

The background is a dark blue-grey color. It features several thin, gold-colored lines that form abstract, geometric shapes. These lines radiate from the central text box, creating a starburst or network-like effect. The lines vary in length and angle, some extending towards the corners of the frame.

Movies

CPSC 392 Final Project

By: Ryan King

Question 1

Of the variables year, gross, votes, budget, runtime, and the various movie genres, which ones have the strongest relationship with a movie's IMDb score?

Variables Used

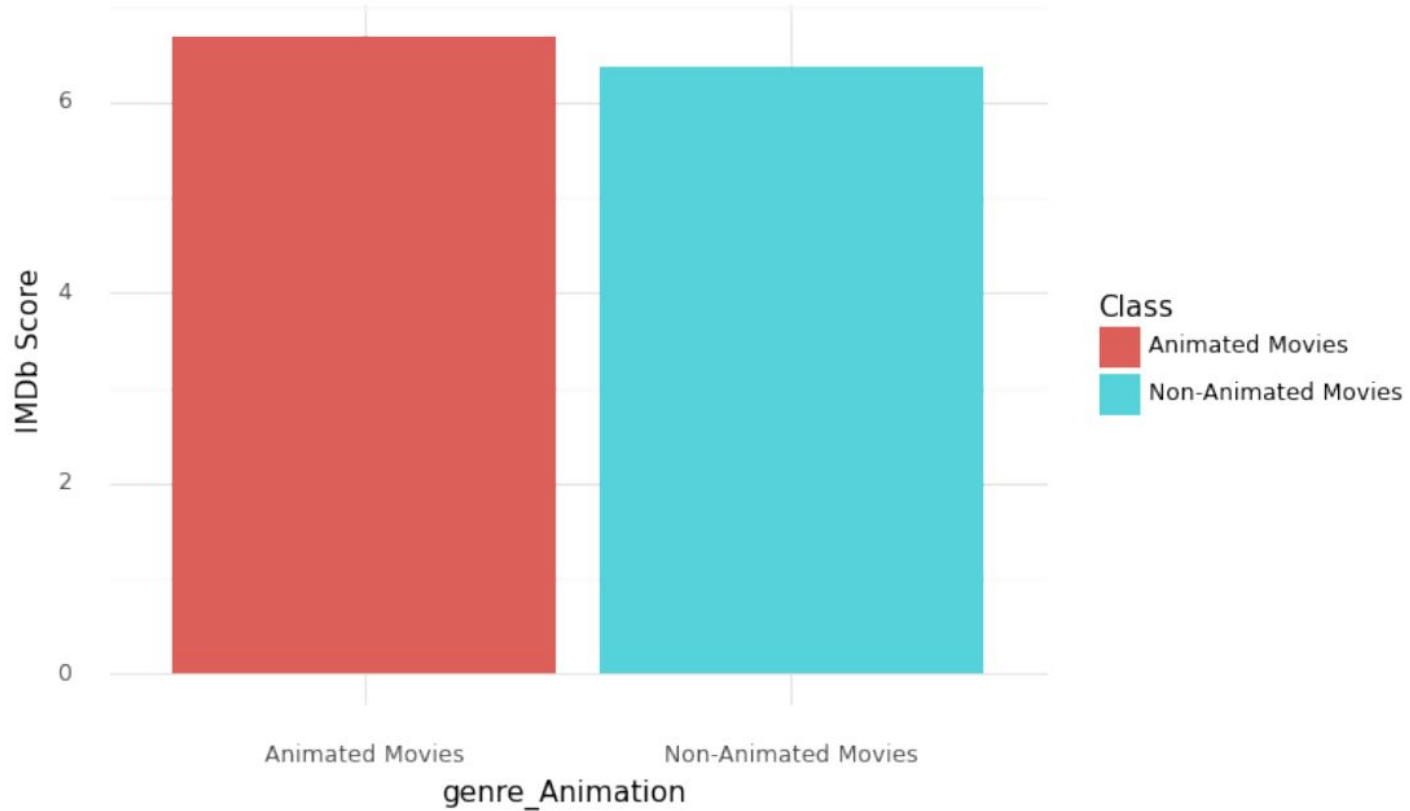
- Year
- Gross
- Votes
- Budget
- Runtime
- genre_Action
- genre_Adventure
- genre_Biography
- genre_Comedy
- genre_Crime
- genre_Drama
- genre_Family
- genre_Fantasy
- genre_Horror
- genre_Mystery
- genre_Romance
- genre_Sci-Fi
- genre_Thriller
- genre_Western

