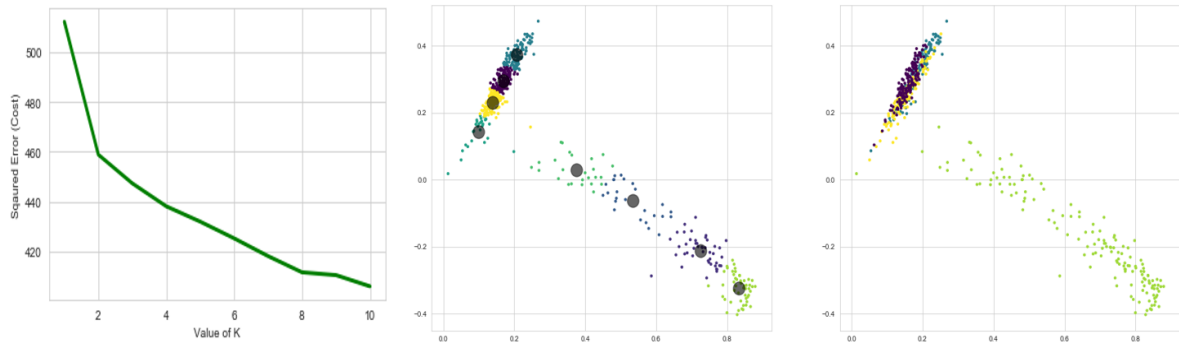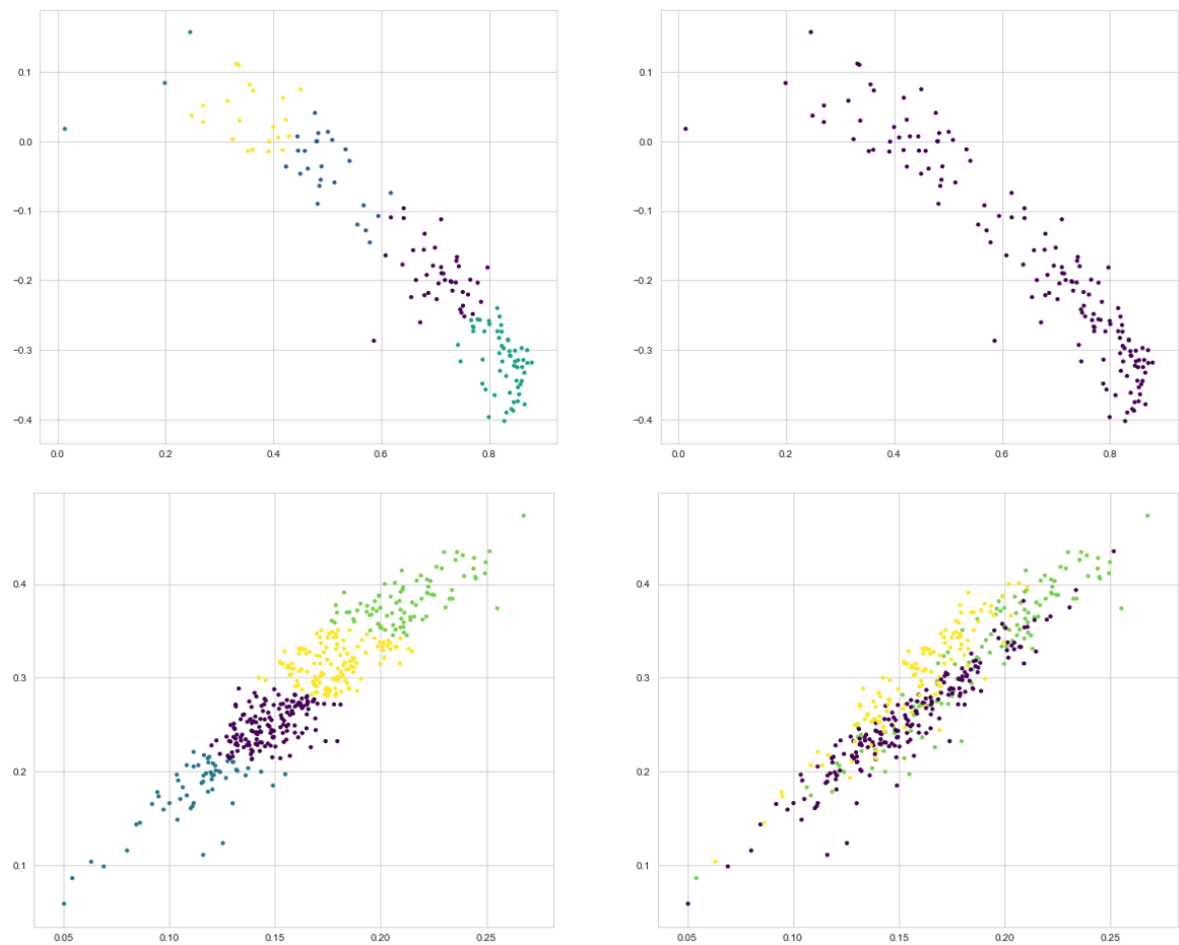# Unsupervised ML-Clustering Analysis

