

Social Network Analysis of Fine Food Reviewers

Dataset Introduction

The Amazon Fine Foods Reviews data set is a text file that contains 568,454 records of user reviews of fine food products from Amazon. Each record contains information on the user ID, product ID, profile name, helpfulness, score, time, summary and text, which can be seen in Figure 1.

Network Structure

This project aims to construct and analyze a network of Amazon fine food reviewers based on the similarity of their product reviews, in other words, the goal is to detect communities of users who have similar tastes or opinions for fine food products on Amazon. Each node will represent a user and edges between the nodes will exist if two users gave the same product the same score in their review. The weights for each edge represent the number of similar reviews each pair of users has. Due to the large size of the dataset, the project will focus on the relationship between the top 50 reviewers, which is based on the number of reviews recorded for each user in the dataset. Figure 2 shows a data frame of the graph data.

Basic Network Statistics

After constructing the network, there are a total of 48 nodes in the network with 461 edges between the nodes. Given that the original dataset had records for 50 users, this indicates that two of the top 50 reviewers had no similar reviews to any of the other top 50 reviewers. On average for each reviewer, about 9 other reviewers in the top 50 have also given a product a review with the same score.

The degree of a node measures the number of neighbours for that node, where neighbours are the other nodes that are connected to the node by an edge. Therefore, in this case, the degree of a node will return the number of reviewers who have given the same product the same score.

Figure 3 in the Appendix shows the degree for each node in the network. We can see that user ‘AKMEY1BSHSDG7’ has the highest degree, with connections to 36 other users who have written at least one similar review to this user. At the bottom of the figure, we can see that user ‘AKZKG2Z7CNV27’ only has one review that is similar to another user in the top 50 reviewers.

Since the network is a weighted graph, we can also compute the weighted degree of each node, where the weighted degree is calculated as the sum of the weighted edges connected to the node (NetworkX 3.2.1 documentation, n.d.). Compared to the unweighted degree, the weighted degree not only factors in the number of connections to the node but also the strength of the connection

measured by the weight of the edge connecting the two nodes. In Figure 4, we can see that the node with the highest weighted degree is now user ‘A1TMAVN4CEM8U8’ with a weighted degree of 461. Given that the weight of a node is the number of similar reviews a user has with other users, we can say that user ‘A1TMAVN4CEM8U8’ has scored a product the same as another user 461 times.

By comparing the two measures, we can say that user ‘AKMEY1BSHSDG7’, who had the highest degree, has the highest number of connections, whereas user ‘A1TMAVN4CEM8U8’, with the highest weighted degree, had a stronger weighted connection to its neighbours. In other words, the second user had a higher number of similar opinions to each of its neighbouring reviewers whereas the first user simply had more similar reviews to a higher number of distinct reviewers.

Shortest Paths

Given this network, we can also find the shortest path between two reviewers. For example, we can first look at the shortest path between the reviewer with the highest simple degree, ‘AKMEY1BSHSDG7’, and the reviewer with the second highest simple degree, ‘AQLL2R1PPR46X’. Using NetworkX’s functions, `shortest_path()` and `shortest_path_length()`, we found that these two users are directly connected by one edge, resulting in a shortest path with a length of 1. On the other hand, the shortest path between the reviewer with the highest simple degree and the reviewer with the lowest simple degree has a length of 2, with user ‘A36MP37DITBU6F’ connecting the two other reviewers.

Node Importance

To better understand the influence of each node in the network, we can look at various measures of node importance such as Degree Centrality, Betweenness Centrality, and Closeness Centrality.

Degree Centrality measures the fraction of total nodes a node is connected to, or how connected the node is to the rest of the network (NetworkX 3.2.1 documentation, n.d.). In Figure 5, we can see that user ‘AKMEY1BSHSDG7’ had the highest degree centrality at about 0.77, which indicates that this user had written at least one similar review to 77% of the top 50 Amazon Users. This is also the same user with the highest degree in Figure II.