Steps

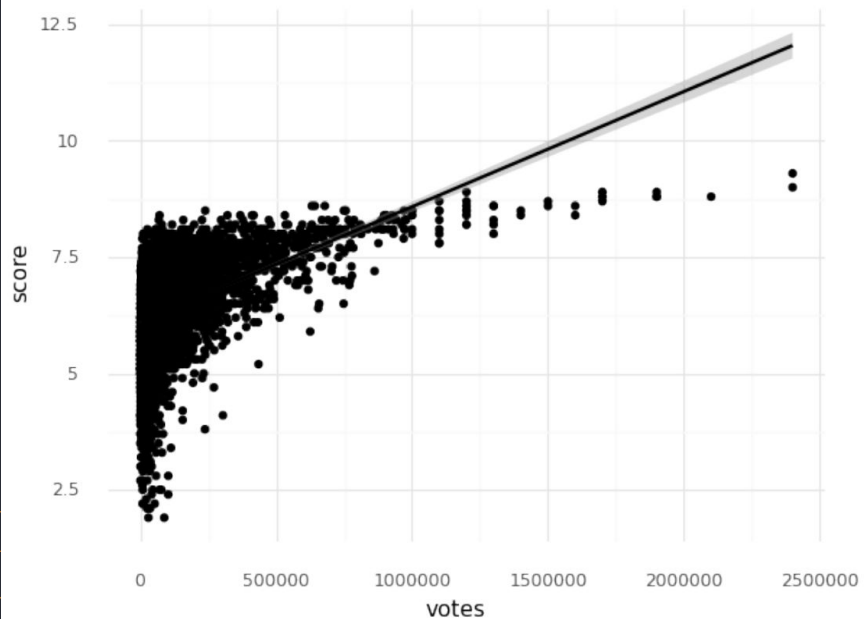
- STEP 1 Train/Test split with an 80/20 split
- STEP 2 Z-Scored variables
- STEP 3 Fit linear regression model
- STEP 4 Stored all variable coefficients in a dataframe
- STEP 5 Created a plot to display those coefficients



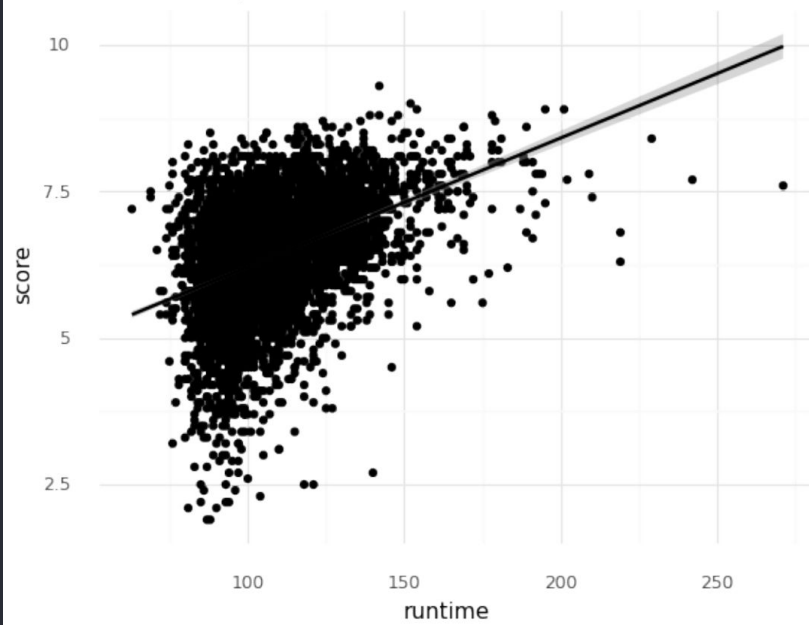
How Animated Films compare to Non-Animated Films



Relationship between IMDb Score and Votes



Relationship between IMDb Score and Movie Runtime



Question 2

When comparing a model using PCA on all the continuous variables other than score, in the dataset and retaining enough PCs to keep 90% of the variance, to a model using all the continuous variables other than score, what is the difference in the mean squared error when predicting the IMDb score of a movie?

