# Wine Cultivar Clustering

Dominic Maranta

# Abstract

- Wine Cultivar data set
- Sourced, cleaned, and modeled the data using Python.
- Unsupervised machine learning problem using clustering to identify cultivars
- Both k-means clustering and agglomerative clustering are highly effective when combined with PCA
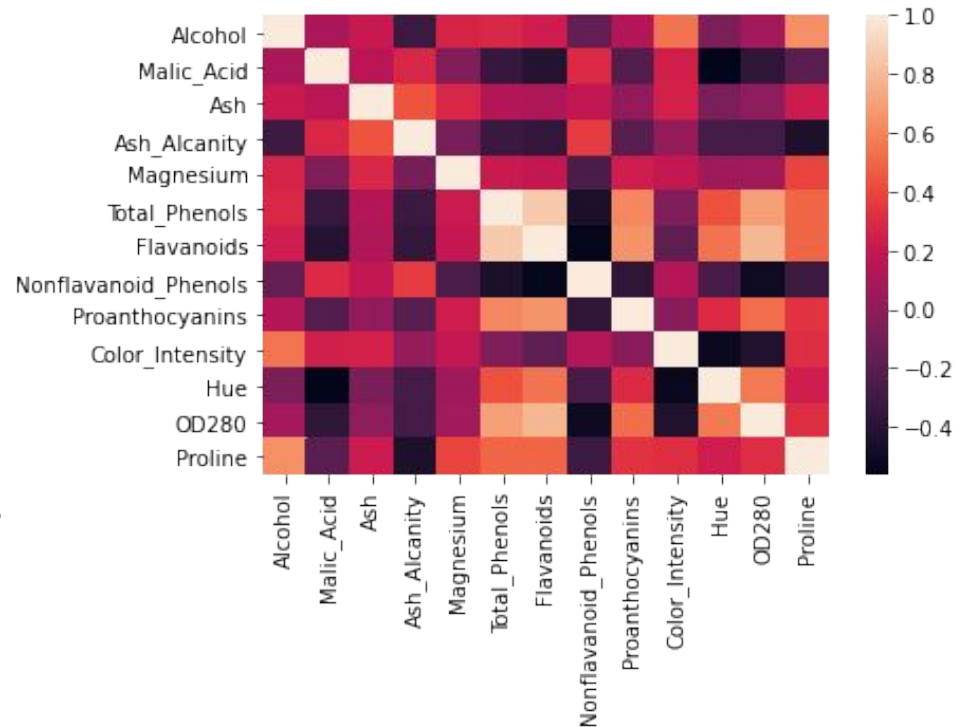
# Problem Statement

- Construct and train a model to cluster wine types
- How can we properly identify the cultivar a wine is from.
- Use EDA to explore and prepare the data.
- Compare and Evaluate machine learning clustering models.

# Data Preprocessing

- Data: 13 variables with 178 instances
- 13 features identify characteristics including alcohol content, hue, and many more.
- No missing values in the 178 observations
- Clean variables and create numeric variables
- Perform EDA before modeling
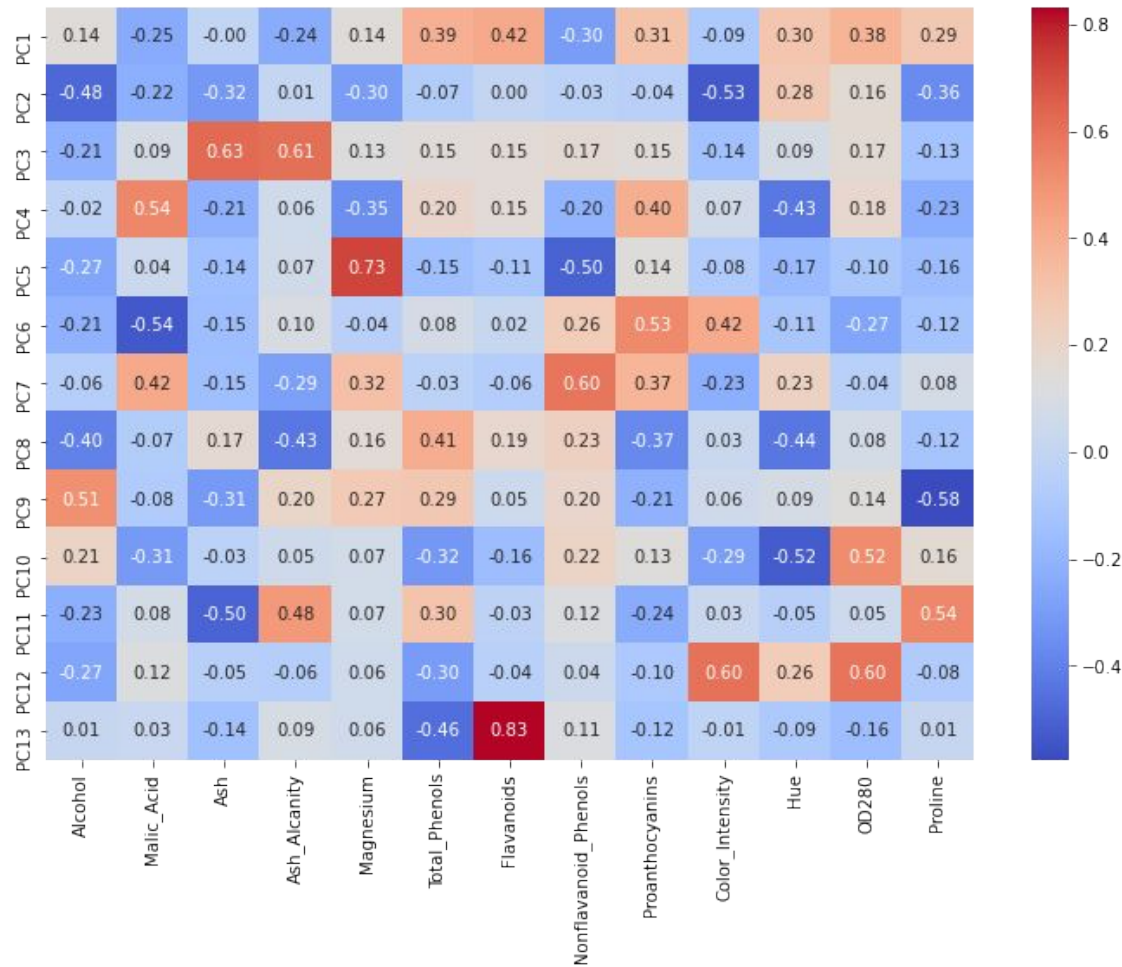
# Data Understanding

