

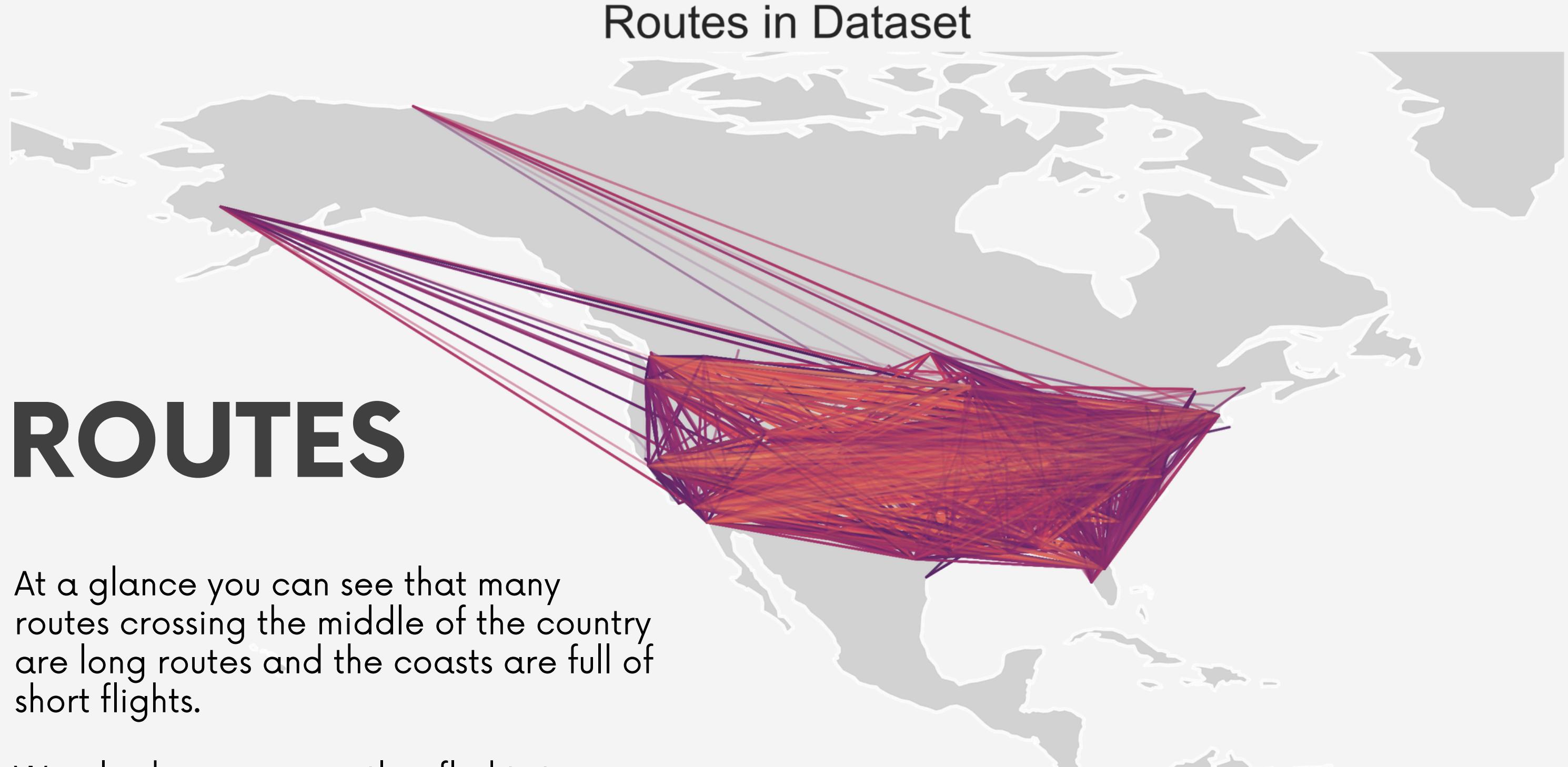
SKY ROUTE

DATE  
09/23/2024

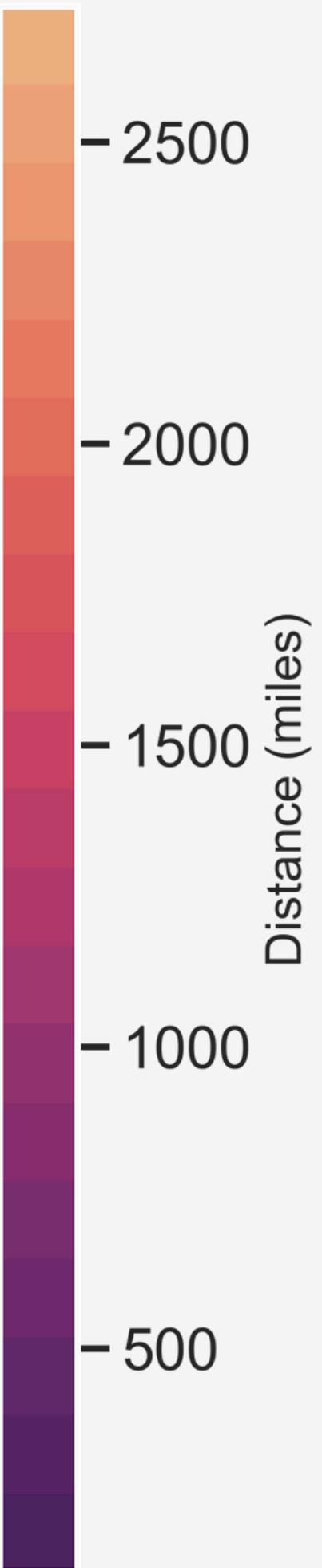
# CUSTERING ANALYSIS

Identifying distinct groups of city pairs with similar passenger patterns

PRESENTED BY  
Byung Jae Bae



We also have some outlier flights to Alaska which can show up in our clustering analysis



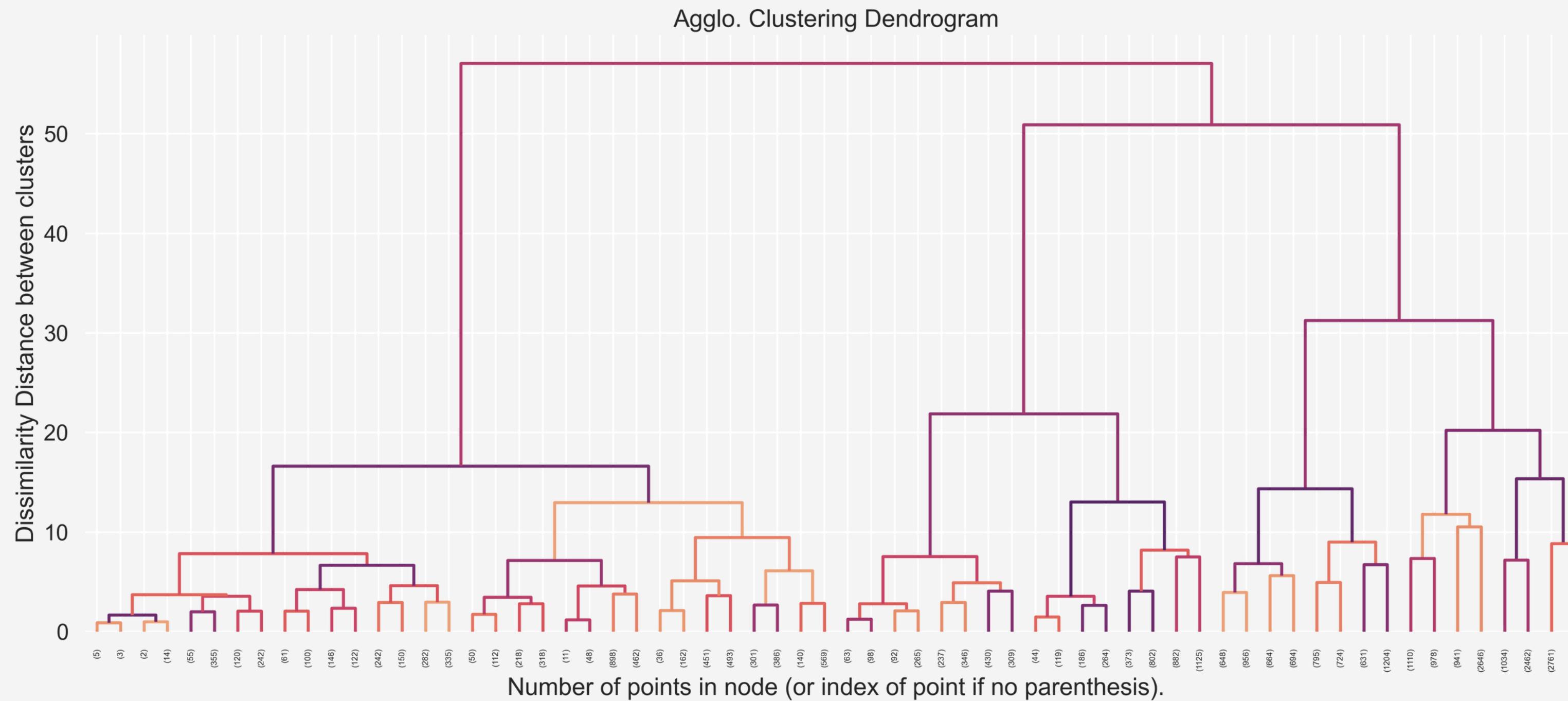
# LIMITATIONS OF ANALYSIS

This clustering analysis does not reflect direct customer behavior and uses metrics like average fares and market share to reflect market demand

# CLUSTERING ON CONTINUOUS FEATURES

Agglomerative Hierarchical Clustering was used to see how well clusters differentiate from different clusters.

We can see that there are about 2-3 clusters that have very high dissimilarity from each other.



# COMPARISON

The top two charts show K-Means with 3 clusters.

