

# Implementation of SciKit Fuzzy

CSE 401: Artificial Intelligence

Utkarsh Gupta

A2305217557

7CSE 8Y

October 21st, 2020

## 1 Fuzzy c-means clustering

Fuzzy logic principles can be used to cluster multidimensional data, assigning each point a membership in each cluster center from 0 to 100 percent. This can be very powerful compared to traditional hard-thresholded clustering where every point is assigned a crisp, exact label.

Fuzzy c-means clustering is accomplished via `skfuzzy.cmeans`, and the output from this function can be repurposed to classify new data according to the calculated clusters (also known as prediction) via `skfuzzy.cmeans_predict`.

### 1.1 Data generation and setup

In this example we will first undertake necessary imports, then define some test data to work with.

```
[1]: from __future__ import division, print_function
import numpy as np
import matplotlib.pyplot as plt
import skfuzzy as fuzz

colors = ['b', 'orange', 'g', 'r', 'c', 'm', 'y', 'k', 'Brown', 'ForestGreen']

# Define three cluster centers
centers = [[4, 2],
           [1, 7],
           [5, 6]]

# Define three cluster sigmas in x and y, respectively
sigmas = [[0.8, 0.3],
           [0.3, 0.5],
           [1.1, 0.7]]

# Generate test data
np.random.seed(42) # Set seed for reproducibility
xpts = np.zeros(1)
ypts = np.zeros(1)
```

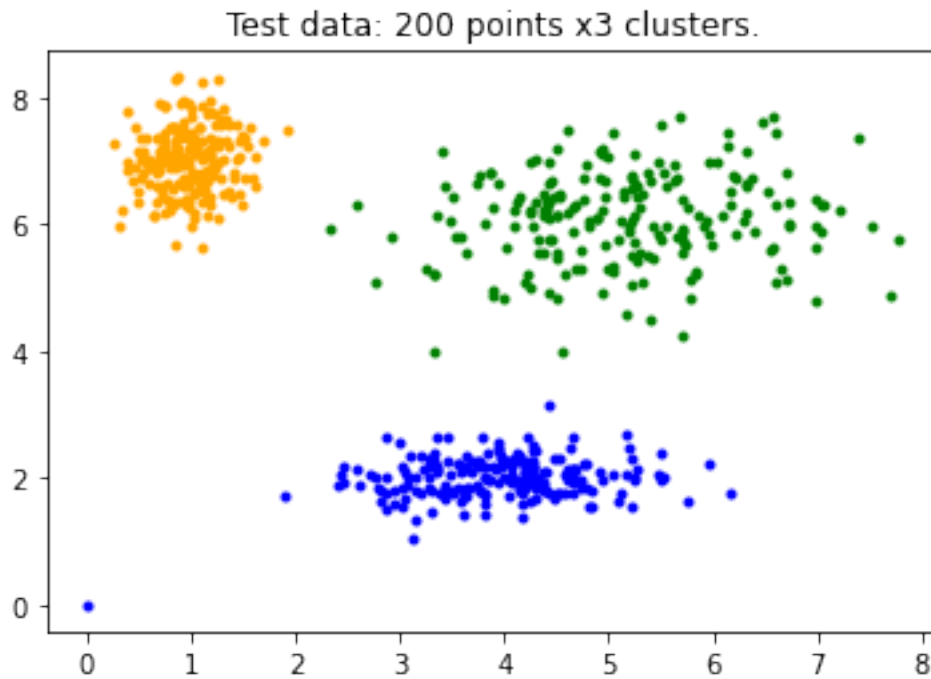
```

labels = np.zeros(1)
for i, ((xmu, ymu), (xsigma, ysigma)) in enumerate(zip(centers, sigmas)):
    xpts = np.hstack((xpts, np.random.standard_normal(200) * xsigma + xmu))
    ypts = np.hstack((ypts, np.random.standard_normal(200) * ysigma + ymu))
    labels = np.hstack((labels, np.ones(200) * i))

# Visualize the test data
fig0, ax0 = plt.subplots()
for label in range(3):
    ax0.plot(xpts[labels == label], ypts[labels == label], '.',
            color=colors[label])
ax0.set_title('Test data: 200 points x3 clusters.')

