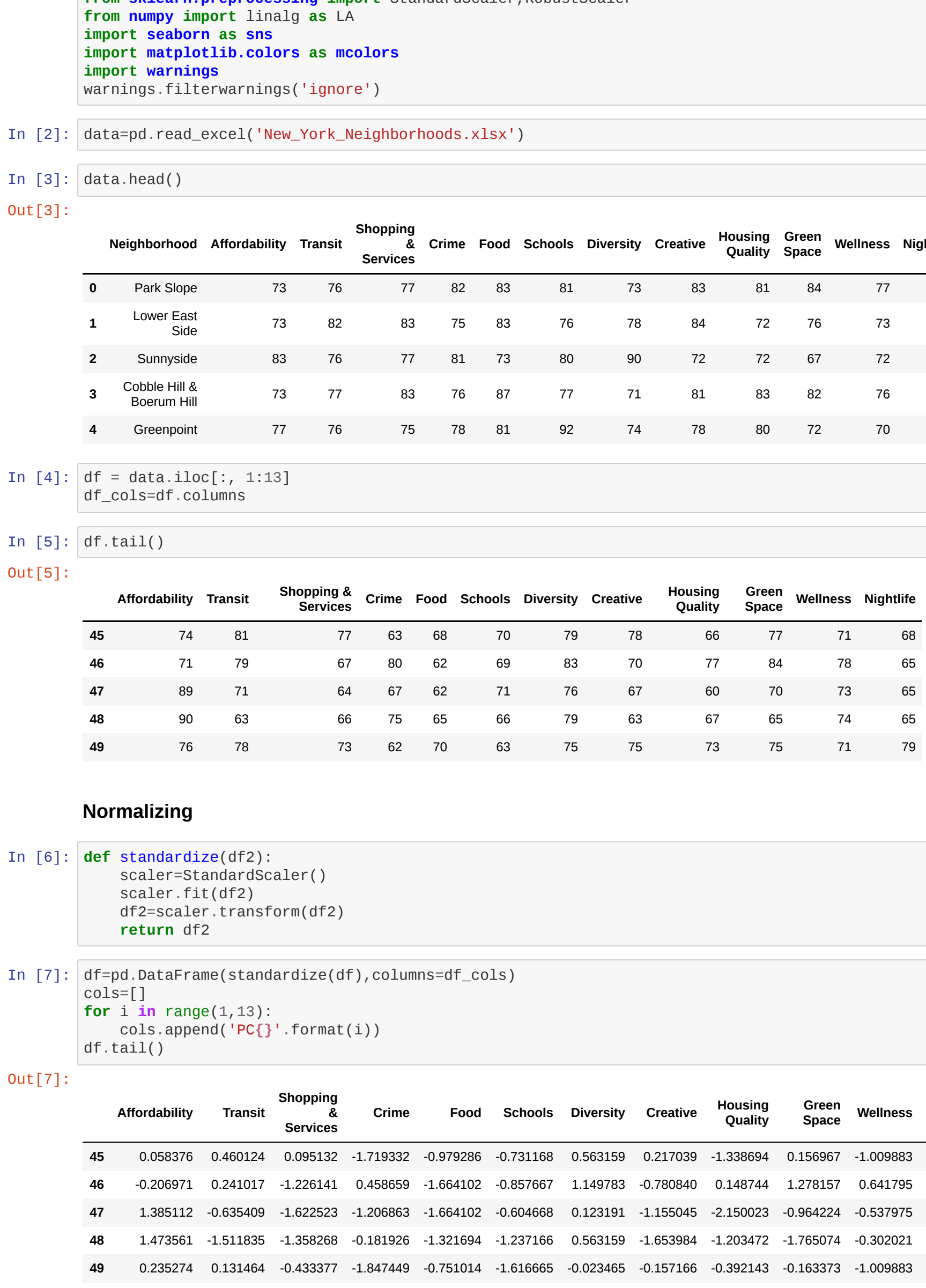


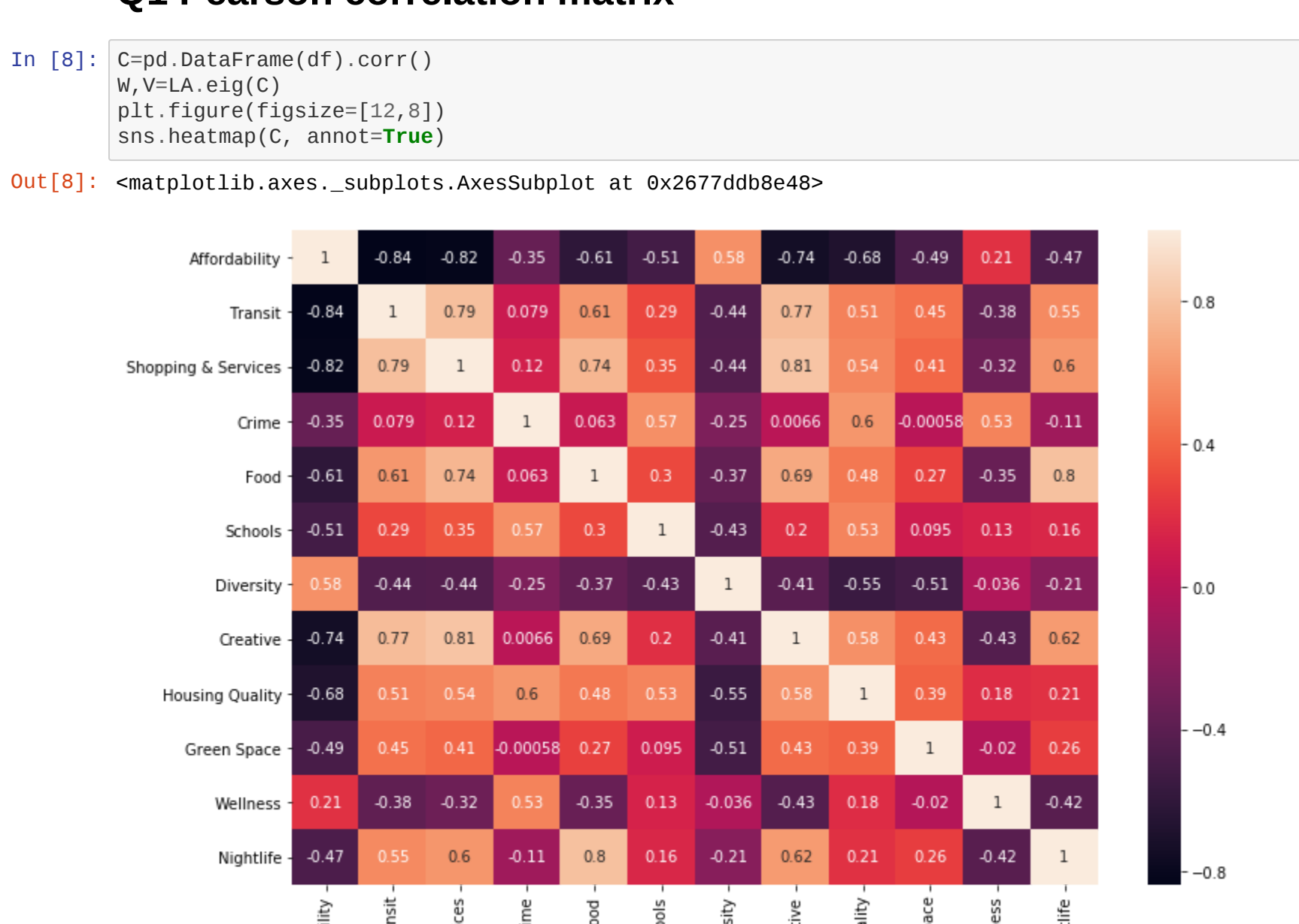
CS306: DATA ANALYSIS AND VISUALIZATION

LAB 6. Principal Component Analysis (PCA) and visualization of multivariate data  
STUDENT ID: 201801407