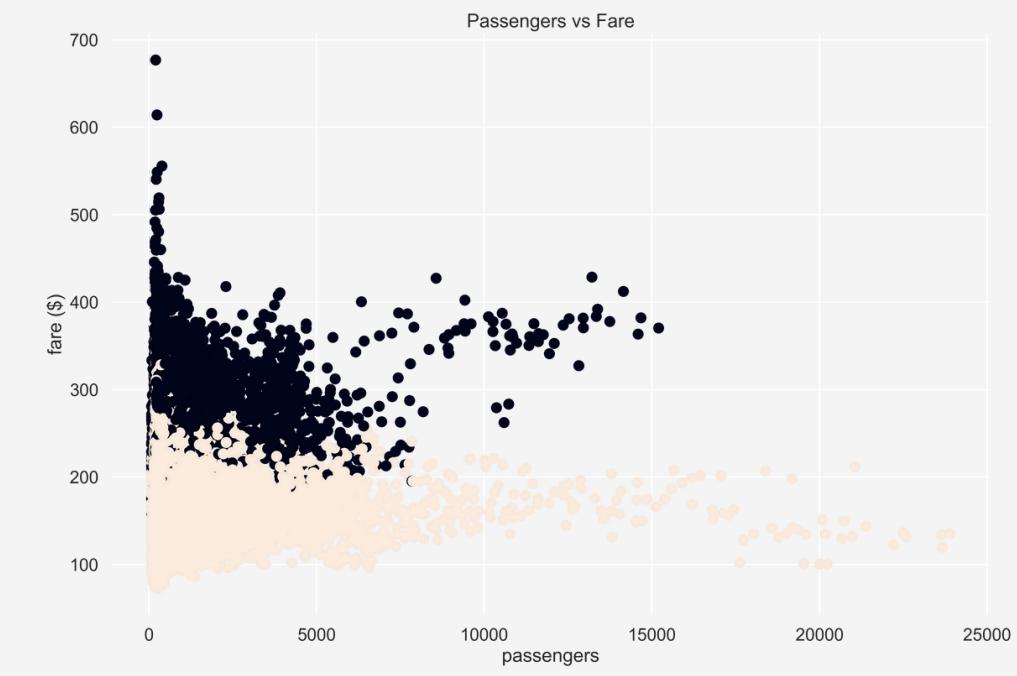
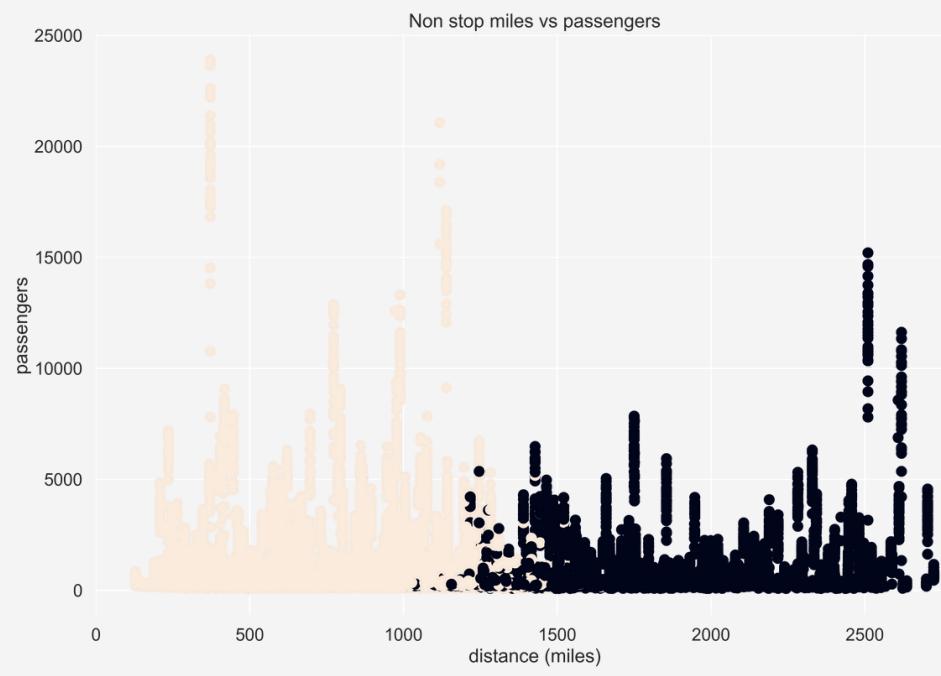
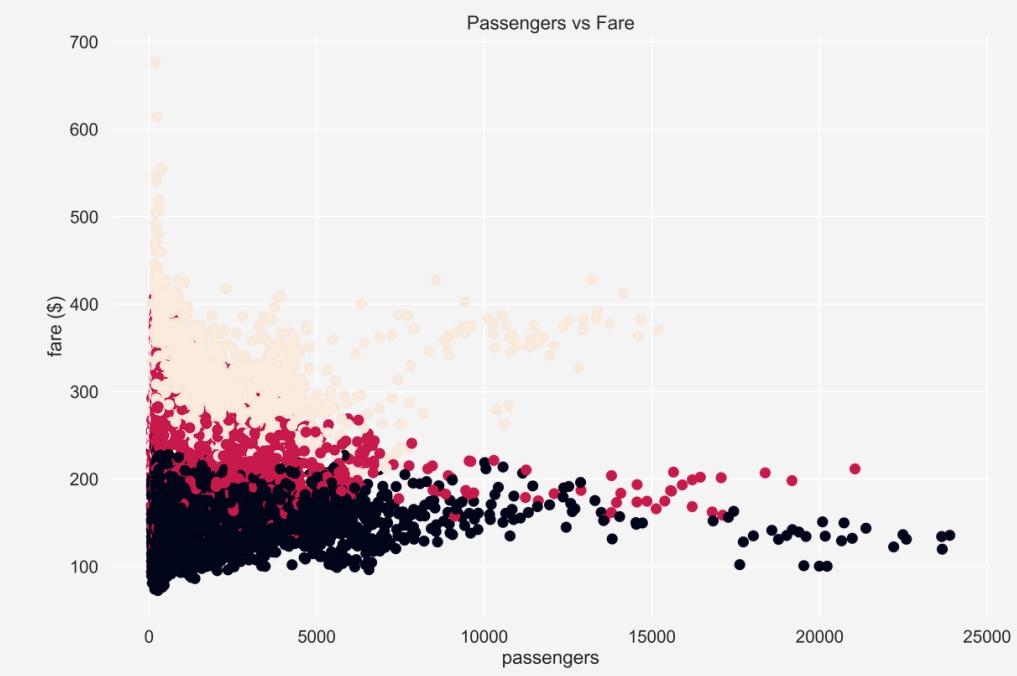
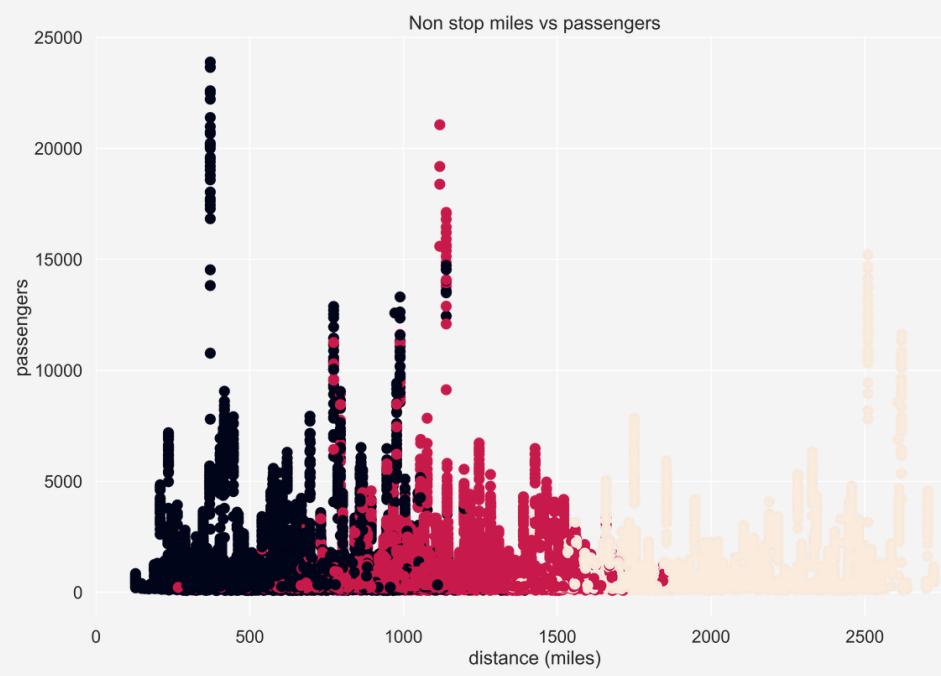
The bottom two charts show K-Means with 2 clusters.

Silhouette score (higher is better):  
0.32236013543110364 (Top)

vs  
0.48941445345090195 (Bottom)

