

```
In [129]: 1 ## importing all necessary packages and libraries
2 import pandas as pd
3 import re
4 import gensim
5 from gensim.utils import simple_preprocess
6 import nltk
7 from nltk.corpus import stopwords
8 from wordcloud import WordCloud
9 import pyLDAvis.gensim
10 import os
11 import pickle
12 import pyLDAvis
13 import warnings
14 warnings.filterwarnings('ignore')
15 import csv
16 from sklearn.metrics.pairwise import cosine_similarity
17 import gensim.corpora as corpora
18 from itertools import chain
19 from collections import defaultdict
20 from scipy.stats import percentileofscore
21 from gensim.parsing.preprocessing import STOPWORDS
22 from nltk.stem import WordNetLemmatizer, SnowballStemmer
23 from nltk.stem.porter import PorterStemmer
24 from nltk.corpus import wordnet
25 import numpy as np
26 np.random.seed(42)
27 from pprint import pprint
28 nltk.download('wordnet')
29 nltk.download('omw-1.4')
30 from gensim.models import CoherenceModel
31 from gensim.corpora import Dictionary
32 from gensim.models import LdaModel
33 from sklearn.feature_extraction.text import TfidfVectorizer
34 from gensim.models import TfidfModel
35 from sklearn.metrics.pairwise import pairwise_distances
36 from sklearn.metrics import pairwise_distances
```

```
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\Manu\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package omw-1.4 to
[nltk_data] C:\Users\Manu\AppData\Roaming\nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!
```

```
In [89]: 1 ## Reading the dataframe into the master dataframe.
2 # File name - 'Contact Center Ticket Description.csv' should be present in the working directory for the code to fetch the data
3 masData = pd.read_csv('Contact Center Ticket Description.csv', sep='|', header = 0, dtype = {0: str}, names = ['Description',
```

```
In [90]: 1 #Creating a copy of the master dataframe in df.
2 df = masData.copy()
```

```
In [91]: 1 ## Filter for 1 source org and define the initial preprocessing steps.
2 print("Kindly choose from the below list of source orgs for the input variable: \n")
3 unique_source_org = masData['Source Org'].drop_duplicates().tolist()
4
5 print(unique_source_org)
```

Kindly choose from the below list of source orgs for the input variable:

```
['TAA APAC', 'EU GBS Log Cases', 'EU NL BaseCone', 'NA Accounting Services', 'NA AVA', 'NA Canada', 'NA TeamMate', 'TAA NA', 'Heal
R', 'GRC ELM', 'GRC Lien', 'LR US', 'LR Belgium', 'LR Netherlands', 'LR Italy', 'LR Poland']
```

```
In [92]: 1 source_org = input("Enter the source org of choice from the list provided: ")
```

Enter the source org of choice from the list provided: NA Canada

```
In [93]: 1 print(f"Topic modelling will proceed for the source org: {source_org}")
```

Topic modelling will proceed for the source org: NA Canada



```

In [99]: 1 #Lemmatizing words to remove verbs and changing to root word.
2 def lemmatize(text):
3     return WordNetLemmatizer().lemmatize(text, pos='v')
4
5 def preprocess(text_list):
6     result = []
7     for text in text_list:
8         processed_text = []
9         for token in gensim.utils.simple_preprocess(' '.join(text)):
10             if token not in gensim.parsing.preprocessing.STOPWORDS and len(token) >= 3:
11                 processed_text.append(lemmatize(token))
12         result.append(processed_text)
13     return result
14
15 data_words=preprocess(data_words)
16
17 #Lemmatizing words to remove nouns and changing to root word.
18 def lemmatize(text):
19     return WordNetLemmatizer().lemmatize(text, pos='n')
20
21 data_words=preprocess(data_words)
22
23 #Lemmatizing words to remove adjectives and changing to root word.
24 def lemmatize(text):
25     return WordNetLemmatizer().lemmatize(text, pos='a')
26
27 data_words=preprocess(data_words)

In [100]: 1 # Create a dictionary to keep track of word frequencies.
2 frequency = defaultdict(int)
3
4 # Iterate over each document in data_words and updating word frequencies.
5 for text in data_words:
6     for word in text:
7         frequency[word] += 1

In [101]: 1 # Creating a List of dictionaries from the the words frequencies of each word.
2 my_dict = frequency
3 my_list = list(my_dict.items())

In [102]: 1 # Total count of all words in the dataset.
2 total_count = sum(count for word, count in my_list)
3
4 # Creating a List of percentages of frequency of each word.
5 word_freq = [(word, round(count/total_count*100, 2)) for word, count in my_list]
6
7 # Calculate percentile ranks of the frequencies.
8 percentiles = [round(percentileofscore([freq for word, freq in word_freq], freq), 2) for word, freq in word_freq]
9
10 # Create a List of tuples with the word, frequency, and percentile rank.
11 word_freq_percentile = [(word, freq, percentile) for (word, freq), percentile in zip(word_freq, percentiles)]
12
13 # Sort the List by percentile rank in descending order.
14 word_freq_percentile_sorted = sorted(word_freq_percentile, key=lambda x: x[2], reverse=True)

In [103]: 1 # Creating a stop word List with words below the 80th Percentile threshold.
2 percentile_stop_words = [word for word, freq, percentile in word_freq_percentile_sorted if percentile < 80]

In [104]: 1 #Retaining data_words before removing percentile based stop words.
2 d1 = data_words
3 v1 = len(sum(d1, []))
4
5 #Removing stop words from data_words using percentile based stop words.
6 stop_words.extend(percentile_stop_words)
7 data_words = remove_stopwords(data_words)
8
9 #Retaining data_words after removing percentile based stop words.
10 d2 = data_words
11 v2 = len(sum(d2, []))
12
13 #Calculating percentage of words retained
14 percent_retained = (v2/v1)* 100
15
16 print("Word count before stopword removal is: ", v1)
17 print("Word count after stopword removal is: ", v2)
18 print("Percentage of words retained is: ", percent_retained)

