

K_in_K_means_Clustering

July 17, 2020

```
[1]: import pandas as pd
import numpy as np
from sklearn import datasets
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
import matplotlib.patches as mpatches
import sklearn.metrics as sm
import math
%matplotlib inline
```

```
[3]: iris = datasets.load_iris()
print(iris.data)
```

```
[[5.1 3.5 1.4 0.2]
 [4.9 3.  1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.5 0.2]
 [5.  3.6 1.4 0.2]
 [5.4 3.9 1.7 0.4]
 [4.6 3.4 1.4 0.3]
 [5.  3.4 1.5 0.2]
 [4.4 2.9 1.4 0.2]
 [4.9 3.1 1.5 0.1]
 [5.4 3.7 1.5 0.2]
 [4.8 3.4 1.6 0.2]
 [4.8 3.  1.4 0.1]
 [4.3 3.  1.1 0.1]
 [5.8 4.  1.2 0.2]
 [5.7 4.4 1.5 0.4]
 [5.4 3.9 1.3 0.4]
 [5.1 3.5 1.4 0.3]
 [5.7 3.8 1.7 0.3]
 [5.1 3.8 1.5 0.3]
 [5.4 3.4 1.7 0.2]
 [5.1 3.7 1.5 0.4]
 [4.6 3.6 1.  0.2]
 [5.1 3.3 1.7 0.5]
 [4.8 3.4 1.9 0.2]]
```

[5. 3. 1.6 0.2]
 [5. 3.4 1.6 0.4]
 [5.2 3.5 1.5 0.2]
 [5.2 3.4 1.4 0.2]
 [4.7 3.2 1.6 0.2]
 [4.8 3.1 1.6 0.2]
 [5.4 3.4 1.5 0.4]
 [5.2 4.1 1.5 0.1]
 [5.5 4.2 1.4 0.2]
 [4.9 3.1 1.5 0.2]
 [5. 3.2 1.2 0.2]
 [5.5 3.5 1.3 0.2]
 [4.9 3.6 1.4 0.1]
 [4.4 3. 1.3 0.2]
 [5.1 3.4 1.5 0.2]
 [5. 3.5 1.3 0.3]
 [4.5 2.3 1.3 0.3]
 [4.4 3.2 1.3 0.2]
 [5. 3.5 1.6 0.6]
 [5.1 3.8 1.9 0.4]
 [4.8 3. 1.4 0.3]
 [5.1 3.8 1.6 0.2]
 [4.6 3.2 1.4 0.2]
 [5.3 3.7 1.5 0.2]
 [5. 3.3 1.4 0.2]
 [7. 3.2 4.7 1.4]
 [6.4 3.2 4.5 1.5]
 [6.9 3.1 4.9 1.5]
 [5.5 2.3 4. 1.3]
 [6.5 2.8 4.6 1.5]
 [5.7 2.8 4.5 1.3]
 [6.3 3.3 4.7 1.6]
 [4.9 2.4 3.3 1.]
 [6.6 2.9 4.6 1.3]
 [5.2 2.7 3.9 1.4]
 [5. 2. 3.5 1.]
 [5.9 3. 4.2 1.5]
 [6. 2.2 4. 1.]
 [6.1 2.9 4.7 1.4]
 [5.6 2.9 3.6 1.3]
 [6.7 3.1 4.4 1.4]
 [5.6 3. 4.5 1.5]
 [5.8 2.7 4.1 1.]
 [6.2 2.2 4.5 1.5]
 [5.6 2.5 3.9 1.1]
 [5.9 3.2 4.8 1.8]
 [6.1 2.8 4. 1.3]
 [6.3 2.5 4.9 1.5]

