

Challenge Report

Angel Salazar - 513236

Master in Applied Artificial Intelligence

TC-4029 - Analytics and Data Science

Data Cleansing, Visualization, Clustering, Classification

Intro

The challenges were designed to put into practice the topics covered in the subject during this academic period.

They are a great example of what a data scientist must face when working with actual data.

This report aims to provide a summary of the challenges and outcomes after applying data scientist and machine learning techniques to a given dataset.

Motivation

Learn and practice the following data scientist and machine learning techniques:

- Data cleansing and data normalization
- Unsupervised learning using a clustering technique (KMeans)
- Supervised learning using multiclass classification (Random Forest / Decision trees)

Research Statements

Find out whether the water quality can be related to its source location

Evaluate and select a multiclass classification model

Data

- Mexico
- Records related to underground water quality
- 1,068 water sources

Methodology



Data cleansing and analytics

Deal with missing data, normalize values, visualize outliers and mean



Clustering

Compute optimal number of clusters, compare clustering result to sempharo coloring result



Model Training

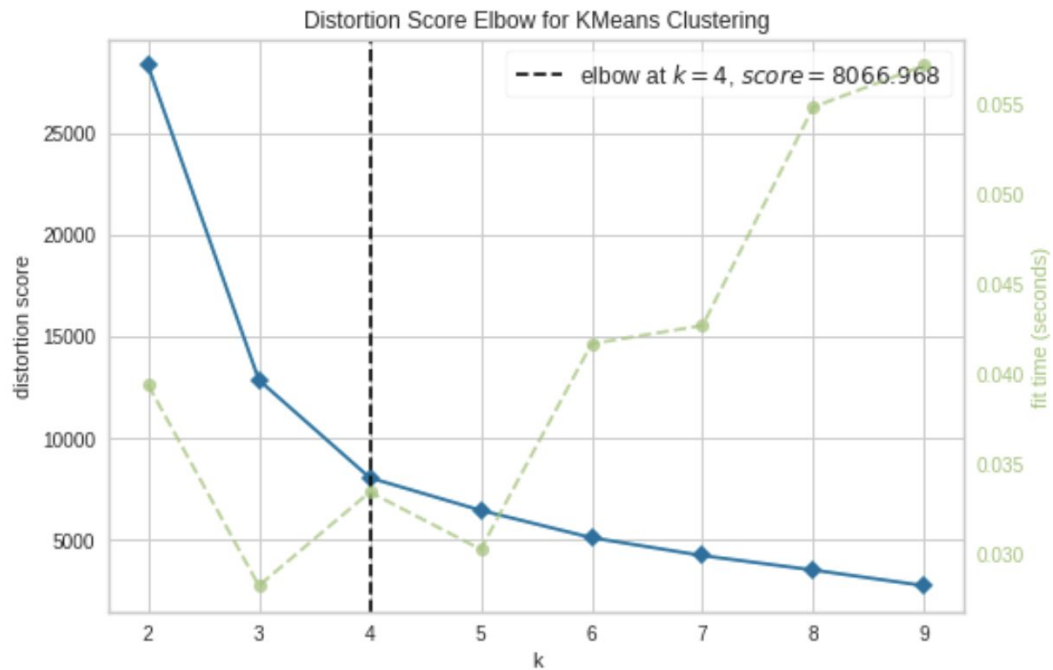
Compute feature importance, match up classification models, select optimal model