```

Word count before stopword removal is: 175774  
Word count after stopword removal is: 165342  
Percentage of words retained is: 94.06510632971884

```
In [105]: 1 # Creating a corpora dictionary of unique tokens for topic modelling.
2 id2word = corpora.Dictionary(data_words)
3
4 # Storing a copy of data_words in texts field
5 texts = data_words
6
7 # Creating a corpus for topic modelling.
8 corpus = [id2word.doc2bow(text) for text in texts]

In [106]: 1 # Choosing the number of topics to in categorize into issue categories
2 num_topics = 7
3
4 # Giving number of topics,bag of words as inputs at an assigned random state 42 to excute our topic modelling
5 lda_model = gensim.models.LdaMulticore(corpus=corpus,
6                                         id2word=id2word,
7                                         num_topics=num_topics, random_state = 42)
8
9 ## The output of this LDA represents the topic number with the weightage of the words contained within the topic
10 pprint(lda_model.print_topics())
11 doc_lda = lda_model[corpus]
12 if not os.path.exists('./results'):
13     os.makedirs('./results')
```

```
[(0,
  '0.099*"payment" + 0.071*"inquiry" + 0.068*"order" + 0.026*"place" + '
  '0.025*"fulfillment" + 0.022*"pay" + 0.020*"download" + 0.020*"account" + '
  '0.018*"cantax" + 0.018*"invoice"'),
 (1,
  '0.088*"cancellation" + 0.056*"invoice" + 0.050*"inquiry" + 0.036*"order" + '
  '0.028*"update" + 0.026*"email" + 0.025*"payment" + 0.024*"cantax" + '
  '0.019*"support" + 0.019*"cancel"'),
 (2,
  '0.105*"support" + 0.093*"transfer" + 0.046*"email" + 0.032*"order" + '
  '0.030*"send" + 0.027*"inquiry" + 0.020*"tax" + 0.018*"cantax" + '
  '0.017*"taxprep" + 0.015*"change"'),
 (3,
  '0.053*"update" + 0.046*"taxprep" + 0.039*"token" + 0.035*"send" + '
  '0.026*"ifirm" + 0.026*"support" + 0.026*"contact" + 0.021*"invoice" + '
  '0.019*"payment" + 0.018*"order"'),
 (4,
  '0.080*"order" + 0.068*"inquiry" + 0.058*"place" + 0.049*"request" + '
  '0.048*"fulfillment" + 0.034*"inv" + 0.031*"account" + 0.027*"payment" + '
  '0.022*"crv" + 0.020*"number"'),
 (5,
  '0.058*"account" + 0.047*"invoice" + 0.047*"ifirm" + 0.041*"access" + '
  '0.035*"inquiry" + 0.031*"send" + 0.026*"taxprep" + 0.022*"pay" + '
  '0.020*"download" + 0.019*"order"'),
 (6,
  '0.201*"order" + 0.142*"inquiry" + 0.140*"place" + 0.133*"fulfillment" + '
  '0.029*"payment" + 0.027*"ifirm" + 0.017*"plcd" + 0.009*"cancel" + '
  '0.008*"refer" + 0.008*"send"')]
```

