

Text Analysis Report

Data:

We use data collected from interviews and online resources like Reddit regarding recycling as the central theme. 20 interview texts without the takeaway section, around 110 reddit posts, and text from an article about recycling from online repositories was used. An important omission is a substantial twitter dataset about recycling and sustainability which could not be used due to twitter api being no longer open source.

Data Preprocessing:

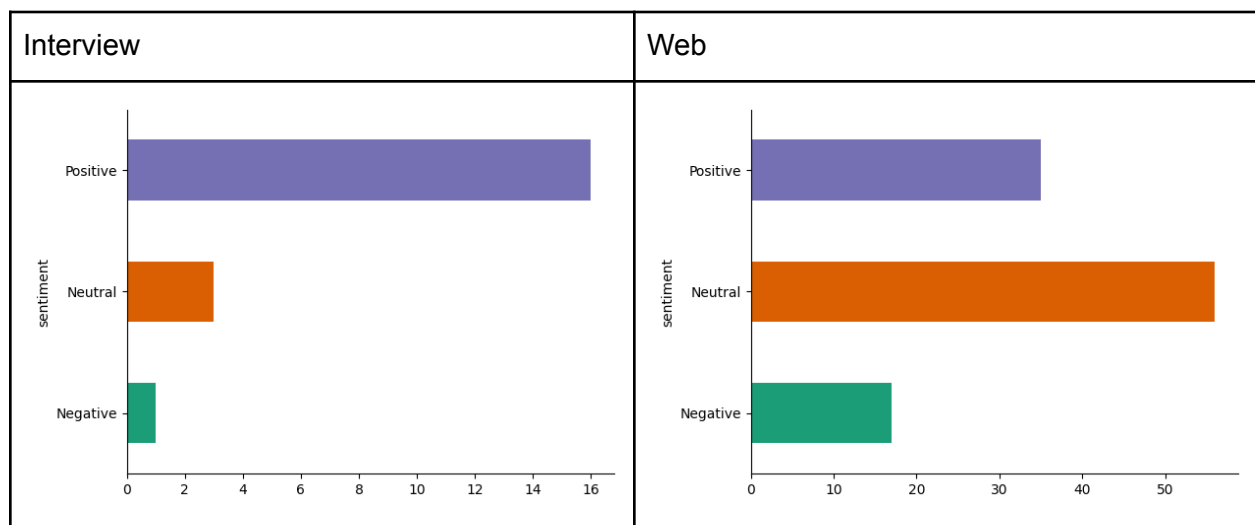
All the data collected is stored as raw text in json format. Interview data and web scrapped data are loaded individually and the following preprocessing steps are performed.

- Removal to URL, hyperlinks, emails, new line characters, quotations
- Conversion to bigram and trigrams
- Lemmatization and POS tagging
- Stop word removal

Sentiment Analysis:

After preprocessing the text data, we use ChatGPT to perform sentiment analysis. Particularly, we use the following prompt to get the sentiment and the associated text which invokes that sentiment:

"You will be provided with a text, and your task is to identify their sentiments about recycling as positive, neutral, or negative and what part of text corresponds to it."



As it can be seen, in case of interviews, there were mostly positive sentiments about recycling and only 1 case of overall negative sentiment. This may be attributed to the fact that we were inquiring about their personal attitudes towards recycling and thus invoking a positive attitude and sense of self awareness. Whereas, in case of web data, the majority of posts have a neutral stance and a substantial amount of posts even have negative sentiment about recycling. This can be attributed to anonymity, since you can express yourself without fear of identification.

a. Positive Sentiments: The following texts from interview reflect positive statement according to ChatGPT

"because she personally thinks that recycling works and cares about the environment"

"Jay is proud of Brooklyn's community-centric approach to recycling and often takes part in local awareness campaigns."

b. Negative Sentiments:

"kaylee says she hates taking out the trash because the trash room is disgusting and smells horrible"

c. Neutral Sentiments:

"honestly with our busy lives recycling is the least of our worries"

"she mentioned her less consistent recycling habits primarily due to confusion regarding labels and the inconvenience associated with the recycling process"

Topic Modeling:

We use gensim based Latent Dirichlet Allocation (LDA) model for topic modeling. Considering that we have restricted ourselves to the theme of recycling, and limited data points, we use 4 topics for interviews and 5 topics for web data to explore. The topics extracted by the LDA model are presented in the table below.