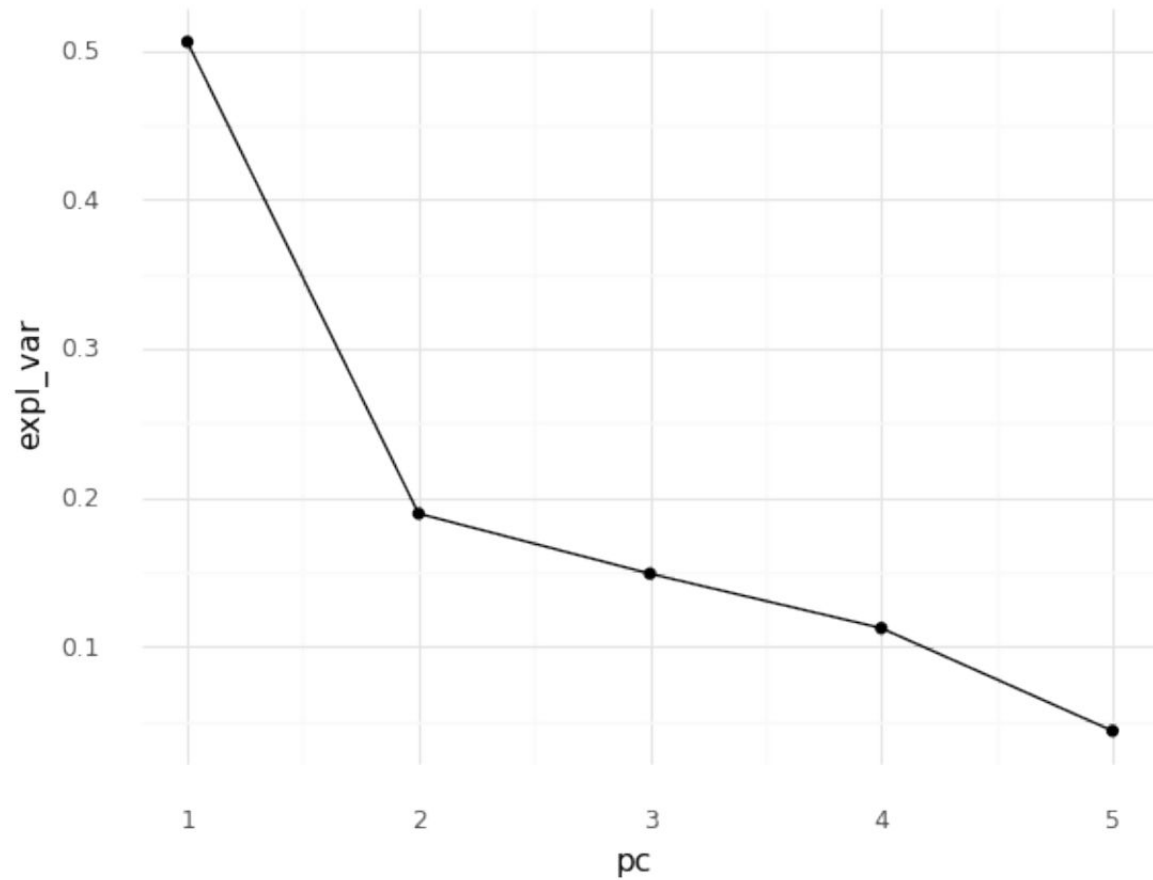
Variables Used

- Year
- Gross
- Votes
- Budget
- Runtime

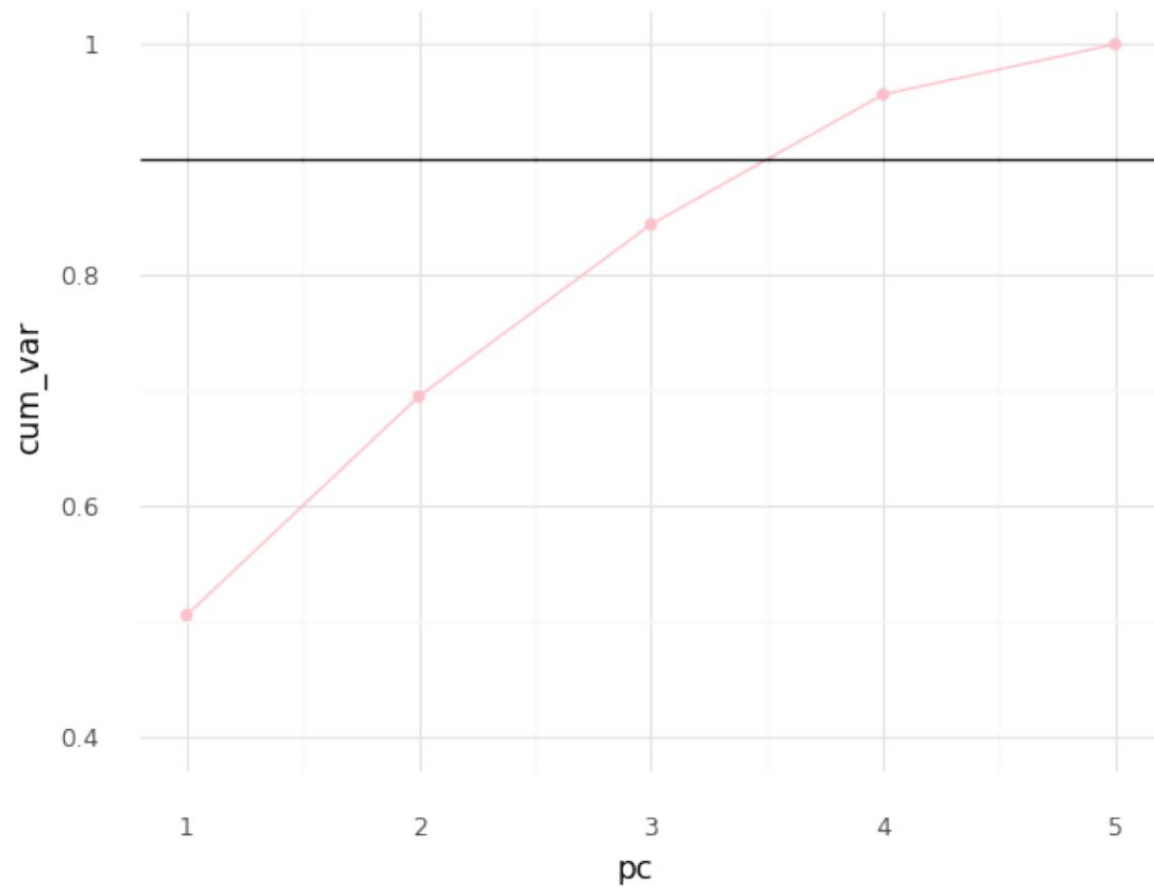
Steps

- STEP 1 Created original model to compare PCA model to
- STEP 2 Calculated MSE values of original model
- STEP 3 Produced Scree and Cumulative Variance plots to get PC values
- STEP 4 Created PCA models using PC values
- STEP 5 Compared MSE values of all models

Scree Plot



Cumulative Variance Plot



MSE Values Produced

Original Model

- Train MSE: 0.611
- Test MSE: 0.619

2 PCs Model

- Train MSE: 0.865
- Test MSE: 0.842

4 PCs Model

- Train MSE: 0.654
- Test MSE: 0.683

Question 3

When considering the variables movie gross, score, and budget, what clusters are shown, and describe what those clusters mean for those groups of movies?