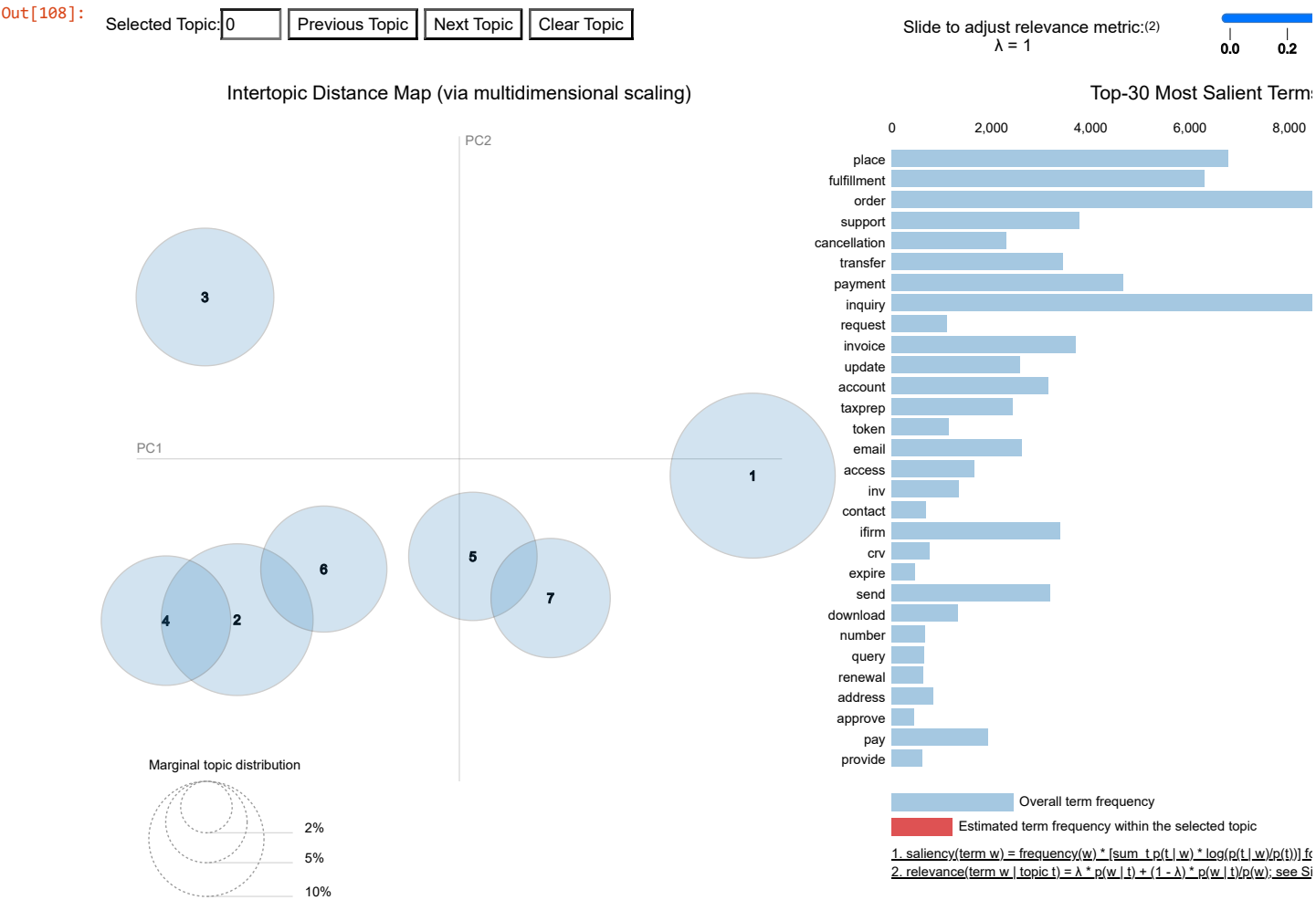
```
In [107]: 1 # Printing topics in a text format for better readability.
2 for topic_idx, topic_terms in lda_model.show_topics(num_topics=num_topics, num_words=40, formatted=False):
3     print("Topic #{}:".format(topic_idx))
4     for term in topic_terms:
5         print(" {} ({:.3f})".format(term[0], term[1]))
6     print("\n")
```

```
research (0.007)
replacement (0.007)
cancel (0.007)
cra (0.006)
create (0.006)
want (0.006)
registration (0.006)
token (0.006)
issue (0.005)
```

```
Topic #1:
cancellation (0.088)
invoice (0.056)
inquiry (0.050)
order (0.036)
update (0.028)
email (0.026)
payment (0.025)
```

In [108]:

```
1 # Initializing pyLDAvis and defining the file path for the prepared data
2 pyLDAvis.enable_notebook()
3 LDAvis_data_filepath = os.path.join('./results/ldavis_prepared_'+str(num_topics))
4
5 # Visualizing the topic modeling results using pyLDAvis
6 if 1 == 1:
7     LDAvis_prepared = pyLDAvis.gensim.prepare(lda_model, corpus, id2word)
8     with open(LDAvis_data_filepath, 'wb') as f:
9         pickle.dump(LDAvis_prepared, f)
10
11 # Saving the HTML visualization
12 with open(LDAvis_data_filepath, 'rb') as f:
13     LDAvis_prepared = pickle.load(f)
14 pyLDAvis.save_html(LDAvis_prepared, './results/ldavis_prepared_'+ str(num_topics) + 'lda.html')
15 LDAvis_prepared
16
17
18 # The above code creates two visualizations.
19
20 # Clusters representing the chosen number of topics in multidimensional scaling, where each access represents principal compo
21
22 # To the right we see top 30 salient terms/words across these topics and the length of the histogram is the representation of
23
24 # When we select one of the topics on the Left, we see the top 30 salient words/terms which represent that particluar topic.
25
26 # The Lambda
```



```
In [109]: 1 #Computing perplexity score
2 print('\nPerplexity: ', lda_model.log_perplexity(corpus))
3
4 # instantiate topic coherence model
5 cm = CoherenceModel(model=lda_model, corpus=corpus, texts=texts, coherence='c_v')
6 # get topic coherence score
7 coherence_lda = cm.get_coherence()
8 print('\nCoherence: ', coherence_lda)
```

