# Finding a suitable location for 'Home from Home'

## Shirley Ng - August 2020

## 1. Introduction

### 1.1 Background

In the past year, Hong Kong has appeared on the world's newspaper headlines more often than the past decade combined. Unprecedented police brutality and violent conflicts in the street have unfortunately become part of everyday life. People in this former vibrant metropolitan city fear the loss of freedoms. Following the Chinese government's decision to impose a new National Security Law, the UK government has committed to open a new immigration route to British National (Overseas) passport holders in Hong Kong.

### 1.2 Problem

Despite the fact that many well-travelled Hong Kongers have been to London for holidays, it is still a challenge to learn everything about life in the UK for an uprooting move. There are numerous videos, social media posts, blogs and forums flowing with varying subjective opinions and information of mixed quality.

### 1.3 Interest

Selecting a suitable location to settle down in a 'home from home' is an important and complex decision. There is a need to analyse more recent and widely trusted data sources to provide an objective data-led view to help them make an informed decision. This project will cover these criteria: a) top common venues, b) property price, c) crime rate, d) employment rate and e) education.

## 2. Data acquisition and cleaning

### 2.1 Data sources

This project analyses 42 cities in England that have data across all criteria. Data from Four Square, government sources like Office of National Statistics and Department of Education as well as some websites are used.

List of data sources

a) List of English cities and the coordinates:
https://www.townscountiespostcodes.co.uk/cities-in-england/ and
https://www.latlong.net/category/cities-235-15.html

b) Top 10 most common venues: data of up to 100 venues within 500m radius of each city are extracted from Four Square data for analysis.

c) Property price: The Annual Price Change by Local Authority for England data is used to get property price in January 2020: https://www.gov.uk/government/publications/uk-house-price-index-england-january-2020/uk-house-price-index-england-january-2020

d) Crime rate: Recorded crime data by Community Safety Partnership area, year ending March 2020 is used: https://www.ons.gov.uk/peoplepopulationandcommunity/crimeandjustice/datasets/recordedcrimedatabycommunitysafetypartnershiparea

e) Employment rate: The employment rate during January to December 2019 from nomis that provides official labour market is used: https://www.nomisweb.co.uk/query/construct/submit.asp?menuopt=201&subcomp=

f) Education: To compare school performance, the 2018-2019 Final Key Stage 2 percentage of pupils meeting expected standard and Final Key Stage 4 Attainment 8 score are used. https://www.compare-school-performance.service.gov.uk/download-data?currentstep=datatypes&regiontype=all&la=0&downloadYear=2018-2019&datatypes=ks2&datatypes=ks4

**2.2 Data cleaning**

Data downloaded and scraped from webpages as listed above. Some of the data are processed by Microsoft Excel VLOOKUP function to extract data of relevant cities. Since several data sources are based on local authority as unit instead of city name, data cleaning involves amending some of the naming like 'Kingston upon Hull' to 'Hull'. The list started with 51 cities but eventually trimmed down to 42 cities to only keep those that have complete data across all criteria.

## 3. Methodology

K-means cluster segmentation is used to group the 42 cities into 5 clusters.

The analysis starts with Four Square data that up to 100 venues within 500m radius of each city are extracted. One Hot Encoding is used to convert venue category data into binary data.

| | City | Accessories Store | American Restaurant | Antique Shop | Argentinian Restaurant | Art Gallery | Art Museum | Asian Restaurant | Auto Garage | BBQ Joint | ... | University | Vegetarian / Vegan Restaurant | Video Game Store | Video Store | Vietnamese Restaurant | Wa St( |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Bath | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | Bath | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 1 | 0 | 0 | 0 | 0 |
| 2 | Bath | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | Bath | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | Bath | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | ... | 0 | 0 | 0 | 0 | 0 | 0 |

5 rows × 208 columns

Property price index of each city is relatively straightforward. Crime and employment data have to take each city's varying population into account, so recorded crime per 1,000 population/household and employment rate are used.