Variables Used

- Gross
- Budget
- Score

Steps

STEP 1

Created a graph of the Sum of Squared Errors for different K values

STEP 2

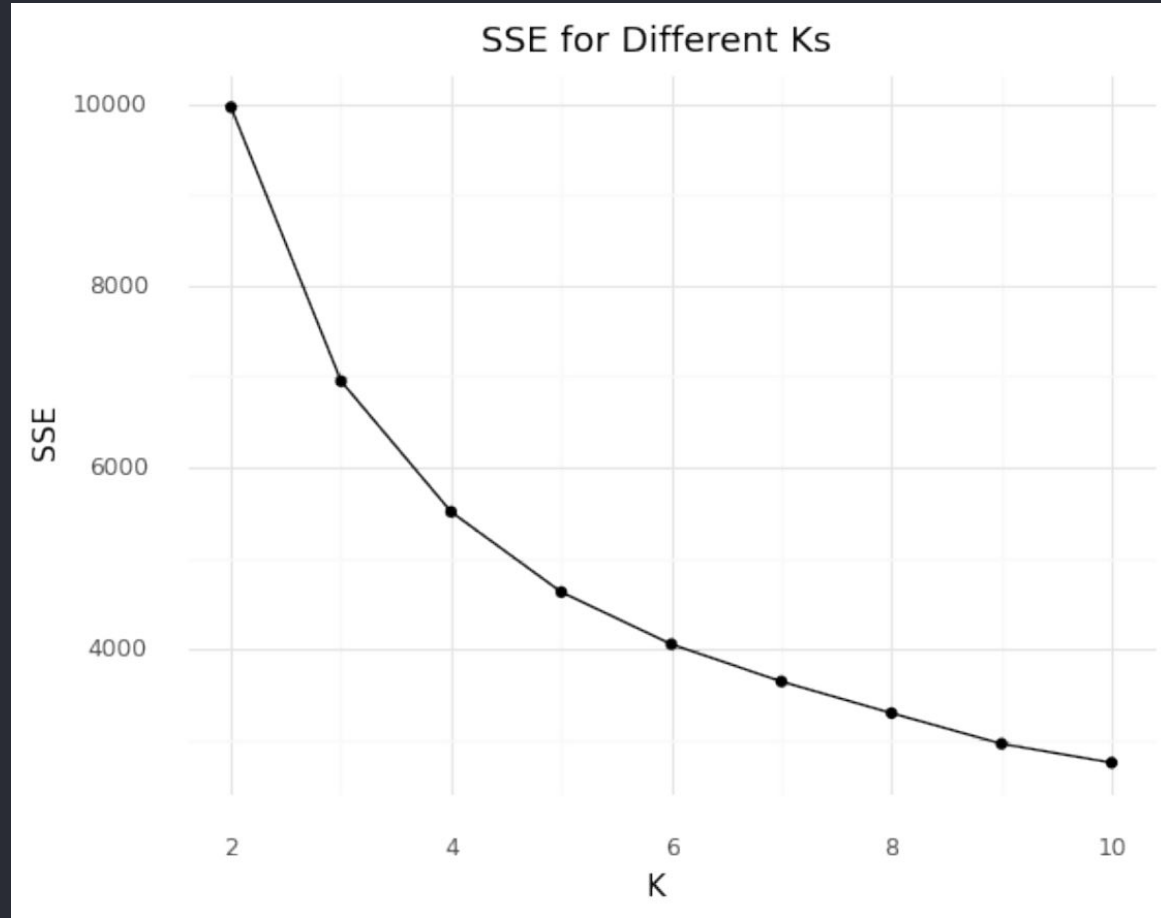
Created a graph for Silhouette Scores for different K values

STEP 3

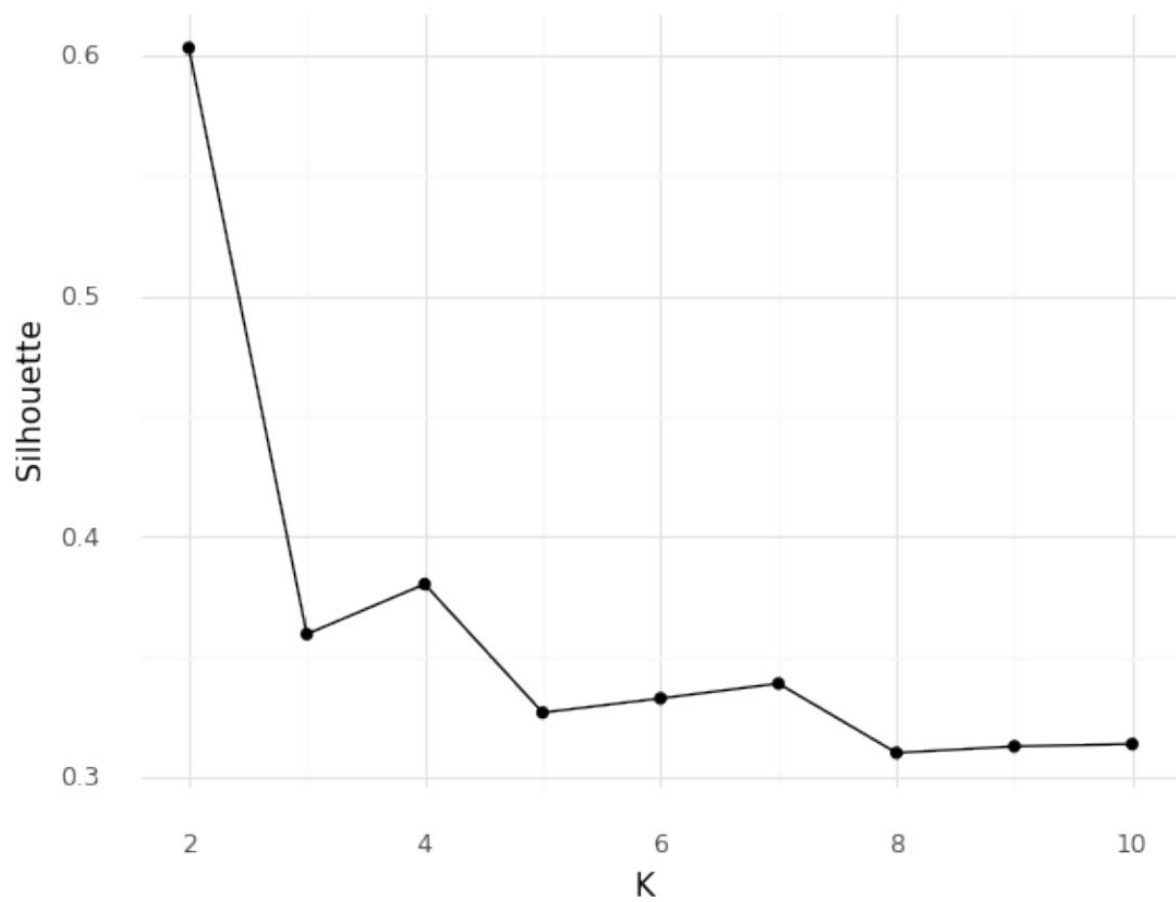
Used graphs to determine there are 3 clusters