Perplexity: -5.336973669946822

Coherence: 0.3485594785242257

```
In [110]: 1 # Initializing LDA using tfidf approach to compute tfidf weight of each word.
2 tfidf = TfidfModel(corpus)
3 corpus_tfidf = tfidf[corpus]
4 lda_model = gensim.models.LdaMulticore(corpus=corpus_tfidf, id2word=id2word, num_topics=num_topics, random_state=42)
5
6 if not os.path.exists('./results'):
7     os.makedirs('./results')
```

```
In [111]: 1 # Printing topics in a text format for better readability for the tfidf approach.
2 for topic_idx, topic_terms in lda_model.show_topics(num_topics=num_topics, num_words=40, formatted=False):
3     print("Topic #{:}:".format(topic_idx))
4     for term in topic_terms:
5         print(" {} ({:.3f})".format(term[0], term[1]))
6     print("\n")
```

Topic #0:  
 payment (0.137)  
 inquiry (0.057)  
 order (0.041)  
 fulfillment (0.027)  
 place (0.026)  
 pay (0.025)  
 reinstate (0.021)  
 invoice (0.018)  
 download (0.017)  
 renewal (0.016)  
 account (0.014)  
 transfer (0.014)  
 cantax (0.013)  
 decline (0.012)  
 step (0.010)  
 monthly (0.009)  
 notify (0.009)  
 ifirm (0.000)

```

In [112]: 1 # Initializing pyLDAvis and defining the file path for the prepared data
2 pyLDAvis.enable_notebook()
3 LDAvis_data_filepath = os.path.join('./results/ldavis_prepared_' + str(num_topics))
4
5 # Visualizing the topic modeling results using pyLDAvis
6 if 1 == 1:
7     LDAvis_prepared = pyLDAvis.gensim.prepare(lda_model, corpus, id2word)
8     with open(LDAvis_data_filepath, 'wb') as f:
9         pickle.dump(LDAvis_prepared, f)
10
11 # Saving the HTML visualization
12 with open(LDAvis_data_filepath, 'rb') as f:
13     LDAvis_prepared = pickle.load(f)
14 pyLDAvis.save_html(LDAvis_prepared, './results/ldavis_prepared_' + str(num_topics) + 'lda.html')
15 LDAvis_prepared

```

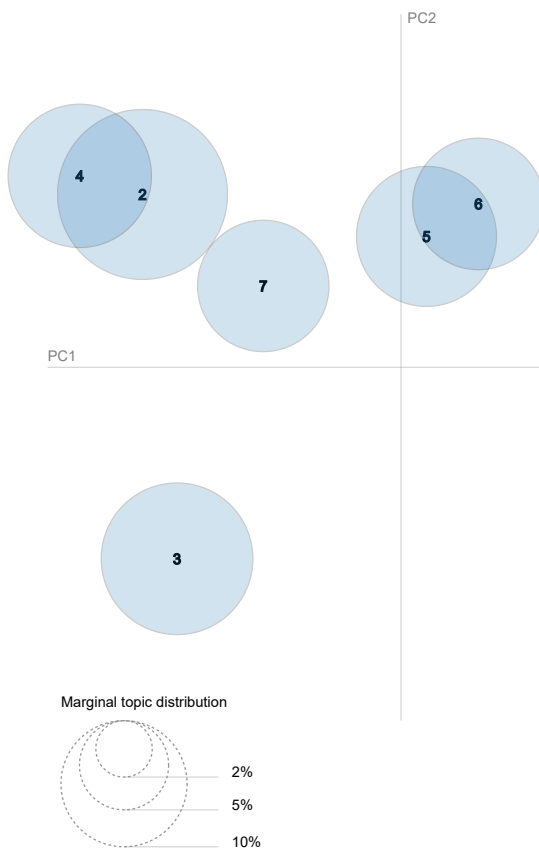
Out[112]:

Selected Topic:

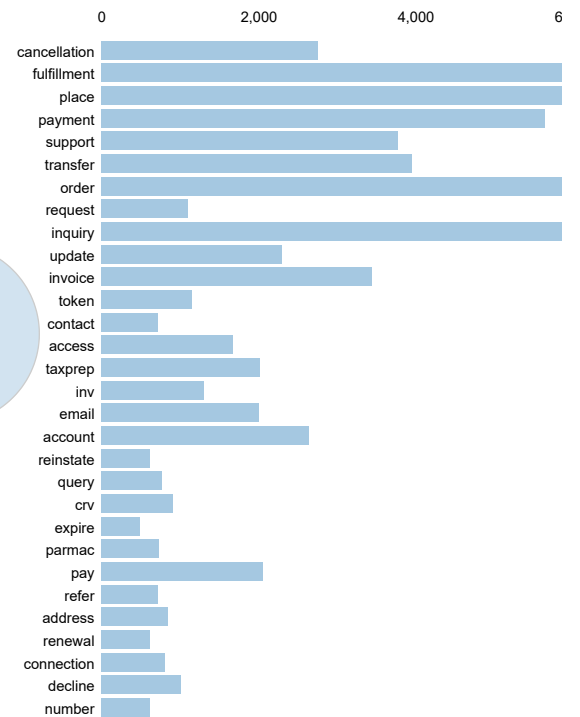
Slide to adjust relevance metric:(2)  
 $\lambda = 1$