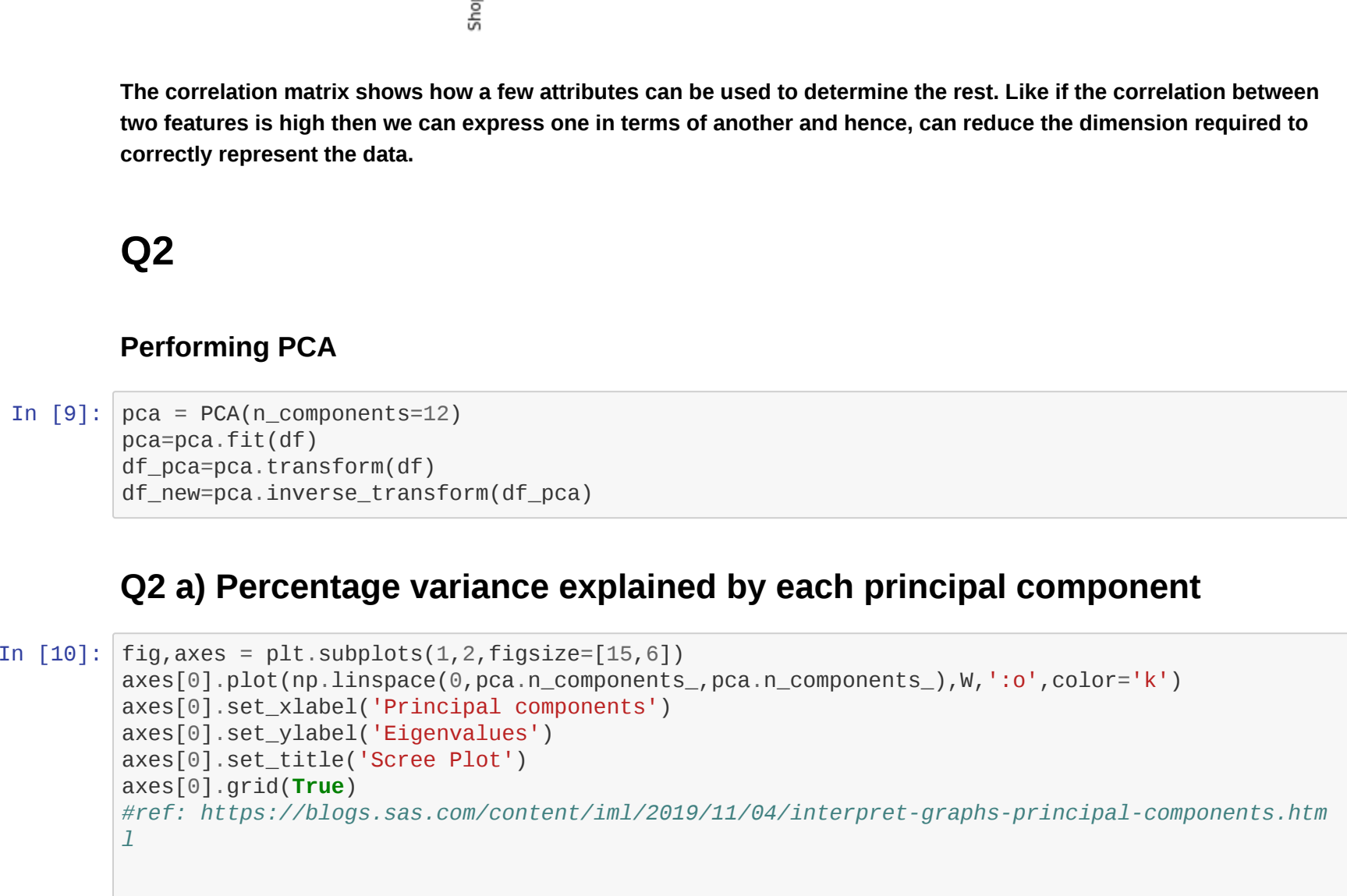
NAME: PRATVI SHAH



Normalizing



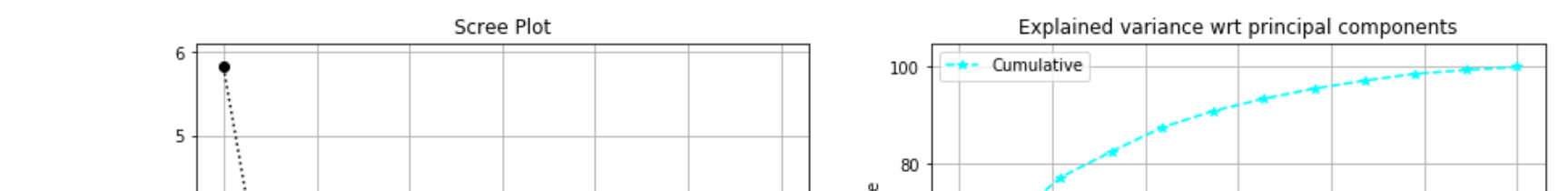
Q1 Pearson correlation matrix



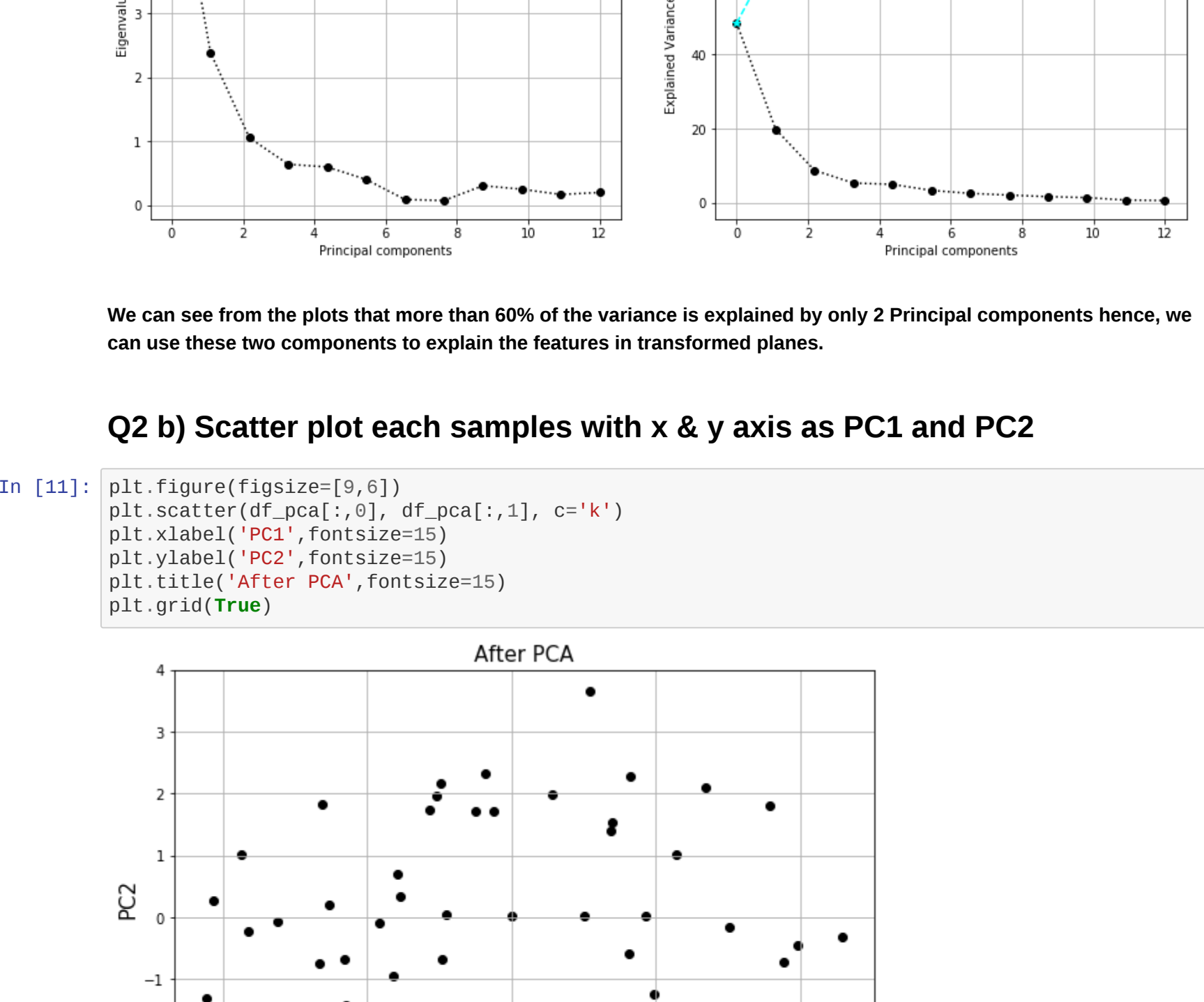
The correlation matrix shows how a few attributes can be used to determine the rest. Like if the correlation between two features is high then we can express one in terms of another and hence, can reduce the dimension required to correctly represent the data.

