



# Vertica ML Python Workshop

## Exercise 10: ML Clustering

Ouali Badr

December 6, 2019

## Executive Summary



*"Science knows no country, because knowledge belongs to humanity, and is the torch which illuminates the world."*

**Louis Pasteur**

VERTICA ML PYTHON allows the users to use Vertica advanced analytics and Machine Learning with a Python front-end Interface. In this exercise, you'll learn some basics to begin your fantastic Data Science Journey with the API. As a summary:

- Create a KMeans Model
- Find the K of the KMeans
- Add the prediction in the Virtual Dataframe
- Evaluate the model

## Contents

<b>1</b>	<b>Presentation</b>	<b>3</b>
<b>2</b>	<b>Functions used during the Exercise</b>	<b>3</b>
2.1	KMeans . . . . .	3
2.2	best_k . . . . .	5
<b>3</b>	<b>Questions</b>	<b>6</b>

## 1 Presentation

Cluster analysis or clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense) to each other than to those in other groups (clusters). We already did some clustering in exercise 2. It was mainly focused on Outliers and Anomaly Detections rather than doing a complete segmentation. We will use exercise 2 functions to segment flowers in the Iris dataset.

## 2 Functions used during the Exercise

### 2.1 KMeans

Create a KMEANS object by using the Vertica Highly Distributed and Scalable KMEANS directly on the data.

#### initialization

**Library:** `vertica_ml_python.learn.cluster`

```
1 class KMEANS (
2     name: str,
3     cursor,
4     n_cluster: int = 8,
5     init = "kmeanspp",
6     max_iter: int = 300,
7     tol: float = 1e-4)
```

#### Parameters

- **name:** *<str>*  
Name of the the model. The model is stored in the DB.
- **cursor:** *<object>*  
DB cursor.
- **n\_cluster:** *<int>*, optional  
Number of clusters.
- **init:** *<str or list>*, optional  
The method used to find the initial cluster centers. The method can be in {random | kmeanspp}. It can be also a list with the initial cluster centers to use.
- **max\_iter:** *<int>*, optional  
The maximum number of iterations the algorithm performs.
- **tol:** *<float>*, optional  
Determines whether the algorithm has converged. The algorithm is considered converged after no center has moved more than a distance of 'tol' from the previous iteration.

#### Methods

The KMeans object has four main methods:

```

1 # Add the cluster prediction in a vDataFrame
  def add_to_vdf(self, vdf, name: str = "")
3
  # Drop the model
5 def drop(self)
7
  # Fit the model with the input columns
  def fit(self, input_relation: str, X: list)
9
  # Plot the model (only possible if the number of columns <= 3)
11 def plot(self)

```

### Attributes

The KMEANS object has two main attributes:

```

1 self.cluster_centers # Position of the cluster centers
  self.metrics # Different metrics to evaluate the model

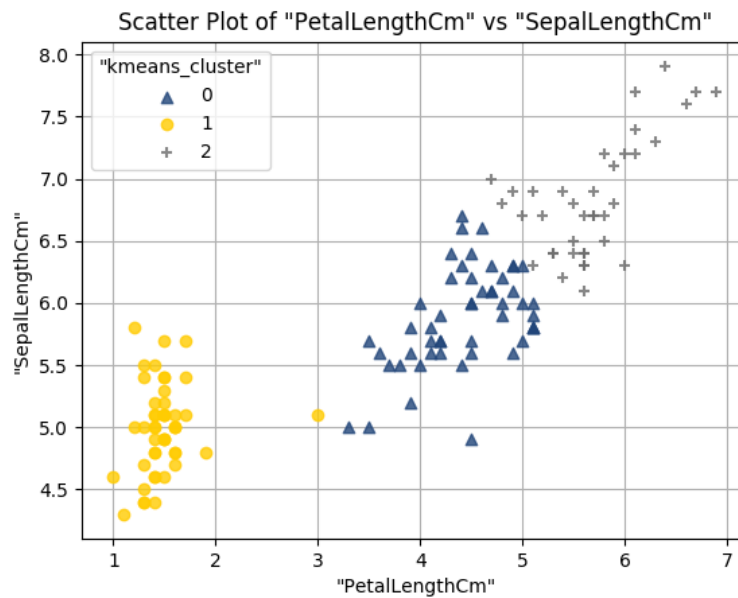
```

### Example

```

from vertica_ml_python.learn.cluster import KMeans
2
# We can build the model
4 model = KMeans("kmeans_iris", cur, n_cluster = 3)
  model.fit("iris", ["PetalLengthCm", "SepalLengthCm"])
6
# We can see the different metrics relative to the model
8 model.metrics
10
# Output
                                     value
12 Between-Cluster Sum of Squares      512.23072
  Total Sum of Squares                  566.03207
14 Total Within-Cluster Sum of Squares  53.801351
  Between-Cluster SS / Total SS         0.9049499969144859
16 converged                           True
18
# And Plot the model
  model.plot()

```



## 2.2 best\_k

**Library:** `vertica_ml_python.learn.cluster.model_selection`

```

1 def best_k(
2     X: list,
3     input_relation: str,
4     cursor,
5     n_cluster = (1, 100),
6     init = "kmeanspp",
7     max_iter: int = 50,
8     tol: float = 1e-4,
9     elbow_score_stop = 0.8)

```

Find the KMeans K by using a score threshold.

### Parameters

- **X:** *<list>*  
List of the predictor columns.
- **input\_relation:** *<str>*  
Relation used to train the model.
- **cursor:** *<object>*  
DB cursor.
- **n\_cluster:** *<tuple>*, optional  
Tuple representing the number of cluster to start with and to end with.

- **init:** *<str or list>*, optional  
The method used to find the initial cluster centers. The method can be in {random | kmeanspp}. It can be also a list with the initial cluster centers to use.
- **max\_iter:** *<int>*, optional  
The maximum number of iterations the algorithm performs.
- **tol:** *<float>*, optional  
Determines whether the algorithm has converged. The algorithm is considered converged after no center has moved more than a distance of 'tol' from the previous iteration.
- **elbow\_score\_stop:** *<float>*, optional  
Stop the Parameters Search when this Elbow score is reached.

### Example

```
1 from vertica_ml_python.learn.model_selection import best_k
3 best_k(["PetalLengthCm", "PetalWidthCm", "SepalLengthCm", "SepalWidthCm"], "
    iris", cur)
5 # Output
3
```

## 3 Questions

Turn on Jupyter with the 'jupyter notebook' command. Start the notebook exercise10.ipynb and answer to the following questions.

- **Question 1:** Find the best K which will segment the data with an elbow score greater than 0.8
- **Question 2:** Use this K to create a KMeans model.
- **Question 3:** Add the prediction in the Virtual Dataframe. Draw an histogram of the Species and the predictions. What do you notice ?