```

[1]: Text(0.5, 1.0, 'Test data: 200 points x3 clusters.')



## 1.2 Clustering

Above is our test data. We see three distinct blobs. However, what would happen if we didn't know how many clusters we should expect? Perhaps if the data were not so clearly clustered?

Let's try clustering our data several times, with between 2 and 9 clusters.

```
[2]: # Set up the loop and plot
fig1, axes1 = plt.subplots(3, 3, figsize=(8, 8))
alldata = np.vstack((xpts, ypts))
fpcs = []

for ncenters, ax in enumerate(axes1.reshape(-1), 2):
    cntr, u, u0, d, jm, p, fpc = fuzz.cluster.cmeans(
        alldata, ncenters, 2, error=0.005, maxiter=1000, init=None)

    # Store fpc values for later
    fpcs.append(fpc)

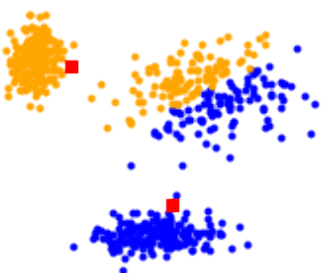
    # Plot assigned clusters, for each data point in training set
    cluster_membership = np.argmax(u, axis=0)
    for j in range(ncenters):
        ax.plot(xpts[cluster_membership == j],
                ypts[cluster_membership == j], '.', color=colors[j])

    # Mark the center of each fuzzy cluster
    for pt in cntr:
        ax.plot(pt[0], pt[1], 'rs')

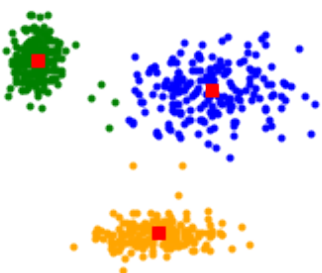
    ax.set_title('Centers = {0}; FPC = {1:.2f}'.format(ncenters, fpc))
    ax.axis('off')

fig1.tight_layout()
```

