# Tema 2 – Bonus: Dog Breeds
## Dumitrescu Gabriel Horia

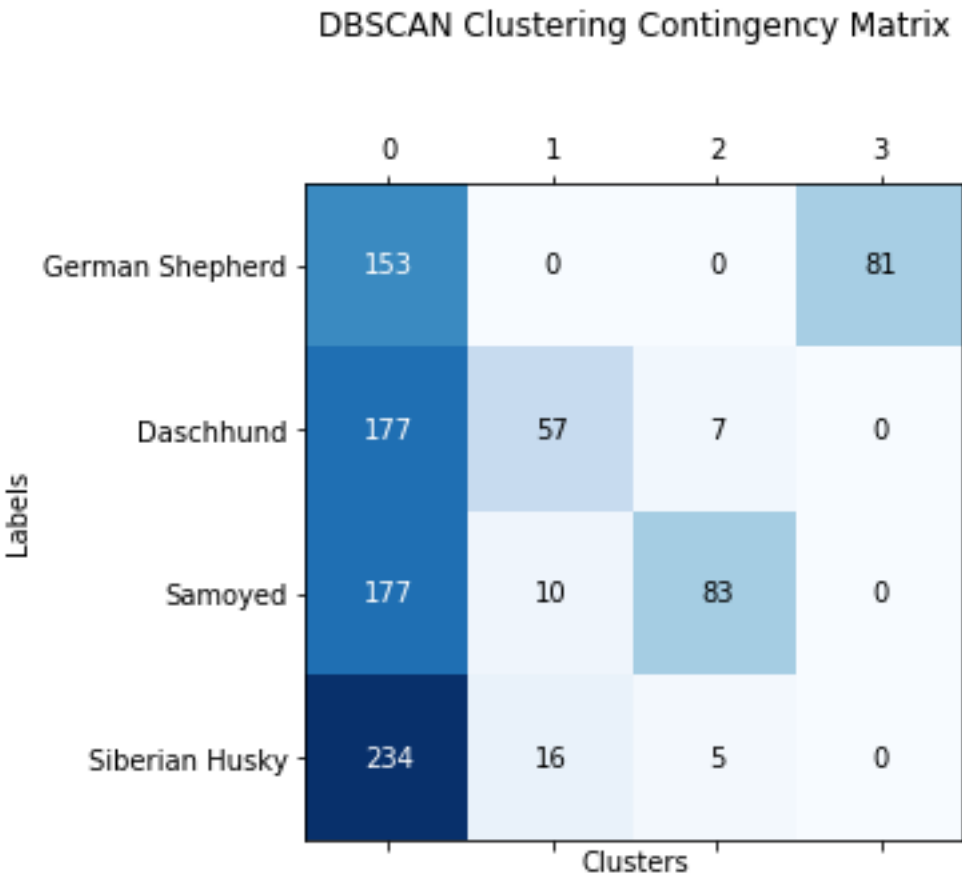## I)Clustering

## a)Kmeans

```
 K-MEANS CLUSTERING
Adjusted random score:  0.3783304148334443
Silhouette score:  0.32348571737418014
Homogeneity score:  0.4502238166890929
Completeness score:  0.5539720433014228
V-measure score:  0.49673856985743325
Fowlkes-Mallows score:  0.5811394558879884
Calinski-Harabaz score:  273.7663824785573
Supervised-like Accuracy:  5.0
Mean Score:  1.1119842882776516
```

### K-means Clustering Contingency Matrix

| Labels | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| German Shepherd | 232 | 2 | 0 | 0 |
| Daschhund | 2 | 216 | 23 | 0 |
| Samoyed | 6 | 232 | 23 | 9 |
| Siberian Husky | 4 | 104 | 25 | 122 |