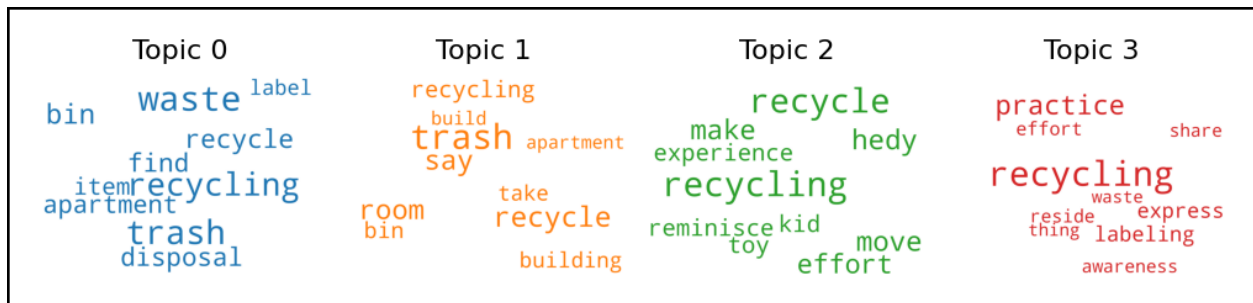
a. Identify Main Topics

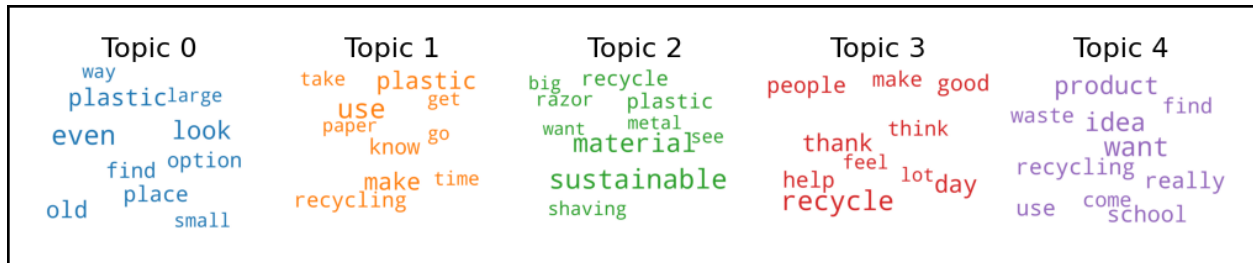
The topics and associated keywords along with the representative sentence (the sentence where contribution of the topic is maximum) are presented in the table. We can see topics like waste disposal bins, trash room in apartment, awareness effort, sustainable materials, plastic recycling.

Topic_Num	Topic_Perc_Contrib	Keywords	Representative Text
0	0	0.9932 waste, trash, recycling, bin, recycle, disposal, find, apartment, item, label	[dispose, trash, use, designate, bin, landfill, recycle, compost, find, apartment, building, sid...
1	1	0.9946 trash, recycle, say, room, recycling, bin, building, take, build, apartment	[recent, graduate, live, year, partner, work, weekend, lot, cleaning, leave, fact, morning, inte...
2	2	0.9906 recycling, recycle, move, hedy, make, effort, kid, toy, experience, reminisce	[move, blend, cultural, experience, influence, perspective, recycle, chat, speak, journey, under...
3	3	0.9893 recycling, practice, labeling, express, effort, reside, awareness, share, thing, waste	[reside, live, year, witness, change, waste, management, practice, recently, family, day, commun...

