## Intuition:

This is a classification problem without ground labels (unsupervised). So the approach was to identify merchants which are similar to each other (using a clustering algorithm – here K means). Ideally, the card transaction price offered to these merchants should be similar. After identifying clusters of similar merchants, if the price offered to them is within 75%-25% of the prices offered to all the merchants in the same group, they will be assigned a label of 'normal' price. On the other hand, if their price is >75 or <25% of the prices offered to the merchants in their group, it will be classified as 'non-competitive' and 'competitive' respectively.

## Limitations of this study:

1. Since 'Accepts card payments' was used for clustering, the column 'Current Provider' was not utilized. In a future iteration of this exercise, it as well.
2. Since the MCC code column has a high granularity and the data does not mention what those codes mean (such as whether a particular code refers to the supermarkets and so on) , it did not make sense to use them to interpret the clusters.
3. The column "Annual Card Turnover" is ambiguous because merchants which currently do not accept card payments have a non-zero value for this column. Is it because they used to accept card payments in the past ?

## Implementation details

1. <u>Converting "Annual Card Turnover" and "Average Transaction Amount" into log scale</u>

The distribution of these metrics was heavily skewed which was impacting clustering performance (even after scaling) and resulting in clusters which were not meaningful.

This was solved by converting them to logarithmic scale.

2. <u>Aligning both parts of fees</u>

All the 'fee' columns had the following structure – "0.55% + 2p" where the first part is the fee applied as a proportion to each transaction and the second part is the flat fee applied to each transaction.

The flat fee was converted to a percentage:

Flat_fee_percentage = ((flat_fee / 100) / Average_transaction_amount) * 100

And summed together

3. <u>Creating a single fee metric</u>

A single fee metric "Fees (%)" was created from the weighted average of all the aligned fee columns based on the following assumptions:

- 40% of card payments are from Mastercard cards, and 60% are from Visa cards

- 90% of card payments are made using Debit cards, 8% are made using credit cards, and 2% are made using business debit cards

Fees (%) =
$(PMD \cdot mastercard\_prob \cdot debit\_prob) + (PVD \cdot visa\_prob \cdot debit\_prob) + (PMC \cdot mastercard\_prob \cdot credit\_prob) + (PVC \cdot visa\_prob \cdot credit\_prob) + (PMBD \cdot mastercard\_prob \cdot business\_debit\_prob) + (PVBD \cdot visa\_prob \cdot business\_debit\_prob)$

Where:

- $PMD P_{MD}$ PMD: Price for "Mastercard Debit"
- $PVD P_{VD}$ PVD: Price for "Visa Debit"
- $PMC P_{MC}$ PMC: Price for "Mastercard Credit"
- $PVC P_{VC}$ PVC: Price for "Visa Credit"
- $PMBD P_{MBD}$ PMBD: Price for "Mastercard Business Debit"
- $PVBD P_{VBD}$ PVBD: Price for "Visa Business Debit"

4. Clustering using K Means

K means clustering was performed on the training dataset (80% samples) with K=7.
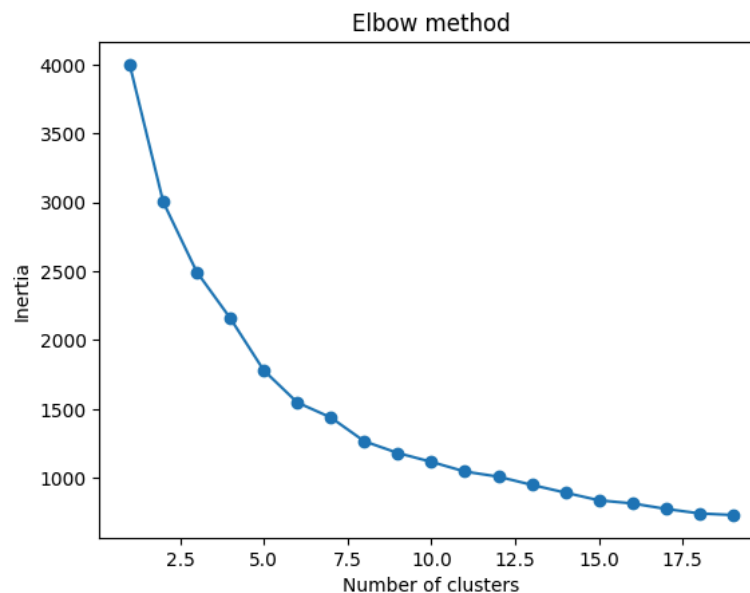


*Figure 1 Inertia vs k*

The clustering performance was evaluated using Silhouette Score and Davies-Bouldin Index.