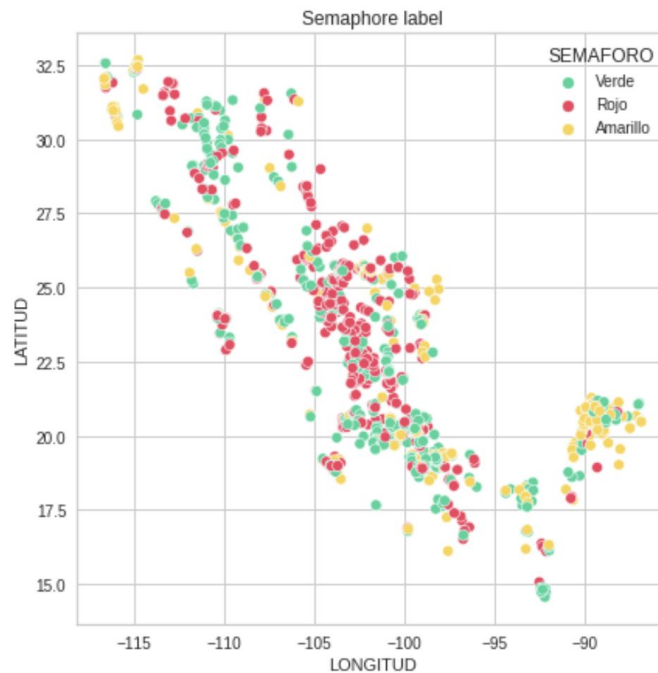
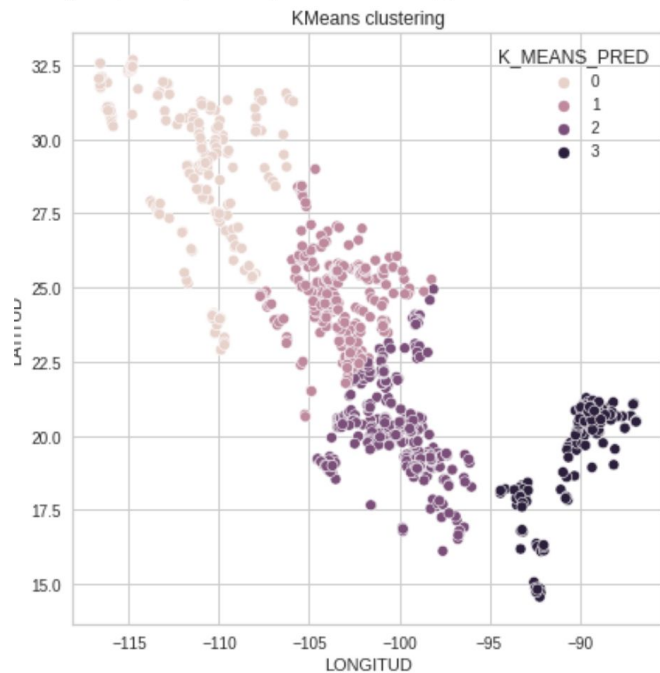
Classify data