The Betweenness Centrality measures how important a node is in terms of connecting other nodes to each other (Neo4j Graph Data Platform, (n.d.-a)). Specifically, the Betweenness Centrality of node V is the fraction of the shortest paths from node A to node B that contain node V over the total number of shortest paths connecting node A to node B (Padmanabhan, 2023). In Figure 6, we can see that all nodes in the network have a very low betweenness centrality score, with user ‘A36MP37DITBU6F’ having the highest betweenness centrality at only 0.059. This is expected as edges are created for each pair of nodes if the two reviewers share a similar review

for a product. Therefore in this network, most pairs of nodes will be directly connected with a shortest path containing only 1 edge, resulting in a very low betweenness centrality for every node.

Finally, the Closeness Centrality measures how far a node is from all the other nodes, with a higher closeness centrality score corresponding to nodes that are closer to all other nodes (Neo4j Graph Data Platform, (n.d.-b)). In Figure 7, we can see that user ‘AKMEY1BSHSDG7’ has the highest closeness centrality score at 0.81; this is the same user with the highest simple degree and degree centrality score. Since a lot of the shortest paths in the network contain only one edge, this user scored the highest in closeness centrality as it had the most neighbours in this network, consequently resulting in shorter distances to all other nodes in the network.

Community Detection

After constructing the network we can try to identify communities within the network using Louvain Modularity. Louvain Modularity identifies communities through optimization based on modularity, where modularity compares the density of connections within clusters and the density of connections between clusters (Padmanabhan, 2023). After running the Louvain Modularity algorithm, the output showed that the algorithm was able to identify 6 communities in this network. Figure 8 shows a plot of the network, with the different colours identifying the different communities.

Community Attributes

To better understand the characteristics of each community, we can analyze the similar reviews between each pair of reviewers in each community. After gathering similar reviews between each pair of reviewers in each community, we can plot the distribution of scores given by the reviewers in each community. Figure 9 shows the average score for each community and Figure 9 shows the bar plots of the number of reviews with each score for each community. From these figures, we can see that the majority of communities have very high scores for their reviews and that Communities 0 and 4 only have reviews with a score of 5. Only Community 2 has reviews with scores ranging from 1-5. This implies that many of the reviewers in each community share positive opinions on products if a high score corresponds to a positive opinion.

Given that the dataset did not include the product name for each review, we can analyze the summary text and review text to better understand what kind of product the user reviewed. To do so, we can create a word cloud of the top 100 words found for each community’s review summaries and review texts. Figures 11-16 show the word clouds generated for each community.

Out of the top words generated for community 0, we can see that this community contains users who shared positive reviews on various snacks such as chips, chocolate, and cereal, describing them with words such as ‘crunch’ or ‘hot’. Community 1 also seems to have reviewers who

mainly reviewed products such as chips, but also other products such as tea, beef, fruit, and juice. Communities 2 seems to focus more on tea reviews, commonly using words such as ‘sweet’ and ‘organic’. Community 3 also focuses on tea reviews, but particularly types like Chamomile and hibiscus. Community 4 has a variety of food reviews including almonds and pasta, with common descriptive words such as ‘bold’ and ‘spicy’. Finally, Community 5 has a strong emphasis on coffee reviews, with descriptive words such as ‘bold’, ‘roast’ and ‘acidic’.

Conclusion

From this project, we were able to construct and analyze reviewers in the Amazon fine foods dataset. To improve on this network in the future, we can try to re-define the edges of the network to better measure the similarity between reviewers. For example, we can take the review text into consideration to measure the similarity of reviews between two users.

References

Betweenness centrality - neo4j graph data science. Neo4j Graph Data Platform. (n.d.-a).
<https://neo4j.com/docs/graph-data-science/current/algorithms/betweenness-centrality/>

Closeness centrality - neo4j graph data science. Neo4j Graph Data Platform. (n.d.-b).
<https://neo4j.com/docs/graph-data-science/current/algorithms/closeness-centrality/#algorithms-closeness-centrality-limitations>

Degree centrality. degree_centrality - NetworkX 3.2.1 documentation. (n.d.).
https://networkx.org/documentation/stable/reference/algorithms/generated/networkx.algorithms.centrality.degree_centrality.html

Graph.degree. Graph.degree - NetworkX 3.2.1 documentation. (n.d.).
<https://networkx.org/documentation/stable/reference/classes/generated/networkx.Graph.degree.html>

Padmanabhan, K. (2023a, December). *Social Network Analysis - Week 1*. Lecture.