## 2D Clustering Analysis of Scatter Plot

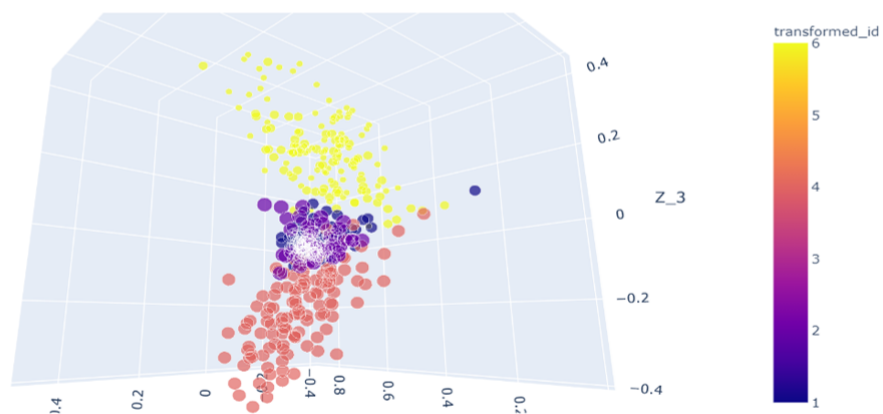### Hook from Mentor's session



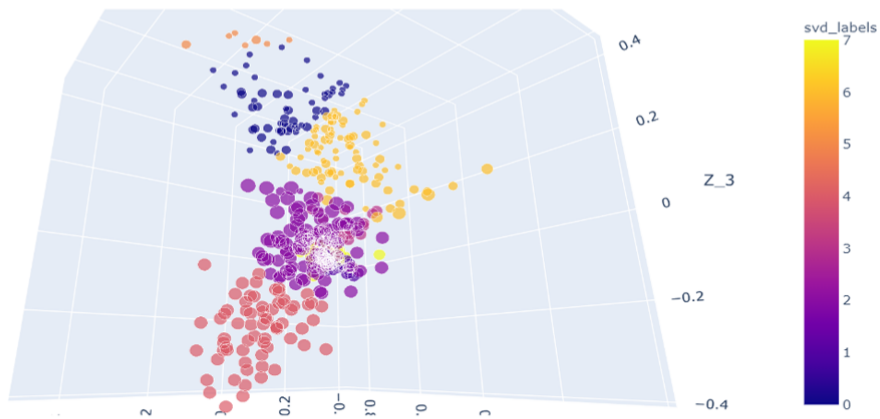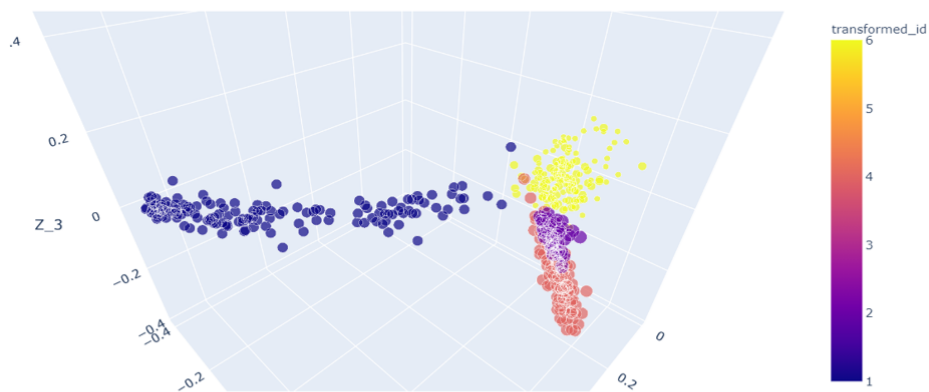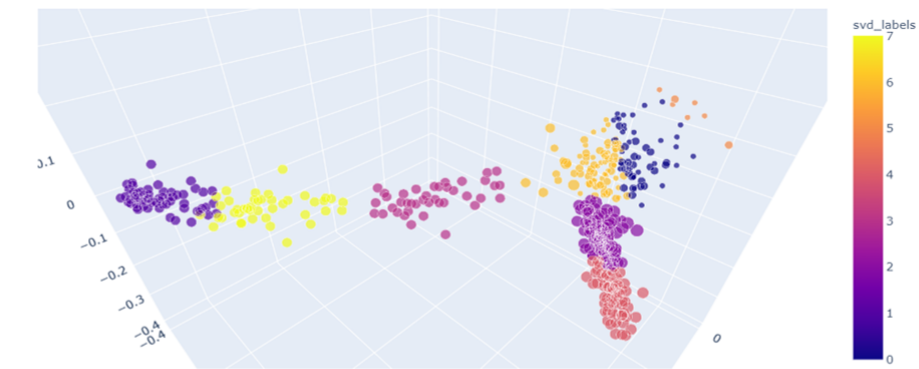### Pick out the special one