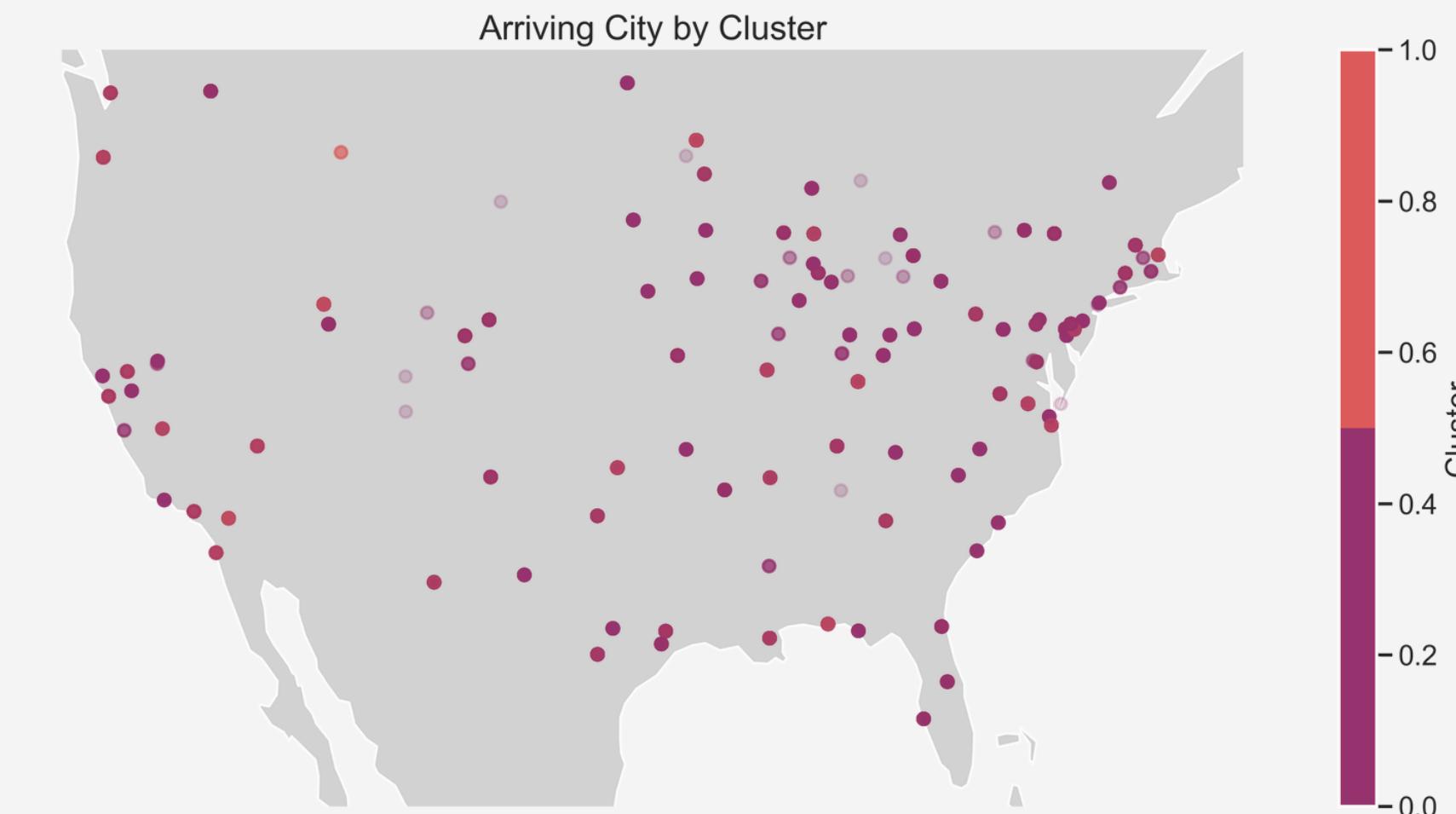
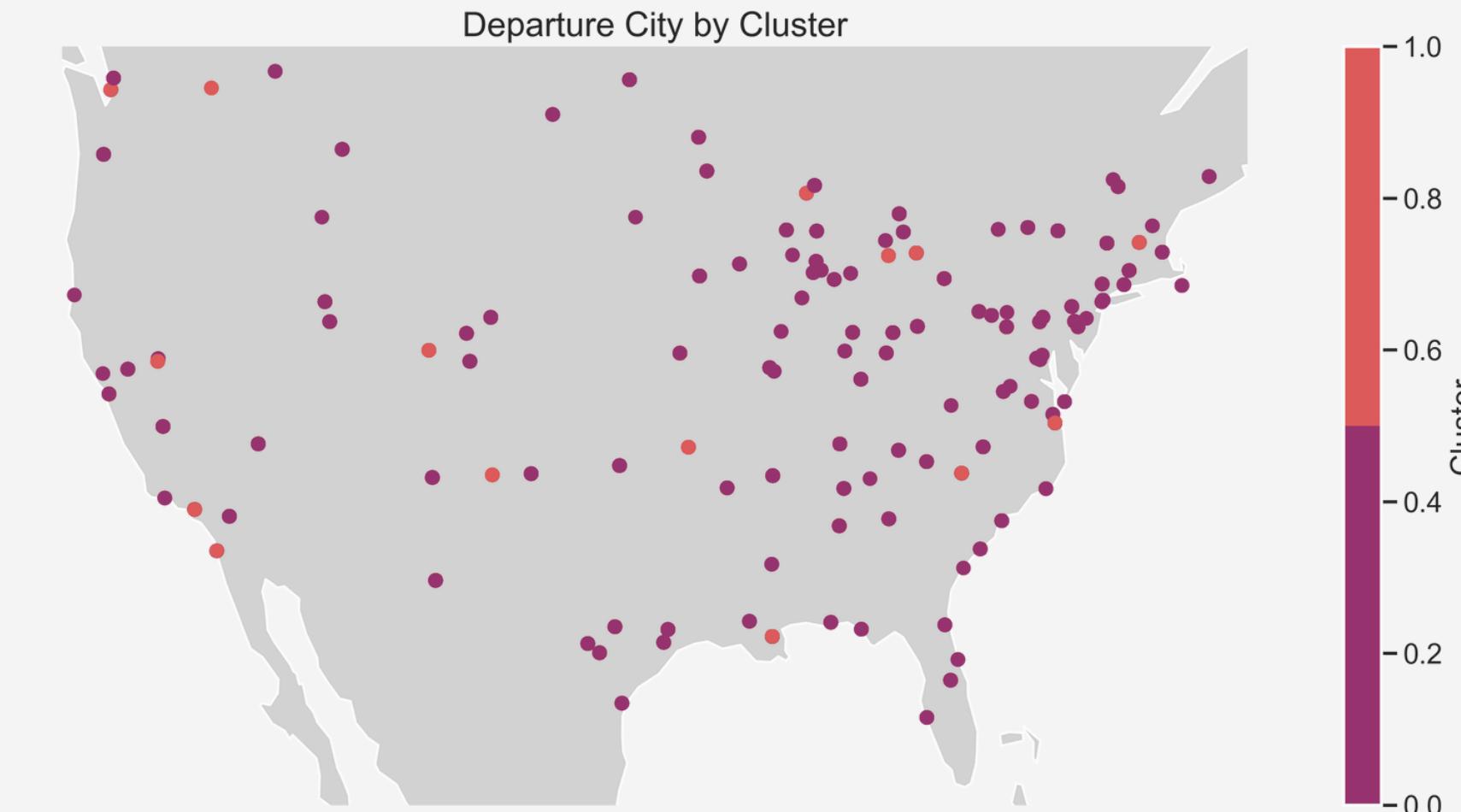
Davies-Bouldin score (lower is better):  
1.1110435866127217 (Top)

vs  
0.7982241607373859 (Bottom)



# CHECKING HOW WELL IT SPLITS GEOGRAPHIC REGIONS

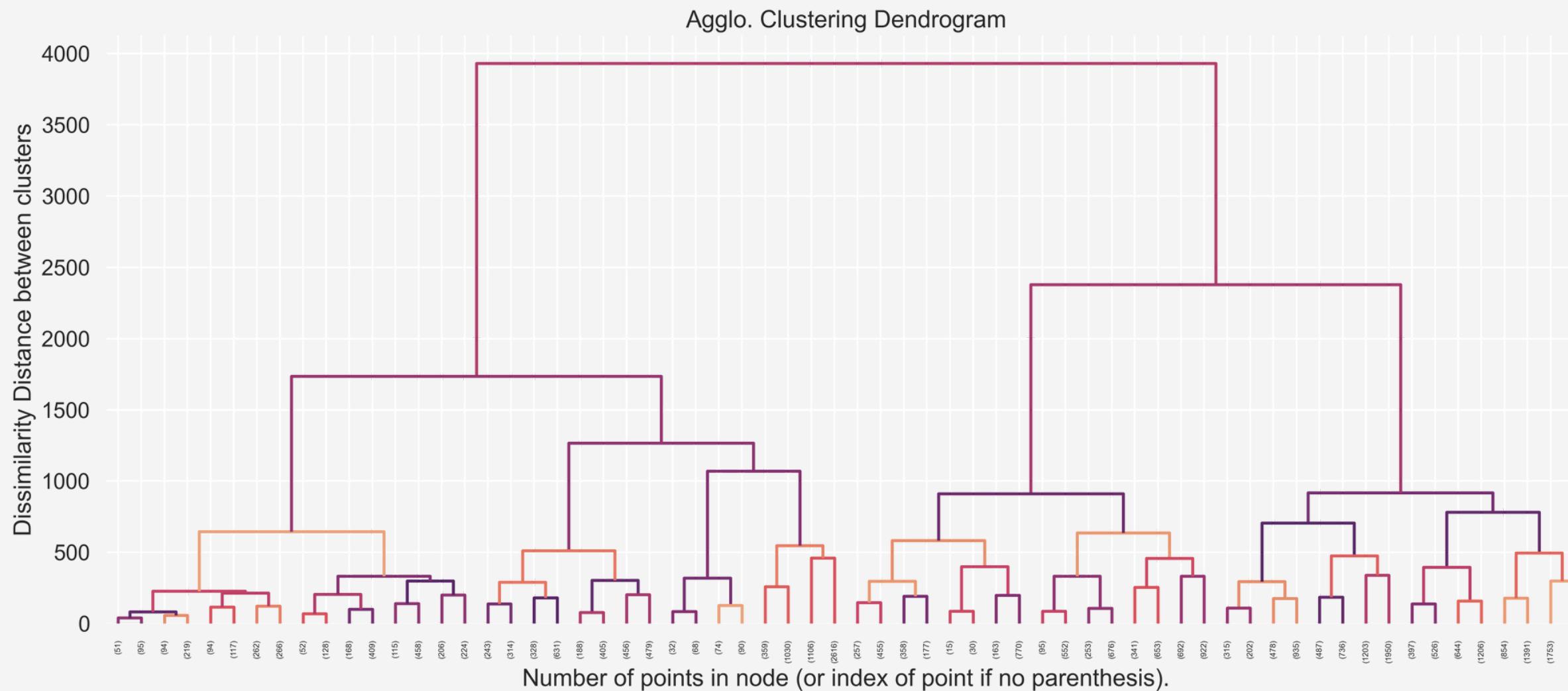
It seems that clustering on non geographic data doesn't reflect much clustering in terms of geographical regions.