Q2

Performing PCA

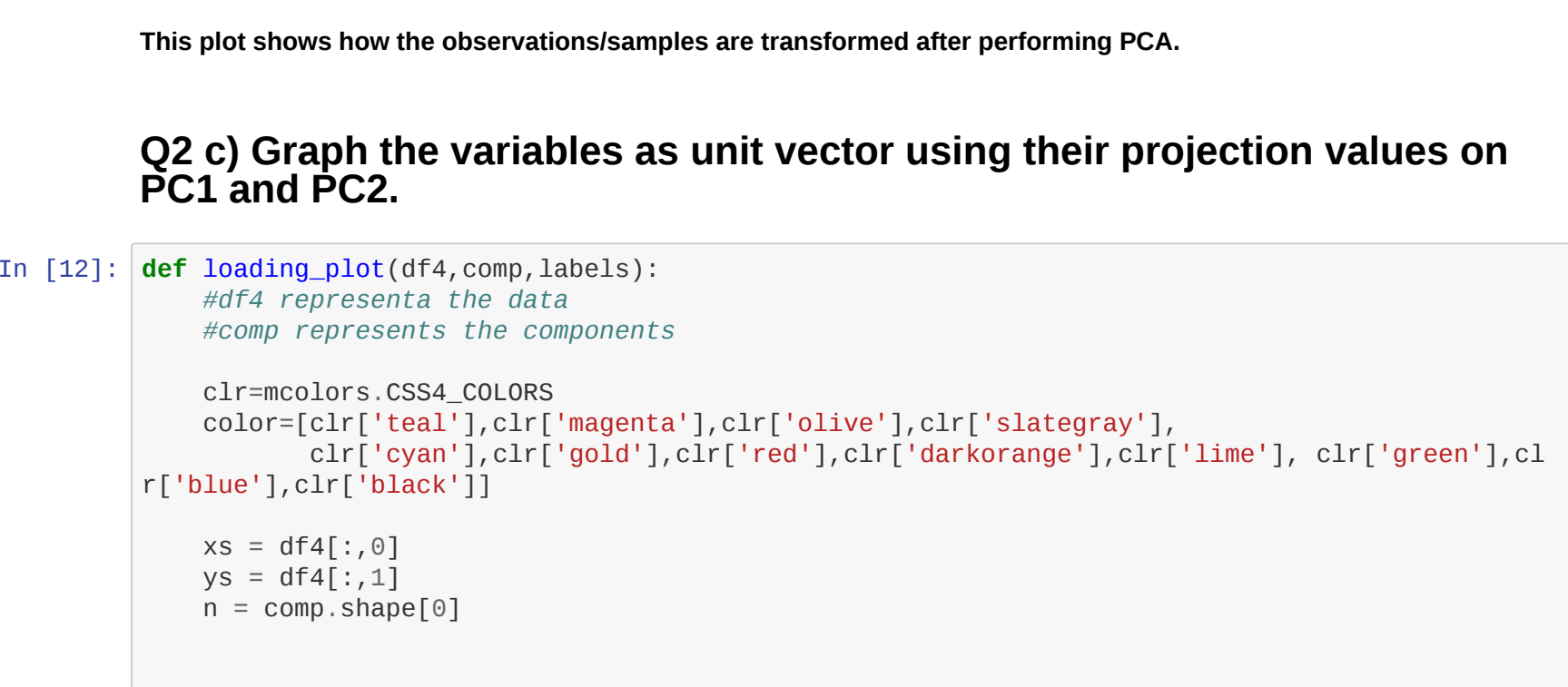


Q2 a) Percentage variance explained by each principal component



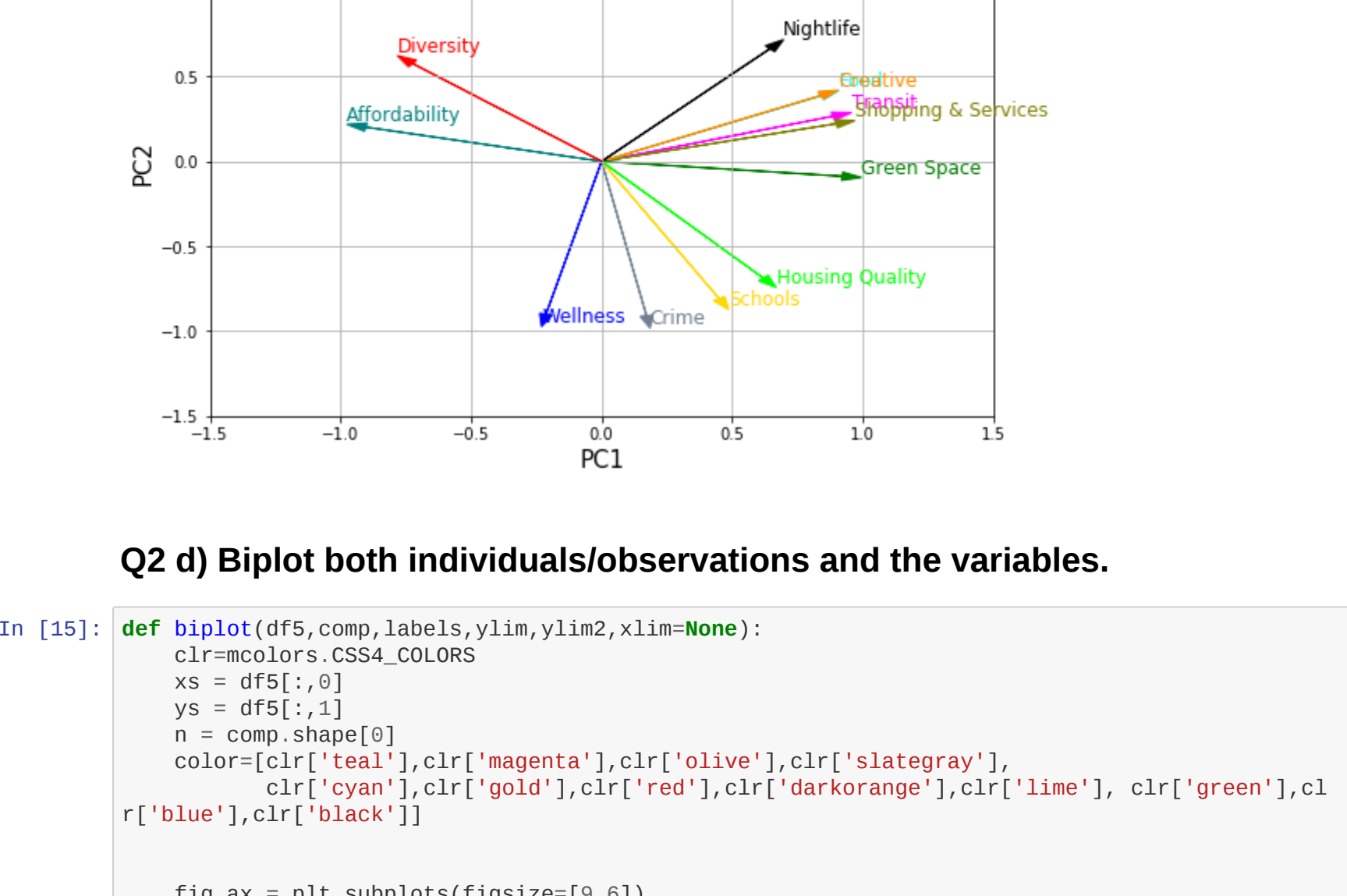
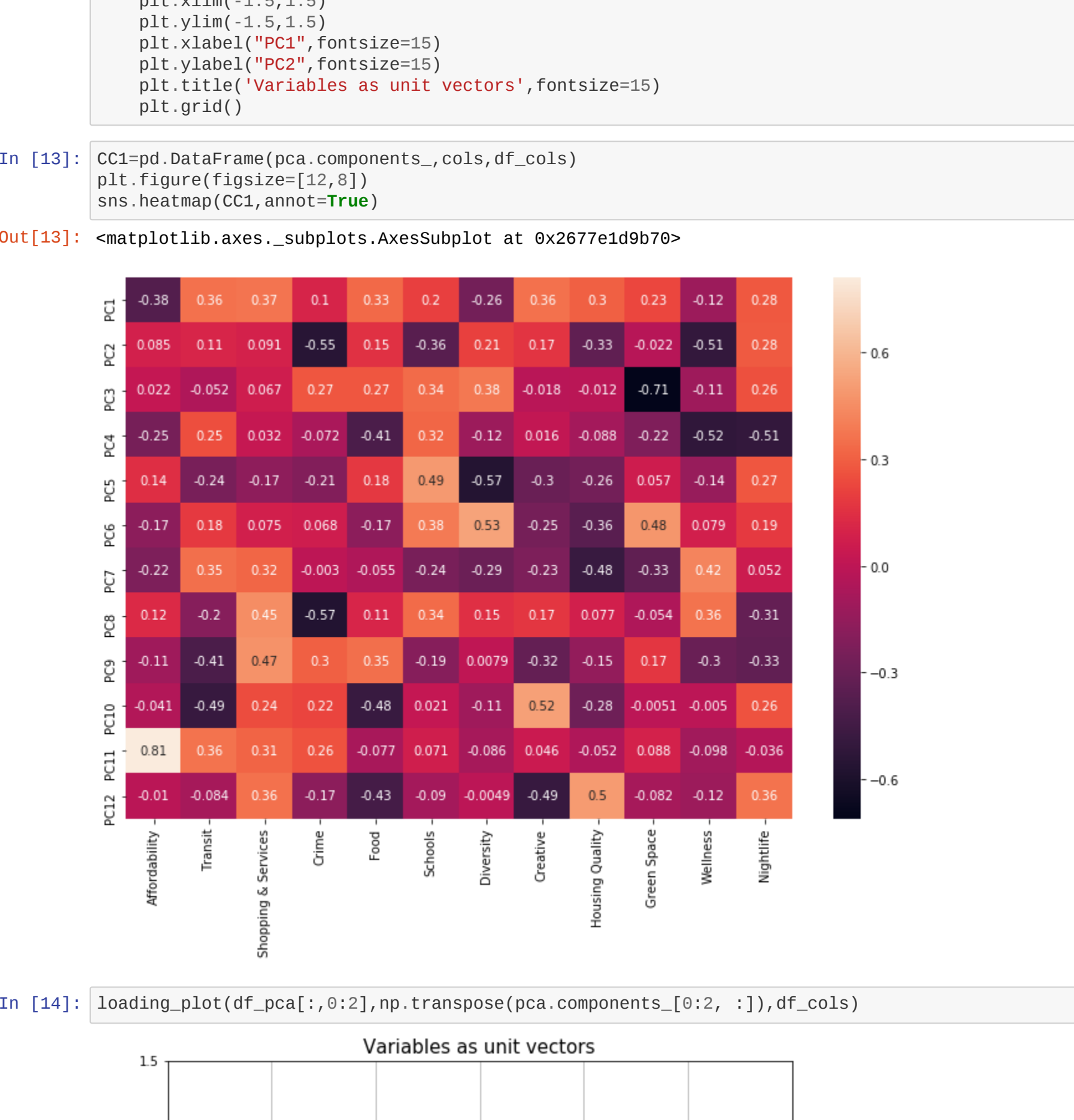
We can see from the plots that more than 60% of the variance is explained by only 2 Principal components hence, we can use these two components to explain the features in transformed planes.