Clusters

# b)DBSCAN

```
 DBSCAN CLUSTERING
Adjusted random score:  0.05347829317180698
Silhouette score:  -0.0239851520353545
Homogeneity score:  0.18320632471622236
Completeness score:  0.2964151767395702
V-measure score:  0.22644996087845962
Fowlkes-Mallows score:  0.4157713255365867
Calinski-Harabaz score:  71.03323976135188
Supervised-like Accuracy:  3.75
Mean Score:  0.7001908470010416
```
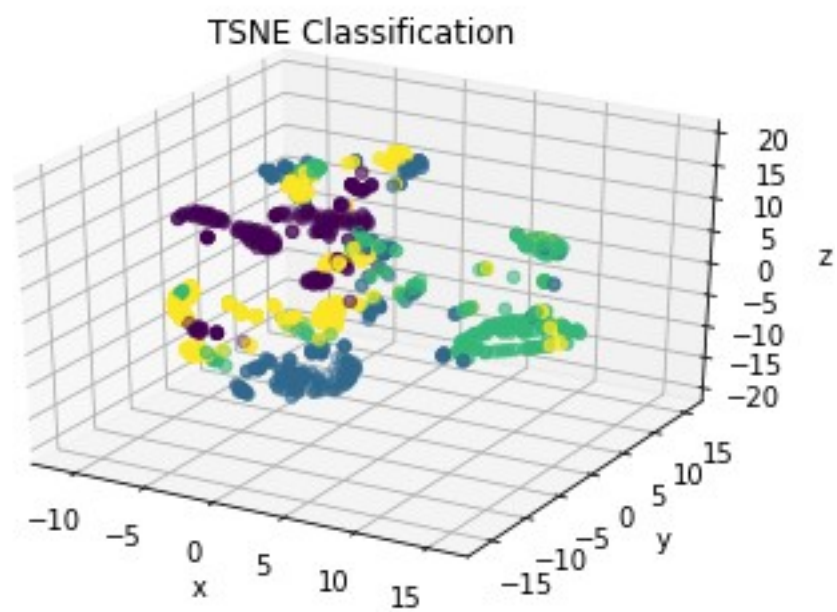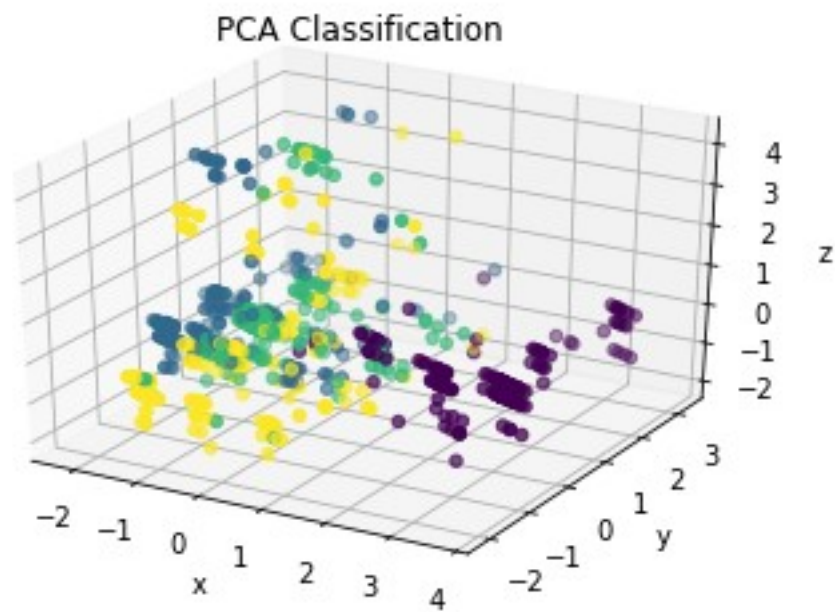
## DBSCAN Clustering Contingency Matrix

| Labels | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| German Shepherd | 153 | 0 | 0 | 81 |
| Daschhund | 177 | 57 | 7 | 0 |
| Samoyed | 177 | 10 | 83 | 0 |
| Siberian Husky | 234 | 16 | 5 | 0 |