Education is relatively tricky because the widely referenced Ofsted ratings only have data by region/local authority that don't match well with the English city list. Also, there is news report questioning its reliability.[1] For Key Stage 2 performance, percentage of pupils meeting expected standard is used to compare primary school performance. For Key Stage 4 performance, Attainment 8 score is used to compare secondary school performance.

```python
# Group the data by 'City', aggregate and calculate the mean for each column
gp_ks2 = df_ed2.groupby('City').agg({'MeetExp':['mean']})
gp_ks2 = gp_ks2.reset_index()
gp_ks2.head()
```

```python
# Group the data by 'City', aggregate and calculate the mean for each column
gp_ks4 = df_ed4.groupby('City').agg({'ATT8SCR':['mean']})
gp_ks4 = gp_ks4.reset_index()
gp_ks4.shape
```

The One Hot Encoding result of nearby venues are combined with the above data regarding property price, crime rate, employment rate and education into one table for running k-means cluster segmentation.

| | City | Accessories Store | American Restaurant | Antique Shop | Argentinian Restaurant | Art Gallery | Art Museum | Asian Restaurant | Auto Garage | BBQ Joint | ... | Women's Store | Total crime | Violent crime | Theft offences | Property_J |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Bath | 0.0 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.000000 | 0.0 | 0.0 | ... | 0.0 | 65.1 | 20.5 | 23.5 | 3 |
| 1 | Birmingham | 0.0 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.000000 | 0.0 | 0.0 | ... | 0.0 | 100.5 | 35.7 | 36.2 | 1 |
| 2 | Bradford | 0.0 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.000000 | 0.0 | 0.0 | ... | 0.0 | 136.1 | 53.2 | 38.8 | 1 |
| 3 | Bristol | 0.0 | 0.010000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.020000 | 0.0 | 0.0 | ... | 0.0 | 114.1 | 34.7 | 40.3 | 2 |
| 4 | Cambridge | 0.0 | 0.011364 | 0.0 | 0.0 | 0.0 | 0.0 | 0.011364 | 0.0 | 0.0 | ... | 0.0 | 127.1 | 30.9 | 63.9 | 4 |

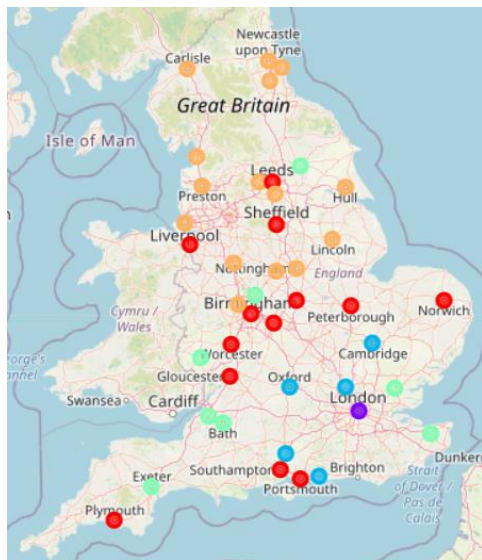Each cluster's numeric data are analysed using the **.describe(exclude=[object])** function.

```python
# Analyse Cluster 1 numeric data
cluster_1 = cities_merged.loc[cities_merged['Cluster Labels'] == 0, cities_merged.columns[[0]] + list(range(1, cities_merged.shape
cluster_1.describe(exclude=[object])
```

| | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) | Cluster Labels |
|---|---|---|---|---|---|---|---|---|
| count | 13.000000 | 13.000000 | 13.000000 | 13.000000 | 13.000000 | 13.000000 | 13.000000 | 13.0 |
| mean | 104.830769 | 37.723077 | 34.223077 | 196070.384615 | 0.745385 | 0.641454 | 44.036201 | 0.0 |
| std | 20.688579 | 7.715045 | 7.542453 | 15003.178977 | 0.049769 | 0.027792 | 2.631869 | 0.0 |
| min | 70.200000 | 25.500000 | 20.900000 | 172357.000000 | 0.640000 | 0.577407 | 40.733333 | 0.0 |
| 25% | 87.800000 | 33.700000 | 31.000000 | 184305.000000 | 0.720000 | 0.635331 | 42.357143 | 0.0 |
| 50% | 104.200000 | 36.500000 | 36.200000 | 199105.000000 | 0.740000 | 0.642636 | 43.989474 | 0.0 |
| 75% | 118.400000 | 46.200000 | 38.400000 | 210018.000000 | 0.770000 | 0.658000 | 44.539024 | 0.0 |
| max | 135.800000 | 48.800000 | 45.000000 | 213279.000000 | 0.840000 | 0.691111 | 50.700000 | 0.0 |

---

[1] BBC News "Ofsted inspection grades challenged" https://www.bbc.co.uk/news/education-47816631

## 4. Results

42 English cities are segmented into 5 clusters as shown in the map and list below:



**Cluster 1**: Birmingham, Chester, Coventry, Gloucester, Leeds, Leicester, Norwich, Peterborough, Plymouth, Portsmouth, Sheffield, Southampton, Worcester.