Topic_Num	Topic_Perc_Contrib	Keywords	Representative Text
0	0	0.9329 even, look, plastic, old, place, option, find, small, way, large	[look, way, tell, unable, blade, attach, also, unable, remove, make, plastic]
1	1	0.9851 use, plastic, make, know, recycling, time, take, go, get, paper	[partner, compost, long, time, recently, move, town, drop, point, available, membership, local, ...
2	2	0.9381 sustainable, material, recycle, plastic, razor, shaving, see, metal, big, want	[energy, require, recycle, limit, amount, want, metal, melt, huge, amount, energy, waste]
3	3	0.9464 recycle, day, thank, people, good, help, make, think, feel, lot	[experience, improved, recycle, make, efficient, stop, recycle, home, public, research, conduct,...
4	4	0.8850 want, idea, product, really, school, use, recycling, waste, find, come	[keen, country, specific, vape, recycling, program]

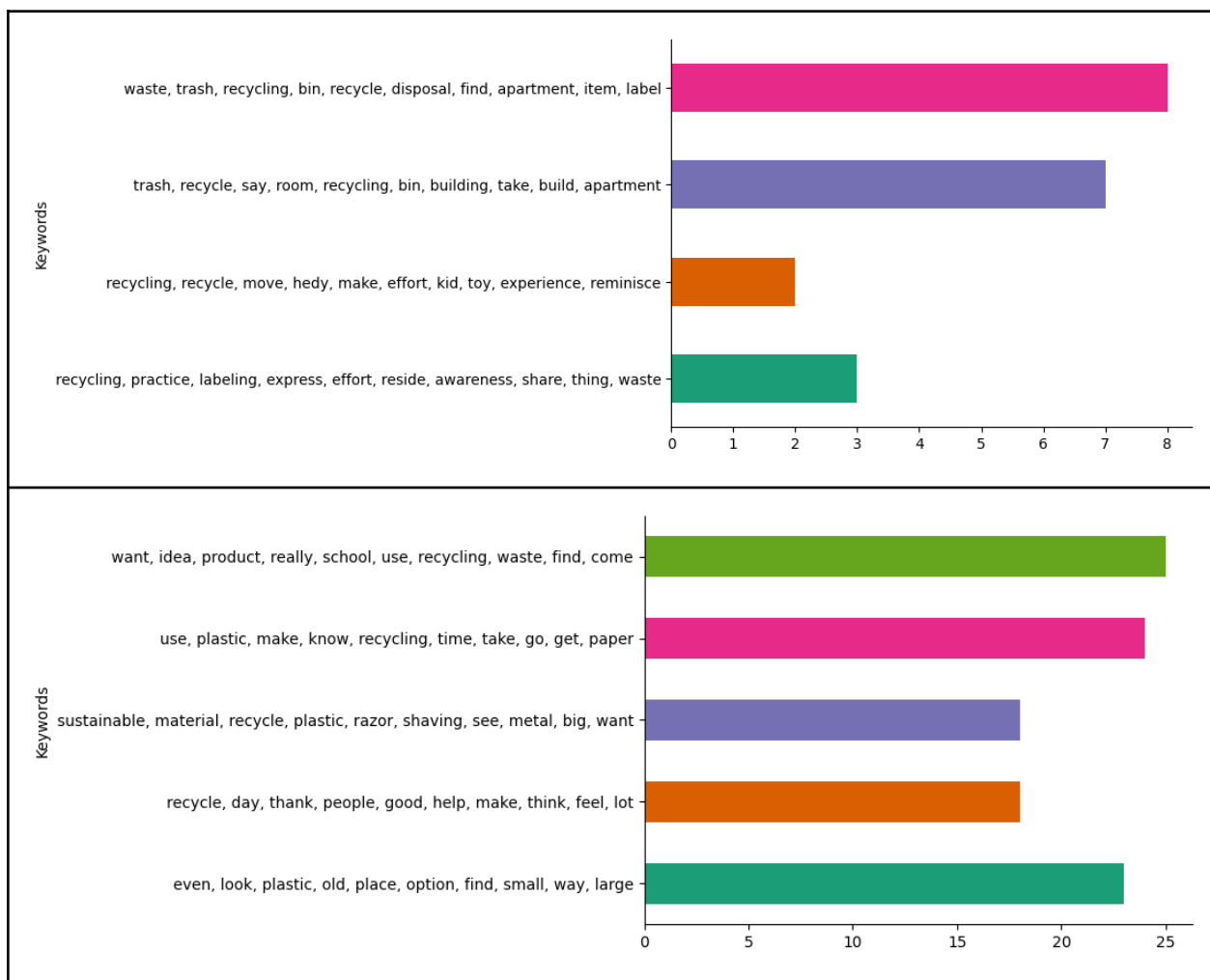
Along with that, the word cloud corresponding to the keywords for each topic is presented below.