Q2 b) Scatter plot each samples with x & y axis as PC1 and PC2

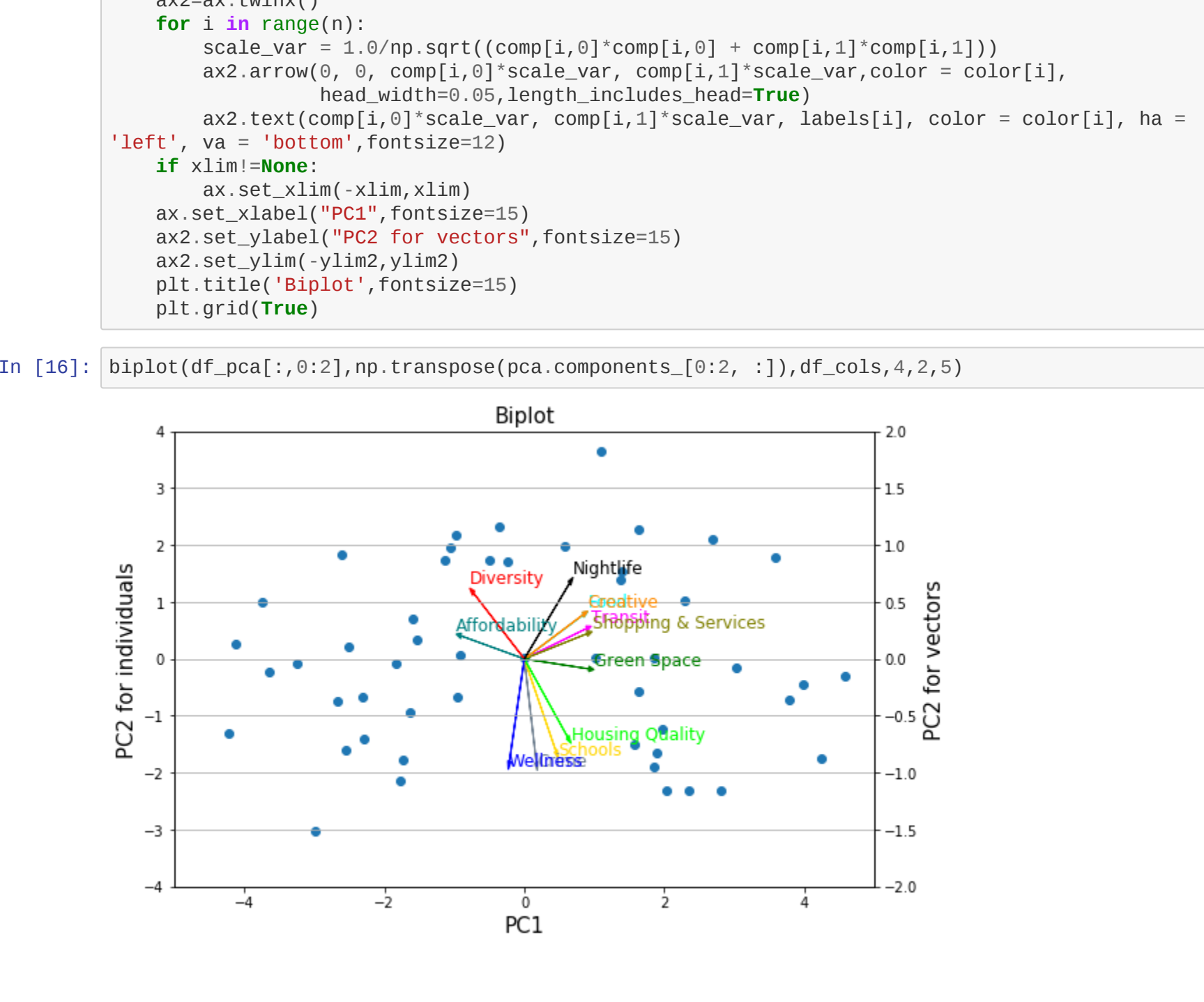


This plot shows how the observations/samples are transformed after performing PCA.

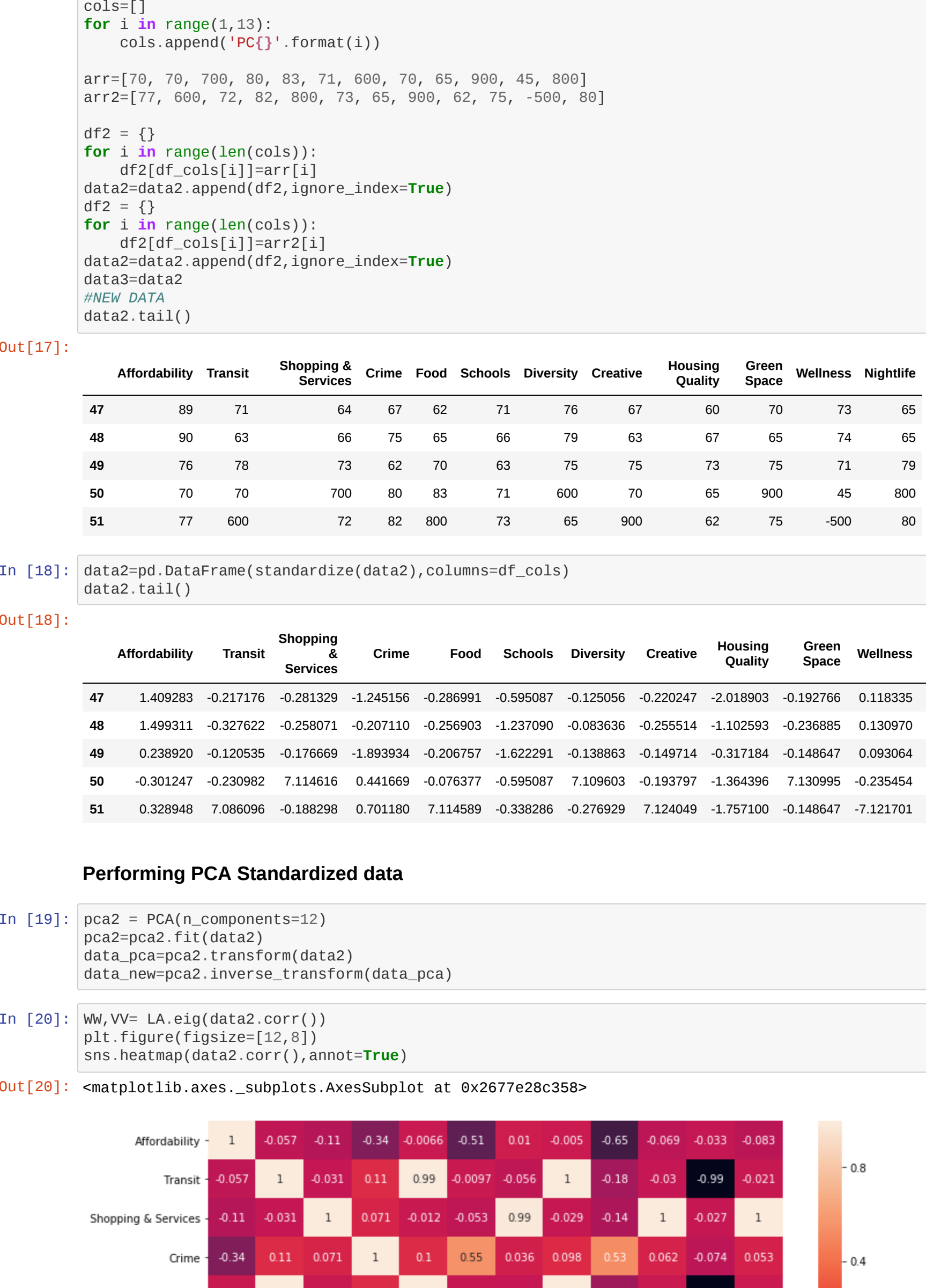
Q2 c) Graph the variables as unit vector using their projection values on PC1 and PC2.