Padmanabhan, K. (2023b, December). *Social Network Analysis - Week 2*. Lecture.

Appendix

Figure 1: Amazon Fine Foods Dataset Sample

ProductId	UserId	ProfileName	Helpfulness	Score	Time	Summary	Text
0 B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1/1	5.0	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...
1 B00813GRG4	A1D87F6ZCVE5NK	dll pa	0/0	1.0	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2 B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1/1	4.0	1219017600	"Delight" says it all	This is a confection that has been around a fe...
3 B000UA0QIQ	A395BORC6FGVXV	Karl	3/3	2.0	1307923200	Cough Medicine	If you are looking for the secret ingredient i...
4 B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0/0	5.0	1350777600	Great taffy	Great taffy at a great price. There was a wid...

Figure 2: Weighted Network Graph Data

	User A	User B	Weight
1	A3FKGKUCI3DG9U	AEC90GPFKLAACW	2
6	A2FRFAQCWZJT3Q	AY1EF0GOH80EK	1
9	A25C2M3QF9G7OQ	A2SZLNSI5KOQJT	5
11	A1IU7S4HCK1XK0	A2FRFAQCWZJT3Q	1
12	A1WX42M589VAMQ	A2DPYMN12HCIOI	1
...
1214	AQLL2R1PPR46X	AY12DBB0U420B	1
1215	A250AXLRBVYKB4	A2FRFAQCWZJT3Q	11
1216	A1X1CEGHTHMBL1	A2R6RA8FRBS608	1
1219	A1P2XYD265YE21	A36MP37DITBU6F	3
1224	A2SZLNSI5KOQJT	A3FKGKUCI3DG9U	2

461 rows × 3 columns

Figure 3: Simple Degree

```
[('AKMEY1BSHSDG7', 36),
 ('A30XHLG6DIBRW8', 32),
 ('AQLL2R1PPR46X', 32),
 ('A3PJZ8TU8FDQ1K', 30),
 ('A2RFAQCWZT3Q', 28),
 ('A25C2M3QF9G7Q0', 28),
 [(('A1TMAVN4CEM8U8', 461),
   ('AD5APY1NKT4', 421),
   ('A1X1CEGHHTHMBL1', 396),
   ('A3PJZ8TU8FDQ1K', 375),
   ('A281NPSIMI1C2R', 256),
   ('A30XHLG6DIBRW8', 249),
   ('A3FY3H6F4249E0', 242),
   ('A2RFAQCWZT3Q', 226),
   ('AKMEY1BSHSDG7', 219),
   ('A1YUL9PCJR3JTY', 208),
   ('A1P2XYD265YE21', 197),
   ('A1LZJZIHULPDV4', 188),
   ('AY1EF0GOH8EK', 167),
   ('A2Y8IDC1FKGNJC', 167),
   ('AQLL2R1PPR46X', 146),
   ('A1WX42M589VAMQ', 143),
   ('A3FKGKUCI3DG9U', 141),
   ('A25C2M3QF9G7Q0', 129),
   ('A2S2LNS15KOQQT', 122),
   ('A1Z54EM24Y40LL', 115),
   ('A1HRYC60VTMYC0', 114),
   ('AQQLWCMRNDFG1', 112),
   ('A1UQBFCERIP7VJ', 103),
   ('A2PNOU7NXB1J4', 96),
   ('A3QNQQKJTL76H0', 95),
   ('A17HMM1M7T9P11', 90),
   ('A250AXLRBVYKB4', 71),
   ('A3HPCRDRX3515', 61),
   ('A31N6KB1600508', 60),
   ('A1IU7S4HCK1XK0', 54),
   ('A2MUGFV2TDQ47K', 52),
   ('A1ZH9LWMX5UCFJ', 51),
   ('AEC98GPFLAAW', 49),
   ('A36MP37DITBU6F', 48),
   ('A2XNOB17796Y6B', 44),
   ('ALSAOZ1V546VT', 40),
   ('A2DPYMN12HCIOI', 35),
   ('AY12DBB0U420B', 35),
   ('A2R6RA8FRBS608', 31),
   ('AZV26LP92E6WU', 30),
   ('A3F3B1JPACN215', 28),
   ('A33AQPJYH7UUXR', 26),
   ('A35R32TA60XD57', 21),
   ('A2M9D9BDHONV3Y', 21),
   ('A36WGRH8T05DKT', 19),
   ('A2GEZJHBV92EVR', 14),
   ('A3D60I36USYOU1', 5),
   ('AKZKG2Z7CNV27', 1)]
```