b. Topic Prevalence

We present the distribution of topics in the interview and web data as histograms. As it can be seen below, topics like waste disposal, trash rooms in apartments are discussed more while topics like awareness efforts, sustainable materials are discussed less.



c. Topic Relationships

In order to identify relationships among the topics, we use two approaches:

1) Sentence topic coloring: In this approach, we color the entire sentence in its representative topic and color each word in the sentence according to the topic which it most likely belongs. This gives us a qualitative understanding of how topics are distributed among the sentences or document texts. As it can be seen, certain sentences are pure i.e. it mostly contains all words from same topic, whereas other sentences have a mix of different topic words.

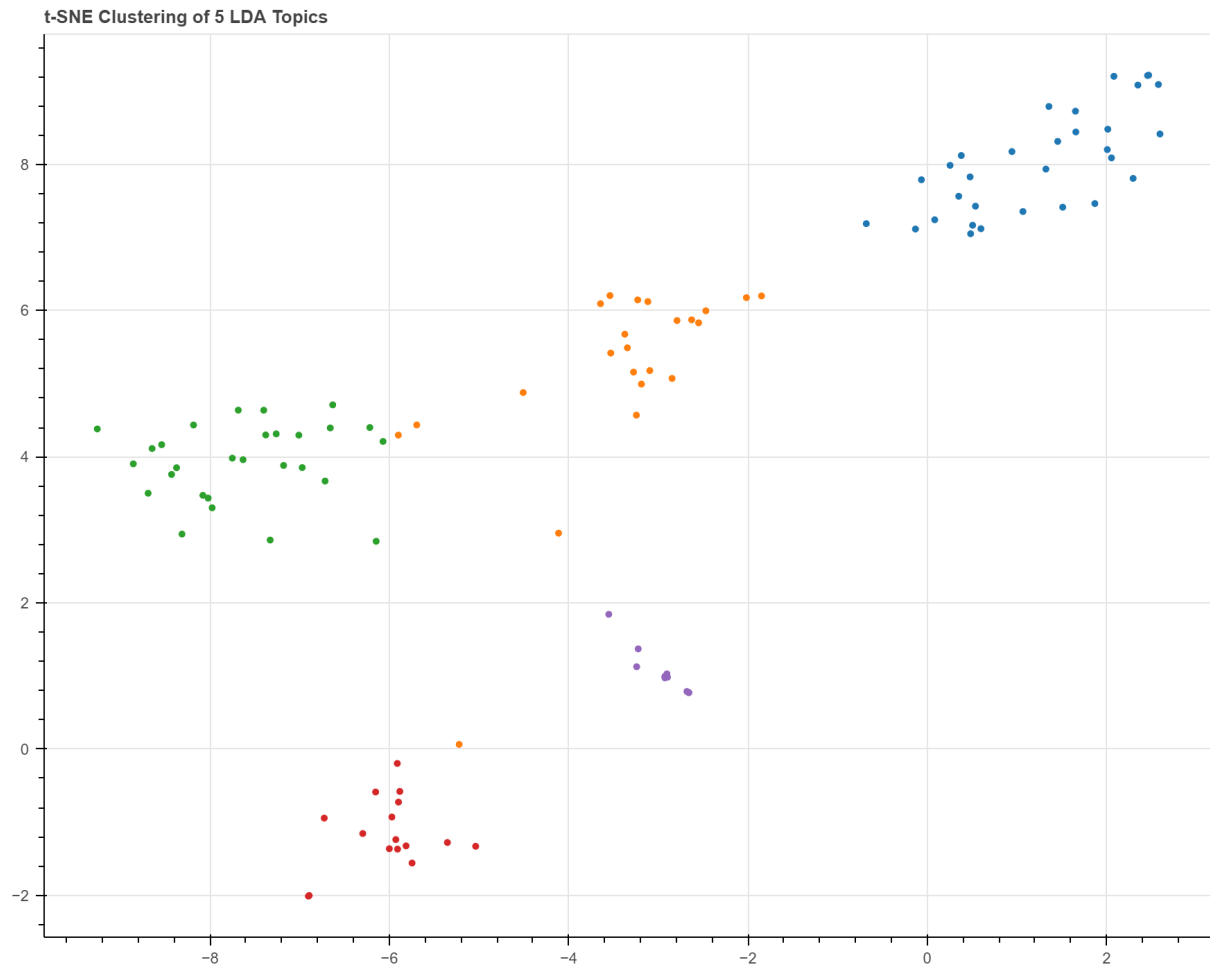
Sentence Topic Coloring for Documents: 0 to 11