A Silhouette score of +1 indicates samples within the same cluster are closer to each other and far away from those in other clusters and a score of -1 represents the opposite end of the spectrum. A score of 0 indicates that some samples lie on the decision boundary between two clusters. The score observed here was 0.33 which shows moderate clustering performance.

The DBI metric measures the ratio of the similarity of each cluster with that of its most similar cluster. A lower score represents a good clustering performance. The score observed for us was 1.18.

The hyperparameter (k=7) was selected based on the elbow-method, if it achieved a reasonable Silhouette Score and DB index score.
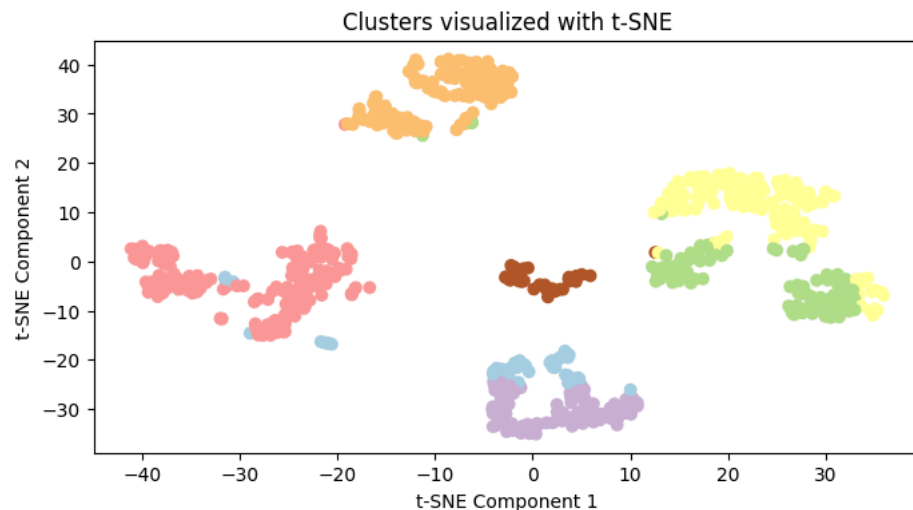


*Figure 2 Visualising the clusters using t-SNE*

## Interpreting Clusters

- Cluster 0
  - o Does not accept Card payments
  - o Is mostly registered
  - o Very high average transaction amount (>£ 1k)
  - o –
  - o **This is an interesting group because the average transaction amount for this group of merchants is very high, but they do not accept card payments. They are probably high value goods (automobiles) or similar merchants.**
- Cluster 1
  - o Accepts card payments
  - o Is registered
  - o Moderate average transaction amount (> £ 100)
  - o Average Annual card turnover (for those merchants which currently accept card payments) < £ 100k (or over 50k)
  - o **Companies which are registered and currently accept card payments. The average transaction size is >£ 100 so not groceries but the average annual card turnover is around £ 100k**
- Cluster 2
  - o Does not accept card payments
  - o Is not registered

- o Low average transaction amount (<£ 100)
- o –
- **o These are companies which are not registered but currently accept card payments and have low ticket size and annual card turnover. So, this group probably represents newly established small stores. It would be interesting to understand why they do not accept card payments yet.**

- Cluster 3
  - o Accepts card payments
  - o Is not registered
  - o Low average transaction amount (<£ 100)
  - o Average Annual card turnover (for those merchants which currently accept card payments) < £ 100k
  - **o These are companies which are not registered but currently accept card payments and have low ticket size and annual card turnover. So, this group probably represents newly established small stores which have not been making much revenue yet.**

- Cluster 4
  - o Does not accept card payments (currently)
  - o Is registered
  - o Low average transaction amount (<£ 100)
  - o –
  - **o These merchants are registered, have a moderate/low ticket size and do not currently accept card payments. It might be because the costs of accepting card payments outweigh the additional revenue, they bring in but we have to look deeper into this. We do not have geographical location about these merchants to understand where they are located.**

- Cluster 5
  - o Accepts card payments (currently)
  - o Is registered
  - o Low average transaction amount (<£ 100)
  - o Average Annual card turnover (for those merchants which currently accept card payments) over £ 100k
  - **o Companies which are registered and currently accept card payments (similar to cluster 1). The average transaction size is <£ 100 so regular stores but the average annual card turnover is over £ 100k. These are the kind of shops where the value of the items are not a lot but people shop more often.**

- Cluster 6
    - o May or may not be registered
    - o May or may not accept card payments
    - o Moderate average transaction amount (> £ 100)
    - o Average Annual card turnover (for those merchants which currently accept card payments) < £ 100k
    - o **Interesting group because it has a mix of merchants who accept and do not accept card payments and may or may not be registered. Otherwise, like cluster 1.**