**Cluster 2**: London.

**Cluster 3**: Cambridge, Chichester, Oxford, St Albans, Winchester.

**Cluster 4**: Bath, Bristol, Canterbury, Chelmsford, Exeter, Hereford, Lichfield, York.

**Cluster 5**: Bradford, Carlisle, Derby, Durham, Hull, Lancaster, Lincoln, Liverpool, Newcastle upon Tyne, Nottingham, Preston, Stoke-on-Trent, Sunderland, Wakefield, Wolverhampton.

See Appendix section for more details about each cluster.
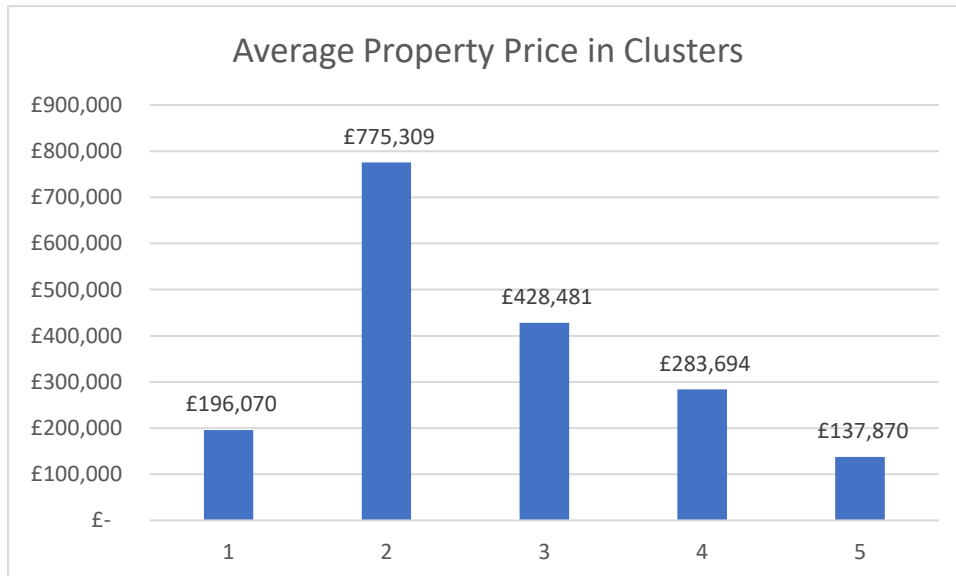
## 5. Discussion

### 5.1 Safety

Crime statistics used in this project is based on recorded crime per 1,000 population/ household. Clusters have mean total crime below 100 are clusters 4 and 3. The lowest total crime recorded cities are Hereford (56) and Lichfield (57), where both are inside cluster 4. The highest total crime recorded cities are Hull (157) and Lincoln (150), where both are inside cluster 5.

