

Using the **"Elbow" method** and the **Silhouette Score** to determine the optimal value of k.

```
1 # Import required packages
2 import pandas as pd
3 from sklearn.preprocessing import MinMaxScaler
4 from sklearn.cluster import KMeans
5 import matplotlib.pyplot as plt
```

```
1 data = pd.read_csv('/wholesale.csv')
2 data.head()
```

↗

	Channel	Region	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
0	2	3	12669	9656	7561	214	2674	1338
1	2	3	7057	9810	9568	1762	3293	1776
2	2	3	6353	8808	7684	2405	3516	7844
3	1	3	13265	1196	4221	6404	507	1788
4	2	3	22615	5410	7198	3915	1777	5185

Determine categorical and continuous variables

```
1 categorical_features = ['Channel', 'Region']
2 continuous_features = ['Fresh', 'Milk', 'Grocery', 'Frozen', 'Detergents_Paper', 'Delicassen']
```

Provide the descriptive statistics on spending items

```
1 data[continuous_features].describe()
```

↗

	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen
count	440.000000	440.000000	440.000000	440.000000	440.000000	440.000000
mean	12000.297727	5796.265909	7951.277273	3071.931818	2881.493182	1524.870455
std	12647.328865	7380.377175	9503.162829	4854.673333	4767.854448	2820.105937
min	3.000000	55.000000	3.000000	25.000000	3.000000	3.000000
25%	3127.750000	1533.000000	2153.000000	742.250000	256.750000	408.250000
50%	8504.000000	3627.000000	4755.500000	1526.000000	816.500000	965.500000
75%	16933.750000	7190.250000	10655.750000	3554.250000	3922.000000	1820.250000
max	112151.000000	73498.000000	92780.000000	60869.000000	40827.000000	47943.000000

To use the categorical features, we need to convert the categorical features to binary using pandas get dummies.

```
1 for col in categorical_features:
2     dummies = pd.get_dummies(data[col], prefix=col)
3     data = pd.concat([data, dummies], axis=1)
4     data.drop(col, axis=1, inplace=True)
5 data.head()
```

↗

	Fresh	Milk	Grocery	Frozen	Detergents_Paper	Delicassen	Channel_1	Channel_2	Region_1	Region_2	Region_3
0	12669	9656	7561	214	2674	1338	0	1	0	0	1
1	7057	9810	9568	1762	3293	1776	0	1	0	0	1
2	6353	8808	7684	2405	3516	7844	0	1	0	0	1
3	13265	1196	4221	6404	507	1788	1	0	0	0	1
4	22615	5410	7198	3915	1777	5185	0	1	0	0	1

To give equal importance to all features, we need to scale the continuous features.

```
1 mms = MinMaxScaler()
2 mms.fit(data)
3 data_transformed = mms.transform(data)
```

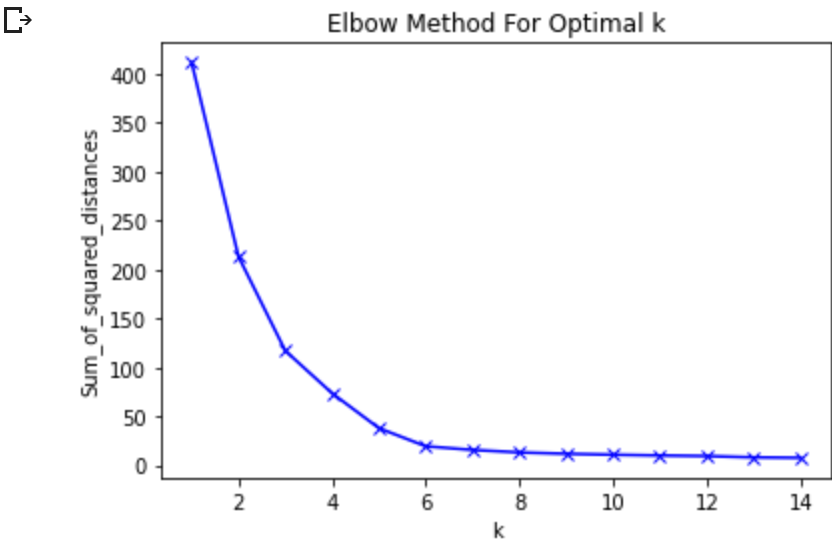
For each k value, we will initialize k-means and use the inertia attribute to identify the sum of squared distances of samples to the nearest cluster center.

```
1 Sum_of_squared_distances = []
2 K = range(1,15)
```

```
3 for k in K:
4     km = KMeans(n_clusters=k)
5     km = km.fit(data_transformed)
6     Sum_of_squared_distances.append(km.inertia_)
```

As k increases, the sum of squared distance tends to zero. Imagine we set k to its maximum value n (where n is number of samples) each sample will form its own cluster meaning sum of squared distances equals zero.

```
1 plt.plot(K, Sum_of_squared_distances, 'bx-')
2 plt.xlabel('k')
3 plt.ylabel('Sum_of_squared_distances')
4 plt.title('Elbow Method For Optimal k')
5 plt.show()
```



Sihouette Method

```
1 import matplotlib.pyplot as plt
2 %matplotlib inline
3 X=data_transformed
4 # Variables to store silhouette score and inertia for different clusters
5 scores = []
6 from sklearn.metrics import silhouette_score
7 # Calculate silhouette scores and inertia for different number of clusters
8 for cluster_number in range(2,10):
9     km = KMeans(n_clusters=k)
10    km = KMeans(n_clusters=cluster_number, random_state=42).fit(X)
11    scores.append(silhouette_score(X,km.labels_))
12 #plot the results
13 plt.clf()
14 plt.figure(figsize=(15,4))
15 #plot Silhouette Score
16 plt.subplot(121)
17 plt.plot(range(2,10), scores, 'bo-')
18 plt.plot(2, scores[0], 'ro-')
19 plt.plot(3, scores[1], 'yo-')
20 plt.title('Silhouette Method For Optimal k')
21 plt.xlabel('Number of Clusters')
22 plt.ylabel('Silhouette Score')
```

```
> Text(0, 0.5, 'Silhouette Score')
<Figure size 432x288 with 0 Axes>
```

