# Instituto Politécnico Nacional

# Escuela Superior de Cómputo

# Evolutionary Computing

# Laboratory session #08: Self Organizing Map (SOM)

## Student: Vargas Romero Erick Efrain

## Professor: Rosas Trigueros Jorge Luis

## Practice completion date: November 10, 2021

## Report delivery date: November 10, 2021

## Theoretical framework

## 1.1 Self Organizing Map

One of the apparent thrends today in all disciplines of engineering sciences has the goal of developing computationally intelligent systems that *"Enable or facilite intelligent behaviour in complex and changing enviroments"*. The concepts, algoriths and models suggested by this new paradign of knoledge-based processing should, have the capability of self-organization, and exhibit an ability to adapt to new situations. One of the methods belonging to the field of computational intelligence which has proven to be a very powerful tool for data analysis is the neurobiologically inspired self-organizing map (SOM). The most common model of SOMs, also known as the Kohonen network, is the topology preserving map proposed by the Finnish researcher Teuvo Kohonen in 1982. The method is a special class of artificial neural networks and is used extensively as a clustering and visualization tool in exploratory data analysis [1].

Among the architectures and algorithms suggested for artificial neural networks, the SOM is trained using unsupervised learning scheme. That means, unlike supervised networks, learning SOM does not rely on predefined target outcomes that would guide the process. Thus, a form of learning by observation, rather than learning by examples takes place to discover the underlying hidden patterns in the data set. In order to learn without a teacher, SOMs apply a competitive learning rule where the output nodes compete among themselves for the opportunity to represent distinct patterns within the input space. During the learning, the feedforward nature of SOMs allows an information flow in only one direction, without looping or cycling, from the input nodes to the output nodes. Note that every node in the input layer is linked (with weights) to every node in the output layer, which makes a SOM a completely connected network. [1]. The formation of a SOM involves three characteristic process: *Competition*, *Cooperation* and *Adaptation*.

### 1.1.1 Competition

The output nodes (neurons) in a self-organizing map compete with each other to best represent the particular input sample. The sucess of representation is measured using a discriminant function, where an input vector is compared with the weight vector of each output node. The particular node

with its connection weights most similar to the input sample is declared winner of the competition, There are a number of different functions to determine the winner, for example: *Best Matching Unit* on the map. The most used one is the *Euclidean Distance*

### 1.1.2  Cooperation

Similar to *"Neurons dealing with closely related pieces of information are close together so that they can interact via short synaptic connections"*, SOM is a topographic organization in which nearby locations in the output space represent inputs with similar properties. This is possible in the presence of neighborhood information. The winning node determines the spatial location of cooperating nodes. These nodes, sharing common features, activate each other to learn something from the same input.

### 1.1.3  Adaptation

The weight vectors of the winner and its neighboring units in the map are adjusted in favor of higher values of their discriminant functions. Through this learning process the relevant nodes become more similar to the input sample. Thus nodes which have a strong response to a particular piece of input data will have an increased chance of responding to similar input data in the future.

# *2*

## Material and equipment

The necessary material for this practice is:

- A computer with the latest *Python* stable version installed

- A text editor

Or is possible to use the google site `https://colab.research.google.com/` that allows us to use a virtual machine with an *Ubuntu* operative system with *Python* installed.

*3*

**Practice development**

## 3.1   Self Organizing Map

To develop this practice we used the Google platform called *Colab* as this platform uses a virtual machine with linux (specifically Ubuntu) we can install some packets and of course we can verify if we have already *Python* installed. To check it we must use the command:

```
1    python --version
```

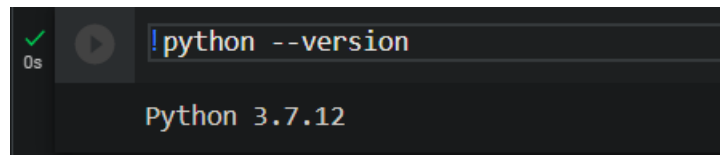If we run this command in our linux terminal using *Colab* we can see the next as the result:

Figure 3.1: Verifying python version

**Instructions**: *Create a SOM using at least 10 indicators about at least 40 countries from the databases in the World Bank. The BMU for each country should be indicated in the final representation of the SOM. The information was obtained from $https://databank.worldbank.org/source/world-development-indicato$*

First of all we got the information from the Worl Bank, from my point of view this step has been the most complicated because a lot of indicators does not exist for some countries (see img.4.1) so we must try to find trying with other period of time or changing the country.
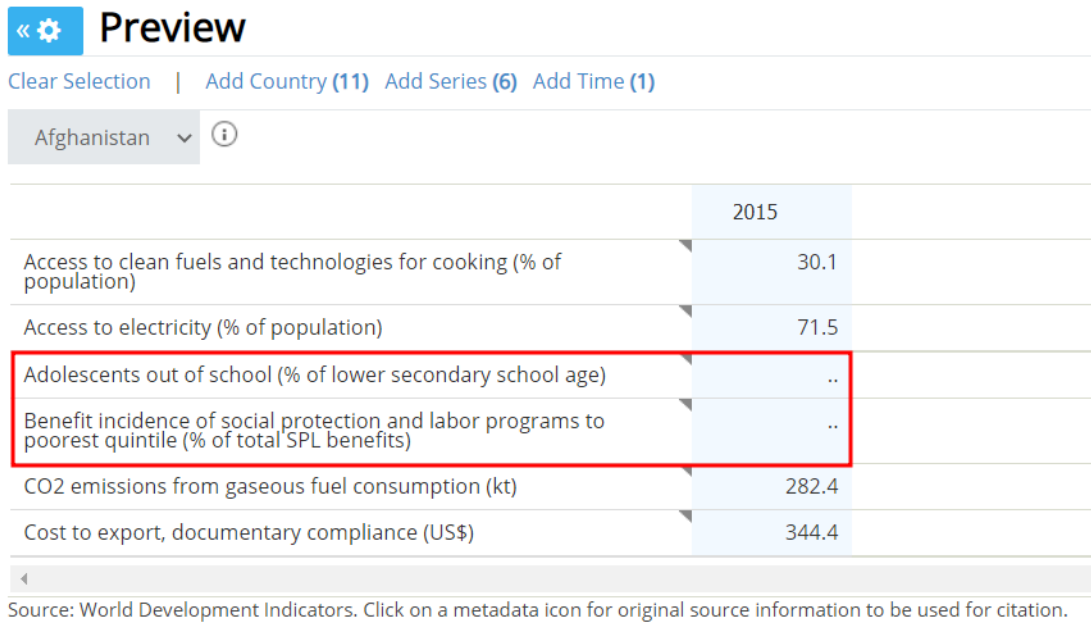
Figure 3.2: Unknown data

Once that we are sure that we have valid data we can get an Excel, CSV or Tabbed TXT files the objective of this is to manipule the file and give it as input of our Self Organizing Map.
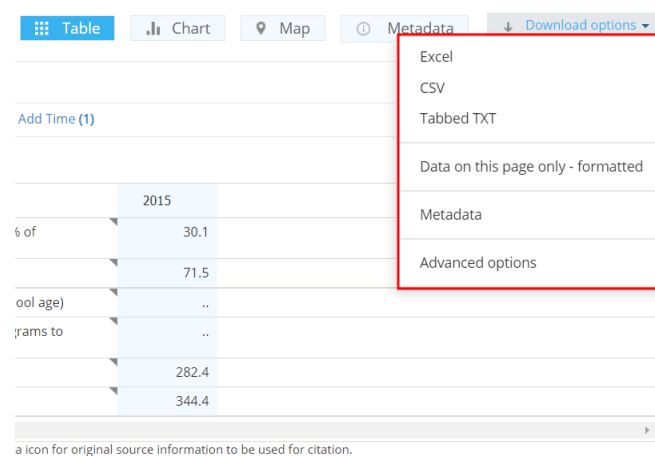