Clusters

c)Agglomerative

```
AGGLOMERATIVE CLUSTERING
Adjusted random score:  0.29773038879125024
Silhouette score:  0.29715029626201056
Homogeneity score:  0.3803601986386643
Completeness score:  0.4842395274060927
V-measure score:  0.4260594522171749
Fowlkes-Mallows score:  0.5331750758120962
Calinski-Harabaz score:  252.39790917335043
Supervised-like Accuracy:  5.0
Mean Score:  1.059816419875327
```
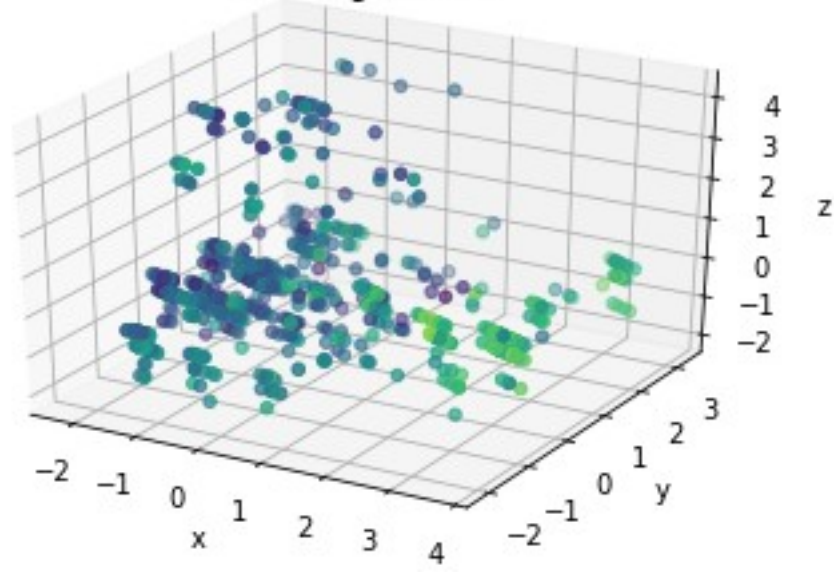
## Agglomerative Clustering Contingency Matrix

| Labels | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| German Shepherd | 209 | 25 | 0 | 0 |
| Daschhund | 1 | 217 | 23 | 0 |
| Samoyed | 8 | 230 | 23 | 9 |
| Siberian Husky | 1 | 117 | 25 | 112 |

Clusters

We can observe all clusterings improved, but this time, kmeans is the best one. Except German Shepherd, everything is hard to discriminate.