Figure 4: Weighted Degree

Figure 5: Degree Centrality

```
[('A36MP37DITBU6F', 0.05861583093828492),
 ('AKMEY1BSHSDG7', 0.05279686198126074),
 ('A30XHLG6DIBRW8', 0.04174389176771267),
 ('A2MUGFV2TDQ47K', 0.03519267337475946),
 ('A3PJZ8T8UDQ1K', 0.029574230918795654),
 ('A2FRFAQCWCJ7TQ', 0.02721618556339733),
 ('A281NPSIMI1C2R', 0.02442729753888083),
 ('A1HRYC60VTMYC0', 0.022889721688026488),
 ('A1UQBFCERIP7VJ', 0.022852991225584762),
 ('AQLL2R1PPR46X', 0.02260657615350786),
 ('A1PZXYD265YE21', 0.01899911986118142),
 ('AY1EF0G0H80EK', 0.017237410025772),
 ('A2Y8IDC1FKGNJC', 0.016922068736142156),
 ('A3HPCRD9RX51S', 0.016243134298401463),
 ('A25CLM9QF967Q', 0.0153066556395754),
 ('ALSAO21V546VT', 0.01466949756237343),
 ('A1IU754HCK1XK0', 0.013553717760456444),
 ('A3FGKGUCI3DG9U', 0.0132970657261896),
 ('A3QNQQKJTL76H0', 0.012890820705938576),
 ('A2SZLNSISIKOQ3T', 0.0126118587236256),
 ('A1Z5E4M24Y40LL', 0.012529069828911898),
 ('A258AXLRVYKB4', 0.0108380641412553),
 ('A2RER8A8PFS668', 0.009677491589670528),
 ('A2DPYMMI1HICIOI', 0.0092458858558776),
 ('A2M9D9BDHONV3Y', 0.00865421352709328),
 ('AY12DBB0U420B', 0.008547595492510206),
 ('A3F3B1PACN215', 0.0083365207796320565),
 ('A2PNOU7NXB1JE4', 0.0080659398169992385),
 ('AEC90GPFKLAAW', 0.00793236760857908),
 ('AZXNOB1T796Y6B', 0.007269410433770456),
 ('A17HMM1M7T9P31', 0.006911972172509754),
 ('AZV26LP92E6WU', 0.006164479698515761),
 ('A1Z9LWMLX5UCFJ', 0.006087788237350189),
 ('A31N6KB1600508', 0.0056411076068800395),
 ('A1LZ3ZIHIPLDV4', 0.005613253866833864),
 ('A2GEZ3HBV92EV', 0.005529005206881639),
 ('A1YUL9PCJR3JTY', 0.005176828925340536),
 ('AQOLWCMRNDFG1', 0.004982098926710231),
 ('A1X1CEGHTHMBL1', 0.004867720173951983),
 ('A1TMAVN4CEM8U8', 0.00458491625122907),
 ('A3A3QJPYH7UXR', 0.003915642708775849),
 ('ADSSPAVIN1KTL4', 0.0038966866292281014),
 ('A3FY3H6F4249E0', 0.003865625936012833),
 ('A1W42M589VAMQ', 0.003899363643859481),
 ('A36WGRH8T05DKT', 0.00216705195834743),
 ('A35R32TA60XD57', 0.0001204922381945755),