# CLUSTERING ON COORDINATES ONLY

Instead of using other features we try clustering on coordinates

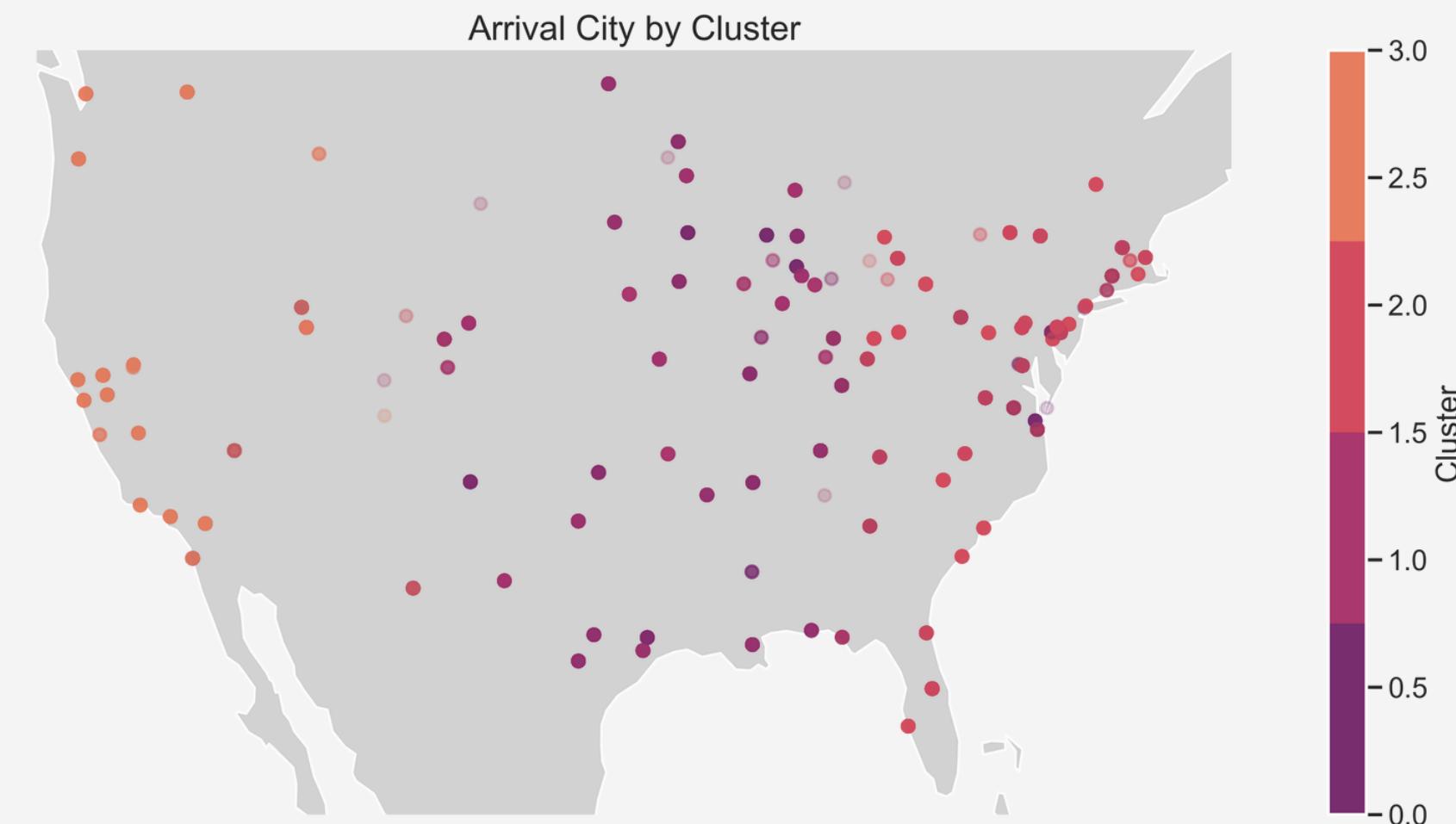
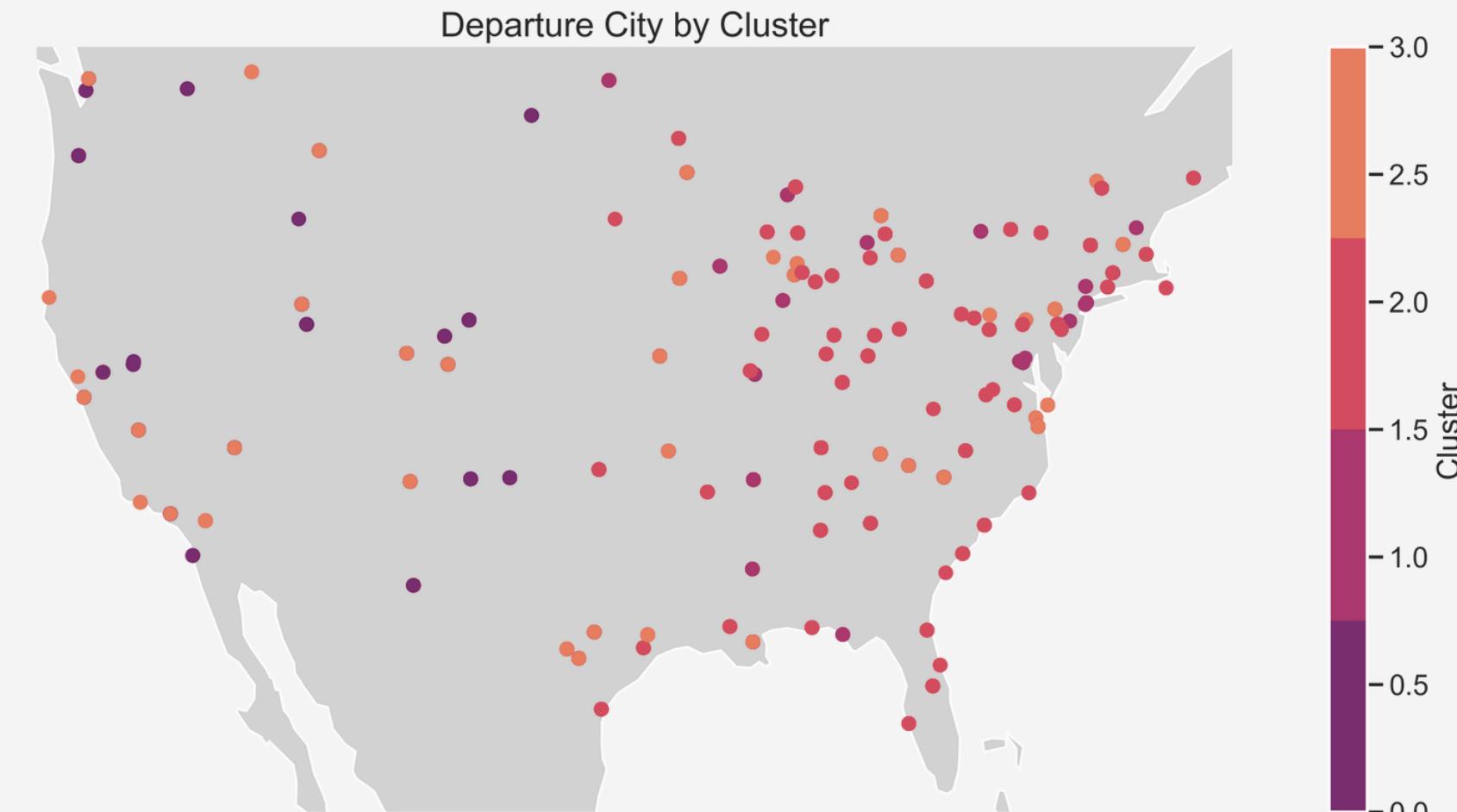
There is 4 nearly equal clusters that seem to appear.



# GEOGRAPHICAL CLUSTERS

Even using coordinates, the clusters are not super defined.

There is a clear east and west coast separation, but the center of the country has some overlap.



# CONCLUSION

There are 2 clear clusters/markets and the more clusters create more overlap, performing less optimally.

Metrics also perform well when market share features are removed from the clustering.

Geographical locations have little impact on good clustering.

We can interpret the data that there is high low-cost and high-cost carrier separation. There are no features that define a middle price tier market. It seems that high-priced low cost carriers and low-priced high cost airlines are both competing in this space. The routes reflects a bimodal distribution.