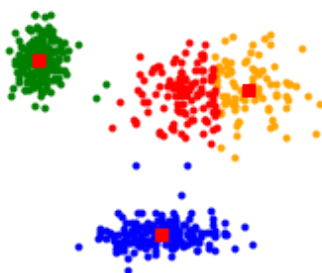
Centers = 2; FPC = 0.79



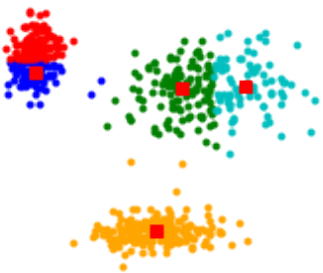
Centers = 3; FPC = 0.88



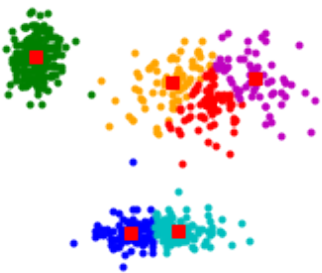
Centers = 4; FPC = 0.81



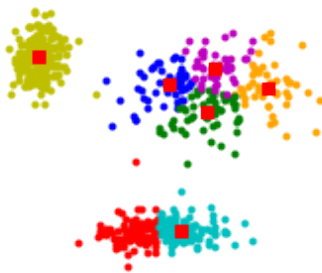
Centers = 5; FPC = 0.72



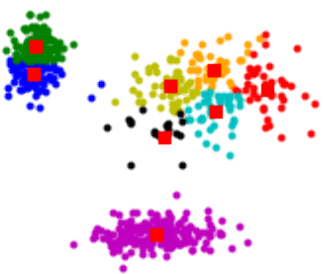
Centers = 6; FPC = 0.71



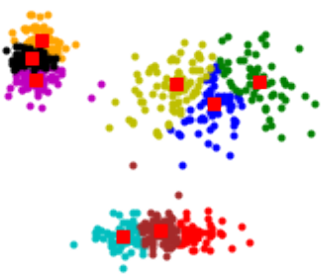
Centers = 7; FPC = 0.69



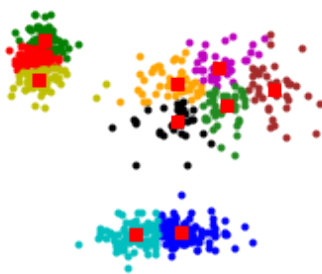
Centers = 8; FPC = 0.64



Centers = 9; FPC = 0.58



Centers = 10; FPC = 0.58

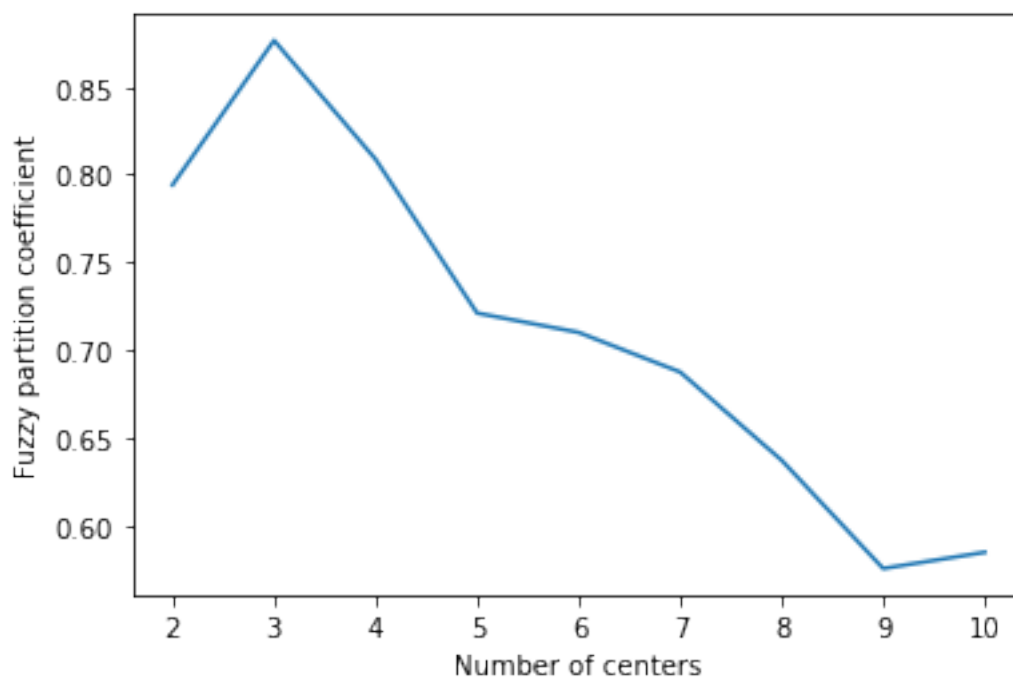


### 1.3 The fuzzy partition coefficient (FPC)

The FPC is defined on the range from 0 to 1, with 1 being best. It is a metric which tells us how cleanly our data is described by a certain model. Next we will cluster our set of data - which we know has three clusters - several times, with between 2 and 9 clusters. We will then show the results of the clustering, and plot the fuzzy partition coefficient. When the FPC is maximized, our data is described best.

```
[3]: fig2, ax2 = plt.subplots()
      ax2.plot(np.r_[2:11], fpcs)
      ax2.set_xlabel("Number of centers")
      ax2.set_ylabel("Fuzzy partition coefficient")
```

```
[3]: Text(0, 0.5, 'Fuzzy partition coefficient')
```



As we can see, the ideal number of centers is 3. This isn't news for our contrived example, but having the FPC available can be very useful when the structure of your data is unclear.

Note that we started with two centers, not one; clustering a dataset with only one cluster center is the trivial solution and will by definition return  $FPC == 1$ .