Figure 3.3: Download data

In our case we choose a CSV file and to get the data needed was mandatory to normalize the data. Normalize the data was a mandatory step because if you do not do that you will get some strange data in the output of the SOM. First of all to manipulate the CSV we used the csv library, it has a lot of different and useful tools to manage CSV files, we get all the data in the CSV file and we store the countries indicators in a multidimensional array (as such as categories) and we get the countries tag too, this in order to make easier to see the countries positions in the map.

```python
import csv

indicators = [[], [], [], [], [], [], [], [], [], []]
countryTrainingVector = []
setCountryTag = set()
countryTag = []

#We will use csv library to get the data from de wold bank
with open('data.csv', mode='r') as csvFile:
    #This object allow us "move" in the CSV (row by row)
    csvReader = csv.DictReader(csvFile)
    lineCount = 0
    for row in csvReader:
        #If we are in the header we can't get data so we ignore it
        if lineCount == 0:
            lineCount += 1
        #We get the value of the current indicator
        currValue = row["2015 [YR2015]"]
        #We get the country code of the current country
        currCountryCode = row["Country Code"]
        #If this one is zero it means that we are in the headers row
        if len(currValue) > 0:
            #We add the current indicator to our indicators matrix
            #First we're "sorting" the indicators by category
            #As we known we have 10 indicators that's why we're using modulo 10
            #And we known that from line 1 we have information about the indicators
            #That's why we must substract 1 from the line count
            #Finally use row["column"] extracts an string that's why we must cast to
    a float
            indicators[(lineCount - 1) % 10].append(float(row["2015 [YR2015]"]))
            #As every row contains the country code, in order to avoid repeated tags
            #we're using a set just to ask if we already visited or added a tag.
            if not currCountryCode in setCountryTag:
                setCountryTag.add(currCountryCode)
                countryTag.append(currCountryCode)
        lineCount += 1
```

The next step is to normalize data, this was done geting the max value of each type of indicator and then divide every indicator value by that maximum.

```python
#We normalize the data getting the maximum of every indicators category
#and dividing by that maximum every indicators element (we'll get a value between 0
    and 1)
for i in range(0, 10):
    maxValue = max(indicators[i])
    for j in range(0, 40):
        indicators[i][j] /= maxValue
```

Then we create our training vector using the information that we normalized but now instead of have the information "sorted" by indicator we assign the indicators that belongs to each contry.

```
#Now ge build the training vector we have 40 elements per category
#The j-th indicator belongs to the i-th country
for i in range(0, 40):
    row = []
    for j in range(0, 10):
        row.append(indicators[j][i])
    #To add elements to our training vector we must add an "id"
    #So the structure is [indicatorsVector, [id]]
    countryTrainingVector.append([row, [i]])
```

Now we can create our SOM and start training it using the vector that we created in the previous step and as you can see we used a SOM of 20 of width and 20 of height this because we need enough space to "move" every node in the SOM. Additionally, we trained the SOM 500 times to make it more accurate.

```
print("Initialization...")
a = SOM(20, 20, 10, 1, False, 0.03)
print("Training...")
a.train(500, countryTrainingVector)
```

Once that we finish with the training we can get the BMU for every country and try to represent them in a graphical way.

```
positions = {}

#We get the BMU (Best Matching Unit)
for i in range(0, 40):
    #We get every country indicators
    currCountry = countryTrainingVector[i][0]
    #We get the BMU it returns the value, and position
    value, x, y = a.predict(currCountry, True)
    #Can exist repeated element's we use the last one
    positions[(x, y)] = countryTag[i]

#We print the number of countries that we get
print(f"He wave {len(positions)} entries.")
#We print a matrix with the countries
printCountries(positions)
```