# Appendix

# 1st Agg. Hierarchical Clustering

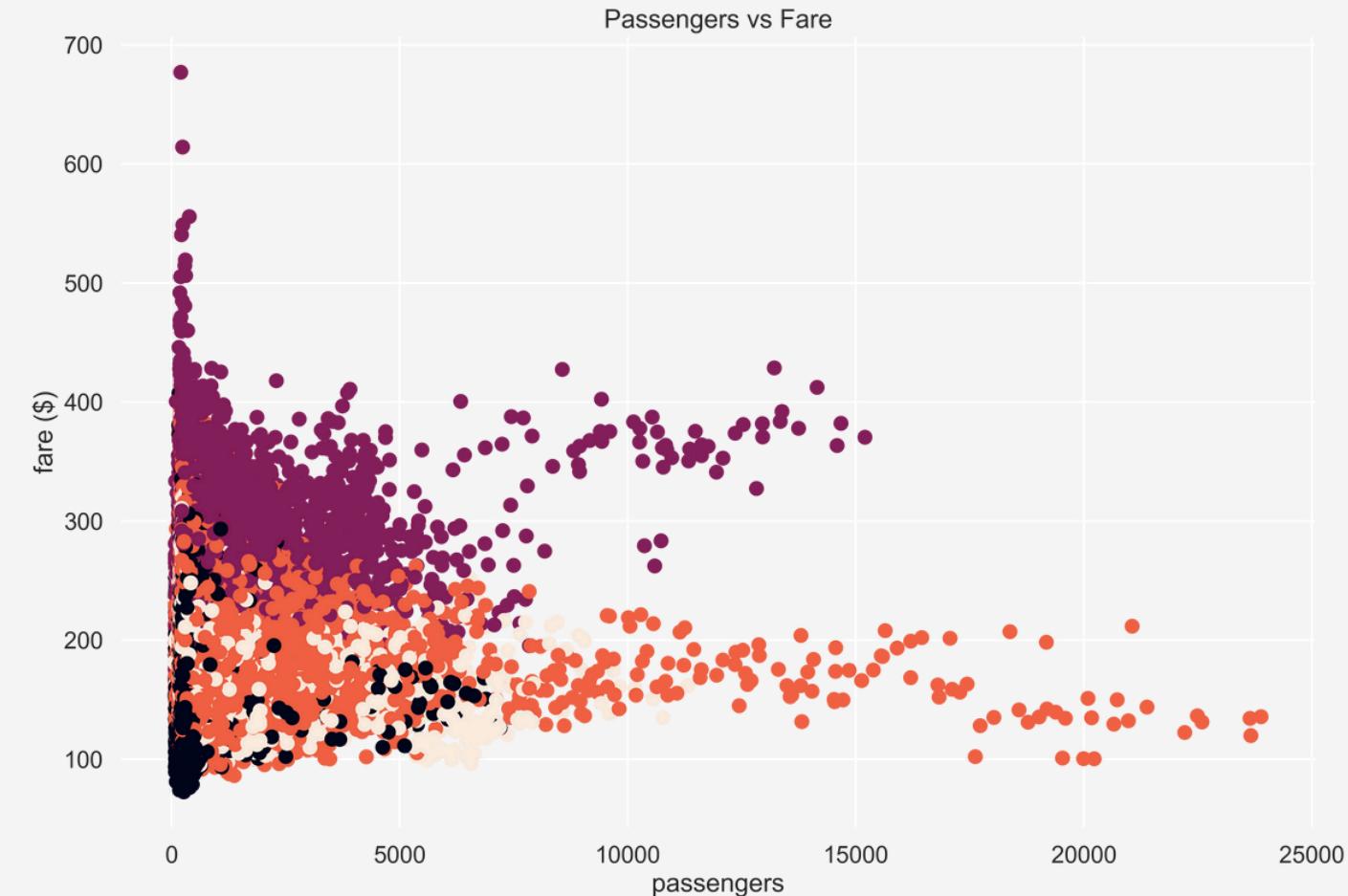
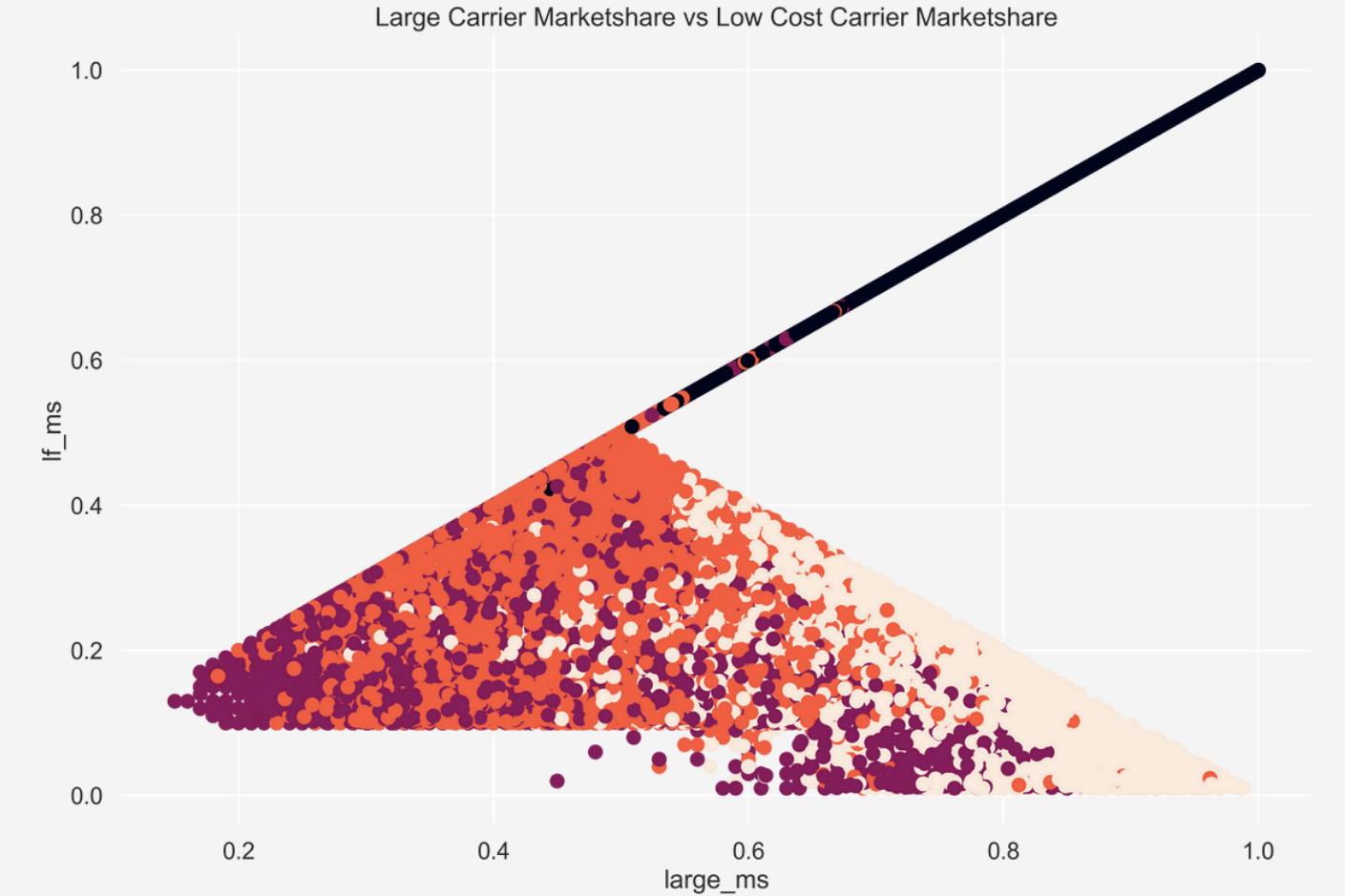
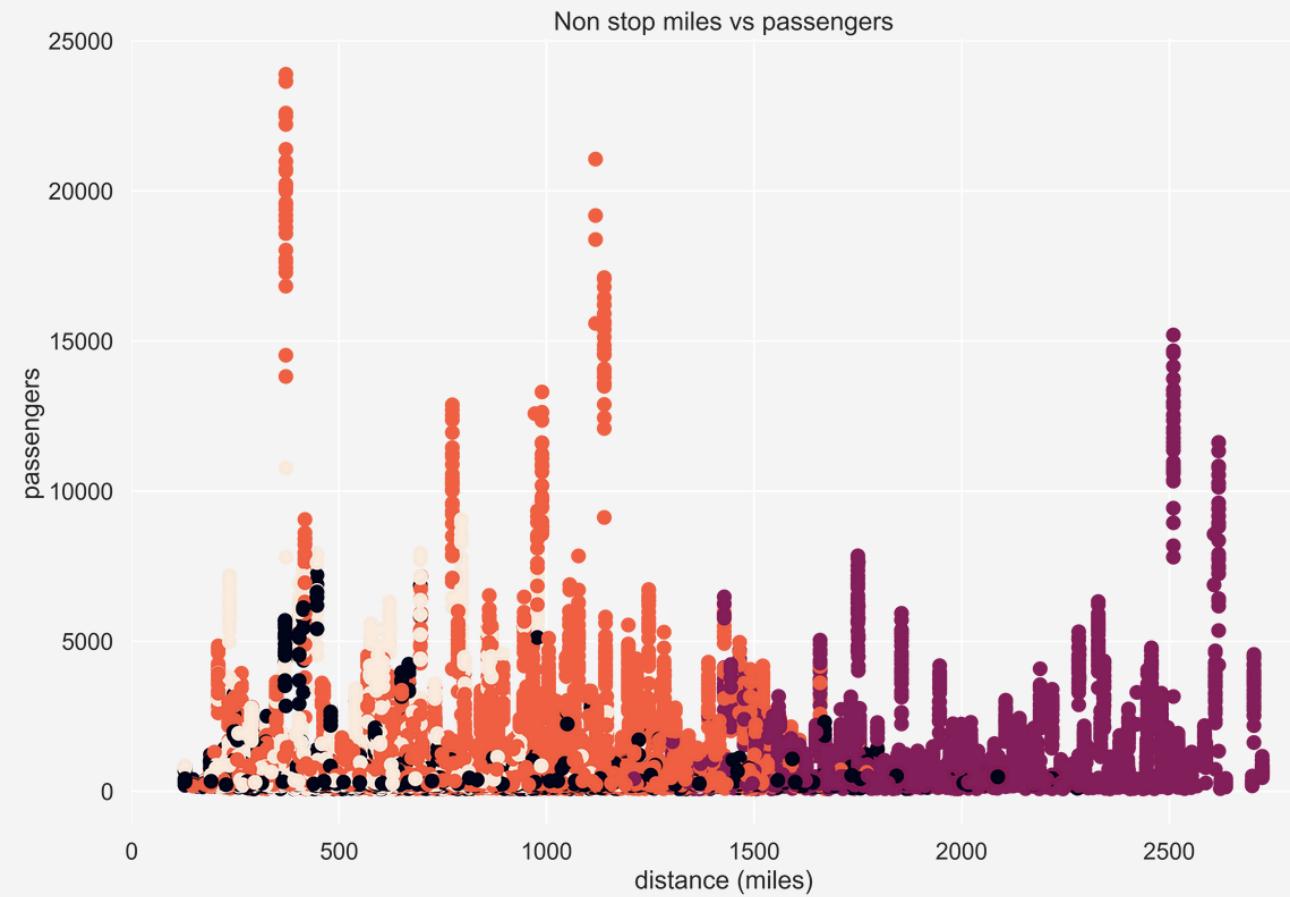
The average silhouette score is:

0.2507565409851805

The Davies-Bouldin score is:

1.2707774363381936

'nsmiles',  
'passengers',  
'fare', 'large\_ms',  
'fare\_lg', 'lf\_ms',  
'fare\_low'

