask_embeddings = embeddings * mask
63 mask_embeddings
64 mask_embeddings.shape
65
66 # Then we sum the remained of the embeddings along axis 1
67 summed = torch.sum(mask_embeddings, 1)
68 summed.shape
69 summed
70
71 # Sum the number of values that must be given attention in each position of the tensor
72 counts = torch.clamp(mask.sum(1), min=1e-9)
73 counts.shape
74 counts
75
76 # Calculate the mean as the sum of the embedding activation's summed divided by the number of values that should be given attention in each position counts
77 mean_pooled = summed / counts
78 mean_pooled.shape
79 mean_pooled
80
81 # The final operations - calculate the cosine similarity between the vectors
82 from sklearn.metrics.pairwise import cosine_similarity
83 mean_pooled = mean_pooled.detach().numpy()
84 data_25g = cosine_similarity(
85     [mean_pooled[0]],
86     [mean_pooled[1 : ]
87     ]
88 )
89 data_25g # data_25g is the final similarities matrix.
90
91 # Print to spreadsheet
92 data_25gT = data_25g.T
93 dfdata_25gT = pd.DataFrame(data_25gT)
94 dfdata_25gT.to_excel("saida.xlsx", index=False)

```

Algorithm 2 Algorithm for the 3D bubble chart

```

1 import plotly.express as px
2 import numpy as np
3 import pandas as pd

```

```

1 df = pd.read_excel('file_3D.xlsx')
2 df.head()

```

	Cluster	Description	Ri + Di	Ri - Di	Cosine Sim	W_cluster
0	COMMUN	Communication	5.42	0.43	1.000000	1
1	INTERF	Component interface	4.69	0.16	0.333691	2
2	SCIENT	Scientific attitude	3.91	0.25	0.332624	2
3	METRIC	Quality metrics	3.69	0.04	0.213171	2
4	GREENSO	Green software development	3.69	0.01	0.212376	2

```

1 fig = px.scatter_3d(df, x='Ri - Di', y='Ri + Di', z='Cosine Sim', size='Ri + Di', color='W_cluster',
2                     hover_data=['Description'])
3 fig.update_layout(scene_zaxis_type="log")
4 fig.show()

```