# II) Data Visualization



PCA Classification
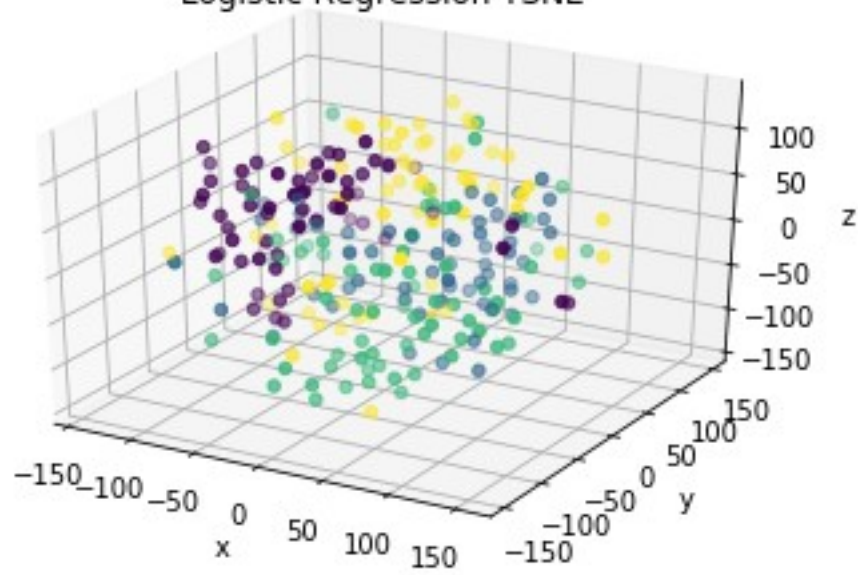


TSNE Classification
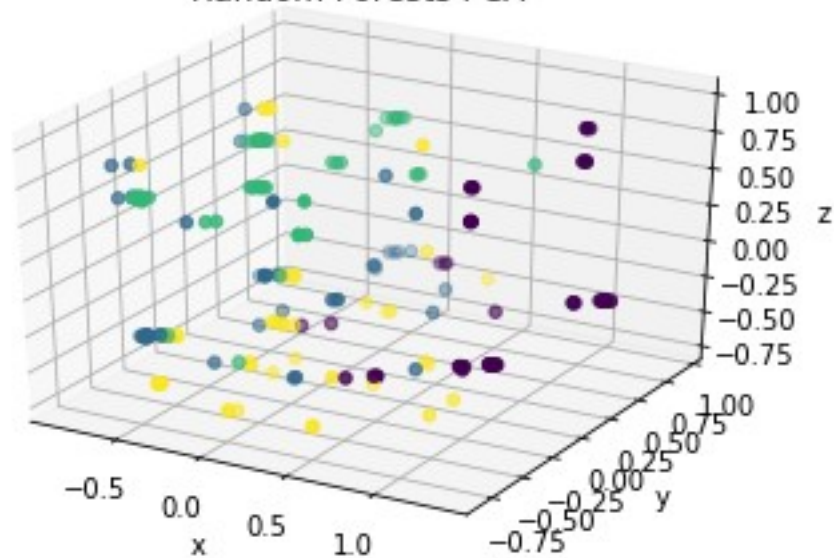
## PCA Regression



