





```
[112]: def create_feature(df, cpc_codes_df):
df['anchor_len'] = df['anchor'].apply(lambda x: len(x.split(' ')))
df['target_len'] = df['target'].apply(lambda x: len(x.split(' ')))

pattern = '[0-9]'
mask = df['anchor'].str.contains(pattern, na=False)
df['num_anchor'] = mask
mask = df['target'].str.contains(pattern, na=False)
df['num_target'] = mask

df['context_desc'] = df['context'].map(cpc_codes_df.set_index('code')['title']).str.lower()

fuzzy_anchor_target_scores = []
fuzzy_anchor_context_scores = []
fuzzy_target_context_scores = []
for index, row in df.iterrows():
    fuzzy_anchor_target_scores.append(fuzz.ratio(row['anchor'], row['target']))
    fuzzy_anchor_context_scores.append(fuzz.ratio(row['anchor'], row['context_desc']))
    fuzzy_target_context_scores.append(fuzz.ratio(row['context_desc'], row['target']))
df['fuzzy_at_score'] = fuzzy_anchor_target_scores
df['fuzzy_ac_score'] = fuzzy_anchor_context_scores
df['fuzzy_tc_score'] = fuzzy_target_context_scores

df['fuzzy_c_score'] = df['fuzzy_at_score'] + df['fuzzy_tc_score']
df['fuzzy_total'] = df['fuzzy_at_score'] + df['fuzzy_c_score']
df['fuzzy_avg_at'] = df['fuzzy_at_score']/df['fuzzy_total']
df['fuzzy_avg_at'] = df['fuzzy_c_score']/df['fuzzy_total']

#df.drop(['fuzzy_ac_score', 'fuzzy_tc_score'], 1, inplace=True)

return df

In [113]: df = create_feature(df.copy(), cpc_codes_df)
df.head()
```

		id	anchor	target	context	score	anchor_len	target_len	section	classes	context_desc	...	spacy_total	spac
0	37d61fd22727659b1	abat	abat pollut	A47	0.50	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.000000	
1	7b9652b7b68b7a4	abat	act abat	A47	0.75	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.388094	
2	36d72442aefd8232	abat	activ catalyst	A47	0.25	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.085102	
3	5296b0c19e1ce60e	abat	elimin process	A47	0.50	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.427271	
4	54c1e3b9184cb5b6	abat	forest region	A47	0.00	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.415349	

5 rows x 40 columns

```
In [114]: all_words = set()

This function is not meant to return anything but to get the words into the set very quickly using a lambda function
```

```
In [115]: def separator_all_words(arr):
all_words.add(arr)
```

```
In [116]: df.head()
```

		id	anchor	target	context	score	anchor_len	target_len	section	classes	context_desc	...	spacy_total	spac
0	37d61fd22727659b1	abat	abat pollut	A47	0.50	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.000000	
1	7b9652b7b68b7a4	abat	act abat	A47	0.75	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.388094	
2	36d72442aefd8232	abat	activ catalyst	A47	0.25	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.085102	
3	5296b0c19e1ce60e	abat	elimin process	A47	0.50	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.427271	
4	54c1e3b9184cb5b6	abat	forest region	A47	0.00	1	2	2	A	47	furniture; domestic articles or appliances; co...	...	0.415349	

5 rows x 40 columns

```
In [117]: df(df.columns[15:])
```

		sent_c_score	sent_total	sent_avg_at	sent_avg_c	bert_at	bert_ac	bert_tc	bert_c_score	bert_total	bert_avg_at
0	0.000000	0.000000	NaN	NaN	0.502910	-0.295080	-0.307132	-0.602212	-0.099302	-5.064444	
1	0.232863	0.232863	0.000000	1.000000	0.366209	-0.295080	-0.299380	-0.594460	-0.228251	-1.604412	
2	0.356796	0.356796	0.000000	1.000000	-0.084905	-0.295080	-0.351655	-0.646735	-0.731640	0.160448	
3	0.321934	0.321934	0.000000	1.000000	0.019158	-0.295080	-0.254250	-0.549330	-0.530172	-0.036136	
4	0.366135	0.366135	0.000000	1.000000	0.172460	-0.295080	-0.083622	-0.378702	-0.206242	-0.836199	
...	...	...	...	...	...	...	...	...	...	...	
36468	0.793261	1.321759	0.399844	0.600156	0.815021	0.068167	0.037802	0.105969	0.920990	0.884940	
36469	0.769052	1.301905	0.409287	0.590713	0.291548	0.068167	0.067315	0.135482	0.427030	0.682733	
36470	0.725379	1.219451	0.405160	0.594840	0.502922	0.068167	-0.000069	0.068098	0.571020	0.880744	
36471	0.793261	1.321759	0.399844	0.600156	0.234150	0.068167	-0.012308	0.055859	0.290009	0.807389	
36472	0.793261	1.321759	0.399844	0.600156	0.179098	0.068167	-0.044824	0.023343	0.202441	0.884691	

