



INSTITUTO POLITÉCNICO NACIONAL  
UNIDAD PROFESIONAL INTERDISCIPLINARIA EN  
INGENIERÍA Y TECNOLOGÍAS AVANZADAS



## PRÁCTICA 1

---

# Inteligencia Artificial

---

**Alumnos:**

Domínguez González Luis Antonio

Munguía Gallegos Mauricio

**Profesor:**

Álvaro Anzueto Ríos

**Grupo:**

3BM8

Fecha de entrega: 29 de Enero del 2020

## Marco teórico

Sistema adaptativo neurodifuso para segmentación automática de imágenes y detección de bordes.

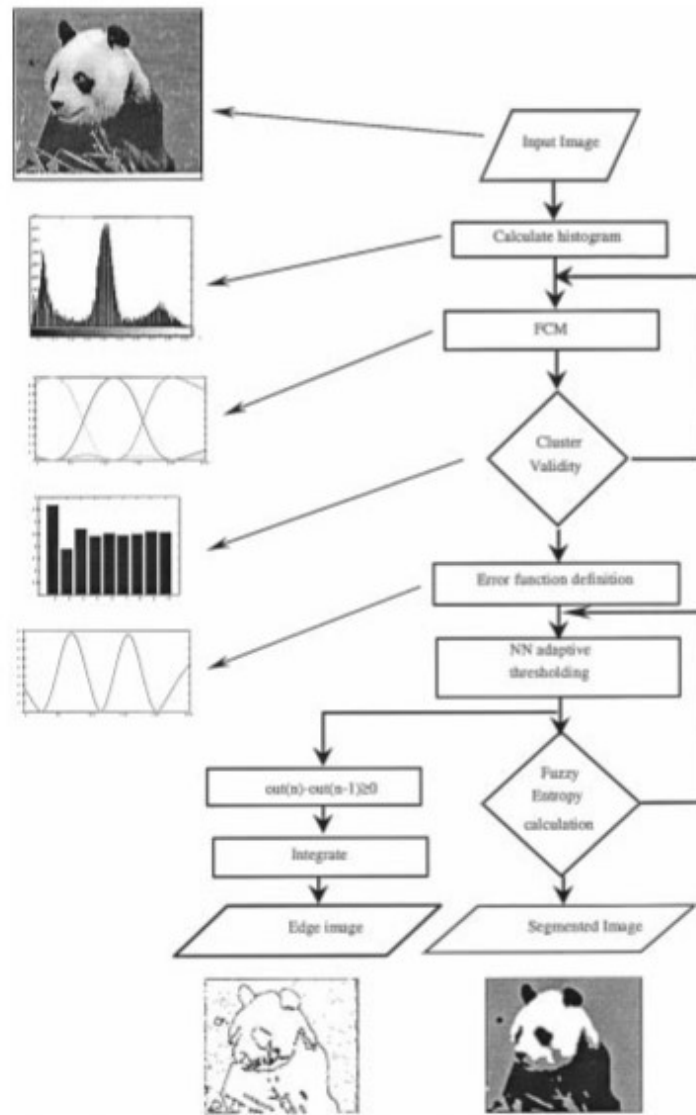


Figura 1 proceso de un sistema adaptativo neurodifuso para segmentación automática de imágenes y detección de bordes

$$PE = -\frac{1}{n \ln(\frac{1}{C})} \sum_{k=1}^n \sum_{i=1}^C [\mu_{ik} \ln(\mu_{ik})]$$

Se ocupa la ecuación como validity cluster de entropía después de haber pasado por todo el proceso que se muestra en la figura 1.

## Algoritmo K-means

“Clustering algorithms” pueden ser aplicados para resolver problemas de segmentación. Consisten en elegir un píxel o región inicial que pertenezca a un objeto de interés, seguido de un proceso interactivo de análisis de vecindades, que decide si cada píxel vecino pertenece o no al mismo objeto. Se utiliza k-means para resolver la detección de masa en mamografías utilizando la información de textura obtenida de los descriptores de Haralick. El algoritmo K-means es uno de los algoritmos de aprendizaje no supervisados más simples. El método sigue los pasos habituales para satisfacer el objetivo principal: agrupar todos los objetos de imagen en K grupos distintos. Primero, se definen K centroides, uno para cada grupo, siendo su posición inicial muy importante para el resultado. Después de eso, se determina una región de propiedad para cada centroide, que agrupa un conjunto de objetos similares. Se inicia la etapa interactiva del algoritmo, en la que se recalcula el centroide de cada grupo para minimizar la función objetivo. Esta función, para K-means, es el método de mínimos cuadrados que se calcula con:

$$J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

Dónde  $\|x_i^{(j)} - c_j\|^2$  es la distancia métrica desde cualquier punto  $x_i^{(j)}$  en el grupo centroide  $c_j$ . Así, la J (función objetiva) representa la medida de similitud de los objetos contenidos en sus respectivos grupos [2].