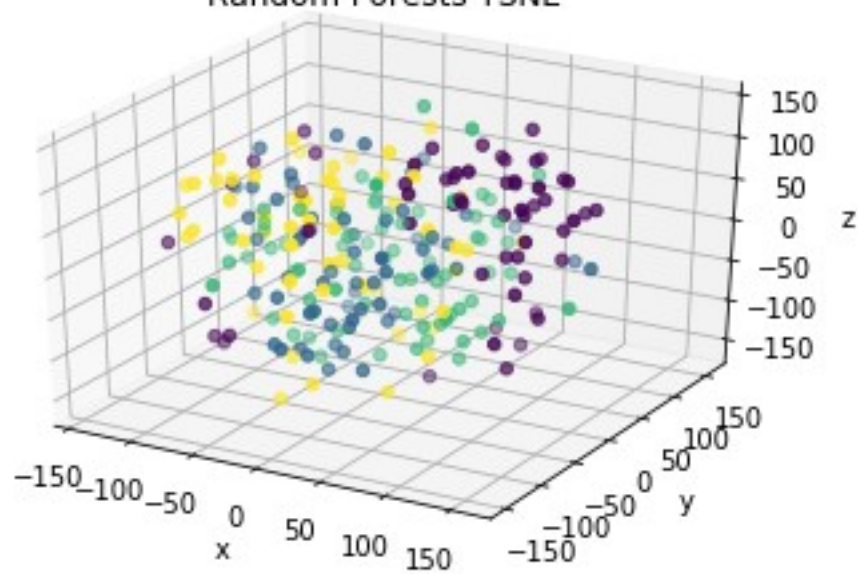## TSNE Regression

Logistic Regression PCA
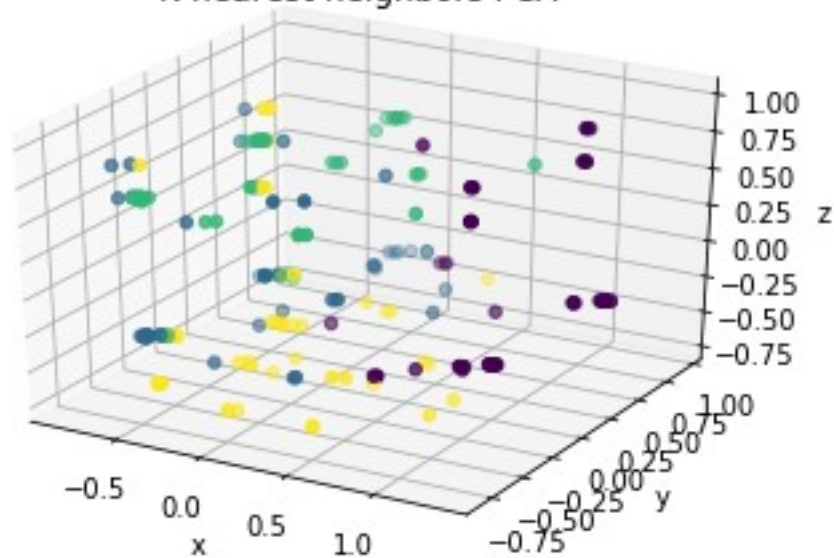


Logistic Regression TSNE
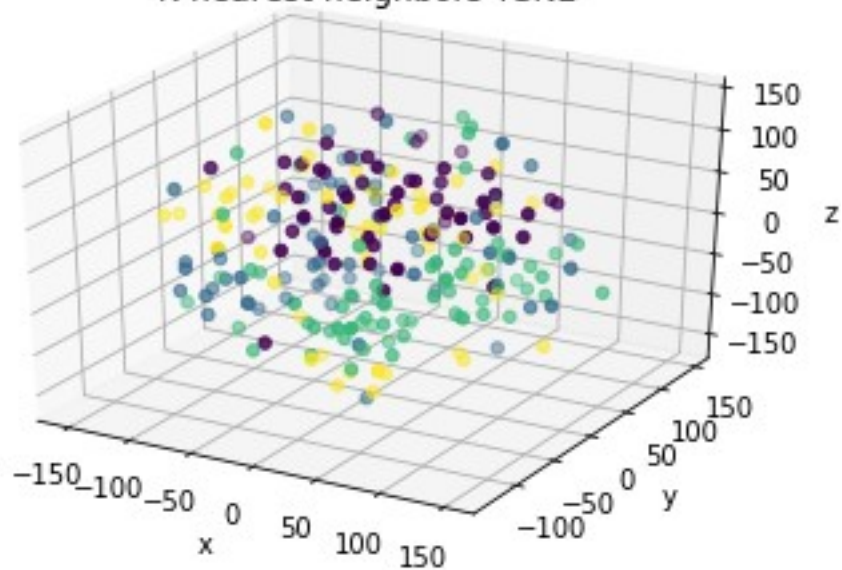