0.0 0.2

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Salient Term:



Overall term frequency

Estimated term frequency within the selected topic

1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) \cdot \left[ \sum_t p(t|w) \cdot \log\left(\frac{p(t|w)}{p(t)}\right) \right]$  fr  
 2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda \cdot p(w|t) + (1 - \lambda) \cdot \frac{p(w|t)p(w)}{p(w)}$ ; see Si

```

In [113]: 1 #Computing perplexity score
2 print('\nPerplexity: ', lda_model.log_perplexity(corpus_tfidf))
3
4 from gensim.models import CoherenceModel
5 # instantiate topic coherence model
6 cm = CoherenceModel(model=lda_model, corpus=corpus_tfidf, texts=texts, coherence='c_v')
7 # get topic coherence score
8 coherence_lda = cm.get_coherence()
9 print('\nCoherence: ', coherence_lda)

```

Perplexity: -6.055344199642336

Coherence: 0.38721797477646663

```
In [114]: 1 #Saving output of the topic modelling in text format in output_csv file in the current working directory.
2 with open('output.csv', 'w', newline='') as file:
3     writer = csv.writer(file)
4     writer.writerow(['Source Org', 'Topic Number', 'Word Name', 'Weightage'])
5     for topic_idx, topic_terms in lda_model.show_topics(num_topics=num_topics, num_words=40, formatted=False):
6         for term in topic_terms:
7             writer.writerow([source_org, topic_idx, term[0], term[1]])
```

```
In [115]: 1 # Using TAA Priority key mapping table as lookup table.
2 lookup_table = pd.read_excel('TAA Priority Key Word Mapping.xlsx')
3
4 #Lemmitizing words in the Lookup table by Lemmitizing nouns, verbs and adjectives.
5 def lemmatize_lookup(text):
6     return WordNetLemmatizer().lemmatize(text, pos='v')
7
8 def preprocess_lookup(df, text_column):
9     result = []
10    for text in df[text_column]:
11        if isinstance(text, float):
12            text = str(float(text)) # convert float to int, then to string
13        processed_text = []
14        for token in gensim.utils.simple_preprocess(text):
15            if token not in gensim.parsing.preprocessing.STOPWORDS and len(token) >= 3:
16                processed_text.append(lemmatize_lookup(token))
17        result.append(' '.join(processed_text))
18    df_processed = df.copy()
19    df_processed[text_column] = result
20    return df_processed
21 lookup_table = preprocess_lookup(lookup_table, 'Keyword')
22
23 def lemmatize_lookup(text):
24     return WordNetLemmatizer().lemmatize(text, pos='n')
25 lookup_table = preprocess_lookup(lookup_table, 'Keyword')
26 def lemmatize_lookup(text):
27     return WordNetLemmatizer().lemmatize(text, pos='a')
28 lookup_table = preprocess_lookup(lookup_table, 'Keyword')
```

```
In [116]: 1 #Doing vlookup of words in Lookup table with words in the topic modelling output.csv
2 def lookup(keyword):
3     result = lookup_table.loc[lookup_table['Keyword'] == keyword]
4     if not result.empty:
5         return result.iloc[0]['Issue Category'], result.iloc[0]['Contact Reason']
6     else:
7         return '***No_Match***', '***No_Match***'
8
9 with open('updated_output.csv', 'w', newline='') as file:
10     writer = csv.writer(file)
11     writer.writerow(['Source Org', 'Topic Number', 'Word Name', 'Weightage', 'Issue Category', 'Contact Reason'])
12
13     for topic_idx, topic_terms in lda_model.show_topics(num_topics=num_topics, num_words=40, formatted=False):
14         for term in topic_terms:
15             issue_category, contact_reason = lookup(term[0])
16             writer.writerow([source_org, topic_idx, term[0], term[1], issue_category, contact_reason])
```

```
In [117]: 1 #####
2 ##The below code automatically generates sentences for each topic and gives the percentage contribution of a particular issue
3 #####
```



```

In [118]: 1 with open('updated_output.csv', 'w', newline='') as file:
2         writer = csv.writer(file)
3         writer.writerow(['Source Org', 'Topic Number', 'Word Name', 'Weightage', 'Issue Category', 'Contact Reason'])
4
5         for topic_idx, topic_terms in lda_model.show_topics(num_topics=num_topics, num_words=40, formatted=False):
6             no_match_count = 0
7             admin_count = 0
8             inquiry_product_support_count = 0
9             orders_payments_count = 0
10            sales_sub_ren_count = 0
11            sme_inq_count = 0
12            unclear_count = 0
13            abandoned_count = 0
14            redacted_count = 0
15            for term in topic_terms:
16                issue_category, contact_reason = lookup(term[0])
17                if issue_category == '**No Match**':
18                    no_match_count += 1
19                elif issue_category == 'Admin':
20                    admin_count += 1
21                elif issue_category == 'Orders and Payments':
22                    orders_payments_count += 1
23                elif issue_category == 'Sales / Subscriptions / Renewals':
24                    sales_sub_ren_count += 1
25                elif issue_category == 'SME Inquiries':
26                    sme_inq_count += 1
27                elif issue_category == 'Unclear description':
28                    unclear_count += 1
29                elif issue_category == 'Abandoned':
30                    abandoned_count += 1
31                elif issue_category == 'Redacted':
32                    redacted_count += 1
33                elif issue_category == 'Inquiries and Product Support':
34                    inquiry_product_support_count += 1
35            writer.writerow([source_org, topic_idx, term[0], term[1], issue_category, contact_reason])
36
37            total_count = no_match_count + admin_count + inquiry_product_support_count + orders_payments_count + sales_sub_ren_co
38            no_match_percent = round((no_match_count / total_count) * 100, 2)
39            admin_percent = round((admin_count / total_count) * 100, 2)
40            inquiry_product_support_percent = round((inquiry_product_support_count / total_count) * 100, 2)
41            orders_payments_percent = round((orders_payments_count / total_count) * 100, 2)
42            sales_sub_ren_percent = round((sales_sub_ren_count / total_count) * 100, 2)
43            sme_inq_percent = round((sme_inq_count / total_count) * 100, 2)
44            unclear_percent = round((unclear_count / total_count) * 100, 2)
45            abandoned_percent = round((abandoned_count / total_count) * 100, 2)
46            redacted_percent = round((redacted_count / total_count) * 100, 2)
47
48            sentence_parts = []
49            if no_match_percent > 0:
50                sentence_parts.append(f"{no_match_percent}% words which are **No Match**")
51            if admin_percent > 0:
52                sentence_parts.append(f"{admin_percent}% words which belong to Admin")
53            if inquiry_product_support_percent > 0:
54                sentence_parts.append(f"{inquiry_product_support_percent}% words which belong to Inquiry and Product Support")
55            if orders_payments_percent > 0:
56                sentence_parts.append(f"{orders_payments_percent}% words which belong to Orders and Payments")
57            if sales_sub_ren_percent > 0:
58                sentence_parts.append(f"{sales_sub_ren_percent}% words which belong to Sales Subscriptions and Renewals")
59            if sme_inq_percent > 0:
60                sentence_parts.append(f"{sme_inq_percent}% words which belong to SME Inquiries")
61            if unclear_percent > 0:
62                sentence_parts.append(f"{unclear_percent}% words which are Unclear")
63            if abandoned_percent > 0:
64                sentence_parts.append(f"{abandoned_percent}% words which are Abandoned")
65            if redacted_percent > 0:
66                sentence_parts.append(f"{redacted_percent}% words which are Redacted")
67
68            sentence = f"Topic {topic_idx} has " + ", ".join(sentence_parts) + "."
69            print(sentence)
70