STEP 4

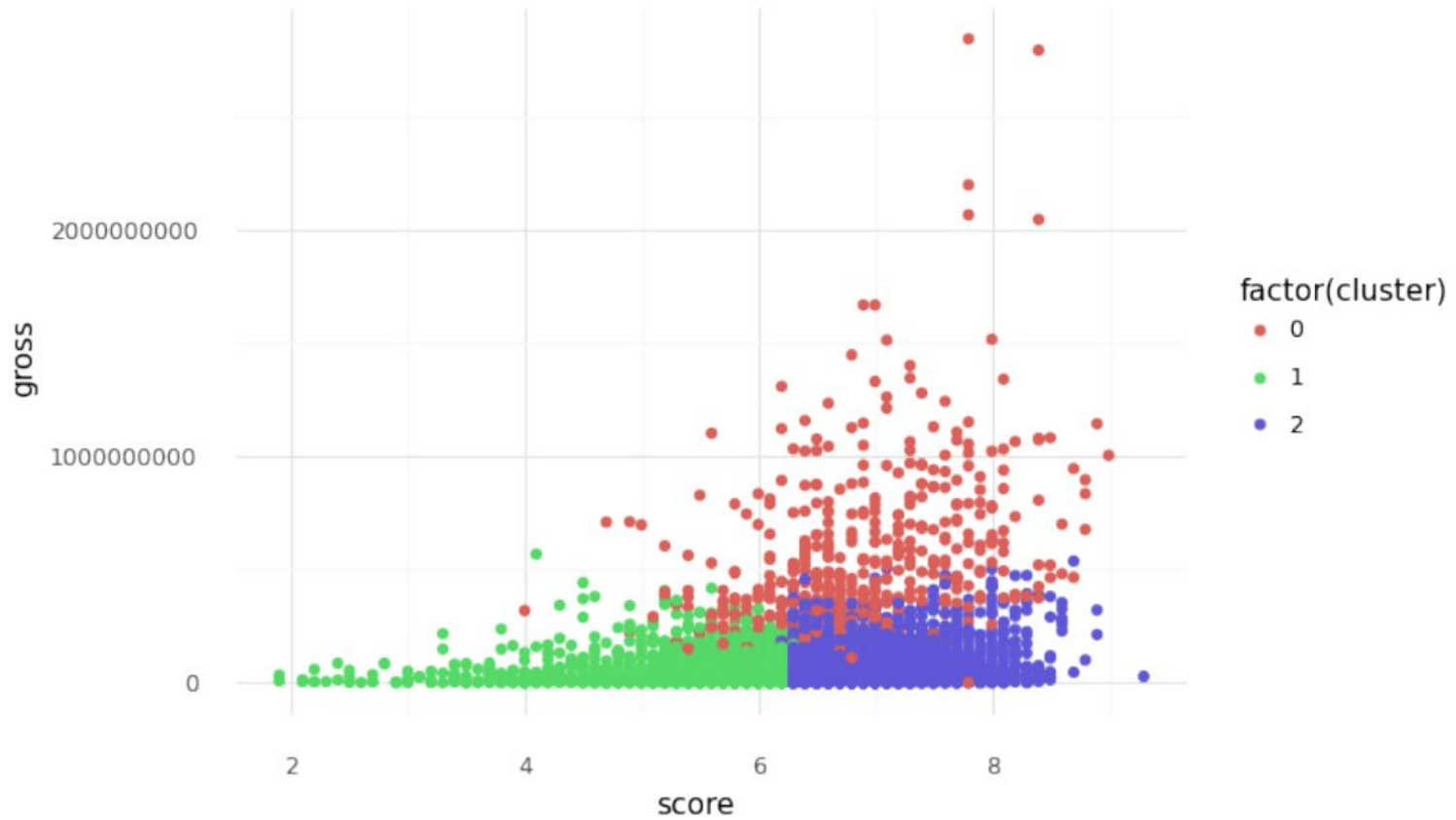
Produced plots to show the clusters within the different variables relationships