Is important to mention that we used the file provided by the professor as a base to develop this practice. Finally we add the whole code in order to make easier to understand all the process that happens during the execution of this program.

```
from random import *
from math import *
import numpy as np
import csv

class Node:
    def __init__(self, FV_size=10, PV_size=10, Y=0, X=0):
        self.FV_size = FV_size
        self.PV_size = PV_size
        self.FV = [0.0] * FV_size  # Feature Vector
        self.PV = [0.0] * PV_size  # Prediction Vector
        self.X = X  # X location
```

```python
13            self.Y = Y  # Y location
14
15            for i in range(FV_size):
16                self.FV[i] = random()  # Assign a random number from 0 to 1
17
18            for i in range(PV_size):
19                self.PV[i] = random()  # Assign a random number from 0 to 1
20
21
22  class SOM:
23
24      # Let radius=False if you want to autocalculate the radis
25      def __init__(
26          self,
27          height=10,
28          width=10,
29          FV_size=10,
30          PV_size=10,
31          radius=False,
32          learning_rate=0.005,
33      ):
34          self.height = height
35          self.width = width
36          self.radius = radius if radius else (height + width) / 2
37          self.total = height * width
38          self.learning_rate = learning_rate
39          self.nodes = [0] * (self.total)
40          self.FV_size = FV_size
41          self.PV_size = PV_size
42          for i in range(self.height):
43              for j in range(self.width):
44                  self.nodes[(i) * (self.width) + j] = Node(FV_size, PV_size, i, j)
45
46      # Train_vector format: [ [FV[0], PV[0]],
47      #                        [FV[1], PV[1]], so on..
48
49      def train(self, iterations=1000, train_vector=[[[0.0], [0.0]]]):
50          time_constant = iterations / log(self.radius)
51          radius_decaying = 0.0
52          learning_rate_decaying = 0.0
53          influence = 0.0
54          stack = []  # Stack for storing best matching unit's index and updated FV
     and PV
55          temp_FV = [0.0] * self.FV_size
56          temp_PV = [0.0] * self.PV_size
57          for i in range(1, iterations + 1):
58              # print "Iteration number:",i
59              radius_decaying = self.radius * exp(-1.0 * i / time_constant)
60              learning_rate_decaying = self.learning_rate * exp(-1.0 * i /
     time_constant)
61              print(i, end=", ")
62              if i % 50 == 0:
63                  print("")
64
65              for j in range(len(train_vector)):
66                  input_FV = train_vector[j][0]
67                  input_PV = train_vector[j][1]
68                  best = self.best_match(input_FV)
```

```python
69                      stack = []
70                      for k in range(self.total):
71                          dist = self.distance(self.nodes[best], self.nodes[k])
72                          if dist < radius_decaying:
73                              temp_FV = [0.0] * self.FV_size
74                              temp_PV = [0.0] * self.PV_size
75                              influence = exp(
76                                  (-1.0 * (dist ** 2)) / (2 * radius_decaying * i)
77                              )
78
79                              for l in range(self.FV_size):
80                                  # Learning
81                                  temp_FV[l] = self.nodes[k].FV[
82                                      l
83                                  ] + influence * learning_rate_decaying * (
84                                      input_FV[l] - self.nodes[k].FV[l]
85                                  )
86
87                              for l in range(self.PV_size):
88                                  # Learning
89                                  temp_PV[l] = self.nodes[k].PV[
90                                      l
91                                  ] + influence * learning_rate_decaying * (
92                                      input_PV[l] - self.nodes[k].PV[l]
93                                  )
94
95                              # Push the unit onto stack to update in next interval
96                              stack[0:0] = [[[k], temp_FV, temp_PV]]
97
98                      for l in range(len(stack)):
99
100                         self.nodes[stack[l][0][0]].FV[:] = stack[l][1][:]
101                         self.nodes[stack[l][0][0]].PV[:] = stack[l][2][:]
102
103     # Returns prediction vector
104     def predict(self, FV=[0.0], get_ij=False):
105         best = self.best_match(FV)
106         if get_ij:
107             return self.nodes[best].PV, self.nodes[best].X, self.nodes[best].Y
108         return self.nodes[best].PV
109
110     # Returns best matching unit's index
111     def best_match(self, target_FV=[0.0]):
112
113         minimum = sqrt(self.FV_size)  # Minimum distance
