# Question 1

```
In [1]:   1  import pandas as pd
          2  from sklearn.preprocessing import StandardScaler
          3  from sklearn.cluster import KMeans
          4  import numpy as np
          5  import matplotlib.pyplot as plt
```

## Import csv into dataframe

```
In [2]:   1  df = pd.read_csv("spenddata.csv")
          2  df.head()
```

Out[2]:

| | Unnamed: 0 | month | var8 | var6 | a.1 | a.2 | a.3 | a.4 | var5 | b.5 | ... | c.276 | c.277 | c.278 | c.279 | c. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 2.0 | NaN | 1 | 5 | 1 | 57 | 34 | 1 | ... | 1 | 0 | 0 | 0 | |
| 1 | 2 | 1 | 2.0 | NaN | 1 | 4 | 1 | 57 | 34 | 2 | ... | 0 | 0 | 0 | 0 | |
| 2 | 3 | 1 | 2.0 | NaN | 1 | 5 | 1 | 57 | 42 | 2 | ... | 0 | 0 | 0 | 0 | |
| 3 | 4 | 1 | 2.0 | NaN | 1 | 6 | 1 | 57 | 34 | 2 | ... | 0 | 0 | 0 | 0 | |
| 4 | 5 | 1 | 2.0 | NaN | 1 | 8 | 1 | 22 | 1 | 1 | ... | 0 | 0 | 0 | 0 | |

5 rows × 301 columns

```
In [3]:   1  df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 18379 entries, 0 to 18378
Columns: 301 entries, Unnamed: 0 to t.158
dtypes: float64(134), int64(165), object(2)
memory usage: 42.2+ MB
```

## Creation of Dummy Variable + Remove Columns not Used

The dummy variables are used to replace the text variables, such that they can be used for computation in the subsequent steps. Columns that will not be used for computation will also be removed. These columns do not add any value in the clustering algorithm.

```
1  def dummy_var9(value):
2      if value == "Mono":
3          return 1
4      elif value == "Multi":
5          return 2
6      else:
7          return 0
8
9  df["var9_int"] = df["var9"].apply(dummy_var9)
10
11 ## remove columns that do not add value
12 df.drop(columns = ["var9", "month", "Unnamed: 0", "year", "respondent.id"],
13         inplace = True)
14
15 df.head()
```

Out[4]:

| | var8 | var6 | a.1 | a.2 | a.3 | a.4 | var5 | b.5 | b.6 | b.7 | ... | c.277 | c.278 | c.279 | c.280 | c.281 | c.28 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2.0 | NaN | 1 | 5 | 1 | 57 | 34 | 1 | NaN | 1 | ... | 0 | 0 | 0 | 1 | 0 | |
| 1 | 2.0 | NaN | 1 | 4 | 1 | 57 | 34 | 2 | 3.0 | 1 | ... | 0 | 0 | 0 | 0 | 0 | |
| 2 | 2.0 | NaN | 1 | 5 | 1 | 57 | 42 | 2 | 1.0 | 1 | ... | 0 | 0 | 0 | 0 | 0 | |
| 3 | 2.0 | NaN | 1 | 6 | 1 | 57 | 34 | 2 | 4.0 | 1 | ... | 0 | 0 | 0 | 1 | 0 | |
| 4 | 2.0 | NaN | 1 | 8 | 1 | 22 | 1 | 1 | NaN | 1 | ... | 0 | 0 | 0 | 0 | 0 | |

5 rows × 297 columns

## Replacing NaN values with 0

This is such that it can be used to calculate the standardised values. In this case, we assume that the columns are not binary. Thus, changing the NaN values to 0 will not "change" the value of the data.

```
1  # replacing all NaN values with 0
2  df.fillna(0, inplace = True)
3  df.head()
```