Predict water quality based on selected features, compute models statistics

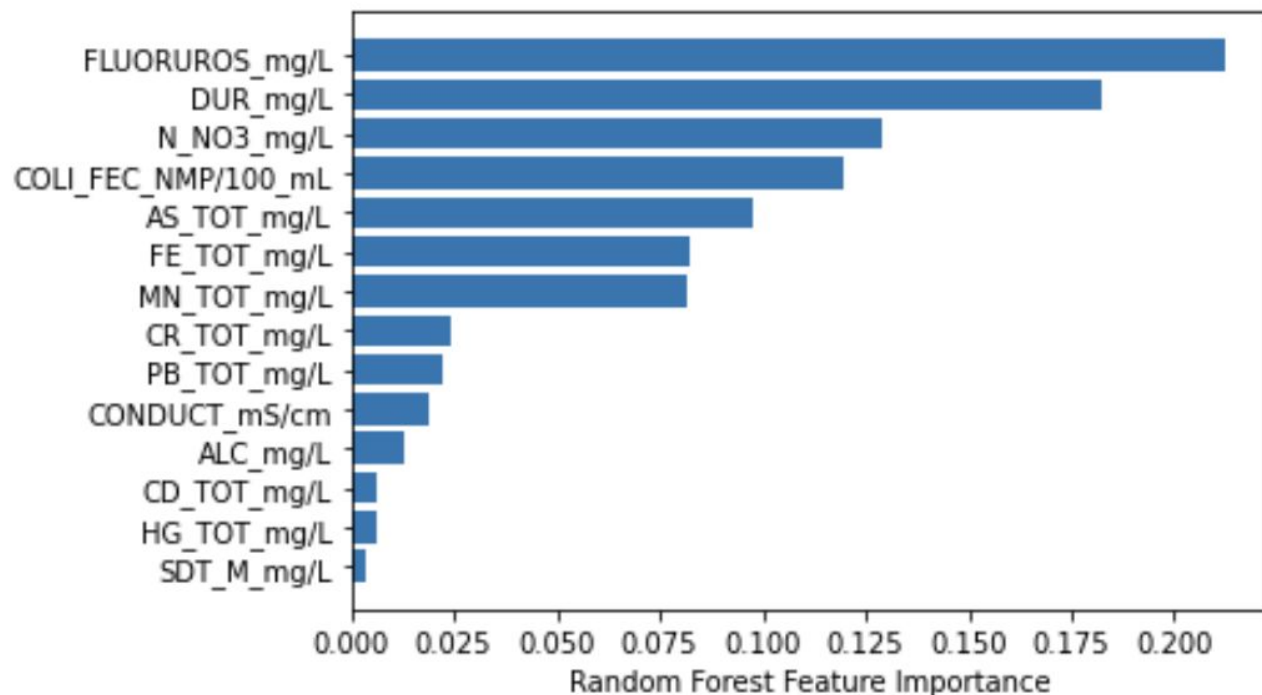
Outcomes



Outcomes



Outcomes



Outcomes

RandomForestClassifier

Regular training

	precision	recall	f1-score	support
Verde	0.78	0.92	0.85	87
Amarillo	0.91	0.58	0.71	55
Rojo	0.98	0.94	0.96	69
micro avg	0.87	0.84	0.86	211
macro avg	0.89	0.81	0.84	211
weighted avg	0.88	0.84	0.85	211
samples avg	0.88	0.84	0.84	211

Cross validation training

accuracy 0.863582443653618

DecisionTreeClassifier

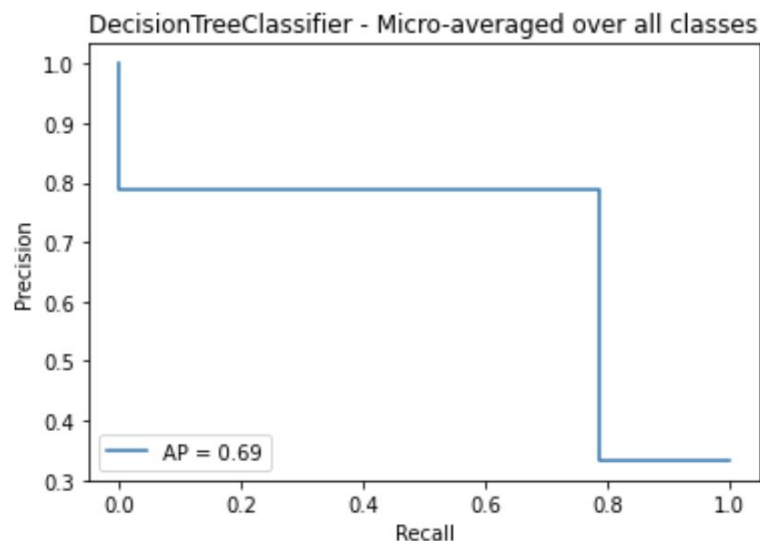
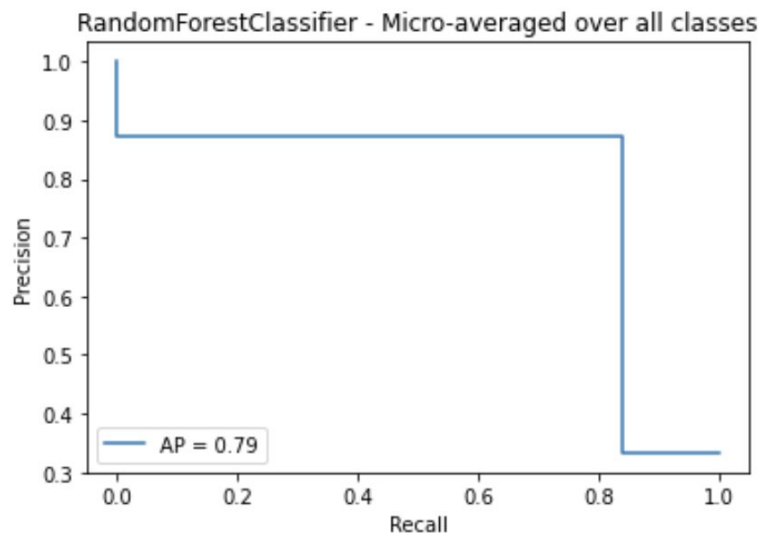
Regular training