## Silhouette Coefficient

Permite que  $D_M = \{\bar{D}_1, \dots, \bar{D}_k\}$  describa un resultado de agrupamiento, es decir, es una partición exhaustiva del conjunto de documentos  $D^*$ . La distancia de un documento  $d \in D^*$  a un agrupamiento  $\bar{D}_i \in D_M$  está dada como

$$dist(d, \bar{D}_i) = \frac{\sum_{p \in \bar{D}_i} dist(d, p)}{|\bar{D}_i|}$$

Permitiendo más adelante ser  $a(d, \bar{D}_i) = dist(d, \bar{D}_i)$  la distancia del documento de a su agrupamiento  $\bar{D}_i (d \in \bar{D}_i)$ , y  $a(d, D_M) = \min_{\bar{D}_i \in D_M, d \notin \bar{D}_i} dist(d, \bar{D}_i)$  la distancia del documento d al agrupamiento de vecindad más cercano.

La silueta  $\vartheta(d, D_M)$  de un documento d está definido cómo

$$s(d, D_M) = \frac{b(d, D_M) - a(d, D_M)}{\max \{a(d, D_M), b(d, D_M)\}}$$

El coeficiente de silueta como

$$SC(D_M) = \frac{\sum_{p \in D^*} s(p, D_M)}{|D^*|}$$

El coeficiente de silueta es una medida de la calidad de agrupamiento, que es bastante independiente de el número de grupos,  $k$ . Experiencias, como las documentadas en (Kaufman y Rousseeuw, 1990), muestran que valores entre 0.7 y 1.0 indican resultados de agrupamiento con una excelente separación entre agrupamientos, *viz.* los puntos de datos están muy cerca del centro de su grupo y lejos del siguiente grupo más cercano. Para el rango de 0.5 a 0.7 se encuentra que los puntos de datos están claramente asignados a los centros de agrupamiento. Valores de 0.25 a 0.5 indican que se pueden encontrar centros de agrupamiento, aunque hay un considerable "ruido", es decir hay muchos puntos de datos que no se pueden asignar claramente a los agrupamientos. Por debajo de un valor de 0.25 se vuelve prácticamente imposible encontrar centros de agrupación significativos y asignar definitivamente la mayoría de puntos de datos[3].

## **DESARROLLO**

Se tiene que desarrollar un programa que indique para un número aleatorio de muestras cuántas clases son las más óptimas para separar todas las muestras usando lo conocido como "Validity Clusters"

Se probaron dos validity clusters siendo Silhouette Coefficient y por entropía los dos utilizados

A continuación se presenta el código donde se uso C-means para el validity cluster de entropía. Se hizo con C-means porque otorga valores basados en lógica difusa que son necesarios para las ecuaciones de entropía.

En el validity cluster de silueta se decidió usar k-means porque no son necesarios los datos basados en lógica difusa y el tiempo de computo era menor

```

1  """Practice 1: Compare Validity Clusters."""
2  import numpy as np
3  from scipy.stats import entropy
4  import matplotlib.pyplot as plt
5  from sklearn.cluster import KMeans
6  from sklearn.datasets import make_blobs
7  from sklearn.metrics import silhouette_score
8
9  # Generating the sample data from make_blobs
10
11  X, Y = make_blobs(centers=4)
12
13  no_of_clusters = [2, 3, 4, 5, 6]
14
15  plt.figure(1)
16  plt.plot(X[:, 0], X[:, 1], 'bo')
17
18  """C-means with Entropy."""
19  n, car = X.shape
20
21  max_data_number = max(max(X[:, 0], X[:, 1], key=lambda x: abs(x[1])))
22
23  for n_clusters in no_of_clusters:
24
25      V = np.random.rand(n_clusters, car) * max_data_number
26      D = np.zeros((n_clusters, 100))
27      U = np.zeros((n_clusters, 100))
28      epsilon = 1
29
30      while epsilon > 10**(-10):
31          Vold = np.array(V)
32          for i in range(n_clusters):
33              for k in range(n):
34                  a = 1**(-10)
35                  for j in range(car):
36                      a = (X[k, j] - V[i, j])**2 + a
37                      D[i, k] = a ** (1 / 2)
38
39              for i in range(n_clusters):
40                  for k in range(n):
41                      b = 0
42                      for j in range(n_clusters):
43                          b = (D[i, k] / D[j, k])**2 + b
44                      U[i, k] = 1 / b
45
46              for i in range(n_clusters):
47                  for j in range(car):
48                      c = 0
49                      d = 0
50                      for k in range(n):
51                          c = ((U[i, k]**2) * X[k, j]) + c
52                          d = ((U[i, k]**2)) + d
53                      V[i, j] = c / d
54          epsilon = sum(sum(abs(Vold[:] - V[:])))
55
56  total_entropy = 0