Out[5]:

| | var8 | var6 | a.1 | a.2 | a.3 | a.4 | var5 | b.5 | b.6 | b.7 | ... | c.277 | c.278 | c.279 | c.280 | c.281 | c.282 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2.0 | 0.0 | 1 | 5 | 1 | 57 | 34 | 1 | 0.0 | 1 | ... | 0 | 0 | 0 | 1 | 0 | ( |
| 1 | 2.0 | 0.0 | 1 | 4 | 1 | 57 | 34 | 2 | 3.0 | 1 | ... | 0 | 0 | 0 | 0 | 0 | |
| 2 | 2.0 | 0.0 | 1 | 5 | 1 | 57 | 42 | 2 | 1.0 | 1 | ... | 0 | 0 | 0 | 0 | 0 | ( |
| 3 | 2.0 | 0.0 | 1 | 6 | 1 | 57 | 34 | 2 | 4.0 | 1 | ... | 0 | 0 | 0 | 1 | 0 | |
| 4 | 2.0 | 0.0 | 1 | 8 | 1 | 22 | 1 | 1 | 0.0 | 1 | ... | 0 | 0 | 0 | 0 | 0 | ( |

5 rows × 297 columns

## Transformation Process

Transformation is done for variables that are highly skewed. Values that are lesser than -1 or greater than 1 will be considered as highly skewed.

```
In [6]:  ▶    1  skewness = df.skew()
              2
              3  var_to_transform = []
              4
              5  for name, value in skewness.iteritems():
              6      if value > 1 or value < -1:
              7          var_to_transform.append(name)
              8  print(var_to_transform)
```

['var6', 'b.7', 'b.9', 'b.10', 'b.11', 'b.13', 'b.14', 'b.15', 'b.16', 'b.17', 'b.18', 'b.20', 'b.21', 'b.22', 'pov6', 'b.23', 'b.24', 'b.25', 'b.26', 'b.27', 'b.28', 'c.29', 'c.30', 'c.31', 'c.32', 'c.34', 'c.35', 'c.36', 'c.37', 'c.39', 'c.40', 'c.41', 'c.42', 'c.43', 'c.44', 'c.45', 'c.46', 'c.47', 'c.48', 'c.49', 'c.50', 'c.51', 'c.52', 'c.53', 'c.54', 'c.55', 'c.56', 'c.57', 'c.58', 'c.59', 'c.60', 'c.61', 'c.62', 'b.63', 'c.65', 'c.66', 'c.67', 'c.68', 'c.69', 'c.70', 'c.71', 'c.72', 'c.73', 'c.74', 'c.75', 'c.76', 'c.77', 'c.78', 'c.79', 'c.80', 'c.81', 'c.82', 'c.83', 'c.84', 'c.85', 'c.86', 'c.87', 'c.88', 'c.89', 'c.90', 'c.91', 'c.92', 'f.105', 'f.106', 'f.107', 'f.108', 'f.109', 'f.110', 'f.111', 'f.112', 'f.113', 'f.114', 'f.115', 'f.116', 'f.117', 'f.118', 'f.119', 'f.120', 'f.121', 'f.122', 'f.123', 'a.124', 'var3', 'var4', 'c.125', 'c.126', 'c.127', 'pp.128', 'pp.129', 'pp.130', 'pp.131', 'pp.132', 'pp.133', 'pp.134', 'c.135', 'c.136', 'c.137', 'c.138', 'c.139', 'c.140', 'c.141', 'c.142', 'c.143', 'c.144', 'c.145', 'c.146', 'c.147', 'c.148', 't7.149', 't7.150', 't7.151', 't7.152', 't7.153', 't7.154', 't7.155', 't7.156', 't7.157', 't7.158', 'c.160', 'c.161', 'c.162', 'c.163', 'c.166', 'c.167', 'c.168', 'c.169', 'c.170', 'c.171', 'c.172', 't7.174', 't7.175', 't7.179', 't7.180', 't7.181', 't7.182', 'a.183', 'a.184', 'a.185', 'a.186', 'var2', 'totshopping.rep', 'var1', 'c.187', 'f.188', 'c.189', 'c.190', 'b.192', 'b.193', 'b.194', 'b.195', 'c.196', 'c.197', 'c.198', 'c.199', 'c.200', 'c.201', 'c.202', 'c.203', 'c.204', 'c.205', 'c.206', 'c.207', 'c.208', 'c.209', 'c.210', 'c.211', 'c.212', 'c.213', 'c.214', 'c.215', 'c.216', 'c.217', 'c.218', 'c.219', 'c.220', 'c.221', 'c.222', 'c.223', 'c.224', 'c.227', 'c.228', 'c.229', 'c.230', 'c.232', 'c.233', 'c.237', 'c.238', 'c.239', 'c.240', 'c.241', 'c.242', 'c.243', 'c.244', 'c.247', 'c.248', 'c.249', 'c.250', 'c.251', 'c.252', 'c.253', 'c.254', 'c.255', 'c.256', 'c.257', 'c.258', 'c.259', 'c.260', 'c.261', 'c.262', 'c.263', 'c.265', 'c.268', 'c.269', 'c.271', 'c.272', 'c.273', 'c.274', 'c.275', 'c.277', 'c.278', 'c.279', 'c.280', 'c.281', 'c.282', 'c.283']

```
In [7]:  ▶    1  for variable in var_to_transform:
              2      new_name = "log_" + variable
              3      df[new_name] = np.log(df[variable] + 1)
              4
              5  # removed non-transformed columns
              6  df.drop(columns = var_to_transform, inplace = True)
              7  df.head()
```