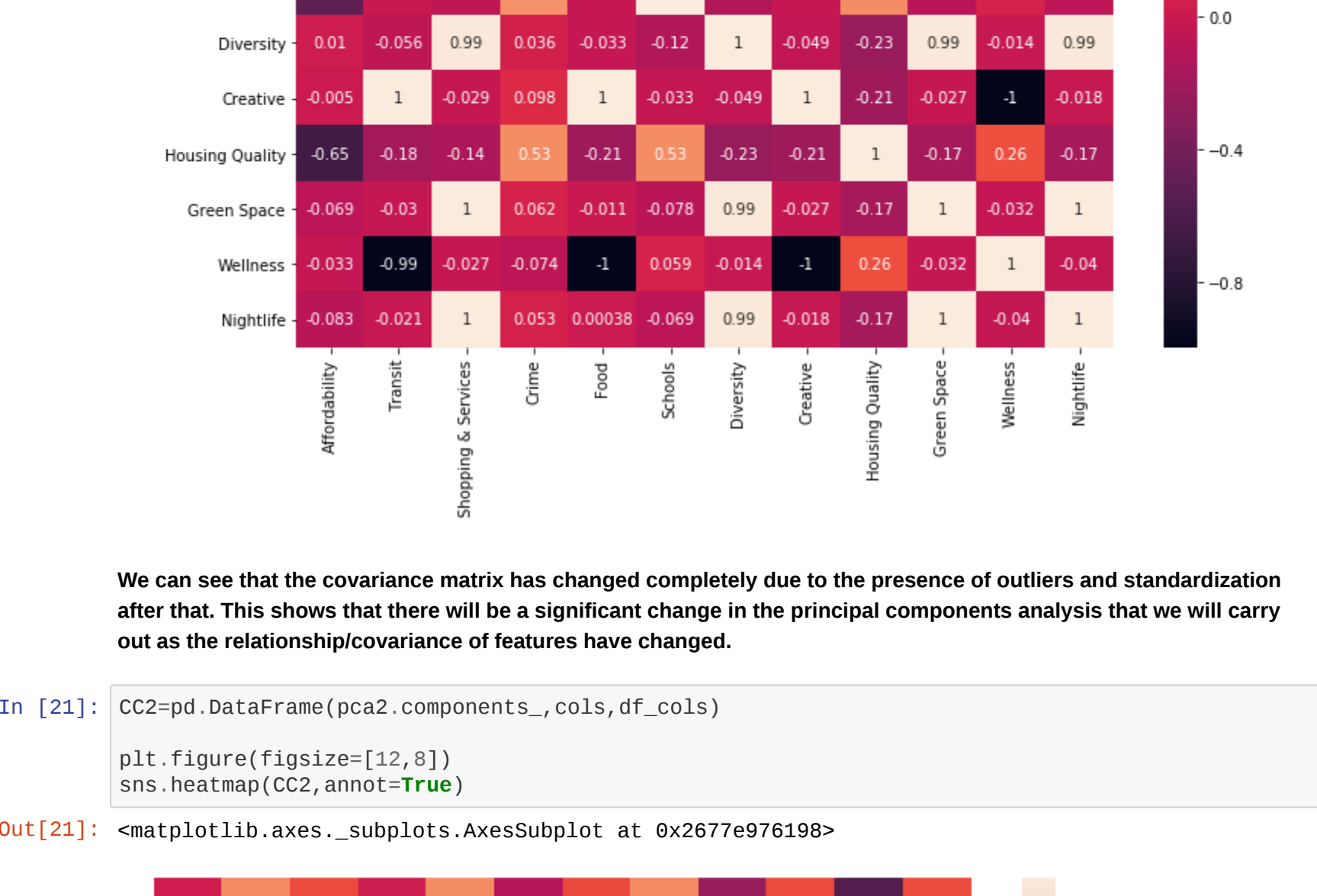
Q2 d) Biplot both individuals/observations and the variables.