 ('A3D60I36USYOU1', 0.0), ('AKZKG2Z7CNV27', 0.0)]
 [('AKMEY1BSHSDG7', 0.8103448275862069),
 ('A30XHLG6DIBRW8', 0.7580645161290323),
 ('AQLL2R1PPR46X', 0.7580645161290323),
 ('A3PJZ8T8UDQ1K', 0.734375),
 ('A2FRFAQCWCJ7TQ', 0.7121212121212122),
 ('A2MUGFV2TDQ47K', 0.7121212121212122),
 ('A281NPSIMI1C2R', 0.7121212121212122),
 ('A25CLM9QF967Q', 0.7014925373134329),
 ('A1P2XD265YE21', 0.7014925373134329),
 ('A281NPSIMI1C2R', 0.7014925373134329),
 ('A1HRYC60VTMYC0', 0.6911764708582353),
 ('A3FKKGUCI3DG9U', 0.6811594282898551),
 ('AY1EF0G0H80EK', 0.6714285714285714),
 ('A1IU754HCK1XK0', 0.6619718309859155),
 ('ALSAO21V546VT', 0.6619718309859155),
 ('A3HPCRD9RX51S', 0.6619718309859155),
 ('A1UQBFCERIP7VJ', 0.6619718309859155),
 ('AEC90GPFKLAAW', 0.6527777777777777),
 ('A36MP37DITBU6F', 0.6527777777777777),
 ('A3QNQQKJTL76H0', 0.6438356164383562),
 ('A17HMM1M7T9P31', 0.6438356164383562),
 ('A1Z5E4M24Y40LL', 0.6351351351351351),
 ('A2SZLNSISIKOQ3T', 0.626666666666666667),
 ('A2DPYMMI1HICIOI', 0.618421052631579),
 ('A2PNOU7NXB1JE4', 0.618421052631579),
 ('A1LZJ2IHIPLDV4', 0.6103896103896104),
 ('A3F3B1PACN215', 0.6103896103896104),
 ('ADSSPAVIN1KTL4', 0.6103896103896104),
 ('AZXNOB1T796Y6B', 0.6103896103896104),
 ('A17HMM1M7T9P31', 0.6025641025641025),
 ('A1Z9LWMLX5UCFJ', 0.6025641025641025),
 ('AY12DBB0U420B', 0.6025641025641025),
 ('A258AXLRVYKB4', 0.6025641025641025),
 ('A1YUL9PCJR3JTY', 0.5949367088607594),
 ('A2R6RA8PRBS608', 0.5875),
 ('AZV26LP92E6WU', 0.5875),
 ('A1TMAVN4CEM8U8', 0.5802469135802469),
 ('A2M9D9BDHONV3Y', 0.5802469135802469),
 ('A3A3QJPYH7UXR', 0.573170731707317),
 ('A3FY3H6F4249E0', 0.573170731707317),
 ('A1W42M589VAMQ', 0.5595238095238095),
 ('A1X1CEGHTHMBL1', 0.5529411764705883),
 ('A2GEZ3HBV92EV', 0.5529411764705883),
 ('AQOLWCMRNDFG1', 0.5402298858574713),
 ('A36WGRH8T05DKT', 0.5222222222222223),
 ('A31N6KB1600508', 0.5222222222222223),
 ('A3D60I36USYOU1', 0.47959183673469385),
 ('A35R32TA60XD57', 0.47474747474747475),
 ('AKZKG2Z7CNV27', 0.3983050847457627)]
```