Out[7]:

| | var8 | a.1 | a.2 | a.3 | a.4 | var5 | b.5 | b.6 | b.8 | b.12 | ... | log_c.273 | log_c.274 | log_c.275 | log_c.2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2.0 | 1 | 5 | 1 | 57 | 34 | 1 | 0.0 | 1.0 | 2 | ... | 0.0 | 0.0 | 0.0 | ( |
| 1 | 2.0 | 1 | 4 | 1 | 57 | 34 | 2 | 3.0 | 1.0 | 2 | ... | 0.0 | 0.0 | 0.0 | ( |
| 2 | 2.0 | 1 | 5 | 1 | 57 | 42 | 2 | 1.0 | 1.0 | 1 | ... | 0.0 | 0.0 | 0.0 | ( |
| 3 | 2.0 | 1 | 6 | 1 | 57 | 34 | 2 | 4.0 | 1.0 | 2 | ... | 0.0 | 0.0 | 0.0 | ( |
| 4 | 2.0 | 1 | 8 | 1 | 22 | 1 | 1 | 0.0 | 2.0 | 1 | ... | 0.0 | 0.0 | 0.0 | ( |

5 rows × 297 columns

## Dataset Standardisation

Data standardisation is peorformed such that all the values will be on the same scale.

```
In [8]:  ▶|   1  scaler = StandardScaler()
             2  scaler.fit(df)
```

Out[8]:  StandardScaler(copy=True, with_mean=True, with_std=True)

```
In [9]:  ▶|   1  df_scaled = scaler.transform(df)
             2  df_scaled
```

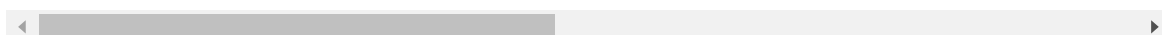Out[9]:  array([[ 0.41797295,  0.        , -0.54669362, ..., -0.19972256,
                 -0.2750927 , -0.13416986],
                [ 0.41797295,  0.        , -0.95477104, ..., -0.19972256,
                  3.6351383 , -0.13416986],
                [ 0.41797295,  0.        , -0.54669362, ..., -0.19972256,
                 -0.2750927 , -0.13416986],
                ...,
                [-1.39358816,  0.        , -0.54669362, ..., -0.19972256,
                 -0.2750927 , -0.13416986],
                [-1.39358816,  0.        , -0.95477104, ..., -0.19972256,
                  3.6351383 , -0.13416986],
                [ 1.32375351,  0.        , -0.13861621, ..., -0.19972256,
                 -0.2750927 , -0.13416986]])

```
In [10]:  ▶|   1  df_prepared = pd.DataFrame(df_scaled, columns = df.columns)
              2  df_prepared.head()
```

Out[10]:

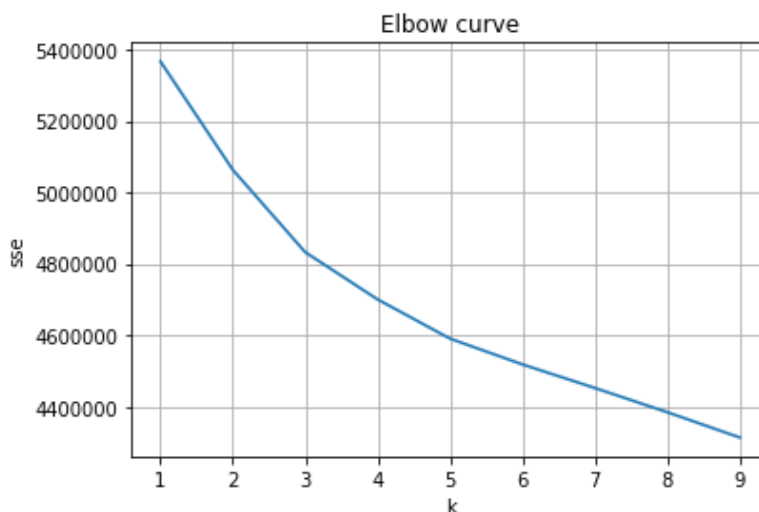|   | var8 | a.1 | a.2 | a.3 | a.4 | var5 | b.5 | b.6 | b.8 | b.12 |
|---|------|-----|-----|-----|-----|------|-----|-----|-----|------|
| 0 | 0.417973 | 0.0 | -0.546694 | 0.0 | 0.008724 | 0.826032 | -1.129119 | -0.879448 | -0.483734 | 0.757060 |
| 1 | 0.417973 | 0.0 | -0.954771 | 0.0 | 0.008724 | 0.826032 | 0.885646 | 1.237100 | -0.483734 | 0.757060 |
| 2 | 0.417973 | 0.0 | -0.546694 | 0.0 | 0.008724 | 1.578506 | 0.885646 | -0.173932 | -0.483734 | -1.320899 |
| 3 | 0.417973 | 0.0 | -0.138616 | 0.0 | 0.008724 | 0.826032 | 0.885646 | 1.942616 | -0.483734 | 0.757060 |
| 4 | 0.417973 | 0.0 | 0.677539 | 0.0 | -0.961049 | -2.277924 | -1.129119 | -0.879448 | 0.239133 | -1.320899 |

5 rows × 297 columns

## KMeans Method

The SSE is calculated and the value of K will be defined by the Elbow Method.

In [11]:  ▶
```
1  sse = []
2
3  for i in range (1, 10):
4      model = KMeans(n_clusters = i, random_state = 0)
5      model.fit(df_prepared)
6      sse.append(model.inertia_)
7
8  sse
```

Out[11]: [5366668.000000001,
5063219.498710746,
4832642.066105606,
4700390.77190287,
4590694.495576666,
4518524.424036916,
4452508.960699122,
4384778.9260005765,
4314006.047367832]

In [12]:  ▶
```
1  k = (range(1, 10))
2  plt.plot(k, sse)
3  plt.title("Elbow curve")
4  plt.xlabel("k")
5  plt.ylabel("sse")
6  plt.grid(True)
7  plt.show()
```



## Performing Clustering

Based on the Elbow curve that is plotted, K can be either 3 or 5. However, the change in SSE when K = 3 is slightly larger, thus the K = 3 will be used in the clustering algorithm.

In [13]:  ▶
```
1  model = KMeans(n_clusters = 3, random_state = 0)
2  model.fit(df_prepared)
```

Out[13]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
        n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto',
        random_state=0, tol=0.0001, verbose=0)

```
In [14]:  ▶  1  print("SSE:", round(model.inertia_, 2))
             2  print("Number of Iterations:", model.n_iter_)
```

SSE: 4832642.07
Number of Iterations: 19

```
In [15]:  ▶  1  df["cluster"] = model.labels_
             2  df["cluster"].value_counts()
```
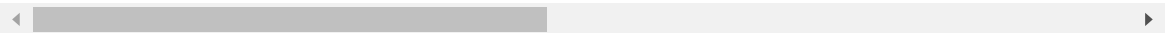
Out[15]: 1    9063
         0    6817
         2    2499
         Name: cluster, dtype: int64

```
In [16]:  ▶  1  cluster_mean = pd.DataFrame(model.cluster_centers_, columns = df_prepared.
             2  cluster_mean
```