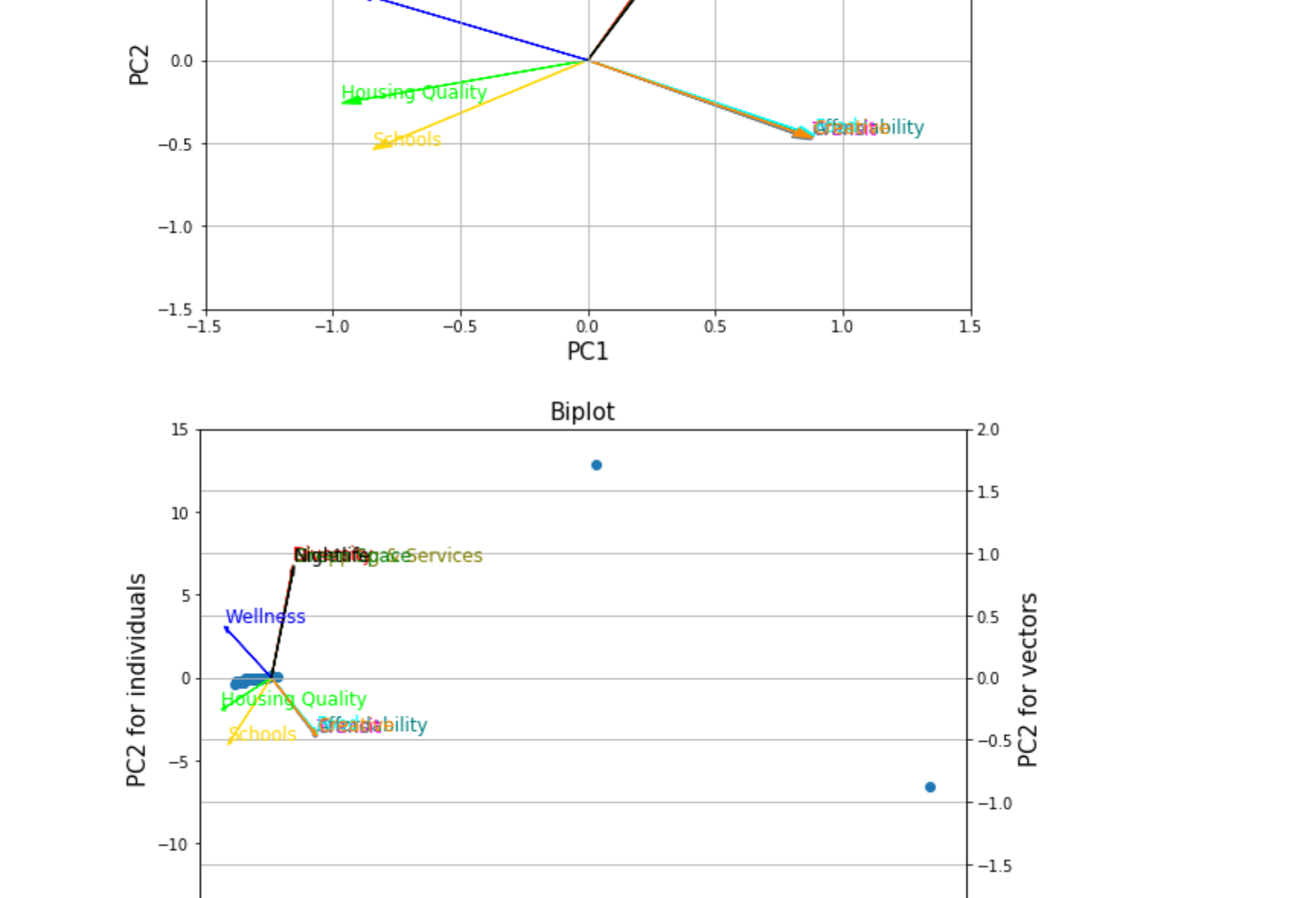
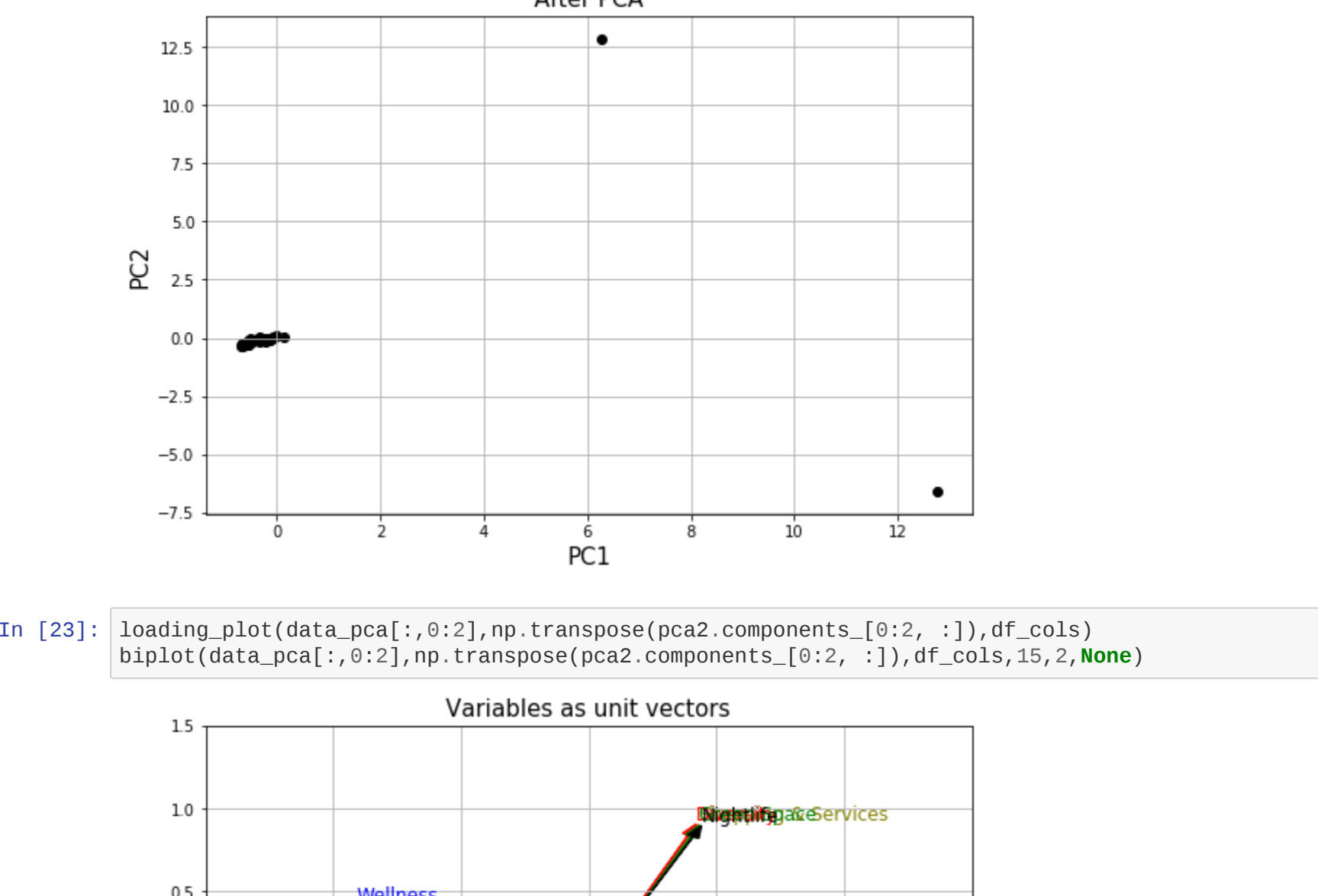