```

```

57
58     for i in range(n_clusters):
59         total_entropy += entropy(U[i])
60
61     total_entropy *= -(1 / 100)
62
63     print("For no of clusters =", n_clusters,
64           " The entropy is :", total_entropy)
65
66     plt.figure('Entropy {}'.format(n_clusters))
67     plt.plot(X[:, 0], X[:, 1], 'bo')
68     plt.plot(V[:, 0], V[:, 1], 'r*')
69
70     """K-means with Silhoutte."""
71
72     for n_clusters in no_of_clusters:
73
74         kmeans = KMeans(n_clusters=n_clusters)
75         # Fitting with inputs
76         kmeans = kmeans.fit(X)
77         # Predicting the clusters
78         labels = kmeans.predict(X)
79         # Getting the cluster centers
80         centers = kmeans.cluster_centers_
81
82         plt.figure('Silhoutte {}'.format(n_clusters))
83         plt.plot(X[:, 0], X[:, 1], 'bo')
84         plt.plot(centers[:, 0], centers[:, 1], 'r*')
85
86         # The silhouette_score gives the average value for all the samples.
87         silhouette_avg = silhouette_score(X, labels)
88
89         print("For no of clusters =", n_clusters,
90               " The average silhouette_score is :", silhouette_avg)
91
92     plt.show()

```

---

## CONCLUSIÓN

El método de entropía se vuelve tardado por el hecho de que se utiliza el algoritmo de c-means para el agrupamiento ya que como se ha mencionado las ecuaciones de entropía necesitan una matriz con valores basados en lógica difusa además de que el tiempo de computo cada vez se vuelve mayor cuando agregamos más clases ya que si queremos llegar el  $\epsilon$  adecuado se tiene que ejecutar más veces el algoritmo para alcanzar el valor de  $\epsilon$  adecuado en cambio con el método de la silueta el tiempo de ejecución fue notoriamente más rápido ya que se utilizó el algoritmo de k-means para el agrupamiento donde no se necesitaron valores de lógica difusa. Se pudo observar que utilizando el método de la silueta los valores que arrojaba hacían más sentido con las gráficas que muestra el programa. Con todo esto se puede concluir que el método de la silueta otorga mejor y más rápida validación sobre el número de clases a usar.

## REFERENCIAS

[1] Boskovitz, V., & Guterman, "An adaptive neuro-fuzzy system for automatic image segmentation and edge detection". *IEEE Transactions on fuzzy systems*, 2002, vol. 10, no 2, p. 247-262. [online], consultado el 29 de enero del 2020 en

<https://ieeexplore.ieee.org/abstract/document/995125>

[2] Martins,O.Braz,G. Aristófanés,S. Paiva, Anselmo. Gattass, Marcelo. "Detection of Masses in Digital Mammograms using K-means and Support Vector Machine". *ELCVIA: electronic letters on computer vision and image analysis*, [online], 2009, Vol. 8, Num. 2, pp. 43, consultado el 29 de enero del 2020 en

<https://www.raco.cat/index.php/ELCVIA/article/view/150156> [View: 30-01-2020]

[3] Hotho, A., Maedche, A., & Staab, S. (2002). "Ontology-based text document clustering". *KI*, 16(4), 48-54.[online], consultado el 29 de enero del 2020 en

[https://www.kde.cs.uni-kassel.de/wp-content/uploads/benz/hotho/pub/Ontology\\_based\\_Text\\_Document\\_Clustering\\_2002.pdf](https://www.kde.cs.uni-kassel.de/wp-content/uploads/benz/hotho/pub/Ontology_based_Text_Document_Clustering_2002.pdf)