## Random Forests PCA
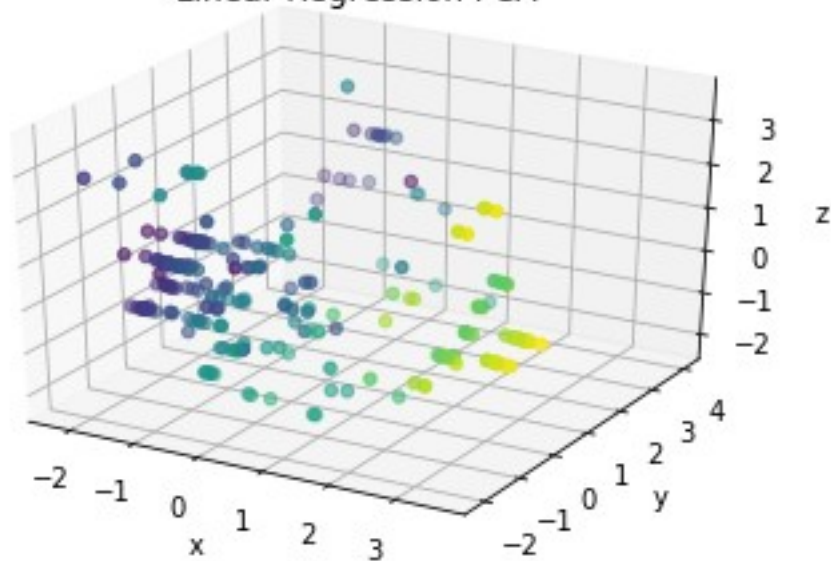


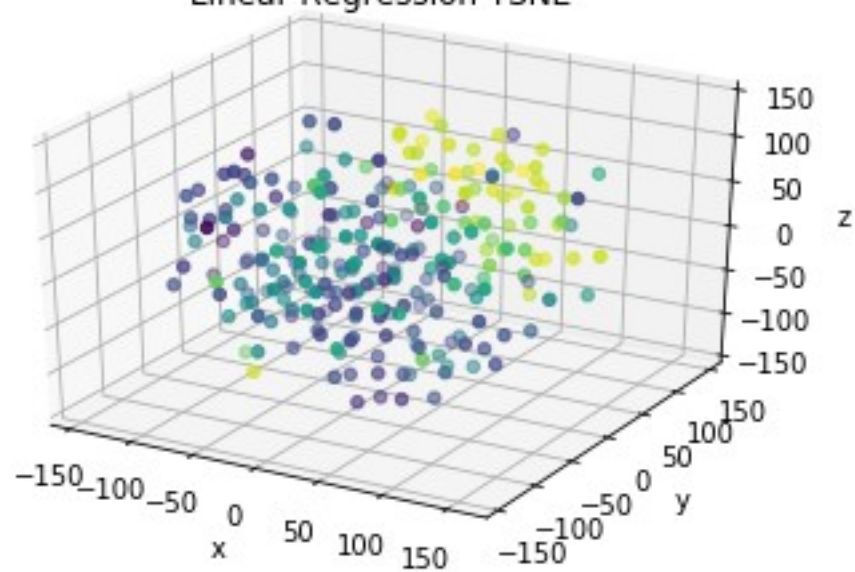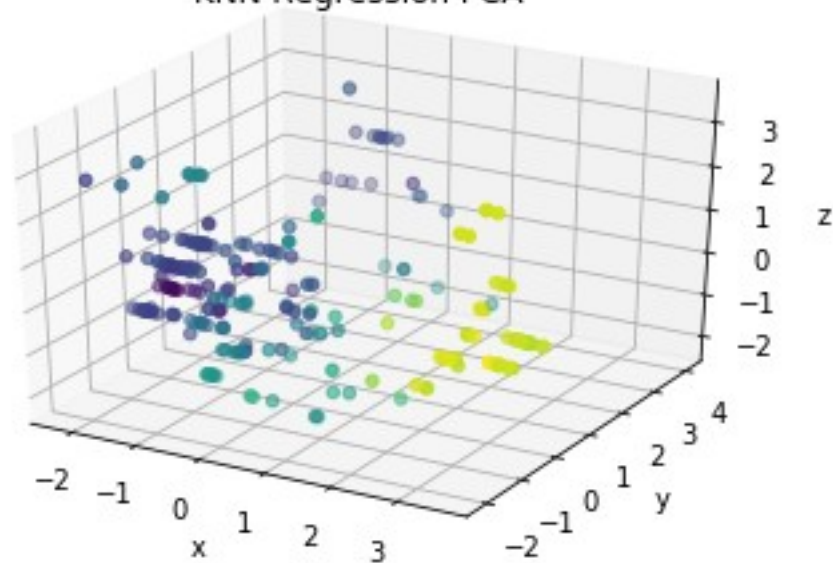## Random Forests TSNE