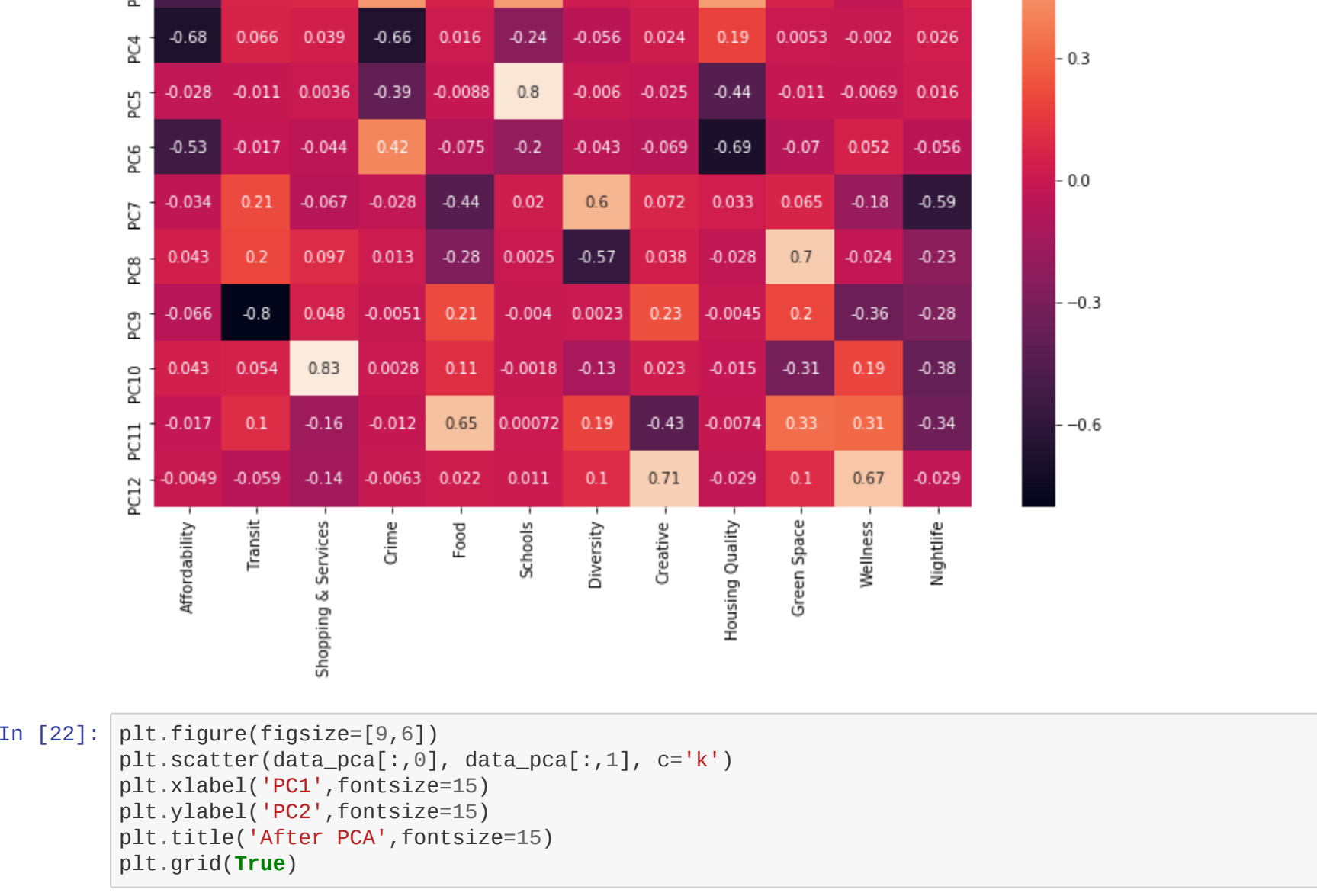
Q3 Introducing outliers