114         minimum_index = 1  # Minimum distance unit
115         temp = 0.0
116         for i in range(self.total):
117             temp = 0.0
118             temp = self.FV_distance(self.nodes[i].FV, target_FV)
119             if temp < minimum:
120                 minimum = temp
121                 minimum_index = i
122
123         return minimum_index
124
125     def FV_distance(self, FV_1=[0.0], FV_2=[0.0]):
126         temp = 0.0
```

```python
127            for j in range(self.FV_size):
128                temp = temp + (FV_1[j] - FV_2[j]) ** 2
129
130            temp = sqrt(temp)
131            return temp
132
133    def distance(self, node1, node2):
134        return sqrt((node1.X - node2.X) ** 2 + (node1.Y - node2.Y) ** 2)
135
136 #Receives the countries info (x, y, and country tag)
137 def printCountries(countries):
138     result = []
139     #Creating a matrix of 20 per 20 elements
140     for i in range(0, 20):
141         row = []
142         for j in range(0, 20):
143             row.append(" X ")
144         result.append(row)
145
146     #In position x, y we replace by the country tag
147     for x, y in countries.keys():
148         result[x][y] = countries[(x, y)]
149
150     #Printing the result
151     for i in range(0, 20):
152         print(f"{result[i]}")
153
154
155 #The training vector contains the next information:
156 #Access to electricity (% of population)
157 #Surface area (sq. km)
158 #Scientific and technical journal articles
159 #Rural population
160 #Population, total
161 #Population, male
162 #Population, female
163 #Military expenditure (% of GDP)
164 #Imports of goods and services (% of GDP)
165 #Armed forces personnel, total
166
167 import csv
168
169 indicators = [[], [], [], [], [], [], [], [], [], []]
170 countryTrainingVector = []
171 setCountryTag = set()
172 countryTag = []
173
174 #We will use csv library to get the data from de wold bank
175 with open('data.csv', mode='r') as csvFile:
176     #This object allow us "move" in the CSV (row by row)
177     csvReader = csv.DictReader(csvFile)
178     lineCount = 0
179     for row in csvReader:
180         #If we are in the header we can't get data so we ignore it
181         if lineCount == 0:
182             lineCount += 1
183         #We get the value of the current indicator
184         currValue = row["2015 [YR2015]"]
```

```
185          #We get the country code of the current country
186          currCountryCode = row["Country Code"]
187          #If this one is zero it means that we are in the headers row
188          if len(currValue) > 0:
189              #We add the current indicator to our indicators matrix
190              #First we're "sorting" the indicators by category
191              #As we known we have 10 indicators that's why we're using modulo 10
192              #And we known that from line 1 we have information about the indicators
193              #That's why we must substract 1 from the line count
194              #Finally use row["column"] extracts an string that's why we must cast to
        a float
195              indicators[(lineCount - 1) % 10].append(float(row["2015 [YR2015]"]))
196              #As every row contains the country code, in order to avoid repeated tags
197              #we're using a set just to ask if we already visited or added a tag.
198              if not currCountryCode in setCountryTag:
199                  setCountryTag.add(currCountryCode)
200                  countryTag.append(currCountryCode)
201          lineCount += 1
202
203 #We normalize the data getting the maximum of every indicators category
204 #and dividing by that maximum every indicators element (we'll get a value between 0
        and 1)
205 for i in range(0, 10):
206     maxValue = max(indicators[i])
207     for j in range(0, 40):
208         indicators[i][j] /= maxValue
209
210 #Now ge build the training vector we have 40 elements per category
211 #The j-th indicator belongs to the i-th country
212 for i in range(0, 40):
213     row = []
214     for j in range(0, 10):
215         row.append(indicators[j][i])
216     #To add elements to our training vector we must add an "id"
217     #So the structure is [indicatorsVector, [id]]
218     countryTrainingVector.append([row, [i]])
219
220 print("Initialization...")
221 a = SOM(20, 20, 10, 1, False, 0.03)
222
223 print("Training...")
224 a.train(500, countryTrainingVector)
225
226 positions = {}
227
228 #We get the BMU (Best Matching Unit)
229 for i in range(0, 40):
230     #We get every country indicators
231     currCountry = countryTrainingVector[i][0]
232     #We get the BMU it returns the value, and position
233     value, x, y = a.predict(currCountry, True)
234     #Can exist repeated element's we use the last one
235     positions[(x, y)] = countryTag[i]
236
237 #We print the number of countries that we get
238 print(f"He wave {len(positions)} entries.")
239 #We print a matrix with the countries
240 printCountries(positions)
```