Doc 0:	actually	also	apartment	ask away bin board build building cardboard catch cleaning come community ...
Doc 1:	actually	also	apartment	ask bin build building cardboard come feel garbage habit last live ...
Doc 2:	also	ask building	cleaning	community different environment help leave live make plastic properly put ...
Doc 3:	also	community dish live	make management recycling	rule take wash year box day set ...
Doc 4:	ask come good habit long	make recycle	recycling	think trash want work get much ...
Doc 5:	community feel good habit leave	make management move recycle	recycling	take information mixed sometimes ...
Doc 6:	ask catch make move properly	recycle	recycling	separate take think even get system throw ...
Doc 7:	bin environment need recycling	trash however information	living sort typically	concern express impact waste ...
Doc 8:	particularly plastic recycling	sort challenge clear	environmental impact paper waste approach	item process voice ...
Doc 9:	also bin habit need recycle	recycling	system challenge clear confusion	grow importance mention share ...
Doc 10:	also apartment bin community	compost feel habit need recycle	recycling	trash however sort clear ...
Doc 11:	apartment compost good management public	recycling	trash set sort clear ensure	environmental waste awareness ...

Sentence Topic Coloring for Documents: 0 to 11

Doc 0:	amount energy huge limit	melt metal recycle	require want waste ...
Doc 1:	actually altogether avoid bag ban basically	bring company default delivery	edit employee food get ...
Doc 2:	want waste food try accept area become bin cause	collective come compost contribute	district ...
Doc 3:	waste food recycling use become come compost recently	advice afford allow also apartment	appreciate ...
Doc 4:	recycle homeconduct efficient experience helpful improved	make observation public research stop thank ...	
Doc 5:	recycle know plastic take throw friend recently	make always answer away clean conserve container ...	
Doc 6:	bag plastic recycling recyclable alternative anywhere aware chip choice foil	garbage happy material mix ...	
Doc 7:	plastic way also make attach blade look remove tell	unable ...	
Doc 8:	come add cover windshield winter ...		
Doc 9:	actually bag food plastic use time recyclable batch entire esp fact layer machine multiple	...	
Doc 10:	recycle actually bring edit plastic put recycling throw try use whole bin	homenew...	
Doc 11:	bag get use much multiple keep place pretty ad break bunch check clothe hole	...	

2) t-SNE plot of documents: We use t-SNE plot to cluster documents according to their topics to identify how varied the documents are in terms of topic separation. As it can be seen, certain topics are well separated, while others are marginally related.