```

Topic 0 has 57.5% words which are \*\*No Match\*\*, 7.5% words which belong to Admin, 17.5% words which belong to Inquiry and Product and Payments, 5.0% words which belong to Sales Subscriptions and Renewals.  
Topic 1 has 62.5% words which are \*\*No Match\*\*, 12.5% words which belong to Admin, 12.5% words which belong to Inquiry and Product s and Payments, 2.5% words which belong to Sales Subscriptions and Renewals.  
Topic 2 has 65.0% words which are \*\*No Match\*\*, 15.0% words which belong to Admin, 10.0% words which belong to Inquiry and Product s and Payments.  
Topic 3 has 62.5% words which are \*\*No Match\*\*, 10.0% words which belong to Admin, 15.0% words which belong to Inquiry and Product s and Payments, 2.5% words which belong to SME Inquiries.  
Topic 4 has 60.0% words which are \*\*No Match\*\*, 12.5% words which belong to Admin, 15.0% words which belong to Inquiry and Product s and Payments.  
Topic 5 has 47.5% words which are \*\*No Match\*\*, 22.5% words which belong to Admin, 17.5% words which belong to Inquiry and Product s and Payments, 2.5% words which belong to SME Inquiries.  
Topic 6 has 57.5% words which are \*\*No Match\*\*, 12.5% words which belong to Admin, 15.0% words which belong to Inquiry and Product s and Payments, 2.5% words which belong to Sales Subscriptions and Renewals.