# 4

# Screens, graphs and diagrams

The result of this practice is a 2-Dimensional matrix that contains information about the positions of the BMU of every country (see img. **??**). Something important to mention are the indicators used in this practice, they were:

1. Access to electricity (% of population)

2. Surface area (sq. km)

3. Scientific and technical journal articles

4. Rural population

5. Population, total

6. Population, male

7. Population, female

8. Military expenditure (% of GDP)

9. Imports of goods and services (% of GDP)

10. Armed forces personnel, total

And bases in that information and based in the result of this test, you can see that are many countries at the same let's say "level" for example China, Russia and the United States we known that they are countries with a high livel of development, and we can see that many different african countries are in the same area, like Madagascar, Uganda, Ghana, etc. Maybe use different indicators can help us to determine more accurate this result because we do not have a clear relation between these indicators.

```
He wave 40 entries.
['CHN', 'X ', 'X ', 'X ', 'X ', 'X ', 'RUS', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'ETH', 'X ', 'MDG', 'X ', 'SLE']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'UGA', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'NGA', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'CAN', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'ZWE']
[' X ', 'X ', 'X ', 'USA', 'X ', 'X ', 'X ', 'BRA', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'CMR', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'ARG', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'SEN', 'X ']
[' X ', 'X ', 'X ', 'X ', 'JPN', 'X ', 'X ', 'MEX', 'X ', 'X ', 'KAZ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
['EGY', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'NPL', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'NZL', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'ITA', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'GHA']
[' X ', 'X ', 'TUR', 'X ', 'X ', 'X ', 'X ', 'X ', 'ESP', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'SWE', 'X ', 'X ', 'X ', 'CHE', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'GBR', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
['PAK', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'NLD', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'THA', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'COL', 'X ', 'X ', 'X ', 'X ', 'CHL', 'X ', 'X ', 'PRT', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
[' X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ', 'X ']
['ISR', 'X ', 'X ', 'X ', 'CUB', 'X ', 'ECU', 'X ', 'X ', 'X ', 'BOL', 'X ', 'X ', 'POL', 'X ', 'X ', 'VNM', 'X ', 'X ', 'SGP']
```

Figure 4.1: Countries BMU positions

*5*

## Conclusions

This practice has been interesting but required a lot of preprocesing and prepare a lot of data to feed our SOM in order to get some results. The behaviour of a SOM is something amazing because as we mentioned in the *Theoretical framework* use the behaviour of living beings for example the competition between the items inside the SOM, the Cooperation between them in order to get a better result and how they addapt to new knoledge. Something additional that we learn in this practice is to manipulate CSV files using Pyhton, at least for me Python is a new tool but very powerful, additionally Python has a lot of libraries that can help us to make easier to process information.

Talking about the result of the practice I really believe that choose more carefully our indicators can help us to get better reslts, I mention this because as you can see my indicators looks like does not have a lot of relation, for example choose only things related about economy, population, etc. I think can give us more information that choose indicators about different topics, but maybe to manage that case we must use something different thatn a SOM or probably modify our base algorithm.

# Bibliography

[1] U. Asan and C. Ercan *An Introduction to Self-Organizing Maps*. Accessed on November 10th, 2021. Available online: `https://www.researchgate.net/publication/263084866_An_Introduction_to_Self-Organizing_Maps/link/0f317539c1430454cf000000/download`