Out[16]:

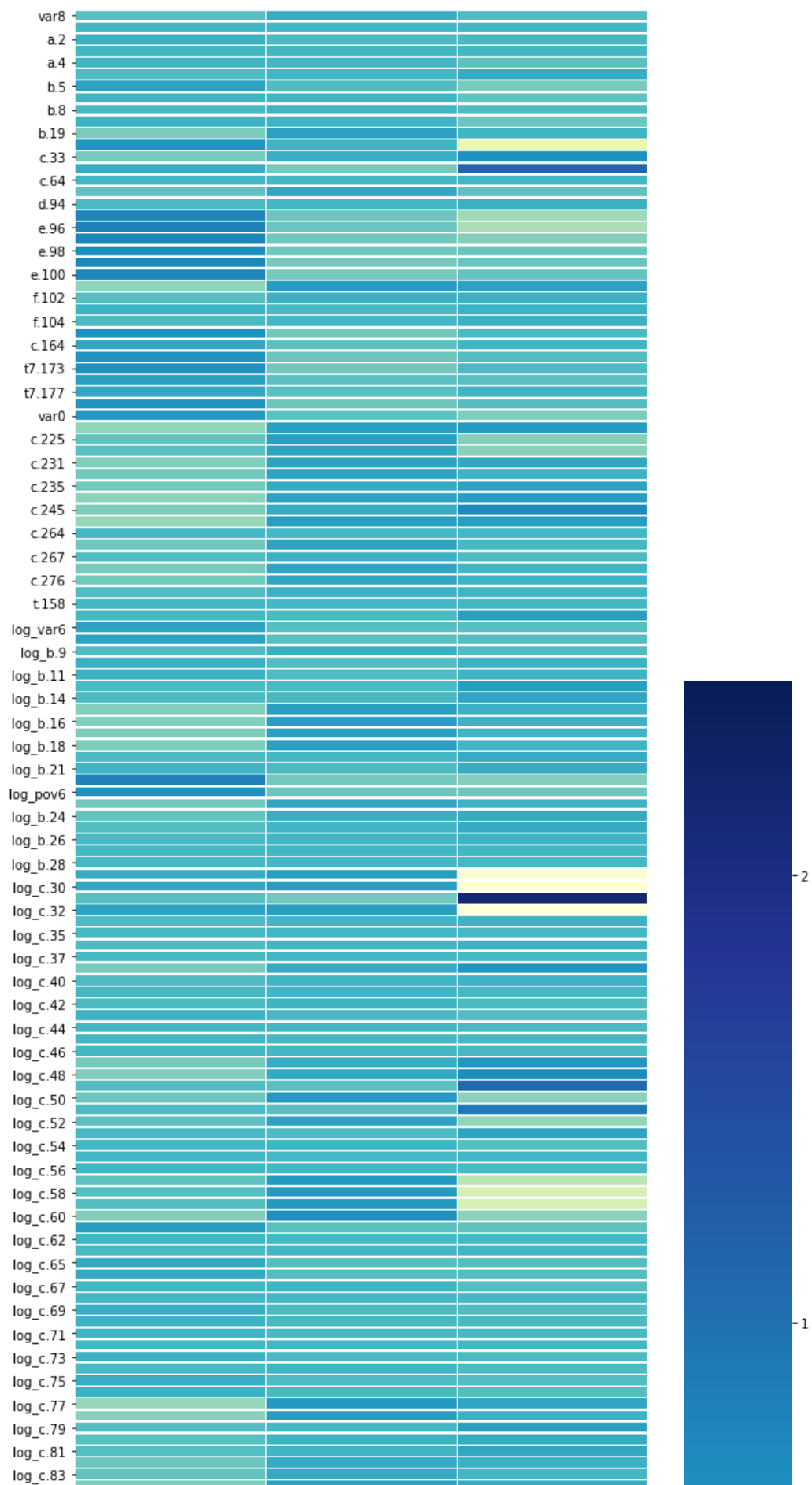|   | var8 | a.1 | a.2 | a.3 | a.4 | var5 | b.5 | b.6 | b.8 | b.12 |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | -0.213164 | 0.0 | 0.119627 | 0.0 | 0.036770 | -0.091919 | 0.369616 | 0.035539 | -0.040173 | 0.068472 |
| **1** | 0.188604 | 0.0 | -0.081252 | 0.0 | 0.027805 | 0.016175 | -0.132074 | 0.044036 | 0.060869 | 0.054778 |
| **2** | -0.102516 | 0.0 | -0.031657 | 0.0 | -0.201141 | 0.192082 | -0.529285 | -0.256651 | -0.111164 | -0.385443 |