# 3D Clustering Analysis of Scatter Plot
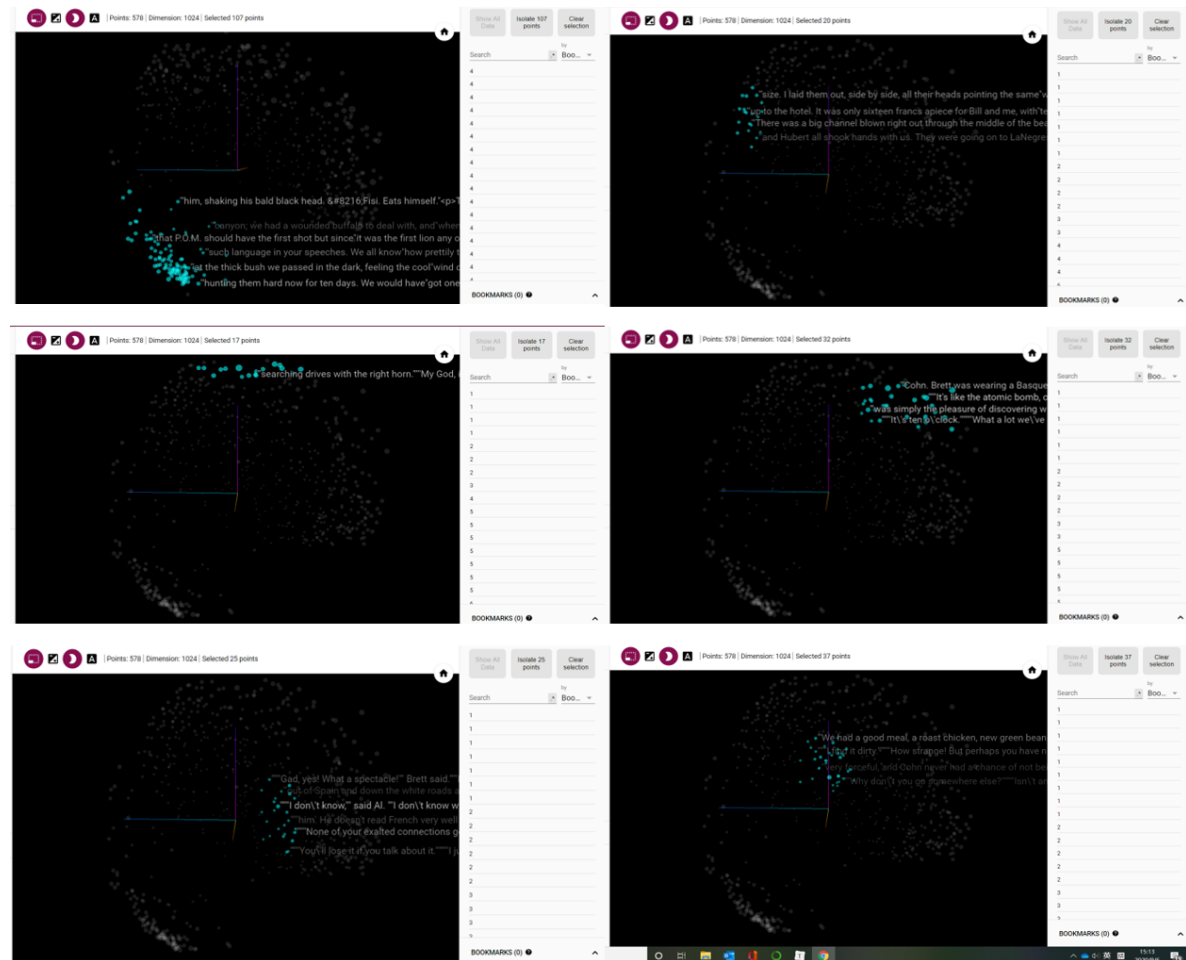
## Plotly Method

# Embedding Projector Method

## Find more about the fourth book --- Green Hills of Africa (1935):



## Randomly select points and analyze

# Summary of Clustering Analysis above

I will introduce and explain all the plots from top to bottom, left to right.

Basically, I want to make a 3D plot in Plotly of the Unsupervised Machine Learning technique, here I determine to apply clustering analysis for texts of six books from Hemingway.

The plots above start from the method in mentor session. After we gaining needed data by TextBlob package, Tf-idf method in sklearn, we determine to set number of clusters in K-means method equal to eight since we can clearly see an 'elbow' in the first graph. I use surjection to get associated id and transform_id, to avoid overlapping.

According to the 2D plots, the left one clearly shows eight clusters with separated centers, which means the machine can learn to pick similar sentences into one cluster no matter which book it belongs to. However, if we look at the right graph, it only has four clusters which means, one transformed_id contains four different clusters and others are combined and hard to be separated.

After plotting each book, I find the fourth book, Green Hills of Africa(id = 3), is the special one, and the texts of it dominate the whole right-bottom part. As for other five books, we find they contain some texts corresponding to the left graph. Additionally, It also illustrates that only four books are clear enough and rational to use K-means Clustering method.