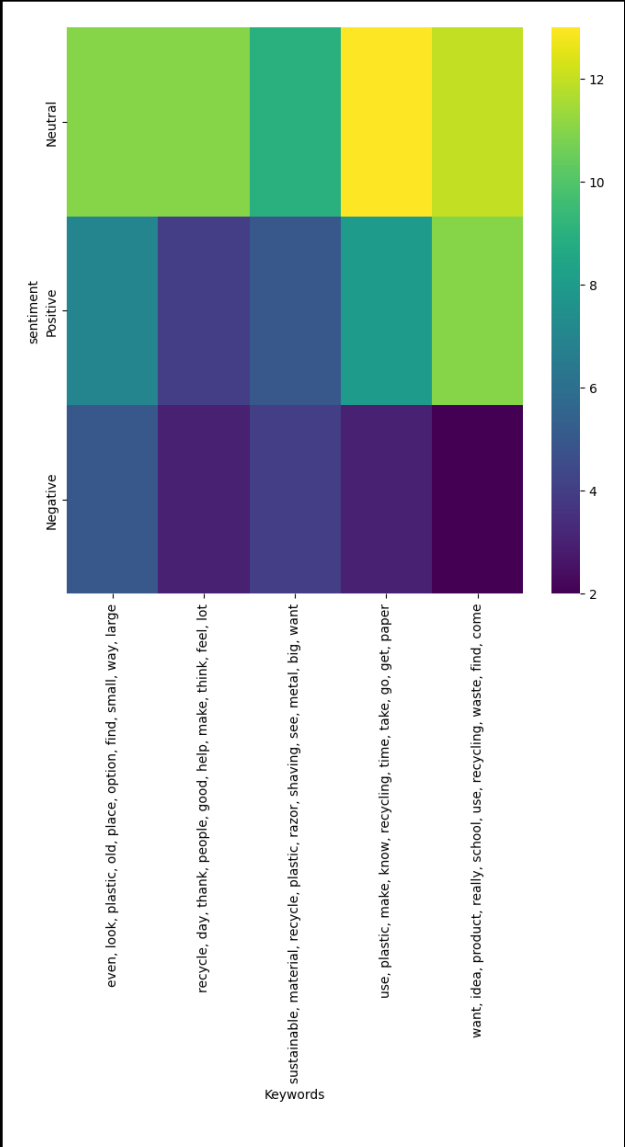
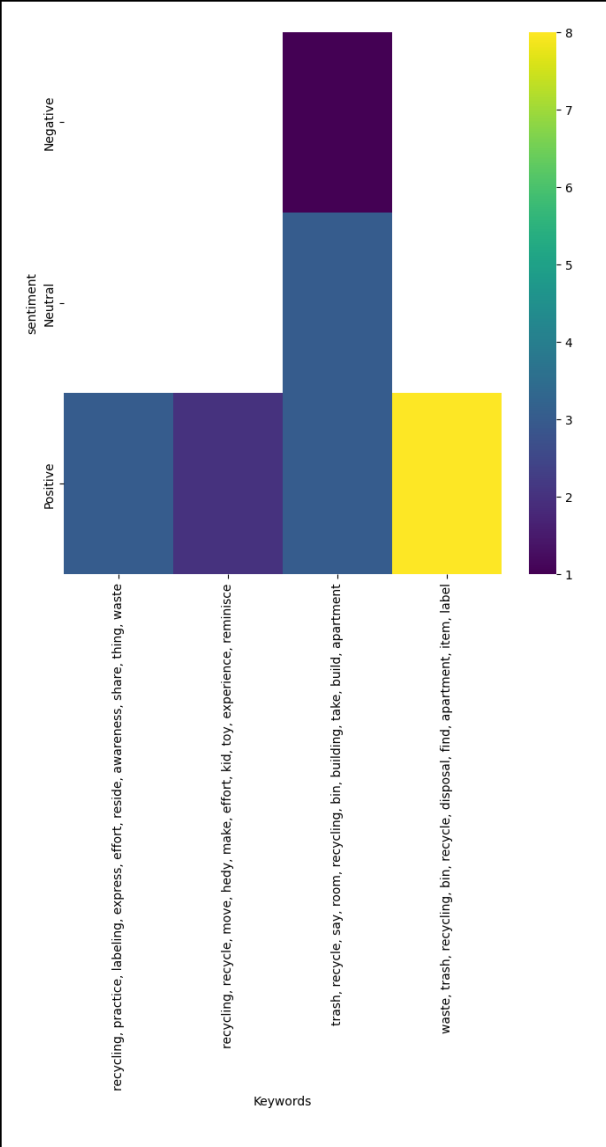
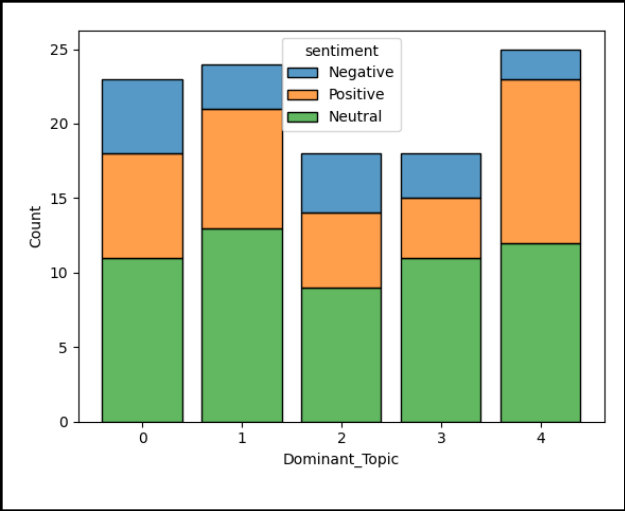
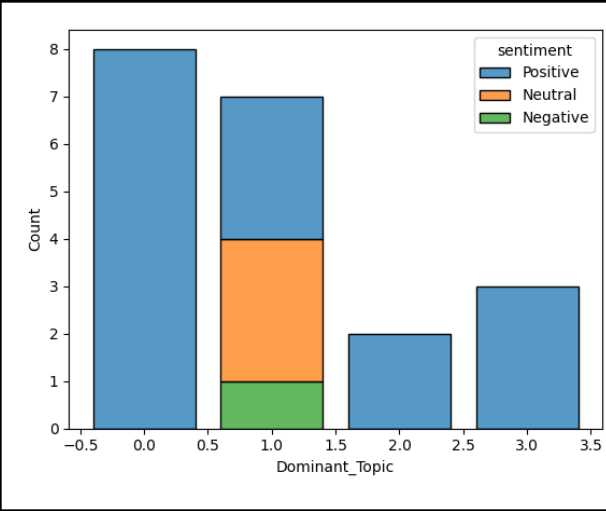


Cross-Referencing Sentiment and Topics

In order to analyze the distribution of sentiments across topics, we plot stacked histograms which are presented below. It can be seen that in case of interviews, only one topic had diversity in terms of sentiments. The topic related to trash rooms in the apartment had mixed sentiments, while other topics had only positive sentiments associated with them.

But in case of web data, the sentiments are equally distributed among the topics and each topic has a fair representation of positive, negative and neutral sentiments.

Furthermore, we also present correlation plots of topics and sentiments to get a more granular understanding of keywords and sentiments.



Limitation:

One limitation of our study is the interview data we had to work with. Because each of us conducted our own interviews, they were all of different lengths and not necessarily focused on the same topics. Additionally, because we didn't record our interviews word for word and instead summarized them, the language we used doesn't necessarily depict what the interviewees actually said, but rather how we chose to recount the interview. Therefore, our own biases may have affected our results. Additionally, we don't have a lot of data, so the trends we see are not strong. Furthermore, the data we collected was limited in quantity and scope, as public comments on recycling are meager on sites like Reddit. And while we did have twitter ids for data on recycling, due to the paid nature of twitter api, we couldn't collect that data. In order to improve our study, access to data from active forums like twitter along with broadening the scope could lead to substantial increase in depth of the insights.