Aligning with the ranking of total crime, clusters 5 and 1 are also top 2 in violent crime. Cluster 2 is the third in total crime, but it is top in theft offences; followed by cluster 3 that is the fourth in total crime but second in theft offences.
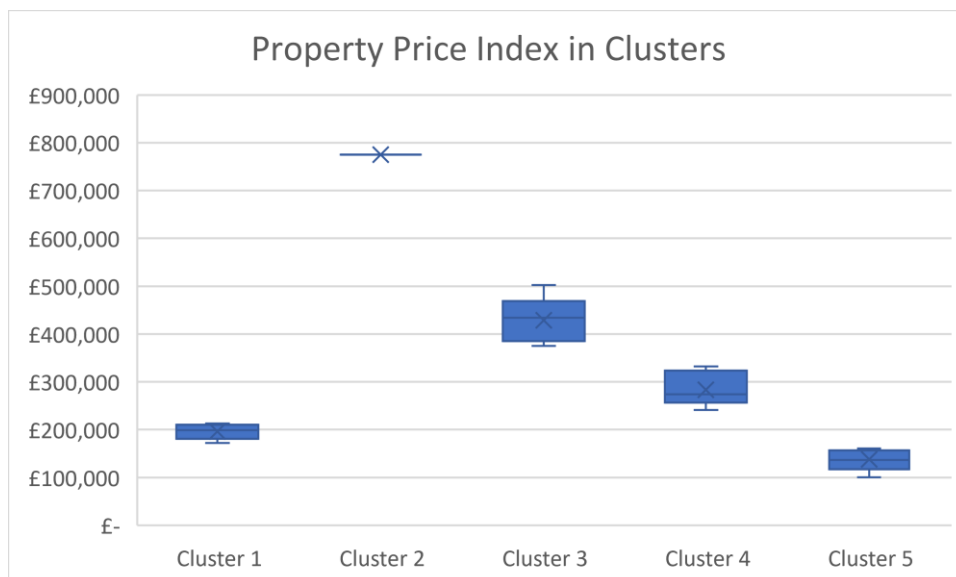
**5.2 Property Price**

Property price index in London (cluster 2) at £775,309 is over 5 times higher than the average of cluster 5, which has the lowest property price index at £137,870; and over 7 times higher than Durham, which is a city with the lowest property price index at £100,643.



The box and whisker plot shows the price range in clusters 3 and 4 are wider than the other clusters. Property price index range from £375k in Chichester to £502k in St. Albans within cluster 3; and range from £241k in Hereford to £332k in Chelmsford within cluster 4.



**5.3 Employment Rate**

Clusters 3 and 4 have the top 2 employment rate at 81% and 78%; while cluster 5 has the lowest employment rate at 71%. City with the highest employment rate is Chichester at 87% in cluster 3 and the lowest employment rate is Bradford at 63% in cluster 5.

Employment Rate in Cluster

### 5.4 Education

Looking at the cluster's mean primary education percentage of pupils meeting expected standard, it ranges from 64% in cluster 1 to 72% in cluster 2. The box and whisker plot reveal a huge disparity in cluster 3 with a mean value of 67% but ranging from Chichester at 60% to Winchester at 76%. The outlier in cluster 1 showing 57% is Portsmouth.
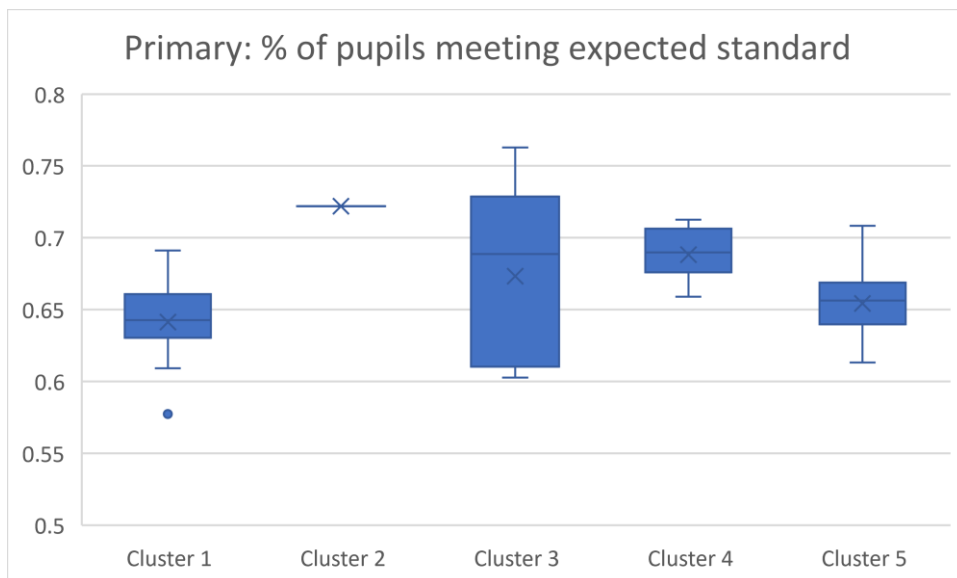


Primary: % of pupils meeting expected standard

For secondary education performance, Attainment 8 score is used. The cluster's mean values range from 43.9 in cluster 5 to 46.6 in cluster 2. Once again, the box and whisker plot reveal the disparity in clusters 4 that score range from 39.4 in Bath to 54.8 in Chelmsford. Cluster 5 has equally long whisker that scores range from the lowest 35.1 in Newcastle upon Tyne to second highest at 51.1 in Lancaster. Gloucester is an outlier in cluster 1 with 50.7.