3 rows × 297 columns

◀ ▬▬▬▬▬▬▬▬▬▬▬▬▬▬ ▶

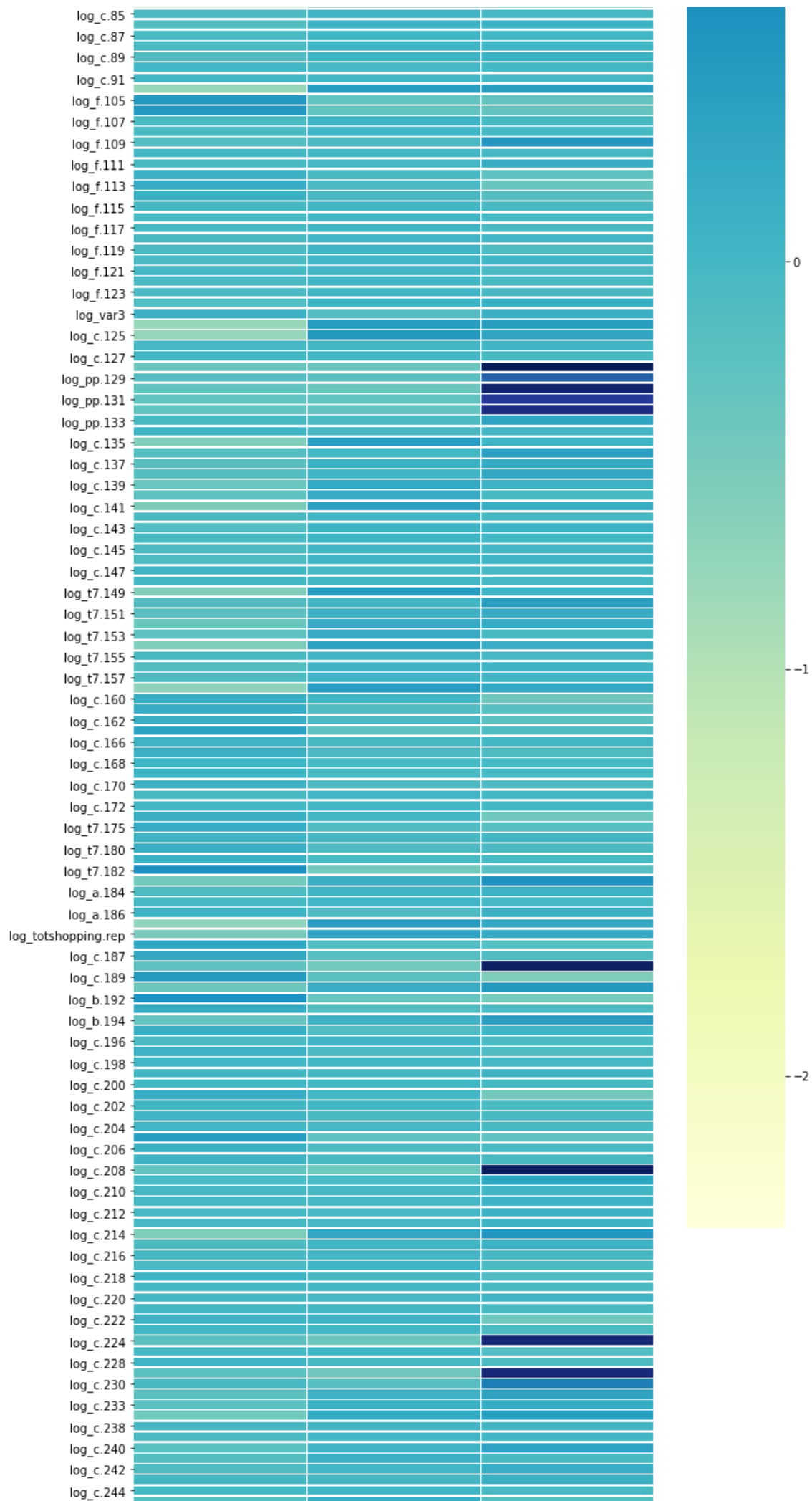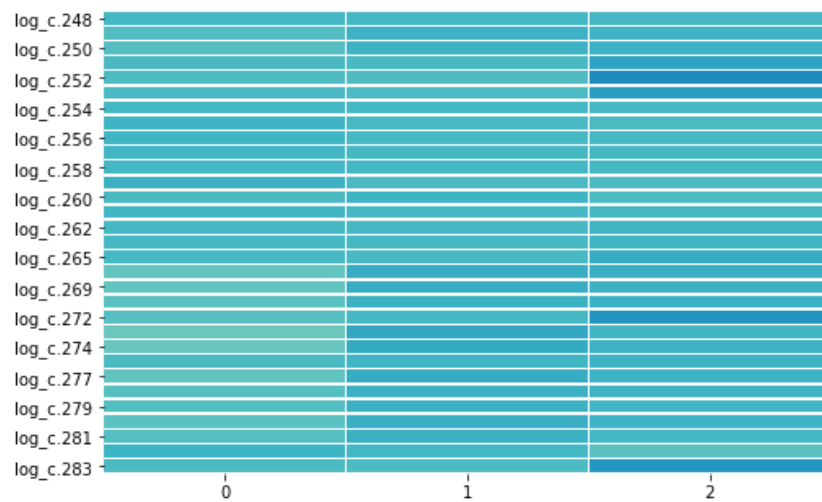## Visualising the Results

```
1  import seaborn as sns
2  plt.figure(figsize = (10,50))
3  sns.heatmap(cluster_mean.T, linewidths = .5, cmap = "YlGnBu")
```
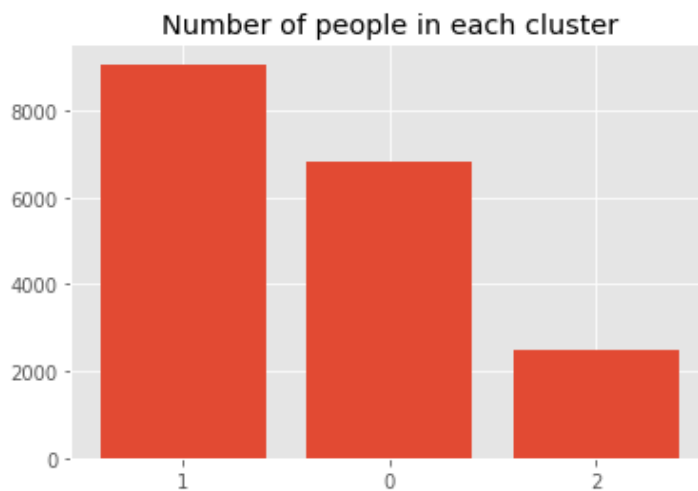
Out[17]: <matplotlib.axes._subplots.AxesSubplot at 0x1c14a550b48>

log_c.248
log_c.250
log_c.252
log_c.254
log_c.256
log_c.258
log_c.260
log_c.262
log_c.265
log_c.269
log_c.272
log_c.274
log_c.277
log_c.279
log_c.281
log_c.283

In [19]:

```python
clusters = df["cluster"].value_counts()

cluster_index = []
cluster_value = []

for index, value in clusters.iteritems():
    cluster_index.append(str(index))
    cluster_value.append(value)

plt.bar(cluster_index, cluster_value)
plt.style.use("ggplot")
plt.title("Number of people in each cluster")
plt.show()
```



Number of people in each cluster

Based on the results, we see that:

- Cluster 0: Generally light colours, which might indicate that they are the lower spenders
- Cluster 1: Colours are in between clusters 0 and 2, thus indicating that they belong to mid range
- Cluster 2: Colours are darker, thus indicating that they are the higher spenders

Most respondents come from cluster 1, followed by cluster 0 and 2.