As I get high dimensional text vectors initially, I consider to use PCA method to transform it to three-dimensional data to avoid losing more information. Hence, 3D plots here are more precise and meaningful. For the whole shape, if we consider the projection of that, we find it kind of similar to 2D plots above. However, the key feature is that it can be divided into six clusters with respect to different transformed_id rather than four. Although each cluster does not match the correct book wholly, it contains parts of texts belonging to the true book.

In addition, to make the 3D plot more accurate and meaningful, I consider the embeddings in the text, since recognizing certain relationship between grams contributes to the Clustering Analysis. Here I use Embedding Projector, a tool in Tensorflow, to get the dynamic 3D plot. I provide some representative screenshots to better understand the Unsupervised ML to text relations among these six books.

From light style, I focus on how bonded the texts of the fourth book are. The two graphs show that no matter 10 nearest points or 50 nearest points from the text of Green Hills of Africa, most of them also belong to this book.

From dark style, we can see the clear two opposite parts of the whole 'ball' shape. Unsurprisingly, as I select several points in the first graph, we see nearly all of them belong to the fourth book. However, as I select different small parts of the other side, other five pictures all show that the highlight points are labeled for several books, not one or two books. I can also claim that, even considering embeddings, the cluster results do not change a lot. Overall, except Green Hills of Africa, ML-generated clusters are not totally corresponding to the clusters of real books.