## Secondary: Attainment 8 Score

**5.5 Top 10 Most Common Venues**

Interesting to see 'Platform' from train station being the first most common venue in Peterborough and Plymouth within cluster 1. Cluster 2 is London that has 'Theatre' as its first most common venue. Pub, coffee shop and café are the first most common venue across clusters 3 and 4. Cluster 5 has 15 cities in it is the biggest, which also has a more diverse profile in its first most common venue. Apart from the popular pub and coffee shop, it has clothing store, supermarket, platform, hotel, discount store and Indian restaurant taking up the top most common venue.

## 6. Conclusion

The k-means cluster segmentation confirms that London is 'one of a kind'. Therefore, the correlation analysis is performed across all clusters, as well as excluding London to reveal any correlations hidden by the unique London data.

As shown in the table below, high violent crime is a strong predictor for low property price, low employment rate and lower secondary school performance. In contrast, employment rate and property price index are positively correlated outside of London. No surprise in the strong positive correlation between primary and secondary school performance as well.

| | Total crime | Violent crime | Theft offences | Property price | Employment rate | Primary Education | Secondary Education |
|---|---|---|---|---|---|---|---|
| **Property** | | | | | | | |
| All | -18% | -76% | 80% | - | 28% | 89% | 84% |
| Excl. LDN | -79% | -96% | 20% | - | 98% | 61% | 79% |
| **Employment** | | | | | | | |
| All | -89% | -82% | -12% | 28% | - | 25% | 56% |
| Excl. LDN | -89% | -99% | 1% | 98% | - | 66% | 86% |
| **Pri. Edu** | | | | | | | |
| All | -35% | -75% | 53% | 89% | 25% | - | 94% |
| Excl. LDN | -86% | -78% | -46% | 61% | 66% | - | 94% |
| **Sec. Edu** | | | | | | | |
| All | -64% | -91% | 36% | 84% | 56% | 94% | - |
| Excl. LDN | -98% | -92% | -39% | 79% | 86% | 94% | - |

## Appendix – Clusters numeric data analysis

**Cluster 1**: Birmingham, Chester, Coventry, Gloucester, Leeds, Leicester, Norwich, Peterborough, Plymouth, Portsmouth, Sheffield, Southampton, Worcester.

| | City | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) | Cluster Labels | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Birmingham | 100.5 | 35.7 | 36.2 | 189161.0 | 0.640017 | 0.635331 | 45.800971 | 0 | Pub | Pizza Place | Brewery | Sandwich Place | |
| 9 | Chester | 85.2 | 33.7 | 20.9 | 209431.0 | 0.800380 | 0.658000 | 46.454545 | 0 | Pub | Hotel | Restaurant | Historic Site | |
| 11 | Coventry | 77.8 | 25.5 | 31.0 | 184305.0 | 0.732214 | 0.638191 | 40.855556 | 0 | Pub | Coffee Shop | Café | Clothing Store | |
| 16 | Gloucester | 104.2 | 37.1 | 36.0 | 210982.0 | 0.840471 | 0.646792 | 50.700000 | 0 | Pub | Coffee Shop | Sandwich Place | Pharmacy | |
| 20 | Leeds | 130.6 | 46.2 | 45.0 | 191578.0 | 0.712668 | 0.625380 | 43.167568 | 0 | Coffee Shop | Café | Hotel | Bar | |
| 21 | Leicester | 115.7 | 40.5 | 36.6 | 177556.0 | 0.720974 | 0.664823 | 44.181818 | 0 | Coffee Shop | Pub | Bar | Indian Restaurant | |
| 28 | Norwich | 135.8 | 48.8 | 40.0 | 202582.0 | 0.774290 | 0.642636 | 42.357143 | 0 | Pub | Coffee Shop | Italian Restaurant | Café | |
| 31 | Peterborough | 109.8 | 36.5 | 38.2 | 199105.0 | 0.732090 | 0.609221 | 42.073684 | 0 | Platform | Clothing Store | Coffee Shop | Pub | |
| 32 | Plymouth | 87.8 | 36.5 | 22.4 | 176390.0 | 0.739620 | 0.640411 | 43.989474 | 0 | Platform | Coffee Shop | Restaurant | Grocery Store | |
| 33 | Portsmouth | 118.4 | 46.3 | 35.1 | 212171.0 | 0.717864 | 0.577407 | 43.130000 | 0 | Fast Food Restaurant | Pub | Sporting Goods Shop | Supermarket | |
| 38 | Sheffield | 99.0 | 30.0 | 38.4 | 172357.0 | 0.738864 | 0.646174 | 44.539024 | 0 | Bar | Café | Hotel | Pub | |
| 39 | Southampton | 127.8 | 46.9 | 41.3 | 210018.0 | 0.745385 | 0.663418 | 44.487500 | 0 | Coffee Shop | Bar | Pub | Grocery Store | |
| 49 | Worcester | 70.2 | 26.7 | 23.8 | 213279.0 | 0.802099 | 0.691111 | 40.733333 | 0 | Pub | Coffee Shop | Bar | Café | |

| | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) |
|---|---|---|---|---|---|---|---|
| count | 13.000000 | 13.000000 | 13.000000 | 13.000000 | 13.000000 | 13.000000 | 13.000000 |
| mean | 104.830769 | 37.723077 | 34.223077 | 196070.384615 | 0.745385 | 0.641454 | 44.036201 |
| std | 20.688579 | 7.715045 | 7.542453 | 15003.178977 | 0.049769 | 0.027792 | 2.631869 |
| min | 70.200000 | 25.500000 | 20.900000 | 172357.000000 | 0.640000 | 0.577407 | 40.733333 |
| 25% | 87.800000 | 33.700000 | 31.000000 | 184305.000000 | 0.720000 | 0.635331 | 42.357143 |
| 50% | 104.200000 | 36.500000 | 36.200000 | 199105.000000 | 0.740000 | 0.642636 | 43.989474 |
| 75% | 118.400000 | 46.200000 | 38.400000 | 210018.000000 | 0.770000 | 0.658000 | 44.539024 |
| max | 135.800000 | 48.800000 | 45.000000 | 213279.000000 | 0.840000 | 0.691111 | 50.700000 |

**Cluster 2**: London. Since there is only one city in this cluster, no analysis is required.

| | City | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) | Cluster Labels | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 25 | London | 101.6 | 24.8 | 50.5 | 775309.0 | 0.74506 | 0.721974 | 46.574713 | 1 | Theater | Burger Joint | Coffee Shop | Pub |

**Cluster 3**: Cambridge, Chichester, Oxford, St Albans, Winchester.

| | City | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) | Cluster Labels | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | Cambridge | 127.1 | 30.9 | 63.9 | 435174.0 | 0.809375 | 0.688630 | 48.104762 | 2 | Pub | Coffee Shop | Clothing Store | Café |
| 10 | Chichester | 61.0 | 18.1 | 25.2 | 375252.0 | 0.874816 | 0.602800 | 43.700000 | 2 | Pub | Italian Restaurant | Coffee Shop | Clothing Store |
| 30 | Oxford | 109.4 | 27.6 | 55.0 | 433918.0 | 0.794011 | 0.617674 | 42.600000 | 2 | Coffee Shop | Pub | Café | Restaurant |
| 40 | St Albans | 61.1 | 17.8 | 27.0 | 502294.0 | 0.790497 | 0.694286 | 46.376923 | 2 | Coffee Shop | Pub | Sandwich Place | French Restaurant |
| 47 | Winchester | 67.7 | 21.2 | 27.4 | 395768.0 | 0.781333 | 0.762857 | 47.280000 | 2 | Pub | Bakery | Coffee Shop | Clothing Store |

| | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) |
|---|---|---|---|---|---|---|---|
| count | 5.00000 | 5.000000 | 5.00000 | 5.000000 | 5.000000 | 5.000000 | 5.000000 |
| mean | 85.26000 | 23.120000 | 39.70000 | 428481.200000 | 0.670000 | 0.673249 | 45.612337 |
| std | 30.87852 | 5.869157 | 18.32048 | 48544.265874 | 0.173781 | 0.064731 | 2.361624 |
| min | 61.00000 | 17.800000 | 25.20000 | 375252.000000 | 0.470000 | 0.602800 | 42.600000 |
| 25% | 61.10000 | 18.100000 | 27.00000 | 395768.000000 | 0.490000 | 0.617674 | 43.700000 |
| 50% | 67.70000 | 21.200000 | 27.40000 | 433918.000000 | 0.790000 | 0.688630 | 46.376923 |
| 75% | 109.40000 | 27.600000 | 55.00000 | 435174.000000 | 0.790000 | 0.694286 | 47.280000 |
| max | 127.10000 | 30.900000 | 63.90000 | 502294.000000 | 0.810000 | 0.762857 | 48.104762 |

**Cluster 4**: Bath, Bristol, Canterbury, Chelmsford, Exeter, Hereford, Lichfield, York.

| | City | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) | Cluster Labels | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Bath | 65.1 | 20.5 | 23.5 | 330975.0 | 0.731638 | 0.692414 | 39.353846 | 3 | Pub | Coffee Shop | Café | Cocktail Bar |
| 4 | Bristol | 114.1 | 34.7 | 40.3 | 285296.0 | 0.778697 | 0.658955 | 43.990566 | 3 | Café | Pub | Coffee Shop | Italian Restaurant |
| 6 | Canterbury | 97.8 | 36.8 | 31.6 | 302100.0 | 0.731429 | 0.711034 | 44.020000 | 3 | Pub | Coffee Shop | Café | Italian Restaurant |
| 8 | Chelmsford | 92.0 | 34.6 | 28.5 | 331973.0 | 0.815789 | 0.676327 | 54.772727 | 3 | Pub | Italian Restaurant | Bar | Department Store |
| 15 | Exeter | 74.7 | 28.4 | 21.3 | 262342.0 | 0.773933 | 0.676000 | 44.990000 | 3 | Pub | Café | Clothing Store | Tea Room |
| 17 | Hereford | 56.3 | 23.1 | 17.5 | 241217.0 | 0.843224 | 0.692381 | 45.530000 | 3 | Pub | Clothing Store | Coffee Shop | Café |
| 22 | Lichfield | 57.0 | 21.2 | 21.6 | 259701.0 | 0.793831 | 0.712500 | 49.180000 | 3 | Pub | Coffee Shop | Bar | Pharmacy |
| 50 | York | 65.3 | 21.8 | 25.7 | 255955.0 | 0.770476 | 0.687097 | 46.160870 | 3 | Pub | Café | Bar | Cocktail Bar |

| | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) |
|---|---|---|---|---|---|---|---|
| count | 8.000000 | 8.000000 | 8.000000 | 8.000000 | 8.000000 | 8.000000 | 8.000000 |
| mean | 77.787500 | 27.637500 | 26.250000 | 283694.875000 | 0.661250 | 0.688338 | 45.999751 |
| std | 21.184121 | 6.865428 | 7.187688 | 34873.284164 | 0.157701 | 0.018100 | 4.479630 |
| min | 56.300000 | 20.500000 | 17.500000 | 241217.000000 | 0.460000 | 0.658955 | 39.353846 |
| 25% | 63.075000 | 21.650000 | 21.525000 | 258764.500000 | 0.485000 | 0.676245 | 44.012642 |
| 50% | 70.000000 | 25.750000 | 24.600000 | 273819.000000 | 0.750000 | 0.689739 | 45.260000 |
| 75% | 93.450000 | 34.625000 | 29.275000 | 309318.750000 | 0.772500 | 0.697069 | 46.915652 |
| max | 114.100000 | 36.800000 | 40.300000 | 331973.000000 | 0.820000 | 0.712500 | 54.772727 |

**Cluster 5**: Bradford, Carlisle, Derby, Durham, Hull, Lancaster, Lincoln, Liverpool, Newcastle upon Tyne, Nottingham, Preston, Stoke-on-Trent, Sunderland, Wakefield, Wolverhampton.