K-nearest neighbors PCA
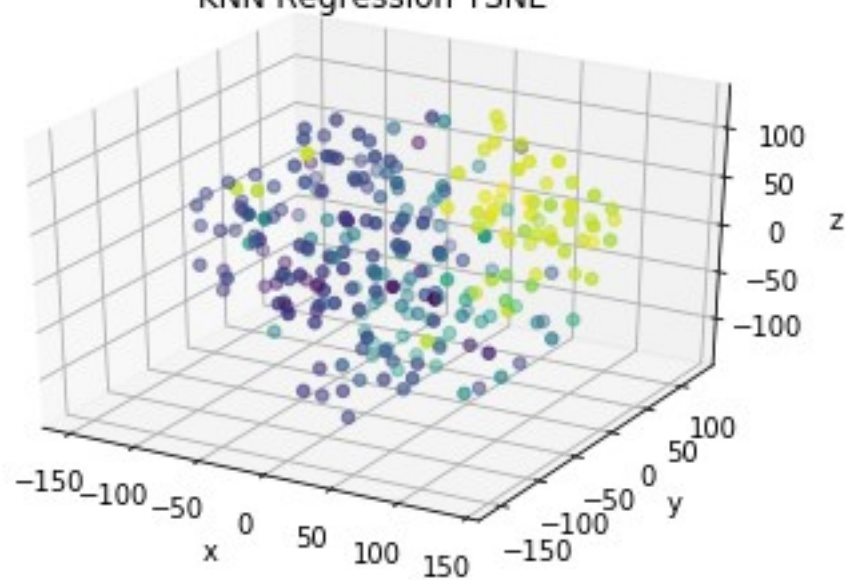

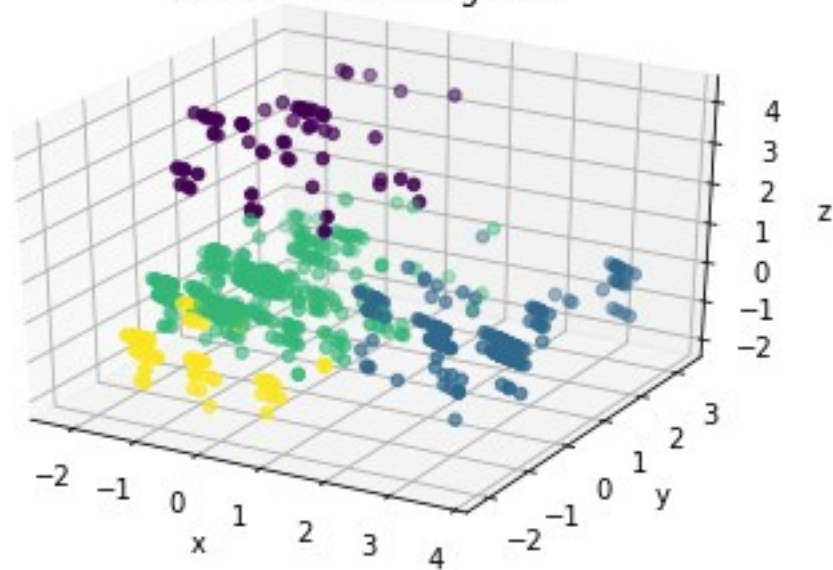K-nearest neighbors TSNE

Linear Regression PCA



Linear Regression TSNE
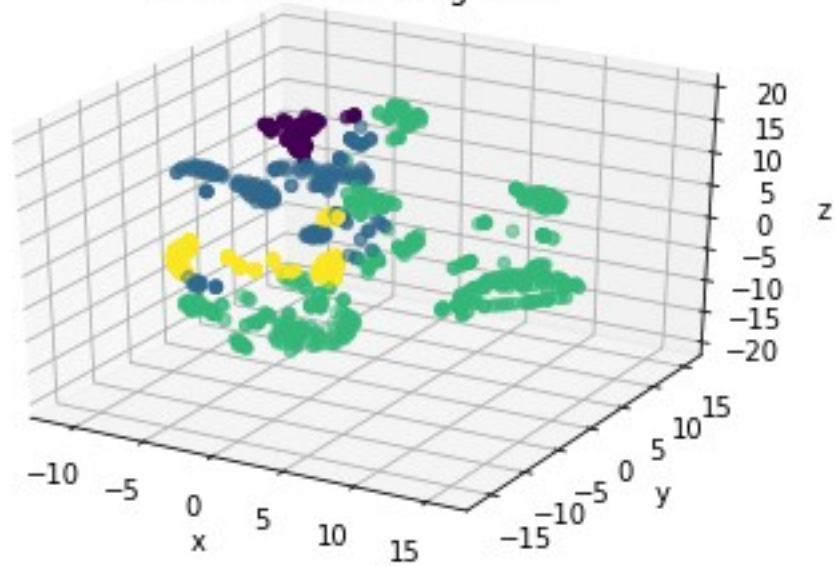
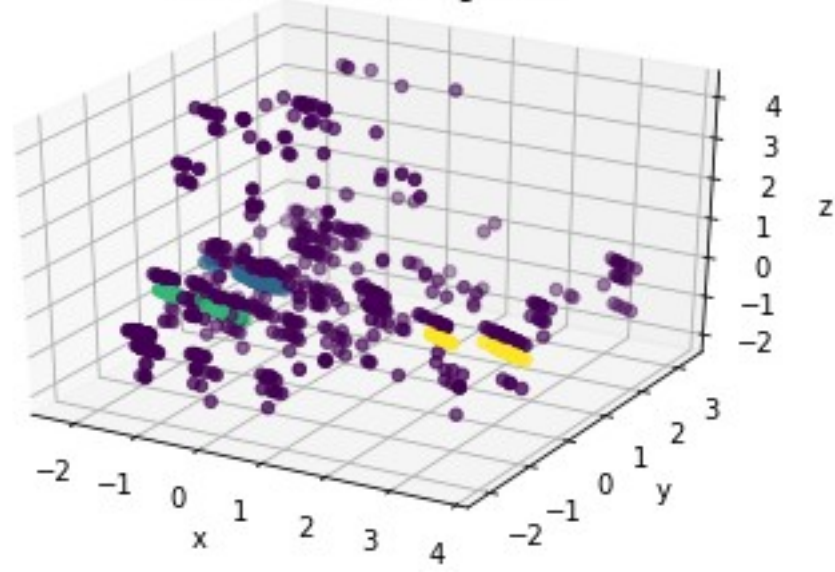## KNN Regression PCA