Figure 7: Closeness Centrality

Figure 6: Degree Betweenness

Figure 8: Community Identification

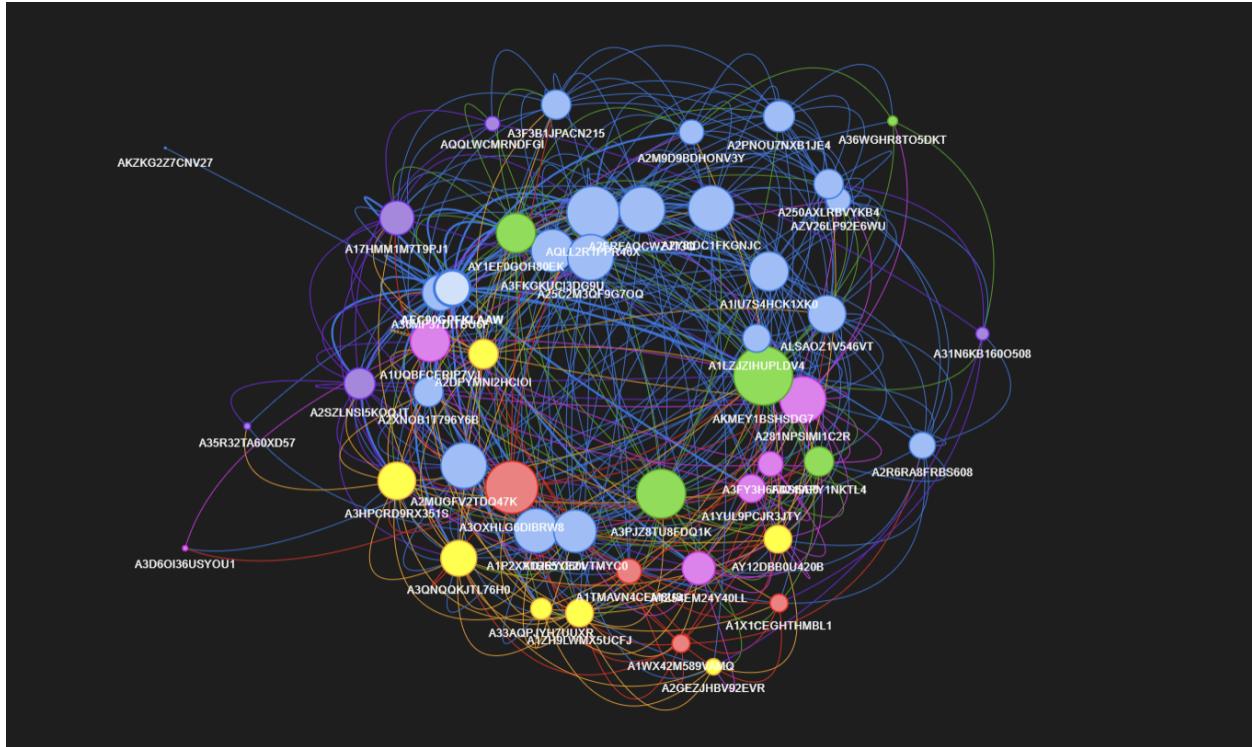


Figure 9: Average Review Score by Community

Community	Average Review Score
0	5.000000
1	4.614035
2	4.554723
3	4.829352
4	5.000000
5	4.148148