In [119]:

```
1 # The below code first removes the basic english stopwords and extended list of stopwords from the priority mapping key table
2 lookup_table.rename(columns={'Issue Category':'Topic Number'},inplace=True)
3 lookup_table=lookup_table.dropna()
4
5 stop_words = stopwords.words('english')
6 stop_words.extend(['wolters', 'kluwer', 'from', 'subject', 're', 'edu', 'use', 'fw', 'Fw', 'FWD', 'fwd', 'hi', 'may', 'know',
7                   'regards', 'iphone', 'android', 'interested', 'australia', 'australian', 'wolterskluwer'])
```

In [120]:

```
1 # Performing pre-processing steps on the updated lookup table.
2 lookup=lookup_table[['Topic Number','Keyword']].drop_duplicates()
3 lookup=lookup.dropna()
4 lookup_sw=lookup[lookup['Keyword'].isin(stop_words)]
5 lookup=lookup[~lookup['Keyword'].isin(stop_words)] ###
6 lookup=lookup.sort_values(by='Keyword',ascending=True)###
7 lookup.head(10)
```

Out[120]:

	Topic Number	Keyword
129	Admin	
126	Orders and Payments	
22	Inquiries and Product Support	
173	SME Inquiries	
219	Abandoned	abandon
214	Abandoned	abandoned
169	Admin	accept
107	Admin	access
109	Admin	account
114	Admin	acct

In [121]:

```
1 # Removing empty keywords from the key mapping table
2 lookup=lookup[lookup['Keyword']!='']
```

In [123]:

```
1 # Creating a pivot table to represent each issue in terms of all the words n the data dictionary.
2 lookup['val']=1
3 lookup_pivot= lookup.pivot_table(index='Topic Number', columns='Keyword', values='val', fill_value=0) ###
4 lookup_pivot=lookup_pivot.reset_index()
5 lookup_pivot.head(10)
```

Out[123]:

Keyword	Topic Number	abandon	abandoned	accept	access	account	acct	ach	add	address	...	troubleshoot	try	unable	update	upload
0	Abandoned	1	1	0	0	0	0	0	0	0	...	0	0	0	0	0
1	Admin	0	0	1	1	1	1	0	1	1	...	0	0	0	0	0
2	Inquiries and Product Support	0	0	0	0	0	0	0	0	0	...	1	1	1	1	1
3	Orders and Payments	0	0	0	0	0	0	1	0	0	...	0	0	0	0	0
4	Redacted	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
5	SME Inquiries	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
6	Sales / Subscriptions / Renewals	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
7	Unclear description	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0

8 rows × 159 columns

```
In [124]: 1 # Retaining issue categories where there was a match in the keywords in the vloop section above.
2 lda_output = pd.read_csv('updated_output.csv')
3 lda_output=lda_output[lda_output['Issue Category']!='***No_Match***']
4 lda_output=lda_output.sort_values(by='Word Name',ascending=True)
5 lda_output=lda_output[['Topic Number','Word Name']].drop_duplicates()
6
7 # Creating a pivot table to represent each topic in terms of all the words n the data dictionary.
8 lda_output['val']=1
9 lda_output_pivot = lda_output.pivot_table(index='Topic Number', columns='Word Name', values='val', fill_value=0)
10 lda_output_pivot=lda_output_pivot.reset_index()
11 lda_output_pivot.head(10)
```

```
Out[124]:
```

	Word Name	Topic Number	access	account	add	address	archive	cch	change	contact	contract	...	renewal	request	reset	return	ship	subscription	sup
	0	0	0	1	0	0	0	1	0	0	0	...	1	0	0	1	0		1
	1	1	0	1	1	1	1	0	0	0	0	...	1	0	0	0	0		0
	2	2	1	1	0	0	1	0	1	0	0	...	0	0	0	0	1		0
	3	3	0	1	0	1	0	0	0	1	0	...	0	0	0	0	0		0
	4	4	1	1	0	1	0	0	0	0	0	...	0	1	0	0	0		0
	5	5	1	1	0	0	1	1	1	0	0	...	0	0	1	0	0		0
	6	6	0	1	1	0	1	1	1	0	1	...	0	1	0	0	0		1

7 rows × 41 columns

```
In [125]: 1 ## The Unions of Keywords - adding missing columns to dataframe lookup_pivot and filling with 0s.
2 lookup_pivot=lookup_pivot.reindex(columns=lookup_pivot.columns.union(lda_output_pivot.columns), fill_value=0)
3 # Adding missing columns to dataframe lda_output_pivot and filling with 0s.
4 lda_output_pivot=lda_output_pivot.reindex(columns=lookup_pivot.columns.union(lda_output_pivot.columns), fill_value=0)
```

```
In [126]: 1 lookup_pivot
```

```
Out[126]:
```

	Topic Number	abandon	abandoned	accept	access	account	acct	ach	add	address	...	troubleshoot	try	unable	update	upload	uptodate
0	Abandoned	1	1	0	0	0	0	0	0	0	...	0	0	0	0	0	
1	Admin	0	0	1	1	1	1	0	1	1	...	0	0	0	0	0	
2	Inquiries and Product Support	0	0	0	0	0	0	0	0	0	...	1	1	1	1	1	
3	Orders and Payments	0	0	0	0	0	0	1	0	0	...	0	0	0	0	0	
4	Redacted	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
5	SME Inquiries	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
6	Sales / Subscriptions / Renewals	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
7	Unclear description	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	