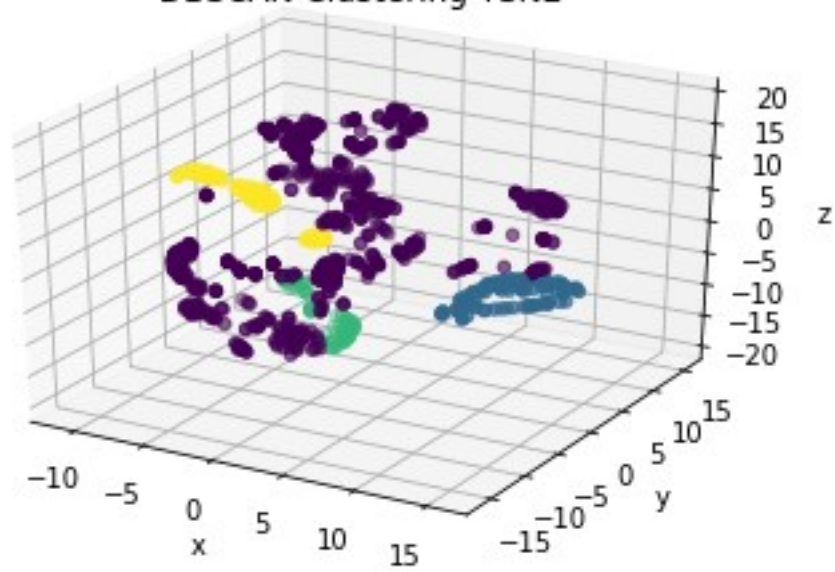## KNN Regression TSNE

## K-means Clustering PCA
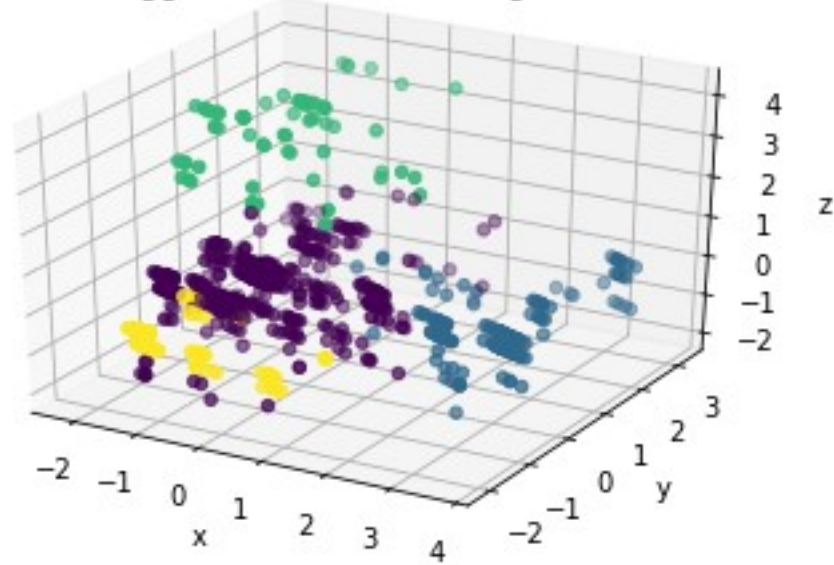


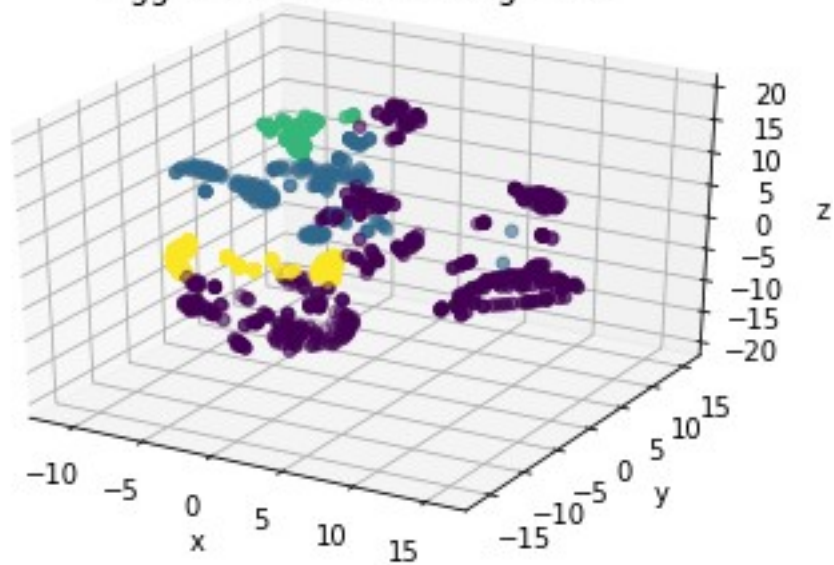## K-means Clustering TSNE

DBSCAN Clustering PCA



DBSCAN Clustering TSNE

Agglomerative Clustering PCA



Agglomerative Clustering TSNE

Needed to run my code:
https://pypi.org/project/imbalanced-learn/
https://pypi.org/project/gender-guesser/