1. **Intro**
   a. To cluster customers into segments to optimize the significance of each customer to the business. This may help modify products according to customers specific needs and help businesses address the concerns of different user cohorts
2. **Data Cleansing Approaches**
   a. Basic EDA - check for missing data (Null), data check to evaluate which columns would be valuable to include
   b. Add any additional variables that would be valuable for the model:
      i. Total ($) Spent
      ii. Cleanse data grouping on user attributes: education, household size, marital status
      iii. Remove unnecessary attributes such as: costcontract, ID, etc.
      iv. Convert categorical variables using LabelEncoder to be consumable for K-Means
   c. Remove outliers:
      i. Remove age > 90
      ii. Income > 75%
3. **Initial Findings (before training):**
   a. income and num of visit per month is negatively correlated
   b. high correlation between wine, fruit, meat
   c. low correlation of complain to any of the purchases (or income)
4. **Data Modeling Approaches**
   a. Standard Scaling
      i. K-Means is distance based algorithms that's affected by the scale of the variable
      ii. If we print - the mean of the variables in the different cluster, we can see that the income range is wider than other variables - if unscaled, our cluster will be more dependent on income variable vs others.
   b. Finding the optimal K:
      i. Elbow Method / Inertia Curve: The basic idea behind this method is that it plots the various values of cost with changing k. As the value of K increases, there will be fewer elements in the cluster. So average distortion will decrease. The lesser number of elements means closer to the centroid. So, the point where this distortion declines the most is the elbow point.
      ii. As seen in the graph, the inertia curve starts to flatten out around k = 3
   c. Success Metrics
      i. Silhouette Score - measures how well the clusters are separated
      ii. Inertia
5. **Findings**
   a. Cluster 0 - Low Intent
      i. Low Spending, across all income bracket
      ii. Low Deals Purchases
      iii. Kids

      b. Cluster 1 - High Intent
          i.     High Spending, High Income
          ii.    Medium Deals Purchases
          iii.   No Kids
          iv.   Slightly Younger Demographic
      c. Cluster 2 - Medium Intent
          i.     Medium Spending, medium bracket
          ii.    High Deals Purchases
          iii.   Kids
          iv.   Older Demographic
          v.    Bigger household size

**6. Next Steps / Limitations:**
      a. Cluster is not as tight and separated as expected, potentially due to the similarities in a lot of the input variables that still have high correlations with income.
      b. Try out different method of unsupervised machine learning and couple this with a supervised learning to model out how to best maximize revenue / business goals.