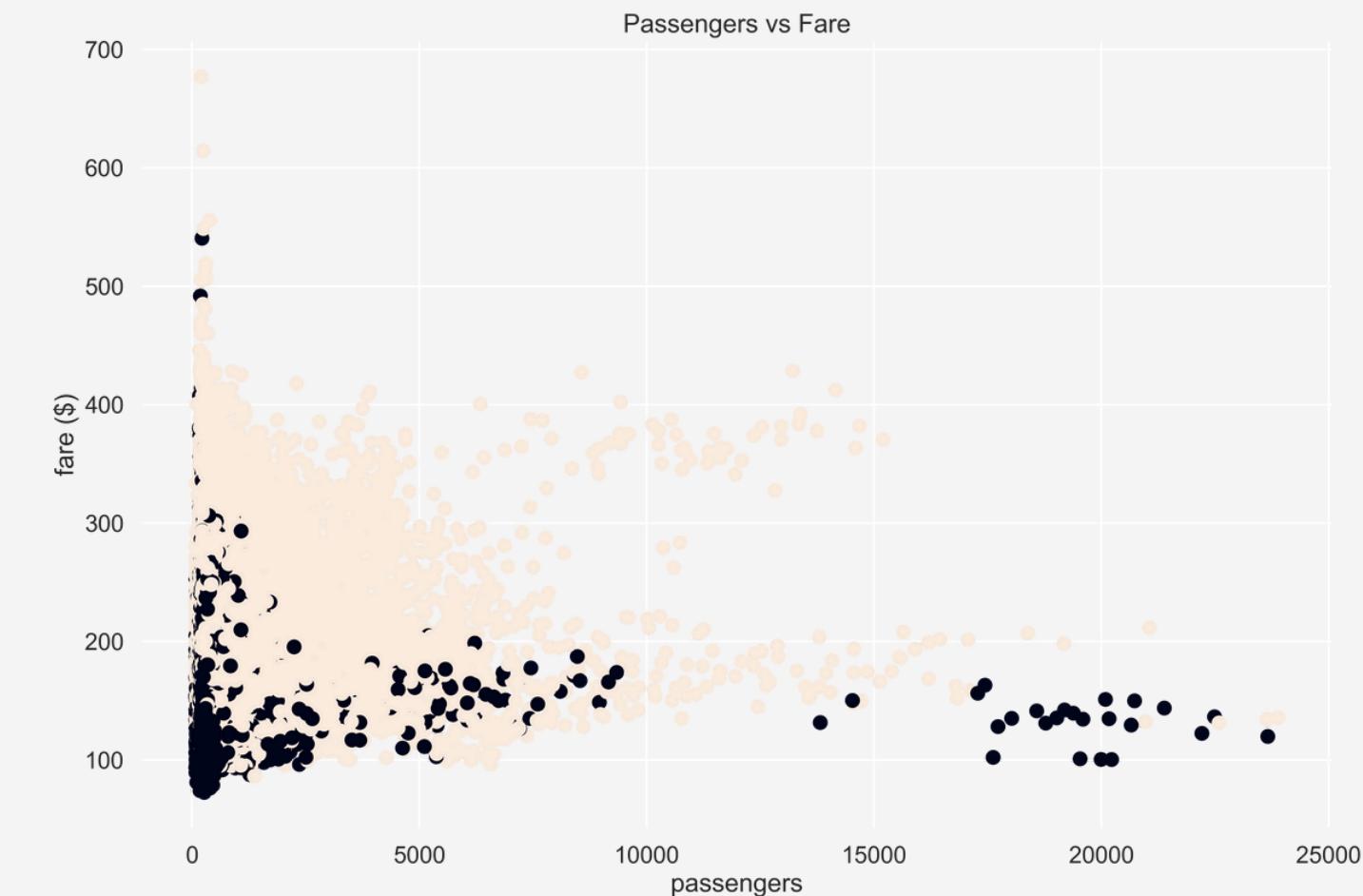
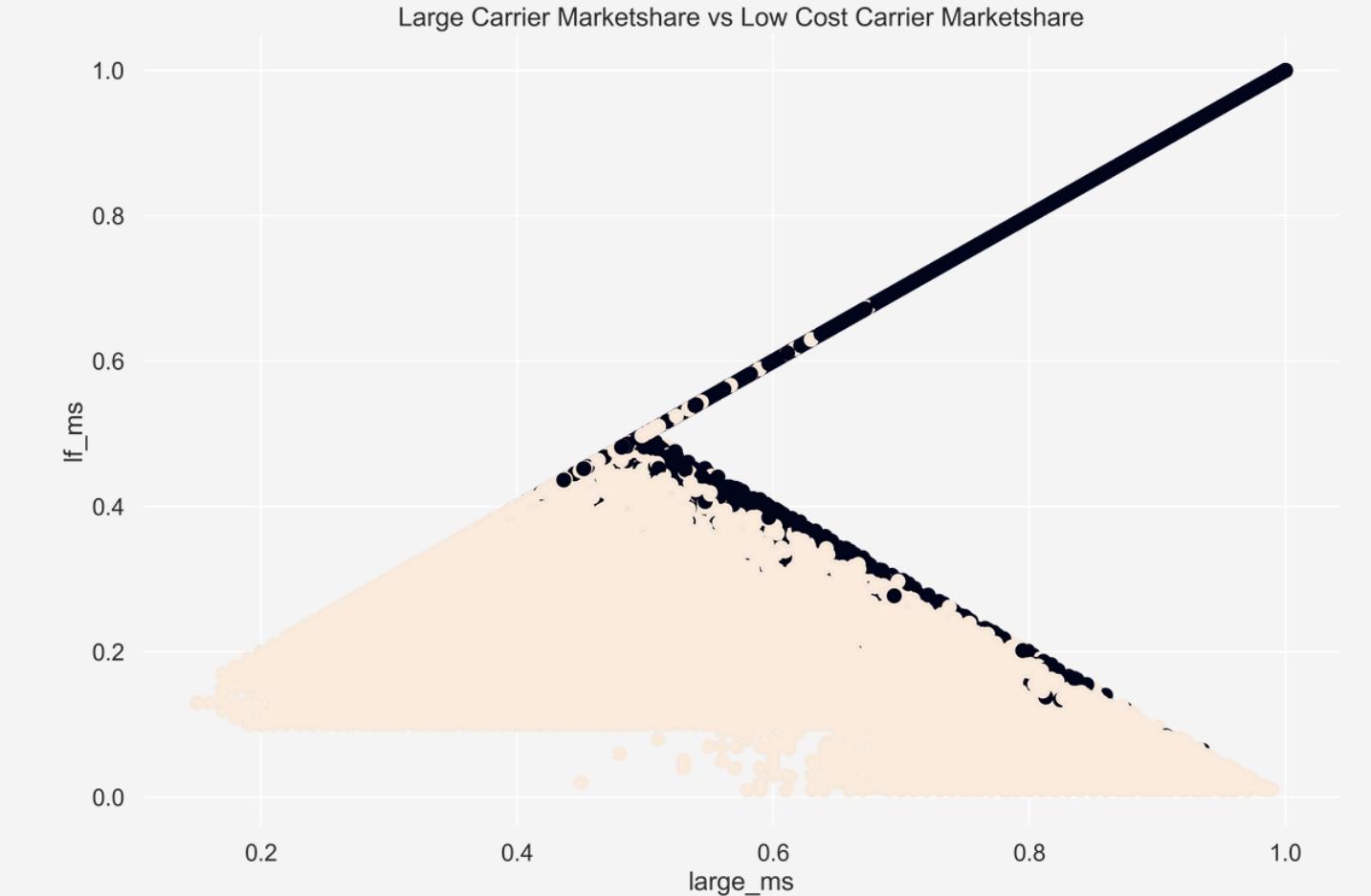
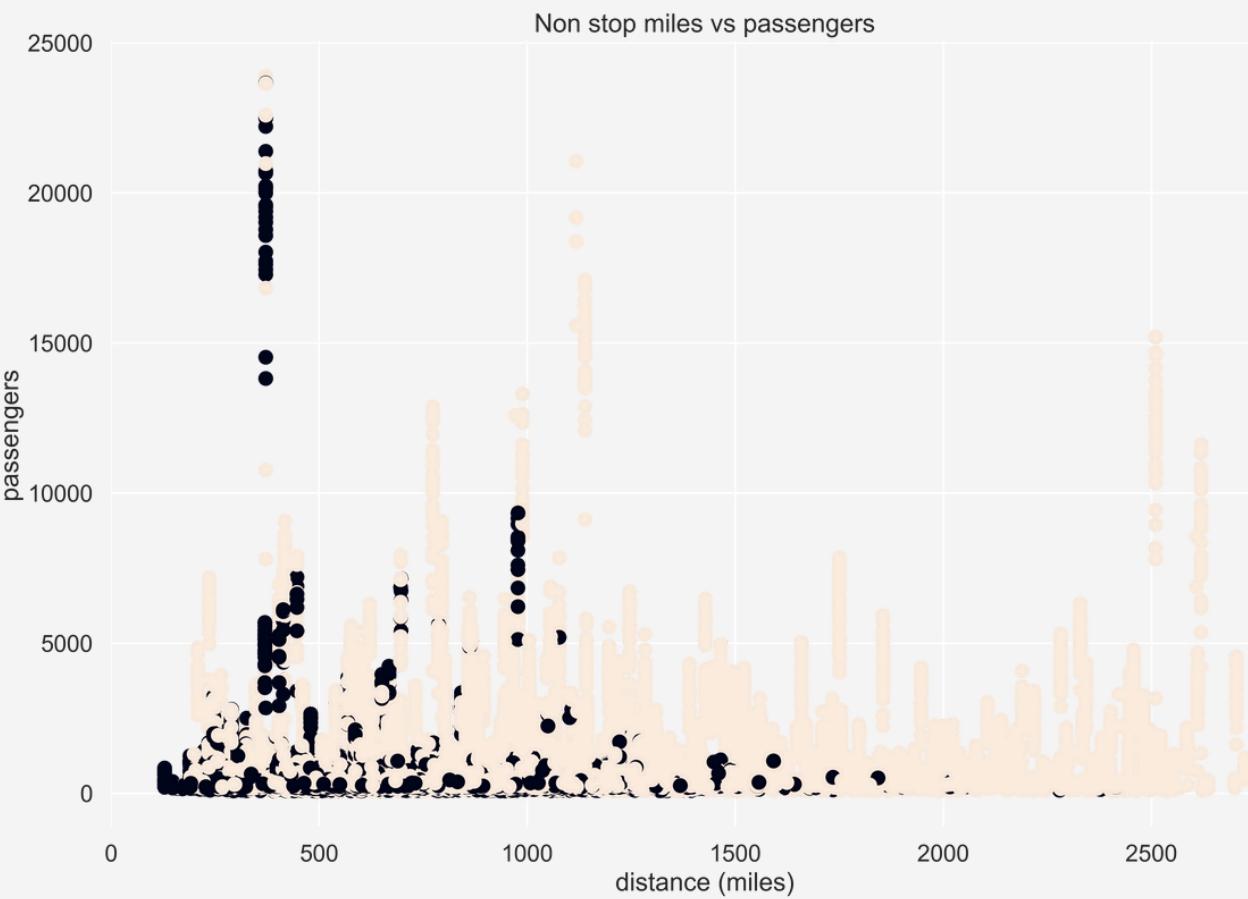