[6.1 2.8 4.7 1.2]
[6.4 2.9 4.3 1.3]
[6.6 3. 4.4 1.4]
[6.8 2.8 4.8 1.4]
[6.7 3. 5. 1.7]
[6. 2.9 4.5 1.5]
[5.7 2.6 3.5 1.]
[5.5 2.4 3.8 1.1]
[5.5 2.4 3.7 1.]
[5.8 2.7 3.9 1.2]
[6. 2.7 5.1 1.6]
[5.4 3. 4.5 1.5]
[6. 3.4 4.5 1.6]
[6.7 3.1 4.7 1.5]
[6.3 2.3 4.4 1.3]
[5.6 3. 4.1 1.3]
[5.5 2.5 4. 1.3]
[5.5 2.6 4.4 1.2]
[6.1 3. 4.6 1.4]
[5.8 2.6 4. 1.2]
[5. 2.3 3.3 1.]
[5.6 2.7 4.2 1.3]
[5.7 3. 4.2 1.2]
[5.7 2.9 4.2 1.3]
[6.2 2.9 4.3 1.3]
[5.1 2.5 3. 1.1]
[5.7 2.8 4.1 1.3]
[6.3 3.3 6. 2.5]
[5.8 2.7 5.1 1.9]
[7.1 3. 5.9 2.1]
[6.3 2.9 5.6 1.8]
[6.5 3. 5.8 2.2]
[7.6 3. 6.6 2.1]
[4.9 2.5 4.5 1.7]
[7.3 2.9 6.3 1.8]
[6.7 2.5 5.8 1.8]
[7.2 3.6 6.1 2.5]
[6.5 3.2 5.1 2.]
[6.4 2.7 5.3 1.9]
[6.8 3. 5.5 2.1]
[5.7 2.5 5. 2.]
[5.8 2.8 5.1 2.4]
[6.4 3.2 5.3 2.3]
[6.5 3. 5.5 1.8]
[7.7 3.8 6.7 2.2]
[7.7 2.6 6.9 2.3]
[6. 2.2 5. 1.5]
[6.9 3.2 5.7 2.3]

```

[5.6 2.8 4.9 2. ]
[7.7 2.8 6.7 2. ]
[6.3 2.7 4.9 1.8]
[6.7 3.3 5.7 2.1]
[7.2 3.2 6.  1.8]
[6.2 2.8 4.8 1.8]
[6.1 3.  4.9 1.8]
[6.4 2.8 5.6 2.1]
[7.2 3.  5.8 1.6]
[7.4 2.8 6.1 1.9]
[7.9 3.8 6.4 2. ]
[6.4 2.8 5.6 2.2]
[6.3 2.8 5.1 1.5]
[6.1 2.6 5.6 1.4]
[7.7 3.  6.1 2.3]
[6.3 3.4 5.6 2.4]
[6.4 3.1 5.5 1.8]
[6.  3.  4.8 1.8]
[6.9 3.1 5.4 2.1]
[6.7 3.1 5.6 2.4]
[6.9 3.1 5.1 2.3]
[5.8 2.7 5.1 1.9]
[6.8 3.2 5.9 2.3]
[6.7 3.3 5.7 2.5]
[6.7 3.  5.2 2.3]
[6.3 2.5 5.  1.9]
[6.5 3.  5.2 2. ]
[6.2 3.4 5.4 2.3]
[5.9 3.  5.1 1.8]]

```

```
[3]: print(iris.target_names)
```

```
['setosa' 'versicolor' 'virginica']
```

```
[4]: print(iris.target)
```

```

[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2
 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 2 2]

```

```

[5]: x = pd.DataFrame(iris.data, columns=['Sepal Length', 'Sepal Width',
                                           'Petal Length', 'Petal Width'])
     y = pd.DataFrame(iris.target, columns=['Target'])

```

```
[6]: x.head()
```

[6]:	Sepal Length	Sepal Width	Petal Length	Petal Width
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2

```
[7]: y.head()
```

[7]:	Target
0	0
1	0
2	0
3	0
4	0

```
[8]: iris_k_mean_model = KMeans(n_clusters=3)
iris_k_mean_model.fit(x)
```

```
[8]: KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
           n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto',
           random_state=None, tol=0.0001, verbose=0)
```

```
[9]: print(iris_k_mean_model.labels_)
```

$$\begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 2 & 0 & 2 & 2 \\ 2 & 2 & 0 & 0 & 2 & 2 & 2 & 2 & 0 & 2 & 0 & 2 & 0 & 2 & 2 & 2 & 2 & 2 & 0 & 2 & 2 & 2 & 0 & 2 \\ 2 & 0 \end{bmatrix}$$

```
[11]: print(iris_k_mean_model.cluster_centers_)
```

```
[[5.9016129  2.7483871  4.39354839 1.43387097]
 [5.006       3.428       1.462       0.246       ]
 [6.85        3.07368421  5.74210526 2.07105263]]
```

```
[12]: predictedY = np.choose(iris_k_mean_model.labels_, [0, 1, 2]).astype(np.int64)
```

```
[13]: sm.accuracy_score(predictedY, y['Target'])
```

```
[13]: 0.24
```

0.1 Interpretation of Confusion Matrix

Correctly identified all 0 classes as 0's correctly classified 48 class 1's but miss-classified 2 class 1's as class 2 correctly classified 36 class 2's but miss-classified 14 class 2's as class 1

```
[14]: sm.confusion_matrix(predictedY, y['Target'])
```

```
[14]: array([[ 0, 48, 14],  
          [50,  0,  0],  
          [ 0,  2, 36]], dtype=int64)
```

```
[15]: x.shape
```

```
[15]: (150, 4)
```