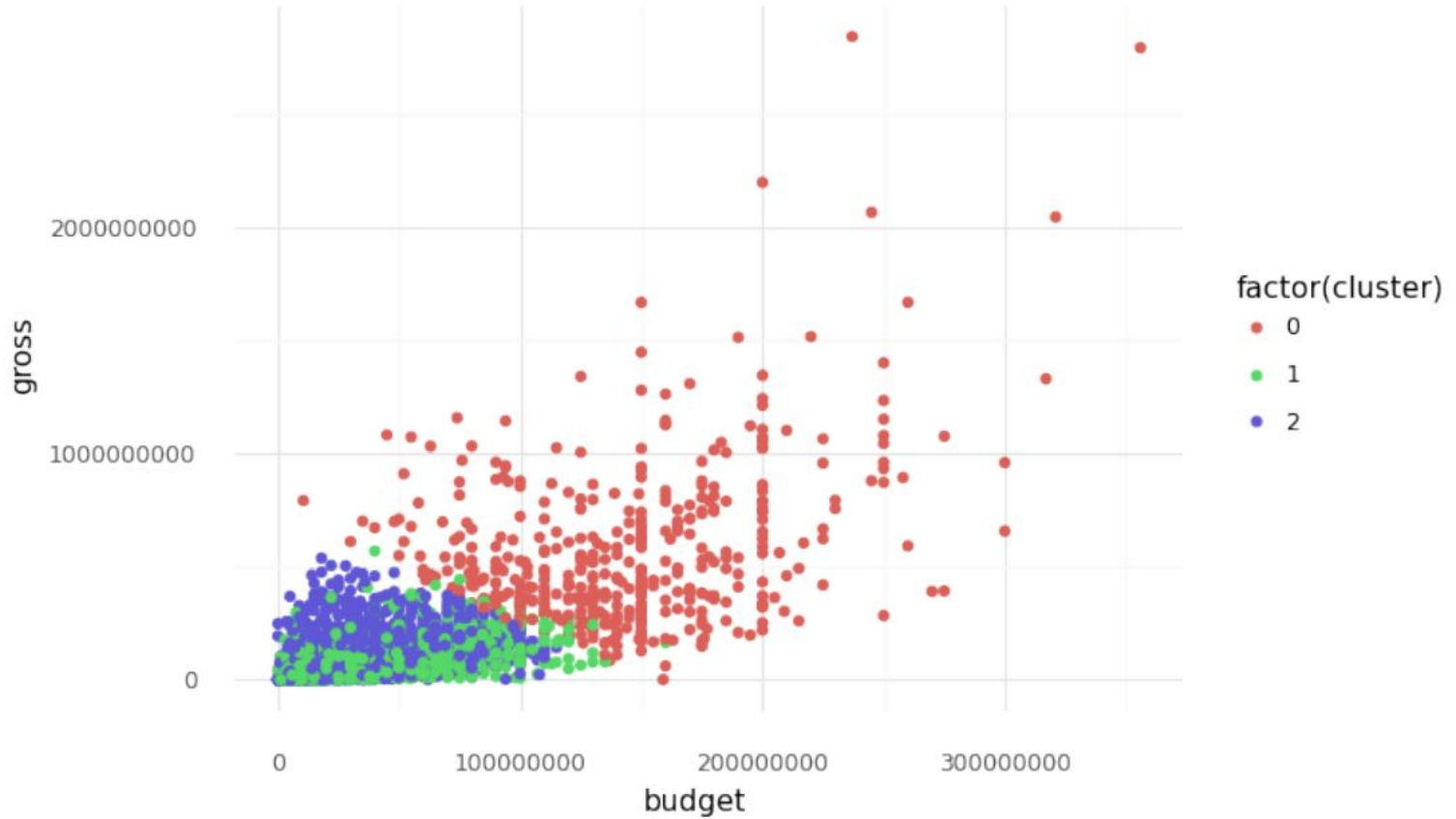
Silhouette Score for Different Ks