Performing PCA Standardized data

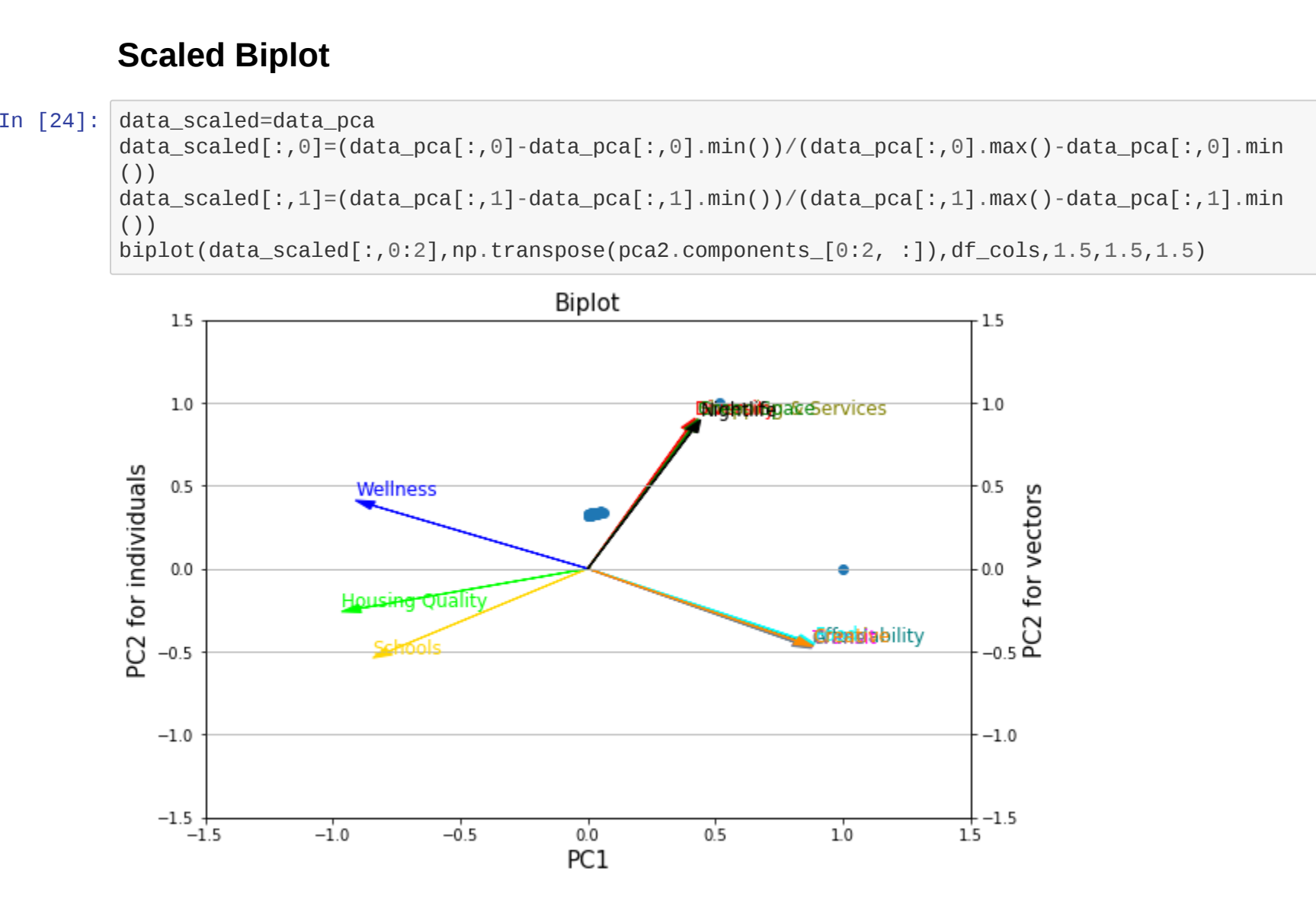


We can see that the covariance matrix has changed completely due to the presence of outliers and standardization after that. This shows that there will be a significant change in the principal components analysis that we will carry out as the relationship/covariance of features have changed.



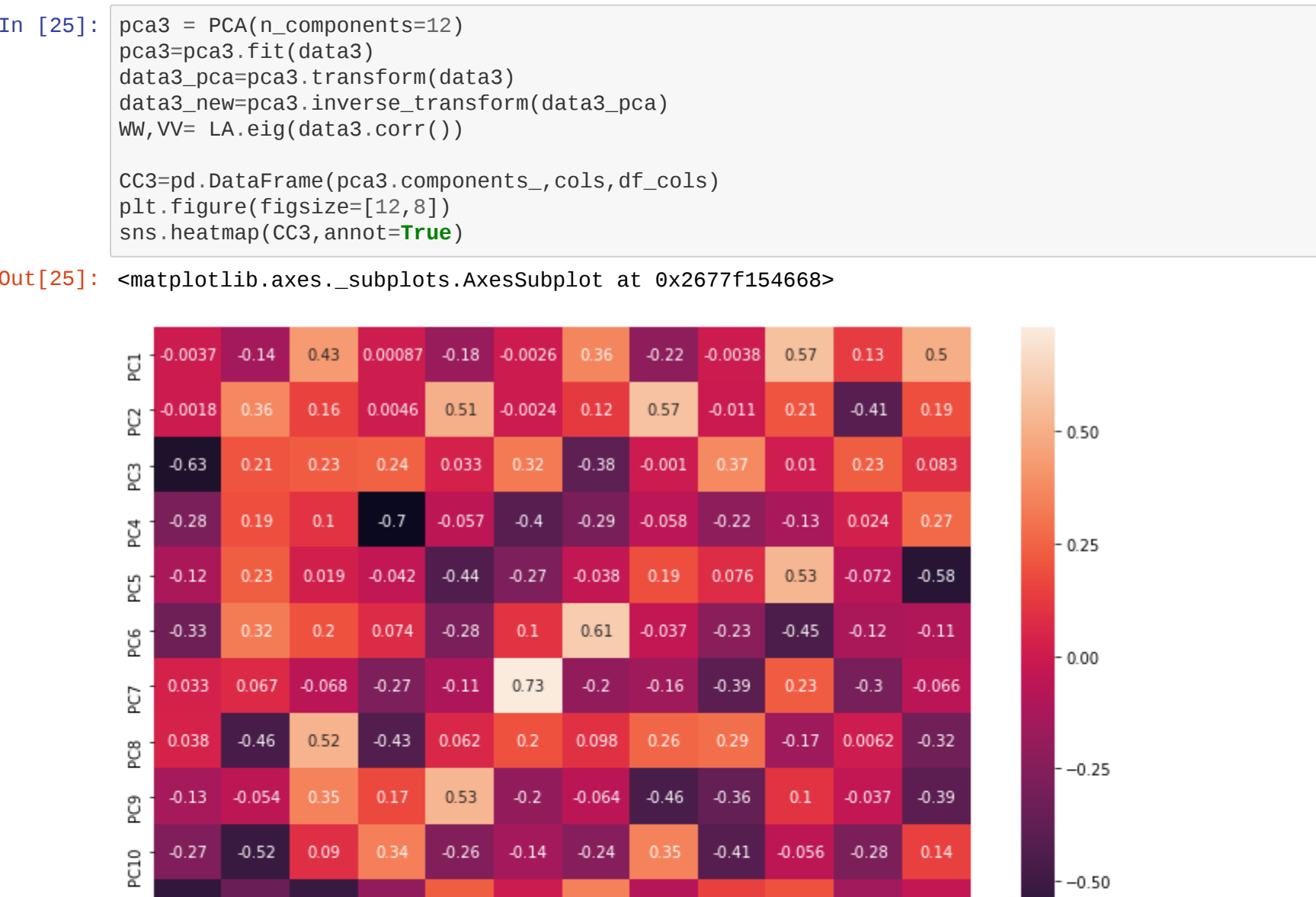
We can see that due to the outliers the exact scatter plot is not visible and the biplot fails to show how the vectors are aligned towards the observations. In order to better visualize this we can scale the scatter plot using minmax scaling

Scaled Biplot



We can see that the scaling of scatter plot gave a better idea as to how the Principal components are changing to accommodate the outliers. The shifting of vectors to point towards the outlier around (0.5,1) in scaled plot shows this.

Performing PCA Unstandardized data



The plots for outliers show that the PCA is highly sensitive to outliers and by introducing a few extreme values can completely change the correlation between the features and hence the principal components also change.

The scaling/standardization also affects the Principal components as the variables are scaled wrt their individual mean and deviation from the mean. This is because the covariance matrix itself changes due to scaling/standardizing. Thus, we can say that any type of scaling or introduction of extreme data points will change the PCs.