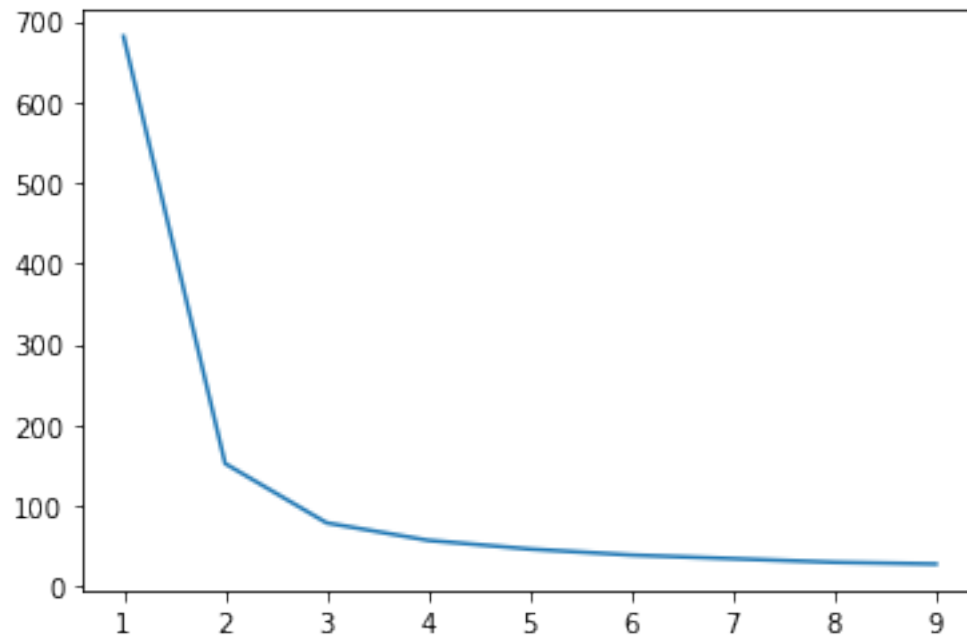
```
[16]: dist_points_from_cluster_center = []  
      K = range(1,10)  
      for no_of_clusters in K:  
          k_model = KMeans(n_clusters=no_of_clusters)  
          k_model.fit(x)  
          dist_points_from_cluster_center.append(k_model.inertia_)
```

```
[17]: dist_points_from_cluster_center
```