| | City | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) | Cluster Labels | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | Bradford | 136.1 | 53.2 | 38.8 | 133763.0 | 0.627727 | 0.617545 | 41.223529 | 4 | Clothing Store | Hotel | Coffee Shop | Bakery |
| 7 | Carlisle | 103.3 | 41.9 | 22.5 | 136912.0 | 0.810924 | 0.661795 | 36.214286 | 4 | Supermarket | Coffee Shop | Rental Car Location | Furniture / Home Store |
| 12 | Derby | 111.0 | 43.0 | 33.9 | 159403.0 | 0.749548 | 0.628706 | 43.481818 | 4 | Platform | Coffee Shop | Pub | Clothing Store |
| 13 | Durham | 97.8 | 39.6 | 27.7 | 100643.0 | 0.737291 | 0.708444 | 47.711111 | 4 | Coffee Shop | Café | Pub | Italian Restaurant |
| 18 | Hull | 158.0 | 55.4 | 51.3 | 113565.0 | 0.717155 | 0.665632 | 41.962500 | 4 | Coffee Shop | Café | Grocery Store | Gay Bar |
| 19 | Lancaster | 82.9 | 31.5 | 27.4 | 151609.0 | 0.693723 | 0.652174 | 51.050000 | 4 | Pub | Clothing Store | Sandwich Place | Coffee Shop |
| 23 | Lincoln | 149.6 | 50.7 | 48.3 | 156640.0 | 0.755014 | 0.613125 | 42.360000 | 4 | Pub | Bar | Café | Hotel |
| 24 | Liverpool | 121.8 | 39.4 | 37.7 | 131051.0 | 0.683695 | 0.639778 | 44.059184 | 4 | Coffee Shop | Hotel | Bar | Café |
| 27 | Newcastle upon Tyne | 122.6 | 36.0 | 43.2 | 160730.0 | 0.667205 | 0.642500 | 35.100000 | 4 | Hotel | Pub | Restaurant | Bar |
| 29 | Nottingham | 129.5 | 39.8 | 46.6 | 151731.0 | 0.639210 | 0.668883 | 46.951064 | 4 | Pub | Bar | Coffee Shop | Café |
| 34 | Preston | 88.6 | 32.4 | 30.5 | 129238.0 | 0.812950 | 0.670174 | 47.270833 | 4 | Discount Store | Hotel | Burger Joint | Pizza Place |
| 41 | Stoke-on-Trent | 109.8 | 42.8 | 34.4 | 113908.0 | 0.740762 | 0.653545 | 44.192000 | 4 | Indian Restaurant | Pub | Sandwich Place | Supermarket |
| 42 | Sunderland | 111.1 | 34.1 | 35.5 | 117032.0 | 0.676728 | 0.664565 | 44.564286 | 4 | Coffee Shop | Clothing Store | Fast Food Restaurant | Pharmacy |
| 44 | Wakefield | 119.0 | 45.6 | 36.9 | 152966.0 | 0.716120 | 0.656275 | 47.769231 | 4 | Clothing Store | Rock Club | Fast Food Restaurant | Brewery |
| 48 | Wolverhampton | 94.9 | 35.7 | 33.6 | 158871.0 | 0.678917 | 0.670000 | 44.419048 | 4 | Supermarket | Women's Store | Fish & Chips Shop | Fast Food Restaurant |

| | Total crime | Violent crime | Theft offences | Property_Jan_2020 | Employment rate | (MeetExp, mean) | (ATT8SCR, mean) |
|---|---|---|---|---|---|---|---|
| count | 15.000000 | 15.000000 | 15.000000 | 15.000000 | 15.000000 | 15.000000 | 15.000000 |
| mean | 115.733333 | 41.406667 | 36.553333 | 137870.800000 | 0.679333 | 0.654209 | 43.888593 |
| std | 21.440804 | 7.308944 | 8.127806 | 19842.520862 | 0.107602 | 0.023894 | 4.243339 |
| min | 82.900000 | 31.500000 | 22.500000 | 100643.000000 | 0.440000 | 0.613125 | 35.100000 |
| 25% | 100.550000 | 35.850000 | 32.050000 | 123135.000000 | 0.655000 | 0.641139 | 42.161250 |
| 50% | 111.100000 | 39.800000 | 35.500000 | 136912.000000 | 0.680000 | 0.656275 | 44.192000 |
| 75% | 126.050000 | 44.300000 | 41.000000 | 154803.000000 | 0.745000 | 0.667258 | 47.110949 |
| max | 158.000000 | 55.400000 | 51.300000 | 160730.000000 | 0.810000 | 0.708444 | 51.050000 |