36473 rows x 25 columns

```
In [118]: df['anchor'].apply(lambda x: separator_all_words(x))
df['target'].apply(lambda x: separator_all_words(x))
df['context_desc'].apply(lambda x: separator_all_words(x))
df['section'].apply(lambda x: separator_all_words(x))
# df['context'].apply(lambda x: separator_all_words(x))

Out[118]: 0      None
1      None
2      None
3      None
4      None
...
36468  None
36469  None
36470  None
36471  None
36472  None
Name: section, Length: 36473, dtype: object

List comprehension of converting a set to a list
```

```
In [119]: entire = [i for i in all_words]
```

Label Encoding

```
In [120]: labelencoder = LabelEncoder()
labelencoder.fit(entire)
```

```
Out[120]: LabelEncoder()
```

```
In [121]: df.to_csv('without_encoding_data_file.csv',index=None)
```

For loop runs a lot faster than the lambda function idk why

```
In [122]: for i in tqdm(np.unique(df['anchor'])):
df['anchor'].replace(i,labelencoder.transform([i])[0],inplace=True)

0% | | 0/733 [00:00<?, ?it/s]
```

```
In [123]: for i in tqdm(np.unique(df['target'])):
df['target'].replace(i,labelencoder.transaform([i])[0],inplace=True)

0% | | 0/26589 [00:00<?, ?it/s]
```

```
In [124]: for i in tqdm(np.unique(df['section'])):
df['section'].replace(i,labelencoder.transform([i])[0],inplace=True)

0% | | 0/8 [00:00<?, ?it/s]
```

```
In [125]: for i in tqdm(np.unique(df['context_desc'])):
df['context_desc'].replace(i,labelencoder.transform([i])[0],inplace=True)

0% | | 0/106 [00:00<?, ?it/s]
```

```
In [126]: df.head()
```

		id	anchor	target	context	score	anchor_len	target_len	section	classes	context_desc	...	spacy_total	spac
0	37d61fd22727659b1	44	46	A47	0.50	1	2	36	47	10103	...	0.000000		
1	7b9652b7b68b7a4	44	257	A47	0.75	1	2	36	47	10103	...	0.388094		
2	36d72442aefd8232	44	276	A47	0.25	1	2	36	47	10103	...	0.085102		
3	5296b0c19e1ce60e	44	8048	A47	0.50	1	2	36	47	10103	...	0.427271		
4	54c1e3b9184cb5b6	44	9734	A47	0.00	1	2	36	47	10103	...	0.415349		

5 rows x 40 columns

```
Out[126]: df
```

		id	anchor	target	context	score	anchor_len	target_len	section	classes	context_desc	...	spacy_total	
0	37d61fd22727659b1	44	46	A47	0.50	1	2	36	47	10103	...	0.000000		
1	7b9652b7b68b7a4	44	257	A47	0.75	1	2	36	47	10103	...	0.388094		
2	36d72442aefd8232	44	276	A47	0.25	1	2	36	47	10103	...	0.085102		
3	5296b0c19e1ce60e	44	8048	A47	0.50	1	2	36	47	10103	...	0.427271		
4	54c1e3b9184cb5b6	44	9734	A47	0.00	1	2	36	47	10103	...	0.415349		
...	...	...	...	...	...	...	...	...	...	...	...	...	...	
36468	8e1386cbefd7245	26847	26860	B44	1.00	2	2	37	44	6191	...	1.779534		
36469	42d9e032d1cd3242	26847	26861	B44	0.50	2	2	37	44	6191	...	1.673599		
36470	208654cc09e14fa3	26847	26865	B44	0.50	2	2	37	44	6191	...	1.483803		
36471	758ec035e69f722b	26847	26868	B44	0.75	2	2	37	44	6191	...	1.779534		
36472	8d135d0c65b06c88	26847	26872	B44	0.50	2	2	37	44	6191	...	1.779534		

36473 rows x 40 columns

```
In [128]: df.to_csv('final_data_file_with_encoding.csv',index=None)
```