Figure 10: Review Scores for Each Community

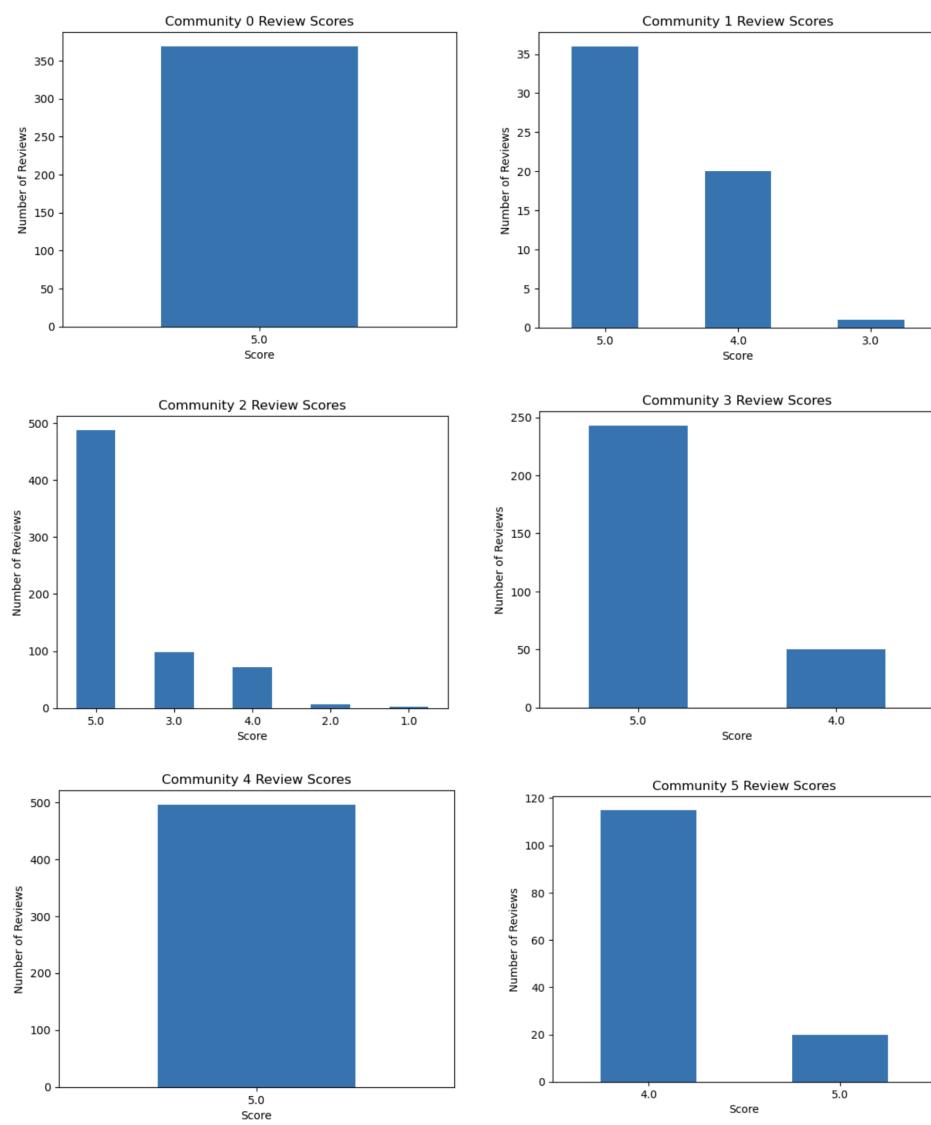


Figure 11: Word Clouds for Cluster 0 (Left: Summary text, Right: Review Text)



Figure 12: Word Clouds for Cluster 1 (Left: Summary text, Right: Review Text)



Figure 13: Word Clouds for Cluster 2 (Left: Summary text, Right: Review Text)



Figure 14: Word Clouds for Cluster 3 (Left: Summary text, Right: Review Text)



Figure 15: Word Clouds for Cluster 4 (Left: Summary text, Right: Review Text)

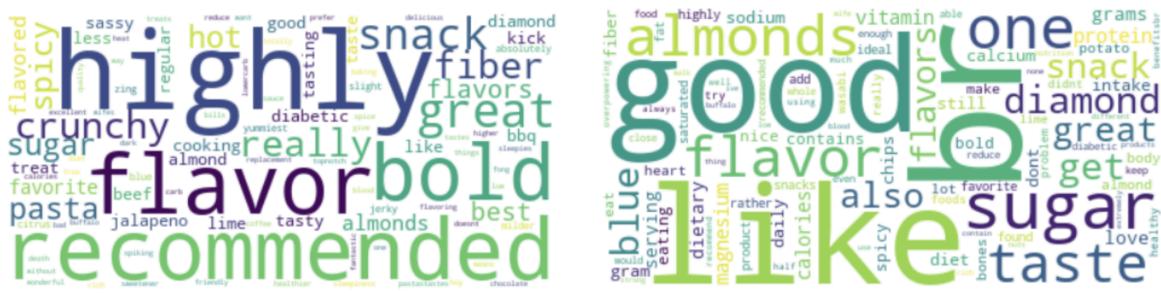


Figure 16: Word Clouds for Cluster 5 (Left: Summary text, Right: Review Text)