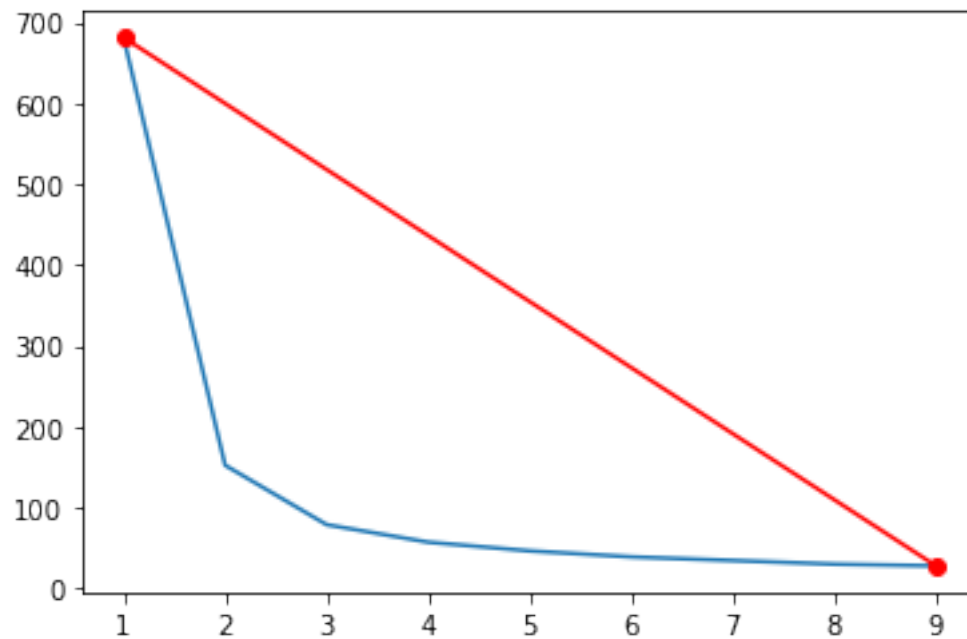
```
[17]: [681.3706,  
      152.34795176035792,  
      78.85144142614601,  
      57.228473214285714,  
      46.44618205128205,  
      39.03998724608725,  
      34.57303082786779,  
      30.063110617452725,  
      27.940751666462198]
```

```
[18]: plt.plot(K, dist_points_from_cluster_center)
```

```
[18]: [<matplotlib.lines.Line2D at 0x1811866d7f0>]
```



```
[19]: plt.plot(K, dist_points_from_cluster_center)
plt.plot([K[0], K[8]], [dist_points_from_cluster_center[0],
                        dist_points_from_cluster_center[8]], 'ro-')
plt.show()
```



```
[20]: x = [K[0], K[8]]
y = [dist_points_from_cluster_center[0], dist_points_from_cluster_center[8]]

# Calculate the coefficients. This line answers the initial question.
coefficients = np.polyfit(x, y, 1)

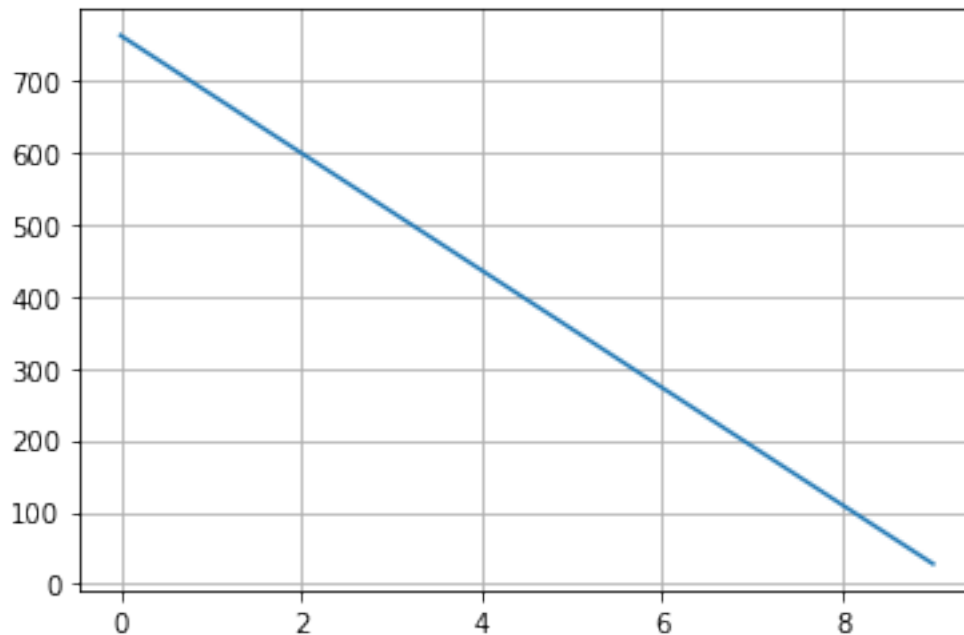
# Print the findings
print('a =', coefficients[0])
print('b =', coefficients[1])

# Let's compute the values of the line...
polynomial = np.poly1d(coefficients)
x_axis = np.linspace(0,9,100)
y_axis = polynomial(x_axis)

# ...and plot the points and the line
plt.plot(x_axis, y_axis)
plt.grid('on')
plt.show()
```

a = -81.67873104169225

b = 763.0493310416921



```
[21]: # Function to find distance
# https://www.geeksforgeeks.org/perpendicular-distance-
# between-a-point-and-a-line-in-2-d/
```



```
def calc_distance(x1, y1, a, b, c):
    d = abs((a * x1 + b * y1 + c)) / (math.sqrt(a * a + b * b))
    return d
```

```
[22]: # (y1 - y2)x + (x2 - x1)y + (x1y2 - x2y1) = 0
# https://bobobobo.wordpress.com/2008/01/07/solving-linear-equations-ax-by-c-0/
a = dist_points_from_cluster_center[0] - dist_points_from_cluster_center[8]
b = K[8] - K[0]
c1 = K[0] * dist_points_from_cluster_center[8]
c2 = K[8] * dist_points_from_cluster_center[0]
c = c1 - c2
```

```
[23]: dist_points_from_cluster_center
```

```
[23]: [681.3706,
152.34795176035792,
78.85144142614601,
57.228473214285714,
46.44618205128205,
39.03998724608725,
34.57303082786779,
30.063110617452725,
27.940751666462198]
```

```
[24]: distance_of_points_from_line = []
for k in range(9):
    distance_of_points_from_line.append(
        calc_distance(K[k], dist_points_from_cluster_center[k], a, b, c))
```

```
[25]: distance_of_points_from_line
```

```
[25]: [0.0,
5.476461109863801,
5.376292957857814,
4.64107998878274,
3.7731535883179514,
2.863896436062676,
1.918656622369384,
0.9739427788034284,
0.0]
```

```
[26]: plt.plot(K, distance_of_points_from_line)
```

```
[26]: [<matplotlib.lines.Line2D at 0x181187d8fd0>]
```