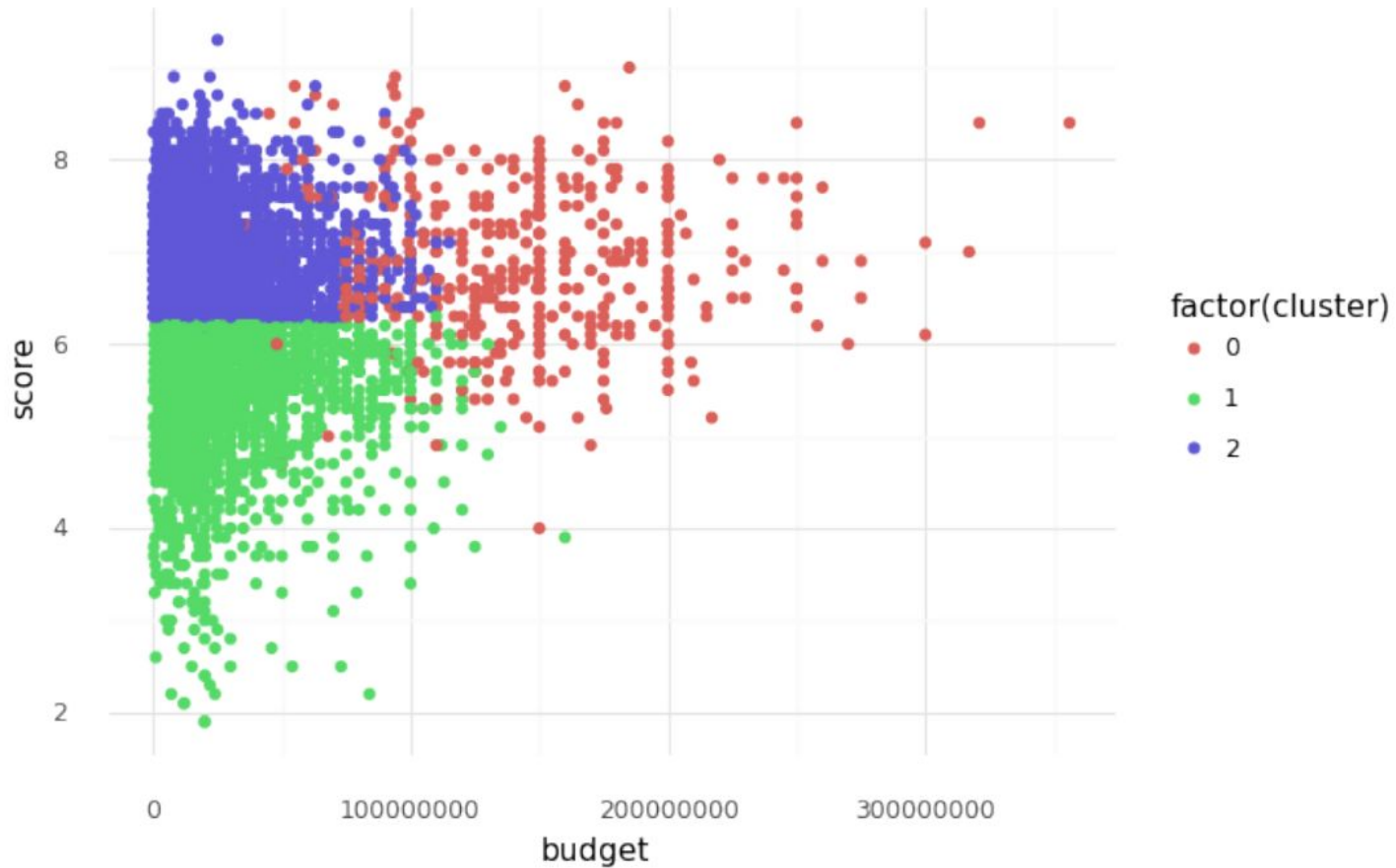
Gross vs Score



Gross vs Budget



Score vs Budget





Thank you!