# 1st KMeans n\_clusters=2

The average silhouette score is:  
0.2901062302893411

The Davies-Bouldin score is:  
1.2307796042861823

'nsmiles',  
'passengers',  
'fare', 'large\_ms',  
'fare\_lg', 'lf\_ms',  
'fare\_low'

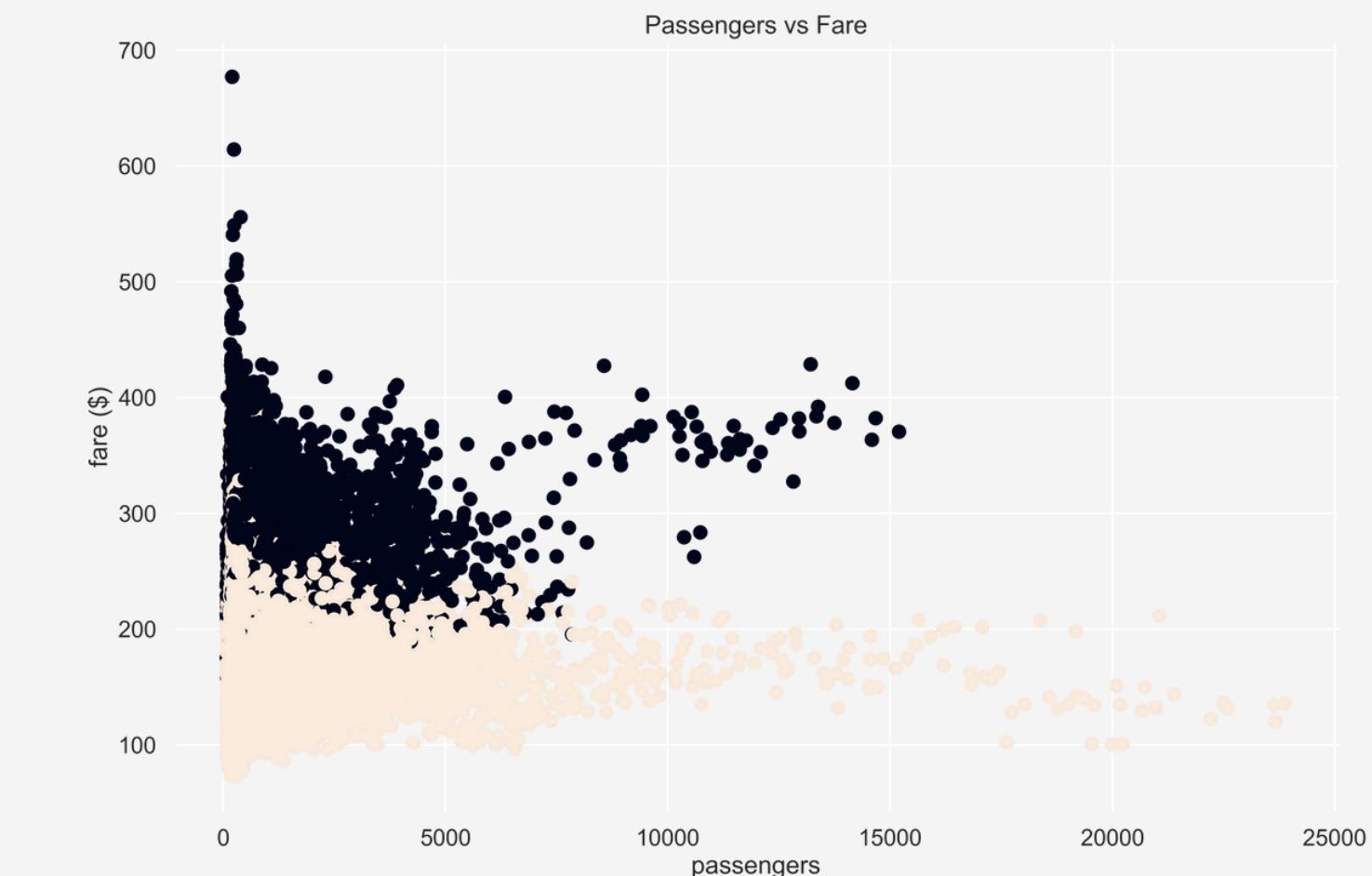
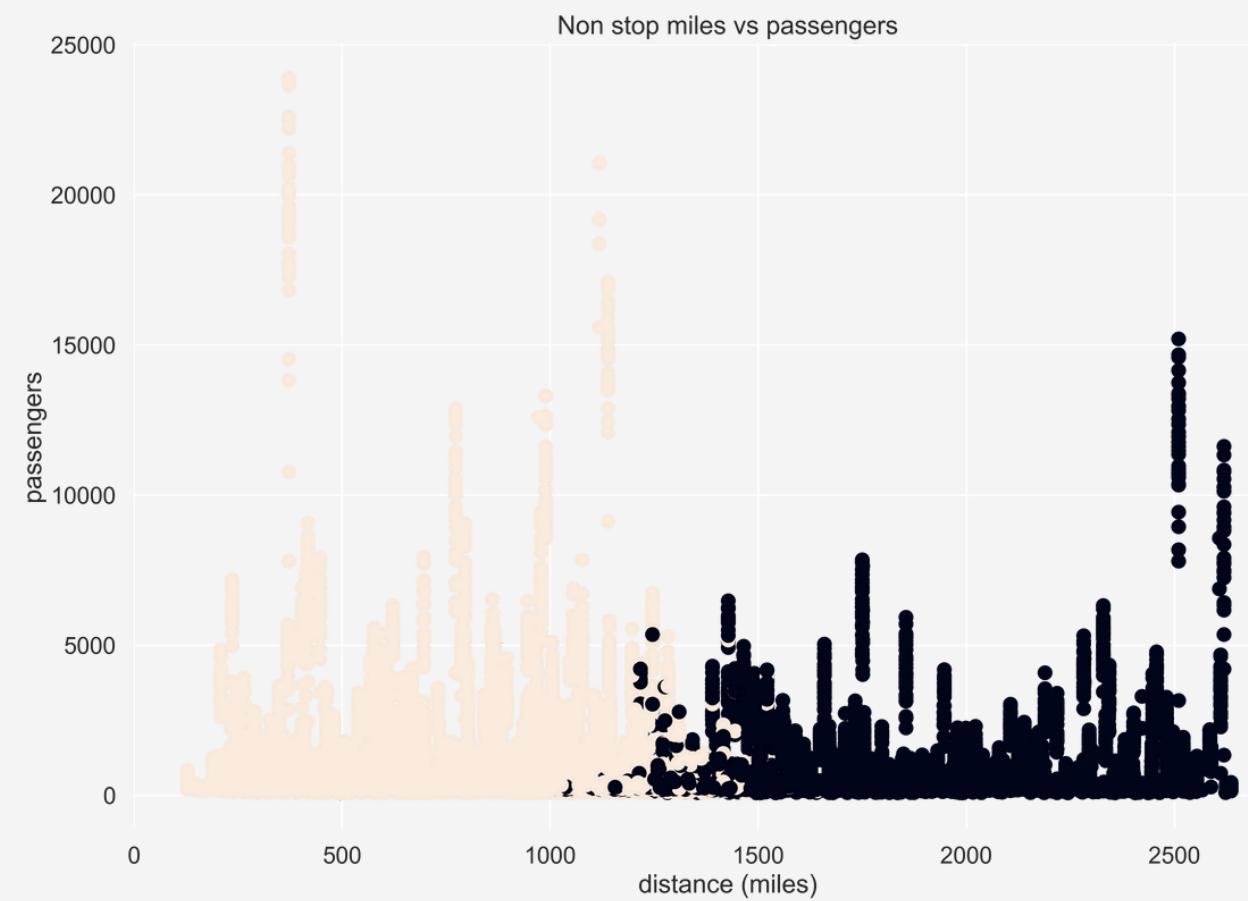


# 2nd KMeans n\_clusters=2 (No market share)

The average silhouette score is: 0.4897933628767886

The Davies-Bouldin score is: 0.797307522739757

'nsmiles',  
'passengers',  
'fare', 'fare\_lg',  
'fare\_low'