- Phenols, flavonoids, nonflavanoids, alcohol and ash alcanity, all produce strong correlation patterns, both inversely and directly
- Examining individual histograms from these 5 variables show distinct humps in the graphs, denoting potential differences in cultivars
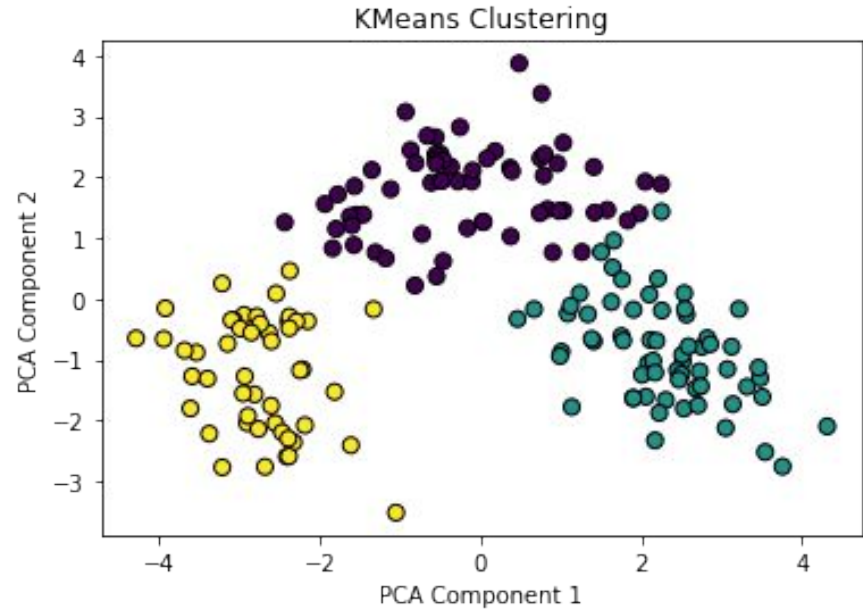
# Modeling

- Before clustering a key step is to standardize all the data.
- Apply a PCA transformation to reduce dimensions, multicollinearity, and assist with visualization.
- First 2 Principal components explain 47% of the variation.
- Phenols, flavonoids, nonflavanoids, and ash alcanity all important in the first PC.
- Alcohol and color intensity very important in the second PC.
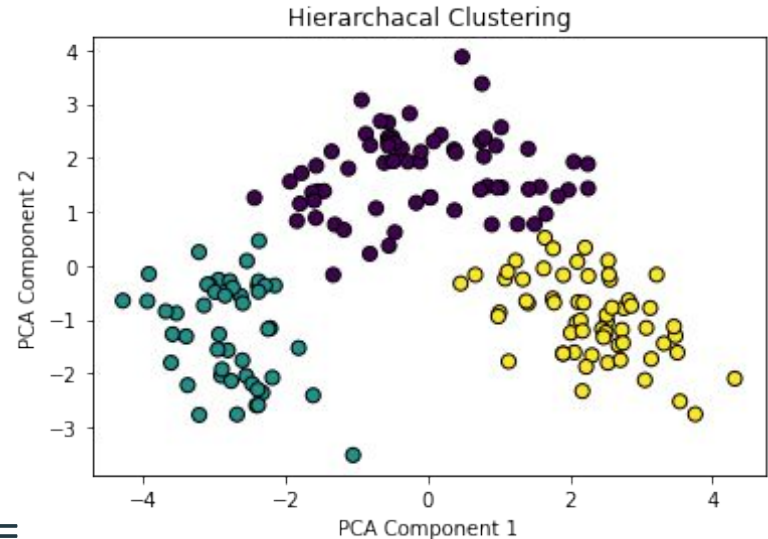
# PCA:

# K Means Clustering

- K-means clustering performed optimally with 2 Principal components, compared to the raw data and an increased number of components
- Evaluated accuracy for three clusters = 96.6%



KMeans Clustering

# Agglomerative Clustering

- Agglomerative clustering performed optimally with 2 Principal components, compared to the raw data and an increased number of components
- Ward linkage methods performed much better than complete, single, and average
- Evaluated accuracy for three clusters = 96.6%



Hierarchacal Clustering

# Outcomes

- Both agglomerative clustering and k-means clustering are effective, as long as PCA is used.
- New wine prediction is made easier with only wine features needed.
- Will allow for expanded understanding of how cultivars differ and which variables are important between cultivars.

# Future Work

- Expand the dataset to include more features, observations, and cultivars.
- Apply supervised machine learning approaches to this data.
- Investigate the data with deep learning approaches.

Thank you!