8 rows × 159 columns

```
In [127]: 1 lda_output_pivot
```

```
Out[127]:
```

	Topic Number	abandon	abandoned	accept	access	account	acct	ach	add	address	...	troubleshoot	try	unable	update	upload	uptodate	user	voice
0	0	0	0	0	0	1	0	0	0	0	...	0	0	0	1	0	0	0	
1	1	0	0	0	0	1	0	0	1	1	...	0	0	0	1	0	0	0	
2	2	0	0	0	1	1	0	0	0	0	...	0	0	0	1	0	0	0	
3	3	0	0	0	0	1	0	0	0	1	...	0	0	0	1	0	0	0	
4	4	0	0	0	1	1	0	0	0	1	...	0	0	0	1	0	0	0	
5	5	0	0	0	1	1	0	0	0	0	...	0	0	0	1	0	0	1	
6	6	0	0	0	0	1	0	0	1	0	...	0	0	0	1	0	0	0	

7 rows × 159 columns

```
In [128]: 1 ## It can be seen from the above two pivot tables, that the number of columns is the same.
```

```
In [130]: 1 # Creating a similarity matrix using Jacard Measure to find relevance between topics and issue categories.
2 similarity_ps = 1 - pairwise_distances(lookup_pivot.values, lda_output_pivot.values, metric='jaccard')
3
4 similarity_df = pd.DataFrame(similarity_ps, index=lookup_pivot.index, columns=lda_output_pivot.index)
5
6 similarity_df.index.name = 'issue category'
7 similarity_df.columns.name = 'Topic Number'
8
9 category_dict = {0: 'Abandoned', 1: 'Admin', 2: 'Inquiries and Product Support', 3: 'Orders and Payments', 4: 'Redacted', 5:
10
11 similarity_df.index = similarity_df.index.to_series().replace(category_dict)
12
13 similarity_df
```

Out[130]:

	Topic Number	0	1	2	3	4	5	6
	issue category							
	Abandoned	0.000000	0.052632	0.055556	0.052632	0.050000	0.040000	0.047619
	Admin	0.051724	0.111111	0.134615	0.090909	0.109091	0.178571	0.107143
	Inquiries and Product Support	0.074468	0.063830	0.053191	0.075269	0.074468	0.081633	0.073684
	Orders and Payments	0.178571	0.185185	0.192308	0.185185	0.222222	0.151515	0.214286
	Redacted	0.000000	0.058824	0.062500	0.058824	0.055556	0.043478	0.052632
	SME Inquiries	0.000000	0.038462	0.040000	0.080000	0.037037	0.064516	0.035714
	Sales / Subscriptions / Renewals	0.095238	0.100000	0.050000	0.047619	0.045455	0.037037	0.090909
	Unclear description	0.000000	0.058824	0.062500	0.058824	0.055556	0.043478	0.052632

```
In [131]: 1 # The below code automatically prints the topic number distribution in each issue category based on the top 3 similarity scor
2 for index, row in similarity_df.iterrows():
3     sorted_topics = row.sort_values(ascending=False).index.values
4     top_topic = sorted_topics[0]
5     second_topic = sorted_topics[1]
6     third_topic = sorted_topics[2]
7     print(f"The issue category '{index}' is most similar to the topic number '{top_topic}' '{second_topic}' and '{third_topic}'")
```

The issue category 'Abandoned' is most similar to the topic number '2' '1' and '3'.

The issue category 'Admin' is most similar to the topic number '5' '2' and '1'.

The issue category 'Inquiries and Product Support' is most similar to the topic number '5' '3' and '0'.

The issue category 'Orders and Payments' is most similar to the topic number '4' '6' and '2'.

The issue category 'Redacted' is most similar to the topic number '2' '1' and '3'.

The issue category 'SME Inquiries' is most similar to the topic number '3' '5' and '2'.

The issue category 'Sales / Subscriptions / Renewals' is most similar to the topic number '1' '0' and '6'.

The issue category 'Unclear description' is most similar to the topic number '2' '1' and '3'.

```
In [132]: 1 # The below code automatically prints the issue category distribution in each topic based on the top 3 similarity scores.
2 for col in similarity_df.columns:
3     sorted_categories = similarity_df[col].sort_values(ascending=False).index.values
4     top_category = sorted_categories[0]
5     second_category = sorted_categories[1]
6     third_category = sorted_categories[2]
7     print(f"The topic number '{col}' is most similar to the issue category '{top_category}' '{second_category}' and '{third_category}'")
```

The topic number '0' is most similar to the issue category 'Orders and Payments' 'Sales / Subscriptions / Renewals' and 'Inquiries

The topic number '1' is most similar to the issue category 'Orders and Payments' 'Admin' and 'Sales / Subscriptions / Renewals'.

The topic number '2' is most similar to the issue category 'Orders and Payments' 'Admin' and 'Redacted'.

The topic number '3' is most similar to the issue category 'Orders and Payments' 'Admin' and 'SME Inquiries'.

The topic number '4' is most similar to the issue category 'Orders and Payments' 'Admin' and 'Inquiries and Product Support'.

The topic number '5' is most similar to the issue category 'Admin' 'Orders and Payments' and 'Inquiries and Product Support'.

The topic number '6' is most similar to the issue category 'Orders and Payments' 'Admin' and 'Sales / Subscriptions / Renewals'.

```
In [133]: 1 ##### END OF CODE #####
```