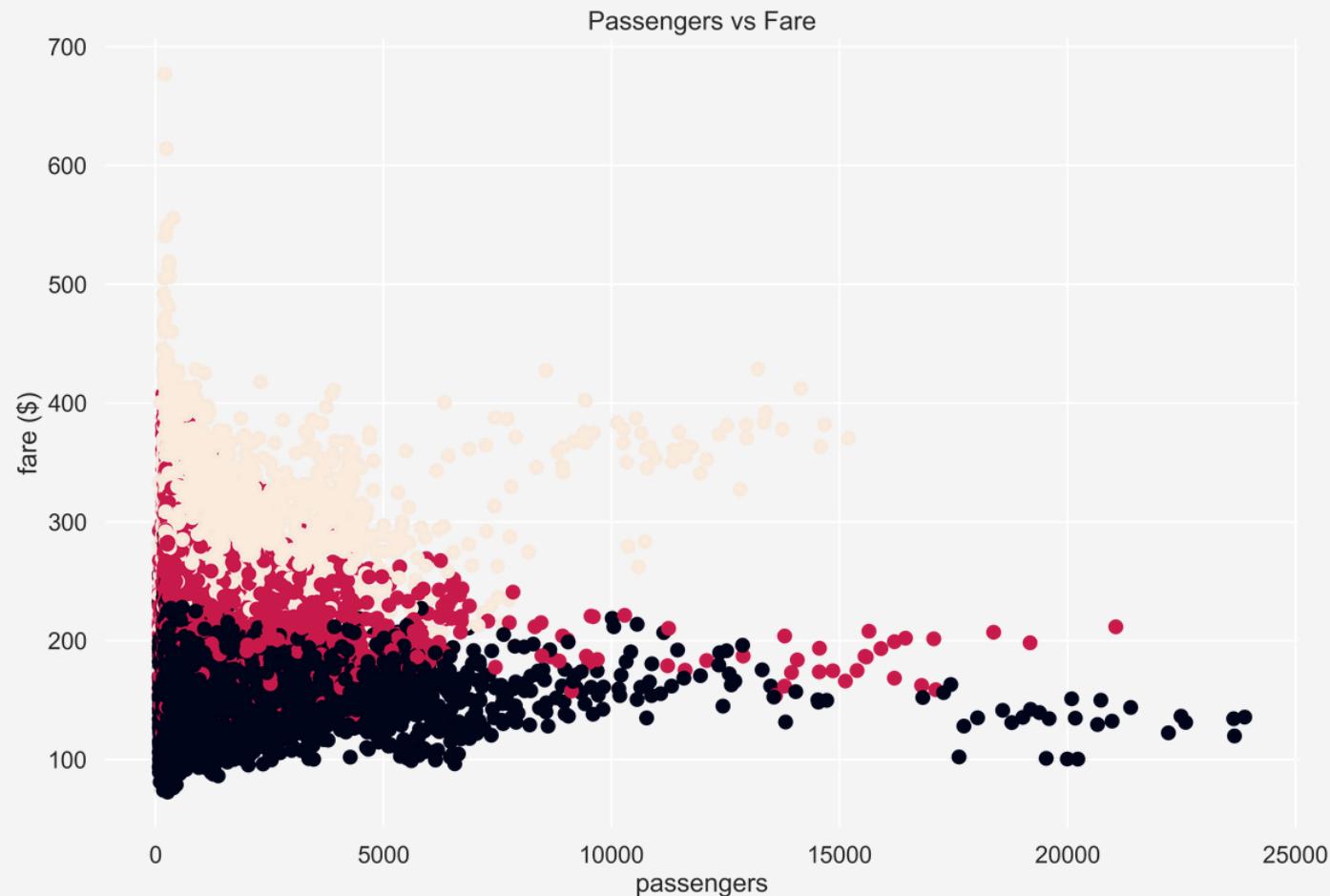
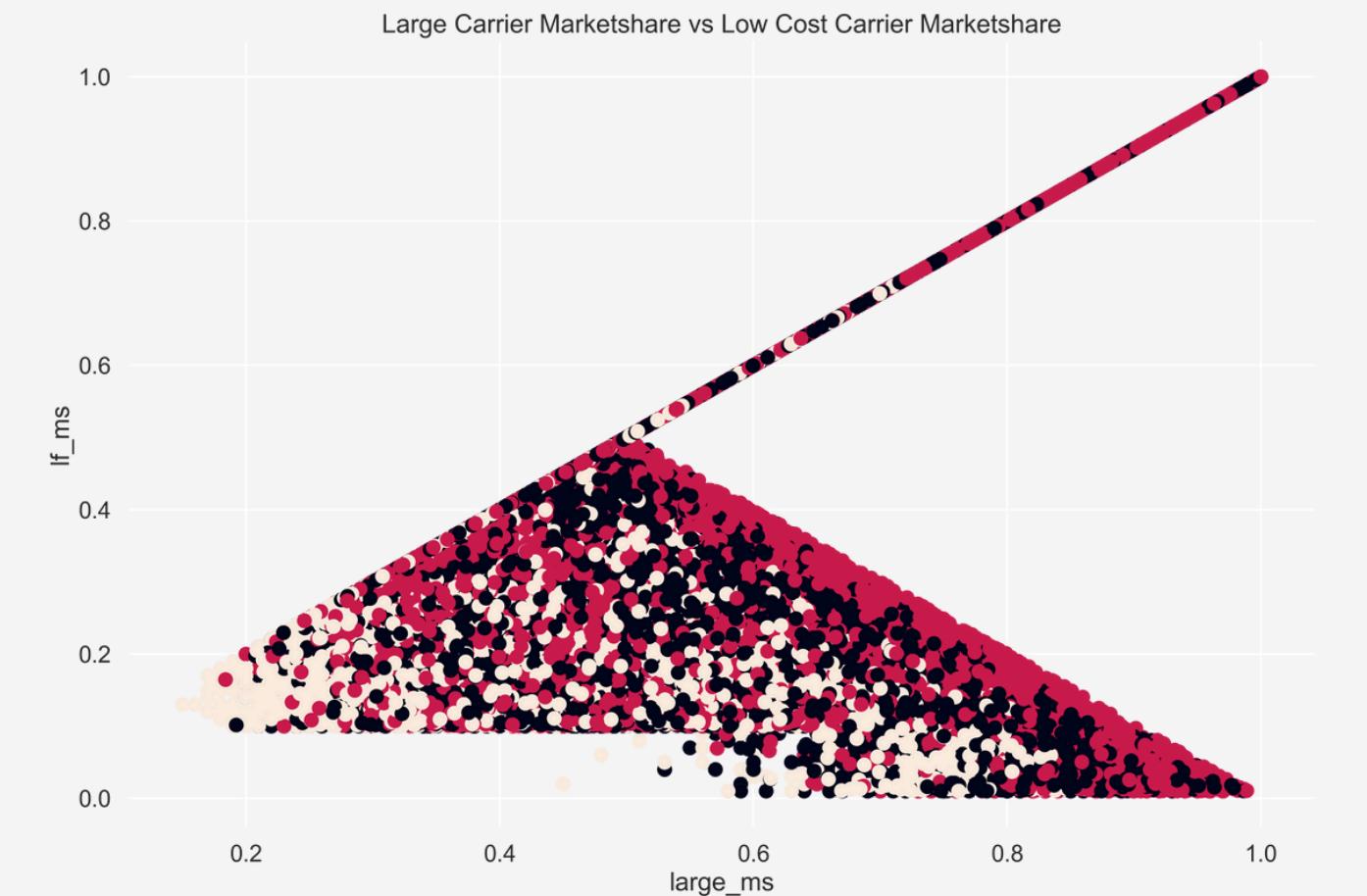
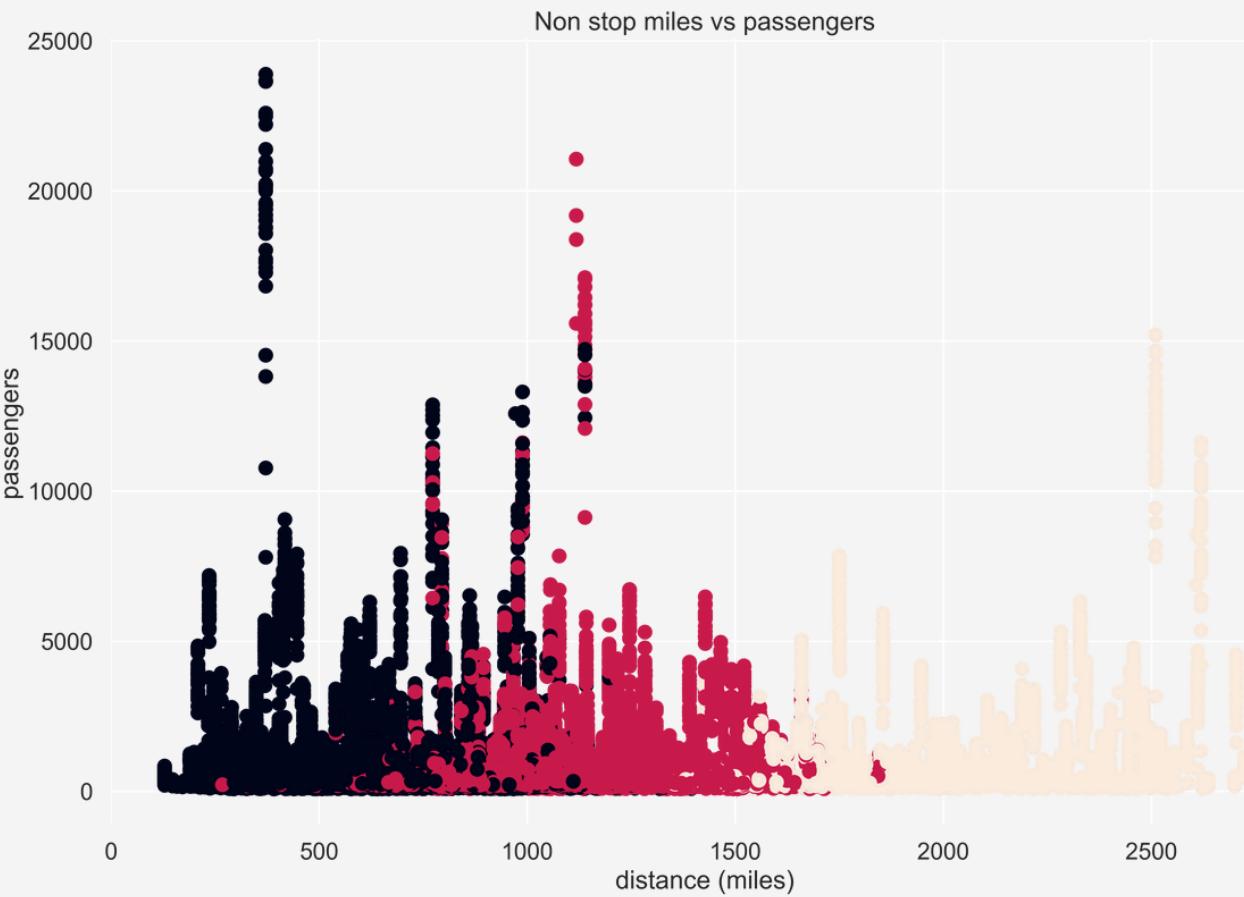


# 3rd KMeans n\_clusters=3 (No market share)

The average silhouette score is:  
0.32238207207898467

The Davies-Bouldin score is:  
1.1109552052058145

'nsmiles',  
'passengers',  
'fare', 'fare\_lg',  
'fare\_low'



# 2nd Agg. Hierarchical Clustering (Coordinates Only)

The average silhouette score is:

0.36295965705198735

The Davies-Bouldin score is:

0.9915806024341887

"city1\_lat",  
"city1\_long",  
"city2\_lat",  
"city2\_long"