	precision	recall	f1-score	support
Verde	0.77	0.77	0.77	87
Amarillo	0.69	0.64	0.66	55
Rojo	0.88	0.93	0.90	69
micro avg	0.79	0.79	0.79	211
macro avg	0.78	0.78	0.78	211
weighted avg	0.78	0.79	0.78	211
samples avg	0.79	0.79	0.79	211

Cross validation training

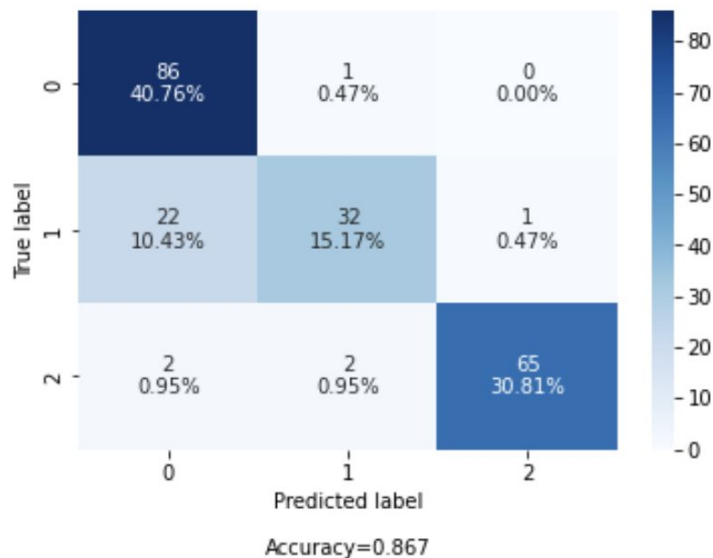
accuracy 0.8315539739027283

Outcomes

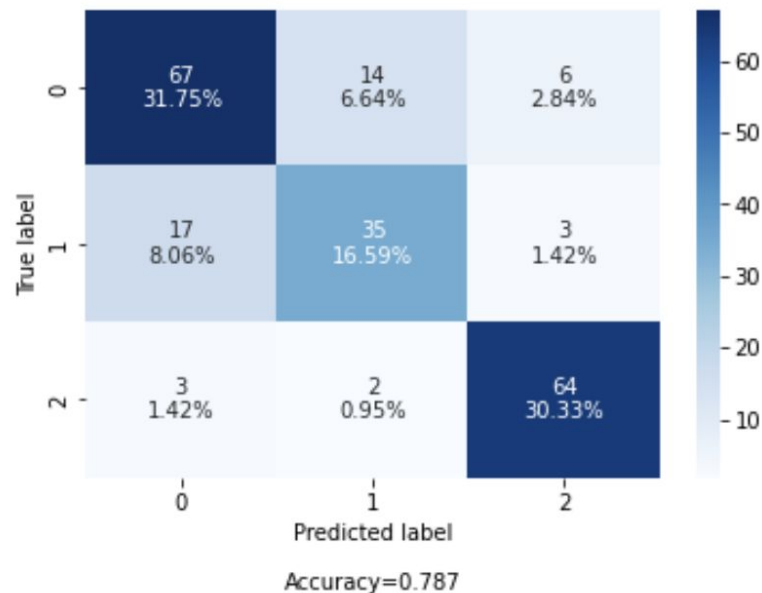


Outcomes

Random Forest Classifier



Decision Tree Classifier



Challenges

- Data cleansing and normalization
- Plotting illustration

Conclusion

Find out whether the water quality is related to its source location

There is no relationship between the entries' coordinates and water quality

Evaluate and select a multiclass classification model

RandomForestClassifier which reported an accuracy reported of 86.3%